

Audio-Visual Glance Network for Efficient Video Recognition

Muhammad Adi Nugroho Sangmin Woo Sumin Lee Changick Kim
 Korea Advanced Institute of Science and Technology (KAIST)
 {madin, smwoo95, suminlee94, changick}@kaist.ac.kr

Abstract

Deep learning has made significant strides in video understanding tasks, but the computation required to classify lengthy and massive videos using clip-level video classifiers remains impractical and prohibitively expensive. To address this issue, we propose Audio-Visual Glance Network (AVGN), which leverages the commonly available audio and visual modalities to efficiently process the spatio-temporally important parts of a video. AVGN firstly divides the video into snippets of image-audio clip pair and employs lightweight unimodal encoders to extract global visual features and audio features. To identify the important temporal segments, we use an Audio-Visual Temporal Saliency Transformer (AV-TeST) that estimates the saliency scores of each frame. To further increase efficiency in the spatial dimension, AVGN processes only the important patches instead of the whole images. We use an Audio-Enhanced Spatial Patch Attention (AESPA) module to produce a set of enhanced coarse visual features, which are fed to a policy network that produces the coordinates of the important patches. This approach enables us to focus only on the most important spatio-temporally parts of the video, leading to more efficient video recognition. Moreover, