



Temporal fidelity enhancement for video action recognition^{*#}

Shaowu XU¹, Xibin JIA^{‡1}, Qianmei SUN², Jing CHANG²

¹Faculty of Information Technology, Beijing University of Technology, Beijing 100124, China

²Beijing Chao-yang Hospital, Capital Medical University, Beijing 100020, China

E-mail: swxu@emails.bjut.edu.cn; jiaxibin@bjut.edu.cn; sunqianmei5825@126.com; cj006006@126.com

Received Mar. 14, 2025; Revision accepted June 4, 2025; Crosschecked July 10, 2025

Abstract: Temporal attention mechanisms are essential for video action recognition, enabling models to focus on semantically informative moments. However, these models frequently exhibit temporal infidelity—misaligned attention weights caused by limited training diversity and the absence of fine-grained temporal supervision. While video-level labels provide coarse-grained action guidance, the lack of detailed constraints allows attention noise to persist, especially in complex scenarios with distracting spatial elements. To address this issue, we propose temporal fidelity enhancement (TFE), a competitive learning paradigm based on the disentangled information bottleneck (DisenIB) theory. TFE mitigates temporal infidelity by decoupling action-relevant semantics from spurious correlations through adversarial feature disentanglement. Using pre-trained representations for initialization, TFE establishes an adversarial process in which segments with elevated temporal attention compete against contexts with diminished action relevance. This mechanism ensures temporal consistency and enhances the fidelity of attention patterns without requiring explicit fine-grained supervision. Extensive studies on UCF101, HMDB-51, and Charades benchmarks validate the effectiveness of our method, with significant improvements in action recognition accuracy.

Key words: Action recognition; Disentangled information bottleneck; Temporal modeling; Temporal fidelity
<https://doi.org/10.1631/FITEE.2500164>

CLC number: TP391.41

1 Introduction

Learning paradigms based on physical systems and neural mechanisms have advanced video-based action recognition by enhancing temporal dependency modeling and discriminative motion pattern extraction (Jiao LC et al., 2024, 2025). While traditional three-dimensional (3D) convolutional networks and two-stream architectures (Liu ZY et al., 2021; Liu Y et al., 2024) have established the

foundation for temporal modeling, their reliance on predefined hierarchies limits the adaptability to long-term dynamics. Multi-scale temporal dependency methods (Yu et al., 2020; Zhou JM et al., 2021; Zhang et al., 2025) partially alleviate this, but remain constrained by fixed parameterization. This has driven the use of temporal attention mechanisms, which dynamically prioritize critical temporal segments and improve motion feature extraction (Wu CY et al., 2022; Jiao JY et al., 2023; Yamazaki et al., 2023; Gao et al., 2024; Zhou JM et al., 2024).

However, these attention-based models for video action recognition (Mondal et al., 2023; Wang H et al., 2024; Wu WH et al., 2024) primarily depend on video-level supervision, providing only coarse-grained supervision despite the inherent fine-grained

[‡] Corresponding author

* Project supported by the National Natural Science Foundation of China (Nos. 62476015, 62171298, and 62476181)

Electronic supplementary materials: The online version of this article (<https://doi.org/10.1631/FITEE.2500164>) contains supplementary materials, which are available to authorized users

ORCID: Shaowu XU, <https://orcid.org/0000-0002-1607-7112>; Xibin JIA, <https://orcid.org/0000-0001-8799-8042>

© Zhejiang University Press 2025

temporal variations in action videos (e.g., eating speed differences or utensil usage patterns). This discrepancy causes temporal infidelity (Liang et al., 2020; Aghaeipoor et al., 2023)—a misalignment between attention distributions and action-relevant temporal segments. For instance, as shown in Fig. 1, for actions such as “eating a sandwich,” a model may mistakenly pay significant attention to adjacent segments that are irrelevant while overlooking key discriminative moments in unseen scenarios. This limitation arises from the training paradigm which depends solely on coarse video-level supervision, and fails to capture the fine-grained temporal dynamics intrinsic to action semantics. Although fine-grained annotations could get alleviated, their acquisition is prohibitively labor-intensive and resource-consuming.

To address this challenge, we propose a novel temporal fidelity enhancement (TFE) framework based on the disentangled information bottleneck (DisenIB) theory (Pan et al., 2021). TFE explicitly optimizes temporal attention fidelity by decoupling the latent video embedding into the action-relevant and semantically redundant components. Unlike conventional attention refinement methods, TFE resolves infidelity through adversarial disentanglement, which maximizes the sufficiency of critical

temporal features while suppressing spurious correlations from noisy contexts. The framework first encodes video segments into temporally discriminative embeddings using pre-trained models, followed by an adversarial process to approximate the DisenIB objective. A disentangler separates embeddings into temporal fidelity-preserving and fidelity-deviating ones, while dual approximators compete to retain action-relevant semantics and discard redundant patterns. This mechanism circumvents the combinatorial complexity of direct optimization, allowing for robust temporal attention alignment without requiring fine-grained supervision. The main contributions of this paper are summarized as follows:

1. We introduce TFE, a novel framework based on the DisenIB theory that explicitly optimizes temporal attention fidelity by disentangling action-relevant semantics from redundant information.
2. We propose an adversarial disentanglement mechanism that approximates the DisenIB-based objective, enabling effective suppression of spurious correlations without requiring fine-grained supervision.
3. Our method achieves outstanding performances on benchmark datasets, including Charades, UCF101, and HMDB-51.

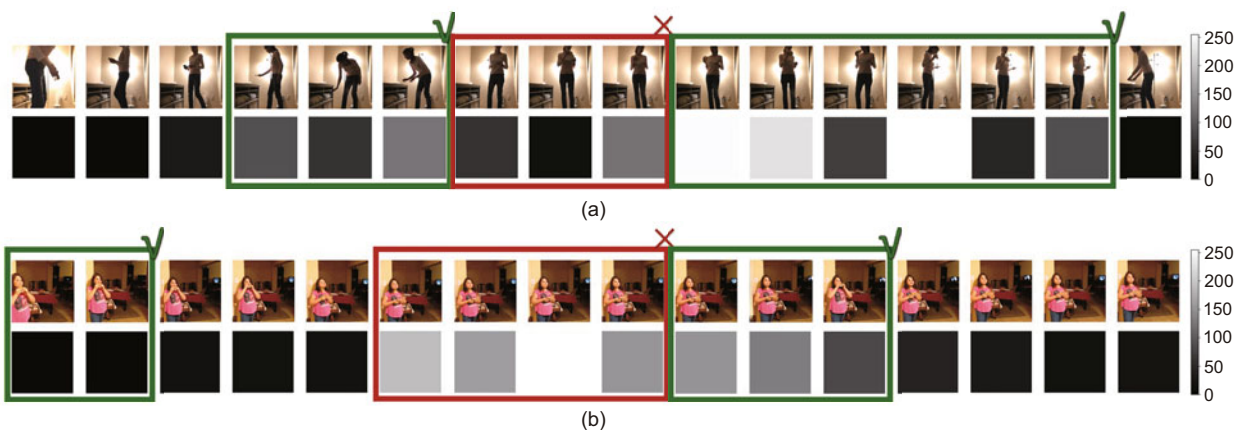


Fig. 1 Temporal attention distribution computed by BIKE ViT-L/14 (Wu WH et al., 2023) on training and testing samples with similar actions. (a) Training sample: the model allocates high attention to adjacent segments related to the action “eating a sandwich.” The whole action is that a person puts a phone onto a shelf, picks up a sandwich, and begins eating it. (b) Testing sample: the model also gives significant attention to adjacent segments, but those including the action “eating a sandwich” receive low attention. The whole action is that a person is standing in front of the pantry smiling and eating a sandwich, while the television is on. In (a) and (b), the first row shows the video segments, and the second row depicts the attention weights, with white indicating higher values and black indicating lower values. Green boxes highlight essential segments based on human cognitive experience, while red boxes indicate non-essential segments. References to color refer to the online version of this figure

2 Related works

Video recognition, unlike image recognition, requires understanding of the temporal evolution of objects. The most common methods use convolutional networks, with many adopting the 3D convolutional neural network (CNN) for action recognition (Carreira and Zisserman, 2017; Tran et al., 2018; Xie et al., 2018). Although advanced 3D CNN-based methods have shown promising results on short-term action benchmarks (Kuehne et al., 2011; Soomro et al., 2012), including CF-IIIH (Liu Y et al., 2024) with counterfactual reasoning and KCMM (Liu Y et al., 2025) with knowledge-driven composition modulation, robust long-term dependency modeling remains crucial for real-world applications (Girdhar et al., 2017; Sigurdsson et al., 2017; Zhou BL et al., 2018; Feichtenhofer et al., 2019; Feichtenhofer, 2020). CSVR (Zhang et al., 2025) enhances dynamic-static feature learning through generative-discriminative synergy, but still shares this limitation. Long-term video recognition (Sigurdsson et al., 2016) requires identifying long-term actions or their subactions (Yu et al., 2020; Zhou JM et al., 2021). Timeception (Hussein et al., 2019a), TRN (Zhou BL et al., 2018), VideoGraph (Hussein et al., 2019b), RhyRNN (Yu et al., 2020), and GHRM (Zhou JM et al., 2021) are methods for addressing long-term dynamics, but they are limited by pre-terminated patterns.

Recently, Vision Transformer models (Vaswani et al., 2017; Fan et al., 2021) and derivatives (Jiao JY et al., 2023; Yamazaki et al., 2023; Gao et al., 2024) have gained attention for their adaptability across video scales. UGPT (Guo et al., 2022) uses self-attention for long-term action dynamics. Wang R et al. (2023) used masked video distillation. TwinFormer (Zhou JM et al., 2024) models fine-to-coarse temporal structures, and MSQNet (Mondal et al., 2023) introduces a multimodal semantic query mechanism. Large-scale pre-trained vision-language models including ViLT-CLIP (Wang H et al., 2024), Text4Vis (Wu WH et al., 2024), and BIKE (Wu WH et al., 2023) have significantly advanced video action recognition through their respective novel cross-modal interaction designs.

However, these attention-based methods largely adhere to the information bottleneck (IB) principle (Tishby et al., 2000; Dimitrov and Miller, 2001)

for video embedding refinement. They use IB-based paradigms (Srivastava et al., 2021; Chi et al., 2022; Cen et al., 2023) to learn compressed yet predictive representations by maximizing mutual information between the video embeddings and target labels. Although this paradigm makes efficient use of coarse-grained (e.g., video-level) supervision, it remains inadequate in separating intra-class dynamic variations when fine-grained (e.g., segment-level) annotations are unavailable, resulting in temporal infidelity.

Our TFE framework addresses temporal infidelity by introducing a novel objective function based on DisenIB (Pan et al., 2021), which resolves temporal infidelity through two key innovations: (1) self-supervised separation of embeddings into action-relevant and redundant components using adversarial optimization, directly targeting the root cause of temporal attention misalignment; (2) dynamic suppression of spurious correlations using competing approximators, outperforming conventional IB's passive compression. TFE achieves temporal fidelity through feature-space purification rather than mere attention redistribution, while maintaining the scalability of coarse-grained supervision paradigms.

3 Method

TFE is a novel learning framework for video action recognition. This section first reviews the general objective function of temporal attention-based models, followed by the formulation of a DisenIB-based learning objective and its associated loss function, subsequently details the neural architecture design, and finally concludes with the complete learning scheme.

3.1 Objective function

3.1.1 Preliminaries

Temporal attention-based models (Guo et al., 2022; Mondal et al., 2023; Wang R et al., 2023; Zhou JM et al., 2024) generally comprise three core components: a visual encoder, a temporal aggregator, and an approximator. Given input video segments, the encoder maps them to embeddings $\mathbf{x} = \{\mathbf{x}_t \mid \mathbf{x}_t \in \mathbb{R}^{1 \times d}, t = 1, 2, \dots, T\}$, where d refers to dimensionality, T is the sequence length, and \mathbb{R} is the set of real numbers. The aggregator computes the temporal attention weights, and aggregates

the embeddings into a video-level representation \mathbf{z}_S , termed “salient embedding” for its focus on high-attention segments. The approximator learns the mapping from \mathbf{z}_S to the target action category distribution. This process emphasizes salient segments while suppressing others; it inherently follows the IB principle (Chen et al., 2018; Watson et al., 2024), minimizing

$$\mathcal{L} = I(\mathbf{x}; \mathbf{z}_S) - I(\mathbf{z}_S; \mathbf{y}), \quad (1)$$

where \mathcal{L} represents the objective function, \mathbf{y} represents the action label, and $I(\cdot; \cdot)$ represents the mutual information.

However, this formulation does not explicitly differentiate action-relevant and action-irrelevant information. Due to substantial fine-grained temporal variations in action videos, models relying on video-level supervision often misallocate attention to redundant segments, resulting in temporal infidelity.

3.1.2 Learning objective

TFE aims to explicitly disentangle video embeddings into two distinct semantic components: an action-relevant embedding that preserves recognition fidelity, and an action-irrelevant counterpart that reduces fidelity. This enables TFE to learn discriminative features for action recognition that are robust to intra-class temporal variations. Based on the DisenIB (Pan et al., 2021) principle, TFE introduces a novel objective function for learning the disentanglement of action semantics and redundant information. Let \mathbf{z}_N denote non-salient video embedding that captures action-irrelevant information, semantically complementary to \mathbf{z}_S . The DisenIB-based objective function (\mathcal{L}_{DIB}) is formalized as follows:

$$\mathcal{L}_{\text{DIB}} = -I(\mathbf{z}_S; \mathbf{y}) + I(\mathbf{z}_N; \mathbf{y}) + I(\mathbf{z}_S; \mathbf{z}_N). \quad (2)$$

Theorem 1 The \mathcal{L}_{DIB} to be minimized is consistent with the maximum compression.

Proof sketch: We first identify the global minimum $\mathcal{L}_{\text{DIB}}^* = -H(\mathbf{y})$. Assuming $\mathcal{L}_{\text{DIB}} - \mathcal{L}_{\text{DIB}}^* < \delta$ (δ is a constant introduced by the existential quantifier during the proof to represent a number greater than 0 that must exist), we derive $H(\mathbf{y}) - I(\mathbf{z}_S; \mathbf{y}) < \delta$, $I(\mathbf{z}_N; \mathbf{y}) < \delta$, and $I(\mathbf{z}_S; \mathbf{z}_N) < \delta$. Given $\mathbf{x} = \mathbf{z}_S \cup \mathbf{z}_N$, it follows $H(\mathbf{x}) - I(\mathbf{x}; \mathbf{z}_N, \mathbf{y}) < 3\delta$. Using the Markov chains $\mathbf{z}_S \leftrightarrow \mathbf{x} \leftrightarrow \mathbf{y}$, $\mathbf{z}_N \leftrightarrow \mathbf{x} \leftrightarrow \mathbf{y}$, and

$\mathbf{z}_S \leftrightarrow \mathbf{x} \leftrightarrow \mathbf{z}_N$, we apply mutual information properties to show $I(\mathbf{x}; \mathbf{y}) = H(\mathbf{y})$. By leveraging inequalities and combining intermediate results, we derive $|I(\mathbf{x}; \mathbf{z}_S) - H(\mathbf{y})| + |I(\mathbf{z}_S; \mathbf{y}) - H(\mathbf{y})| < 25\delta$. Thus, setting $\delta = \epsilon/25$ (ϵ is a constant introduced by the universal quantifier during the proof to represent any number greater than zero), we show that \mathcal{L}_{DIB} is consistent on the maximum compression according to Definition 1 in Pan et al. (2021). Detailed proof is shown in the supplementary materials.

Theorem 1 shows that the \mathcal{L}_{DIB} adheres to the maximum compression principle, which ensures that the learned salient and non-salient video embeddings play distinct roles to optimally balance information preservation and compression, thereby preserving recognition fidelity.

3.1.3 Training strategy

We propose Theorem 2 to derive an optimizable training objective of Eq. (2).

Theorem 2 The global optimum for minimizing \mathcal{L}_{DIB} satisfies

$$D^* = \arg \min_D \mathbb{E}[-\log p(\mathbf{y}|\mathbf{z}_S) + \log p(\mathbf{y}|\mathbf{z}_N)], \quad (3)$$

where $\mathbb{E}(\cdot)$ is the mathematical expectation, D denotes the disentangler, and D^* denotes the optimal disentangler in Eq. (2) that divides input segment embeddings \mathbf{x} into \mathbf{z}_S and \mathbf{z}_N .

Proof sketch: We divide \mathcal{L}_{DIB} into two sections: $\mathcal{L}_1 = -I(\mathbf{z}_S; \mathbf{y}) + I(\mathbf{z}_N; \mathbf{y})$ and $\mathcal{L}_2 = I(\mathbf{z}_S; \mathbf{z}_N)$. First, we show that D^* minimizes \mathcal{L}_1 by showing $\mathbb{E}[-\log p(\mathbf{y}|\mathbf{z}_S) + \log p(\mathbf{y}|\mathbf{z}_N)] \geq \mathbb{E}[-\log p(\mathbf{y}|\mathbf{z}_S^*) + \log p(\mathbf{y}|\mathbf{z}_N^*)]$ for any pair $(\mathbf{z}_S, \mathbf{z}_N)$, leading to $-I(\mathbf{z}_S; \mathbf{y}) + I(\mathbf{z}_N; \mathbf{y}) \geq -I(\mathbf{z}_S^*; \mathbf{y}) + I(\mathbf{z}_N^*; \mathbf{y})$. Second, using a proof by contradiction, we demonstrate that minimizing \mathcal{L}_1 also minimizes \mathcal{L}_2 , ensuring that the $I(\mathbf{z}_S; \mathbf{z}_N)$ is minimized. Thus, D^* provides the global optimum for minimizing \mathcal{L}_{DIB} . Detailed proof can be found in the supplementary materials.

Notably, the optimal solution D^* is based on the true conditional distributions $p(\mathbf{y}|\mathbf{z}_S)$ and $p(\mathbf{y}|\mathbf{z}_N)$, which are analytically intractable, making direct optimization infeasible. Inspired by Liang et al. (2020), TFE uses a self-adversarial training strategy to approximate D^* . Specifically, two approximators, the salient approximator A_S and the non-salient approximator A_N , estimate the variational distributions

$q_S(\mathbf{y}|\mathbf{z}_S)$ and $q_N(\mathbf{y}|\mathbf{z}_N)$. An adversarial optimization framework then alternates between maximizing $\log q_S(\mathbf{y}|\mathbf{z}_S)$ and minimizing $\log q_N(\mathbf{y}|\mathbf{z}_N)$, driving D to the optimal state. This strategy is formalized as follows:

$$D^* = \arg \min_D \mathbb{E}[-\log q_S(\mathbf{y}|\mathbf{z}_S) + \log q_N(\mathbf{y}|\mathbf{z}_N)], \tag{4}$$

$$(A_S^*, A_N^*) = \arg \min_{A_S, A_N} \mathbb{E}[-\log q_S(\mathbf{y}|\mathbf{z}_S) - \log q_N(\mathbf{y}|\mathbf{z}_N)]. \tag{5}$$

3.2 Network architecture

3.2.1 Overview

The TFE framework contains three components: a visual encoder that extracts visual features and contextual semantics from raw video segments, resulting in segment embeddings, a disentangler that divides them into salient and non-salient video embeddings, and a pair of approximators that process respective embeddings using the training strategy outlined in Section 3.1.

Fig. 2 shows our model architecture. First, the visual encoder maps the input segment sequence to embeddings $\mathbf{x} = \{\mathbf{x}_t \mid \mathbf{x}_t \in \mathbb{R}^d, t = 1, 2, \dots, T\}$, where T represents the sequence length. Our design decouples the encoder from the subsequent disentangler–approximator modules, allowing for seamless integration into state-of-the-art (SOTA) methods by simply replacing their final layer with our components.

Subsequently, the disentangler processes input \mathbf{x} to produce salient and non-salient video embeddings using H parallel temporal disentanglement (TD) heads. Each head h generates semantically distinct embeddings $\mathbf{z}_S^h, \mathbf{z}_N^h \in \mathbb{R}^{T \times d}$. Concatenating all heads results in $\mathbf{z}_S = [\mathbf{z}_S^1, \mathbf{z}_S^2, \dots, \mathbf{z}_S^H]$ and

$\mathbf{z}_N = [\mathbf{z}_N^1, \mathbf{z}_N^2, \dots, \mathbf{z}_N^H]$. Two feedforward networks f_{DS} and f_{DN} are then merged to form the final embeddings:

$$(\mathbf{z}'_S, \mathbf{z}'_N) = D(\mathbf{x}) = (f_{DS}(\mathbf{z}_S) + \mathbf{x}, f_{DN}(\mathbf{z}_N) + \mathbf{x}). \tag{6}$$

Finally, TFE incorporates two approximators: A_S and A_N . A_S predicts action labels based on salient video embeddings, while A_N predicts action labels based on non-salient embeddings. In practice, custom-designed networks can be simple fully connected layers; when integrated into existing models, the original approximator structure can be reused. To compute the losses, these approximators use salient embeddings \mathbf{z}'_S and non-salient embeddings \mathbf{z}'_N as inputs and output action predictions:

$$\mathcal{L}_S = \ell(\mathbf{y}, A_S(\mathbf{z}'_S)), \quad \mathcal{L}_N = \ell(\mathbf{y}, A_N(\mathbf{z}'_N)), \tag{7}$$

where $\ell(\cdot, \cdot)$ represents the cross-entropy loss.

The detailed descriptions of the TD module and learning scheme are elaborated in the following subsections.

3.2.2 Temporal disentanglement module

TD divides segment embeddings into two distinct representations: salient video embedding which preserves temporal fidelity for action recognition, and non-salient embedding which encodes redundant temporal contexts. As shown in Fig. 3, our module implements this using a temporal attention mechanism that disentangles input sequences into semantically divergent embeddings.

The module accepts segment embeddings \mathbf{x} as input. It first computes a temporal attention mask \mathbf{m} to probabilistically separate salient and non-salient temporal segments during feature transformation:

$$\mathbf{m} = \text{Gumbel-softmax}(\alpha(\mathbf{x}), k_i) \in \mathbb{R}^{T \times T}, \tag{8}$$

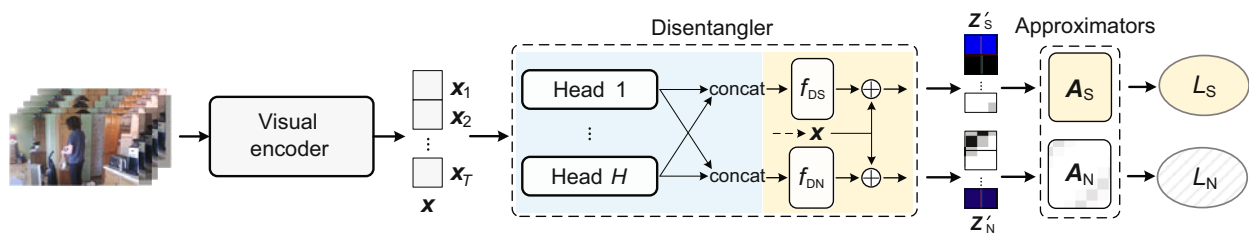


Fig. 2 Overview of TFE. We guide our model to learn temporal attention with fidelity for action recognition with the DisenIB objective. The model consists of a visual encoder, a disentangler, and a pair of approximators

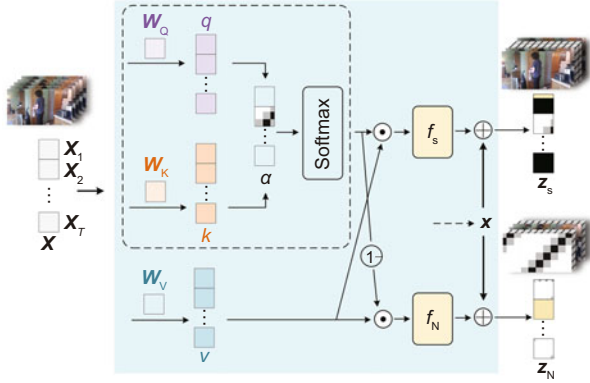


Fig. 3 Architecture of the TD module. Given a segment embedding sequence x , this module outputs a pair of semantically distinct video embeddings (z_S, z_N). W_Q, W_K , and W_V represent the mapping matrices of query, key, and value, respectively. q, k , and v represent the results obtained from the input x after mapping transformations

$$\alpha(x) = \frac{(xW_Q)(xW_K)^T}{\sqrt{d}}, \quad (9)$$

where k_i is an integer satisfying $0 < k_i \leq T$, and $W_Q, W_K \in \mathbb{R}^{d \times d}$ are projection matrices. This differentiable mask m quantifies temporal saliency, where higher values indicate action-critical segments while lower values correspond to non-essential patterns. We adopt Gumbel-softmax reparameterization (Liang et al., 2020) to enable efficient k -hot sampling for combinatorial optimization of feature disentanglement.

The proposed TD module computes two mutually distinct video embeddings as follows:

$$z_S = f_S(mxW_V + x) \in \mathbb{R}^{T \times d}, \quad (10)$$

$$z_N = f_N((1 - m)xW_V + x) \in \mathbb{R}^{T \times d}, \quad (11)$$

where $W_V \in \mathbb{R}^{d \times d}$ represents the value projection matrix, f_S and f_N are feedforward neural networks with the same structure, z_S is the salient video embedding, and z_N is the non-salient video embedding. The complementary mask $1 - m$ ensures that z_N captures temporal features neglected by m . Since m quantifies saliency, $1 - m$ inversely weights segments with lower saliency, which generally encode residual motion patterns (e.g., background hand movements in “eating” actions) or redundant temporal contexts. This explicit inversion ensures mutual exclusivity between z_S and z_N , which corresponds to DisenIB’s requirement for complementary latent disentanglement.

3.2.3 Learning scheme

The proposed TFE adopts a joint objective in expression (12), which is derived from the theoretical bounds in Eqs. (4)–(5), to materialize the DisenIB principle through adversarial coordination:

$$\min_{(A_S, A_N)} (\mathcal{L}_S + \mathcal{L}_N), \quad \min_D (\mathcal{L}_S - \mathcal{L}_N). \quad (12)$$

This is executed via iterative two-phase optimization:

1. Approximator phase: update A_S and A_N using $\nabla(\mathcal{L}_S + \mathcal{L}_N)$ (∇ denotes the gradient) to maximize feature utility from both salient and non-salient embeddings.

2. Disentangler phase: update D with $\nabla(\mathcal{L}_S - \mathcal{L}_N)$ to enforce orthogonality between z_S and z_N .

This coordinated optimization addresses temporal infidelity by enhancing temporal discriminability in z_S by \mathcal{L}_S maximization and suppressing action-irrelevant temporal variations in z_N via \mathcal{L}_N minimization.

4 Experiments

4.1 Datasets

UCF101 (Soomro et al., 2012) is an action recognition dataset of realistic action videos collected from YouTube, with 101 action categories and 13 320 videos. The videos in 101 action categories are grouped into 25 groups, with each group consisting of 4–7 videos of an action. All videos in this dataset are single-label, with an average duration of 5 s, ranging from 4 to 7 s in length. We split the dataset into 9537 videos for training and 3783 videos for testing.

HMDB-51 (Kuehne et al., 2011) is a collection of realistic videos from various sources, including movies and web videos. The dataset comprises 5100 video clips from 51 action categories. This dataset contains only single-label videos, with an average duration of 7 s and video lengths ranging from 4 to 22 s. We split the dataset into 3570 videos for training and 1530 videos for testing.

Charades (Sigurdsson et al., 2016) is a dataset composed of 9848 videos of daily indoor activities. The dataset contains 66 500 temporal annotations for 157 action classes, 41 104 labels for 46 object classes, and 27 847 textual descriptions of the videos. Each video in this dataset has multiple labels, and

the average duration is 29 s, with video lengths ranging from 2 to 187 s. We divide the dataset into 7985 videos for training and 1863 videos for testing.

4.2 Implementation details

In experiments, the TFE framework integrates the established architectures from SOTA methods (Mondal et al., 2023; Wu WH et al., 2023, 2024) as backbone components, preserving their visual encoders while replacing final classification layers with our disentangler–approximator modules. Hyperparameters align with the original implementations, with 16 uniformly sampled segments per video with a frame resolution of 336×336 (Wu WH et al., 2023, 2024) or 224×224 (Mondal et al., 2023). The network is trained using AdamW (Loshchilov and Hutter, 2019) with hyperparameters of $\beta_1 = 0.9$ and $\beta_2 = 0.999$, an initial learning rate $5e-5$ (cosine decay with 5-epoch warmup), and a weight decay of 0.2 over 100 epochs. All the experiments are implemented in PyTorch (Paszke et al., 2019) on an NVIDIA RTX 3090 GPU with 24 GB memory.

4.3 Main results

In temporal attention models, temporal fidelity quantifies how effectively attention mechanisms prioritize action-critical segments. Enhancing temporal fidelity directly correlates with improved recognition accuracy. To validate TFE’s fidelity enhancement capability, we integrate it into three SOTA baselines: MSQNet (Mondal et al., 2023), BIKE (Wu WH et al., 2023), and Text4Vis (Wu WH et al., 2024). Experiments cover UCF101 (Soomro et al., 2012), HMDB-51 (Kuehne et al., 2011), and Charades (Sigurdsson et al., 2016). For multi-label Charades, evaluation uses mean average precision (mAP), while single-label datasets (UCF101/HMDB-51) are assessed via classification accuracy.

4.3.1 Long-term action recognition

Table 1 compares the SOTA approaches on the Charades dataset for long-term video understanding. While 3D convolution-based methods (e.g., ResNet-152 (He et al., 2016)) and short-term temporal approaches (e.g., temporal fields (Sigurdsson et al., 2017)) show limited effectiveness in capturing long-range dependencies, architectures modeling complex temporal relationships achieve superior

performance. Specifically, GHRM (Zhou JM et al., 2021) uses graph-based reasoning, while SlowFast R101 (Feichtenhofer et al., 2019) uses dual-pathway modeling, though both remain constrained by fixed architectural priors. Temporal attention mechanisms address this limitation through dynamic importance weighting of video segments, thereby simultaneously enhancing local action discrimination and global temporal context modeling. This explains the superior performance of TwinFormer (Zhou JM et al., 2024), ActionCLIP (Wang MM et al., 2023), and BIKE (Wu WH et al., 2023) over non-attentive counterparts.

Table 1 Comparisons of the state-of-the-art approaches on Charades in terms of mAP values of sub-actions within long-term actions

Method*	mAP (%)
Temporal Fields (Sigurdsson et al., 2017)	22.4
ResNet-152 (He et al., 2016)	22.8
TRN (Zhou BL et al., 2018)	25.2
RhyRNN (Yu et al., 2020)	25.4
I3D (Carreira and Zisserman, 2017)	32.9
STM (Jiang et al., 2019)	35.3
Timeception (Hussein et al., 2019a)	37.2
VideoGraph (Hussein et al., 2019b)	37.8
GHRM (Zhou JM et al., 2021)	38.3
SlowFast R101 (Feichtenhofer et al., 2019)	43.4
X3D-XL (Feichtenhofer, 2020)	44.3
Method**	mAP (%)
UGPT (Guo et al., 2022)	42.4
TwinFormer (Zhou JM et al., 2024)	43.6
MViT-B (Fan et al., 2021)	43.9
ActionCLIP (Wang MM et al., 2023)	44.3
AdaFocus (Li XH et al., 2023)	47.8
MSQNet (Mondal et al., 2023)	48.5
MSQNet+TFE	49.5
BIKE (Wu WH et al., 2023)	49.4
BIKE+TFE	50.1

The best result is in bold. * Without temporal attention; ** with temporal attention

Our proposed TFE framework achieves the best results, 50.1% mAP with the BIKE backbone and 49.5% with the MSQNet backbone. Notably, TFE shows smaller gains with BIKE than with MSQNet. This is due to their supervision granularity: BIKE leverages segment-level labels by splitting long videos into clips during training, while MSQNet uses untrimmed videos with only video-level supervision. This indicates TFE’s superior capability to approximate DisenIB’s optimal solution

through self-adversarial learning, particularly under coarse-grained supervision without explicit temporal fidelity constraints.

4.3.2 Short-term action recognition

Table 2 compares the recognition results of different models on the short-term video action datasets UCF101 and HMDB-51. While some models not relying on temporal attention, such as TDN (Wang LM et al., 2021), CF-IIH (Liu Y et al., 2024), and STANet (Li XH et al., 2023), achieve good recognition accuracy, attention-based models such as ViLT-CLIP (Wang H et al., 2024), Text4Vis (Wu WH et al., 2024), ZeroI2V (Li Z et al., 2023), and BIKE (Wu WH et al., 2023) generally achieve higher recognition rates. The temporal attention mechanism enhances the model’s sensitivity to short-term

Table 2 Comparisons of the recognition results of different models on short-term action datasets UCF101 and HMDB-51, where the recognition accuracy of single-label videos is used as the evaluation metric

Method*	RA (%)	
	UCF	HMDB
ARTNet (Wang LM et al., 2018)	94.3	74.8
CSVR (Zhang et al., 2025)	94.5	65.5
I3D (Carreira and Zisserman, 2017)	95.6	70.9
CoViFocus (Zheng et al., 2024)	95.8	74.8
TSM (Lin et al., 2019)	95.9	73.5
STM (Jiang et al., 2019)	96.2	72.2
R(2+1)D (Tran et al., 2018)	96.8	74.5
MVFNet (Wu WH et al., 2021)	96.6	75.7
S3D-G (Xie et al., 2018)	96.8	75.9
CF-IIH (Liu Y et al., 2024)	96.9	76.7
TDN (Wang LM et al., 2021)	97.4	76.4
KCMM (Liu Y et al., 2025)	97.4	77.1
STANet (Li XH et al., 2023)	97.6	77.7

Method**	RA (%)	
	UCF	HMDB
LCVE (Ishikawa et al., 2024)	95.6	70.0
VideoMAE (Tong et al., 2022)	96.1	73.3
ActionCLIP (Wang MM et al., 2023)	97.1	76.2
ViLT-CLIP (Wang H et al., 2024)	97.5	73.3
MVD-B (Wang R et al., 2023)	97.5	79.7
ZeroI2V (Li XH et al., 2023)	98.6	83.4
MSQNet (Mondal et al., 2023)	96.0	72.8
MSQNet+TFE	96.5	73.3
Text4Vis (Wu WH et al., 2024)	98.1	81.3
Text4Vis+TFE	98.3	81.4
BIKE (Wu WH et al., 2023)	98.9	83.1
BIKE+TFE	99.1	84.3

The best result is in bold. RA: recognition accuracy. * Without temporal attention; ** with temporal attention

actions and optimizes information extraction, enabling the model to extract action-relevant features.

The TFE framework consistently enhances short-term action recognition across backbones by suppressing non-salient temporal embeddings. The limited gains with Text4Vis originate from its classification-layer text–visual alignment, which diminishes temporal attention reliance when compared to BIKE’s segment-level optimization. MSQNet’s smaller improvement versus BIKE is due to the short-term videos’ inherent low intra-class variation, which reduces the need for segment-level supervision. This performance ceiling highlights that TFE’s attention refinement is secondary to base feature quality in scenarios where BIKE’s visual–textual interaction dominates temporal modeling.

4.4 Ablation studies

We evaluate three key components: TD, optimization via Eq. (4), and optimization via Eq. (5). Table 3 shows ablation results using MSQNet, selected based on findings in Section 4.3 where it exhibits more discernible performance variations than other backbones, enabling clearer ablation insights.

Variants 1 and 2 show that using dual-stream architectures with identical optimization objectives (Eq. (4) or (5)) is insufficient for effective embedding differentiation, and structural complexity degrades overall accuracy. Variant 3 also shows that simply combining these two objectives without TD leads to even lower performance, suggesting that without TD’s structural guidance, the interaction between the two objectives may introduce conflicting gradients that hinder effective representation learning. Variant 5 also shows that using Eq. (5) alone—focusing on approximator distribution matching—fails to induce mutually exclusive salient/non-salient segment embeddings, resulting in

Table 3 Ablation studies of TFE on Charades, UCF101, and HMDB-51 datasets

Method	TD	Opt. via Eq. (4)	Opt. via Eq. (5)	mAP (%)
Variant 1	–	–	✓	48.3
Variant 2	–	✓	–	48.2
Variant 3	–	✓	✓	48.1
Variant 4	✓	✓	–	48.7
Variant 5	✓	–	✓	48.2
TFE	✓	✓	✓	49.5

Opt.: optimization

no significant improvement and even slight performance degradation compared to Variant 1. This indicates that Eq. (5) is insufficient on its own without the complementary effect of adversarially driven embedding separation. In contrast, the opposing prediction requirement of Eq. (4) (i.e., Variant 4) shows a clearer improvement when combined with TD, indicating its effectiveness in leveraging TD to enforce discriminative learning, though gains remain limited due to its non-adversarial formulation.

The comprehensive comparison in Table 3 shows that only the disentangler–approximator framework with adversarial learning enables TFE, thereby improving recognition accuracy.

4.5 Visualization

As shown in Fig. 4, we visualize the temporal attention calculated by the vanilla model and TFE for comparison. This visualization highlights a benefit of TFE: removal of non-essential temporal variation. For the videos shown in Fig. 4, TFE assigns high attention to the first two segments containing the essential information “eating” and “opening television,” avoiding interference from irrelevant information that could disrupt action recognition.

4.6 Computational cost

While the proposed TFE framework introduces additional computational overhead due to its dual-stream structure within a single attention layer, it achieves overall computational efficiency by reducing the total number of attention layers in the network, resulting in improved performance without increasing resource consumption. As shown in Table 4, taking MSQNet as an example, the vanilla model uses 12 Transformer layers for segment embedding refinement, following the front-end vision-language model. Our approach reduces the total number of layers to 7 by replacing the final layer with a disentangler and a pair of approximators to implement TFE. This results in a reduction of 2.5×10^9 floating point operations (FLOPs) and 4.2×10^6 parameters compared to the vanilla model.

4.7 Discussion

Although TFE improves temporal attention models, there are several limitations to consider. The framework’s TD module specifically targets segment-level infidelity, successfully aligning attention learning with segment saliency to improve the overall performance. However, this approach fails to address finer-grained infidelity at the patch level within

Table 4 Performance and efficiency comparison with or without TFE on the Charades dataset

Method	Layer count	FLOPs ($\times 10^9$)	Number of parameters ($\times 10^6$)	mAP (%)
MSQNet (Mondal et al., 2023)	12	1303.9	220.4	48.5
MSQNet+TFE	7	1301.4	216.2	49.5
BIKE (Wu WH et al., 2023)	6	1864.7	230.4	49.4
BIKE+TFE	2	1864.4	211.5	50.1

FLOPs and the number of parameters are measured for a single inference pass with the same input. Layer count refers exclusively to segment embedding refinement modules, excluding the visual encoder

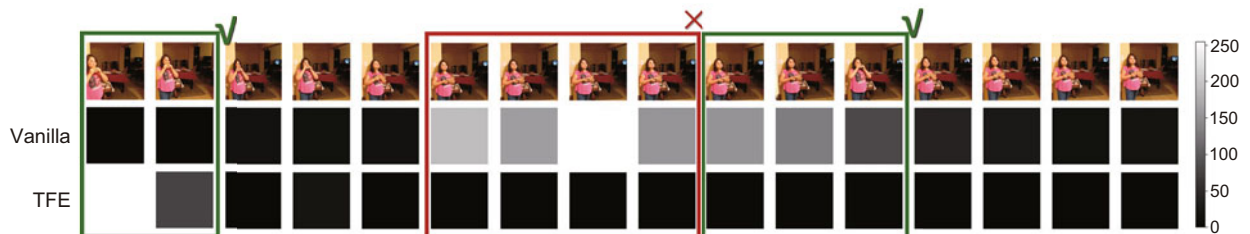


Fig. 4 Visualization of temporal attention calculated by the vanilla model and TFE. In each subfigure, the first row represents the original video segments showing that a person is standing in front of the pantry smiling and eating a sandwich, while the television is on. The second and third rows are the temporal attention computed by the vanilla model and TFE, respectively, with white indicating higher values and black indicating lower values. The green box highlights essential segments based on human cognitive experience, while the red box indicates non-essential ones. References to color refer to the online version of this figure

segments. This limitation becomes particularly evident in short-term action recognition, where TFE shows more modest improvements on MSQNet compared to BIKE across short-term datasets. The performance gap stems from rapid motion changes, creating semantically similar but visually heterogeneous patch dynamics. While BIKE's stronger visual-textual interaction reduces sensitivity to this dynamics, MSQNet's visual-only segment embedding refinement remains affected. As a result, TFE's visual self-adversarial mechanism, without considering patch-level infidelity, achieves suboptimal results for short-term action recognition compared to long-term scenarios. Future work should incorporate patch-level fidelity considerations to better handle highly dynamic scenes.

5 Conclusions

We propose the TFE framework for video action recognition, which approximates the DisenIB through self-adversarial learning. TFE uses a temporal disentangler to divide input videos into two semantically distinct embeddings: an action-relevant one preserving temporal fidelity, and an action-irrelevant counterpart that encodes residual temporal contexts. TFE enforces strict separation between these components by adversarial coordination of dual approximators, enabling robust recognition against intra-class temporal variations without requiring fine-grained supervision. Experiments using Charades, UCF101, and HMDB-51 validate TFE's effectiveness.

Contributors

Shaowu XU proposed the idea, conceived the study, implemented the source code, and drafted the paper. Xibin JIA supervised the research and revised the paper. Qianmei SUN and Jing CHANG contributed to the funding acquisition. Shaowu XU finalized the paper.

Conflict of interest

All the authors declare that they have no conflict of interest.

Data availability

The data that support the findings of this study are available from the corresponding author upon reasonable request.

References

- Aghaeipoor F, Sabokrou M, Fernández A, 2023. Fuzzy rule-based explainer systems for deep neural networks: from local explainability to global understanding. *IEEE Trans Fuzzy Syst*, 31(9):3069-3080. <https://doi.org/10.1109/TFUZZ.2023.3243935>
- Carreira J, Zisserman A, 2017. Quo vadis, action recognition? A new model and the kinetics dataset. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.4724-4733. <https://doi.org/10.1109/CVPR.2017.502>
- Cen J, Zhang SW, Wang X, et al., 2023. Enlarging instance-specific and class-specific information for open-set action recognition. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.15295-15304. <https://doi.org/10.1109/CVPR52729.2023.01468>
- Chen JB, Song L, Wainwright MJ, et al., 2018. Learning to explain: an information-theoretic perspective on model interpretation. *Proc 35th Int Conf on Machine Learning*, p.882-891.
- Chi HG, Ha MH, Chi S, et al., 2022. InfoGCN: representation learning for human skeleton-based action recognition. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.20154-20164. <https://doi.org/10.1109/CVPR52688.2022.01955>
- Dimitrov AG, Miller JP, 2001. Neural coding and decoding: communication channels and quantization. *Netw Comput Neur Syst*, 12(4):441-472. <https://doi.org/10.1080/net.12.4.441.472>
- Fan HQ, Xiong B, Mangalam K, et al., 2021. Multiscale vision Transformers. *Proc IEEE/CVF Int Conf on Computer Vision*, p.6804-6815. <https://doi.org/10.1109/ICCV48922.2021.00675>
- Feichtenhofer C, 2020. X3D: expanding architectures for efficient video recognition. *Proc IEEE Conf on Computer Vision and Pattern Recognition*, p.200-210. <https://doi.org/10.1109/CVPR42600.2020.00028>
- Feichtenhofer C, Fan HQ, Malik J, et al., 2019. SlowFast networks for video recognition. *Proc IEEE/CVF Int Conf on Computer Vision*, p.6201-6210. <https://doi.org/10.1109/ICCV.2019.00630>
- Gao SY, Chen Z, Chen G, et al., 2024. AVSegFormer: audio-visual segmentation with Transformer. *Proc 38th AAAI Conf on Artificial Intelligence*, p.12155-12163. <https://doi.org/10.1609/aaai.v38i11.29104>
- Girdhar R, Ramanan D, Gupta A, et al., 2017. Action-VLAD: learning spatio-temporal aggregation for action classification. *Proc IEEE Conf on Computer Vision and Pattern Recognition*, p.3165-3174. <https://doi.org/10.1109/CVPR.2017.337>
- Guo HJ, Wang HJ, Ji Q, 2022. Uncertainty-guided probabilistic Transformer for complex action recognition. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.20020-20029. <https://doi.org/10.1109/CVPR52688.2022.01942>
- He KM, Zhang XY, Ren SQ, et al., 2016. Deep residual learning for image recognition. *Proc IEEE Conf on Computer Vision and Pattern Recognition*, p.770-778. <https://doi.org/10.1109/CVPR.2016.90>
- Hussein N, Gavves E, Smeulders AWM, 2019a. Timeception for complex action recognition. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.254-263. <https://doi.org/10.1109/CVPR.2019.00034>

- Hussein N, Gavves E, Smeulders AWM, 2019b. VideoGraph: recognizing minutes-long human activities in videos. <https://arxiv.org/abs/1905.05143>
- Ishikawa Y, Kondo M, Kataoka H, 2024. Learnable cube-based video encryption for privacy-preserving action recognition. Proc IEEE/CVF Winter Conf on Applications of Computer Vision, p.6988-6998. <https://doi.org/10.1109/WACV57701.2024.00685>
- Jiang BY, Wang MM, Gan WH, et al., 2019. STM: spatiotemporal and motion encoding for action recognition. Proc IEEE/CVF Int Conf on Computer Vision, p.2000-2009. <https://doi.org/10.1109/ICCV.2019.00209>
- Jiao JY, Tang YM, Lin KY, et al., 2023. DilateFormer: multi-scale dilated Transformer for visual recognition. *IEEE Trans Multim*, 25:8906-8919. <https://doi.org/10.1109/TMM.2023.3243616>
- Jiao LC, Song X, You C, et al., 2024. AI meets physics: a comprehensive survey. *Artif Intell Rev*, 57(9):256. <https://doi.org/10.1007/s10462-024-10874-4>
- Jiao LC, Ma MR, He P, et al., 2025. Brain-inspired learning, perception, and cognition: a comprehensive review. *IEEE Trans Neur Netw Learn Syst*, 36(4):5921-5941. <https://doi.org/10.1109/TNNLS.2024.3401711>
- Kuehne H, Jhuang H, Garrote E, et al., 2011. HMDB: a large video database for human motion recognition. Proc Int Conf on Computer Vision, p.2556-2563. <https://doi.org/10.1109/ICCV.2011.6126543>
- Li XH, Zhu YH, Wang LM, 2023. Zeroi2V: zero-cost adaptation of pre-trained Transformers from image to video. Proc 18th European Conf on Computer Vision, p.425-443. https://doi.org/10.1007/978-3-031-73010-8_25
- Li Z, Zhang RQ, Zou DQ, et al., 2023. Robin: a novel method to produce robust interpreters for deep learning-based code classifiers. Proc 38th IEEE/ACM Int Conf on Automated Software Engineering, p.27-39. <https://doi.org/10.1109/ASE56229.2023.00164>
- Liang J, Bai B, Cao YR, et al., 2020. Adversarial infidelity learning for model interpretation. Proc 26th ACM SIGKDD Int Conf on Knowledge Discovery & Data Mining, p.286-296. <https://doi.org/10.1145/3394486.3403071>
- Lin J, Gan C, Han S, 2019. TSM: temporal shift module for efficient video understanding. Proc IEEE/CVF Int Conf on Computer Vision, p.7082-7092. <https://doi.org/10.1109/ICCV.2019.00718>
- Liu Y, Liu F, Jiao LC, et al., 2024. A knowledge-based hierarchical causal inference network for video action recognition. *IEEE Trans Multim*, 26:9135-9149. <https://doi.org/10.1109/TMM.2024.3386339>
- Liu Y, Liu F, Jiao LC, et al., 2025. Knowledge-driven compositional action recognition. *Patt Recogn*, 163:111452. <https://doi.org/10.1016/j.patcog.2025.111452>
- Liu ZY, Wang LM, Wu W, et al., 2021. TAM: temporal adaptive module for video recognition. Proc IEEE/CVF Int Conf on Computer Vision, p.13688-13698. <https://doi.org/10.1109/ICCV48922.2021.01345>
- Loshchilov I, Hutter F, 2019. Decoupled weight decay regularization. Proc 7th Int Conf on Learning Representations. <https://arxiv.org/abs/1711.05101>
- Mondal A, Nag S, Prada JM, et al., 2023. Actor-agnostic multi-label action recognition with multi-modal query. Proc IEEE/CVF Int Conf on Computer Vision Workshops, p.784-794. <https://doi.org/10.1109/ICCVW60793.2023.00086>
- Pan ZQ, Niu L, Zhang JF, et al., 2021. Disentangled information bottleneck. Proc 35th AAAI Conf on Artificial Intelligence, p.9285-9293. <https://doi.org/10.1609/aaai.v35i10.17120>
- Paszke A, Gross S, Massa F, et al., 2019. PyTorch: an imperative style, high-performance deep learning library. Proc 33rd Int Conf on Neural Information Processing Systems, Article 721.
- Sigurdsson GA, Varol G, Wang XL, et al., 2016. Hollywood in homes: crowdsourcing data collection for activity understanding. Proc 14th European Conf on Computer Vision, p.510-526. https://doi.org/10.1007/978-3-319-46448-0_31
- Sigurdsson GA, Divvala S, Farhadi A, et al., 2017. Asynchronous temporal fields for action recognition. Proc IEEE Conf on Computer Vision and Pattern Recognition, p.5650-5659. <https://doi.org/10.1109/CVPR.2017.599>
- Soomro K, Zamir AR, Shah M, 2012. UCF101: a dataset of 101 human actions classes from videos in the wild. <https://arxiv.org/abs/1212.0402>
- Srivastava A, Dutta O, Gupta J, et al., 2021. A variational information bottleneck based method to compress sequential networks for human action recognition. Proc IEEE/CVF Winter Conf on Applications of Computer Vision, p.2744-2753. <https://doi.org/10.1109/WACV48630.2021.00279>
- Tishby N, Pereira FC, Bialek W, 2000. The information bottleneck method. <https://arxiv.org/abs/physics/0004057>
- Tong Z, Song YB, Wang J, et al., 2022. VideoMAE: masked autoencoders are data-efficient learners for self-supervised video pre-training. Proc 36th Int Conf on Neural Information Processing Systems, Article 732.
- Tran D, Wang H, Torresani L, et al., 2018. A closer look at spatiotemporal convolutions for action recognition. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.6450-6459. <https://doi.org/10.1109/CVPR.2018.00675>
- Vaswani A, Shazeer N, Parmar N, et al., 2017. Attention is all you need. Proc 31st Int Conf on Neural Information Processing Systems, p.6000-6010.
- Wang H, Liu F, Jiao LC, et al., 2024. ViLT-CLIP: video and language tuning clip with multimodal prompt learning and scenario-guided optimization. Proc 38th AAAI Conf on Artificial Intelligence, p.5390-5400. <https://doi.org/10.1609/aaai.v38i6.28347>
- Wang LM, Li W, Li W, et al., 2018. Appearance-and-relation networks for video classification. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.1430-1439. <https://doi.org/10.1109/CVPR.2018.00155>
- Wang LM, Tong Z, Ji B, et al., 2021. TDN: temporal difference networks for efficient action recognition. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.1895-1904. <https://doi.org/10.1109/CVPR46437.2021.00193>
- Wang MM, Xing JZ, Mei JB, et al., 2023. ActionCLIP: adapting language-image pretrained models for video action recognition. *IEEE Trans Neur Netw Learn Syst*,

- 36(1):625-637.
<https://doi.org/10.1109/TNNLS.2023.3331841>
- Wang R, Chen DD, Wu ZX, et al., 2023. Masked video distillation: rethinking masked feature modeling for self-supervised video representation learning. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.6312-6322.
<https://doi.org/10.1109/CVPR52729.2023.00611>
- Watson DS, O'Hara J, Tax N, et al., 2024. Explaining predictive uncertainty with information theoretic Shapley values. Proc 37th Int Conf on Neural Information Processing Systems, Article 320.
- Wu CY, Li YH, Mangalam K, et al., 2022. MeMViT: memory-augmented multiscale vision Transformer for efficient long-term video recognition. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.13577-13587.
<https://doi.org/10.1109/CVPR52688.2022.01322>
- Wu WH, He DL, Lin TW, et al., 2021. MVFNet: multi-view fusion network for efficient video recognition. Proc 35th AAAI Conf on Artificial Intelligence, p.2943-2951.
<https://doi.org/10.1609/aaai.v35i4.16401>
- Wu WH, Wang XH, Luo HP, et al., 2023. Bidirectional cross-modal knowledge exploration for video recognition with pre-trained vision-language models. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.6620-6630.
<https://doi.org/10.1109/CVPR52729.2023.00640>
- Wu WH, Sun Z, Song YX, et al., 2024. Transferring vision-language models for visual recognition: a classifier perspective. *Int J Comput Vis*, 132(2):392-409.
<https://doi.org/10.1007/s11263-023-01876-w>
- Xie SN, Sun C, Huang J, et al., 2018. Rethinking spatiotemporal feature learning: speed-accuracy trade-offs in video classification. Proc 15th European Conf on Computer Vision, p.318-335.
https://doi.org/10.1007/978-3-030-01267-0_19
- Yamazaki K, Vo K, Truong QS, et al., 2023. VLTinT: visual-linguistic Transformer-in-Transformer for coherent video paragraph captioning. Proc 37th AAAI Conf on Artificial Intelligence, p.3081-3090.
<https://doi.org/10.1609/aaai.v37i3.25412>
- Yu TS, Li YK, Li BX, 2020. RhyRNN: rhythmic RNN for recognizing events in long and complex videos. Proc 16th European Conf on Computer Vision, p.127-144.
https://doi.org/10.1007/978-3-030-58607-2_8
- Zhang J, Wan ZF, Hu LQ, et al., 2025. Collaboratively self-supervised video representation learning for action recognition. *IEEE Trans Inform Forens Secur*, 20:1895-1907. <https://doi.org/10.1109/TIFS.2025.3531772>
- Zheng ZW, Yang L, Wang YL, et al., 2024. Dynamic spatial focus for efficient compressed video action recognition. *IEEE Trans Circ Syst Video Technol*, 34(2):695-708.
<https://doi.org/10.1109/TCSVT.2023.3287201>
- Zhou BL, Andonian A, Oliva A, et al., 2018. Temporal relational reasoning in videos. Proc 15th European Conf on Computer Vision, p.831-846.
https://doi.org/10.1007/978-3-030-01246-5_49
- Zhou JM, Lin KY, Li HX, et al., 2021. Graph-based high-order relation modeling for long-term action recognition. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.8980-8989.
<https://doi.org/10.1109/CVPR46437.2021.00887>
- Zhou JM, Lin KY, Qiu YK, et al., 2024. TwinFormer: fine-to-coarse temporal modeling for long-term action recognition. *IEEE Trans Multim*, 26:2715-2728.
<https://doi.org/10.1109/TMM.2023.3302471>

List of supplementary materials

- 1 Proof of Theorem 1
- 2 Proof of Theorem 2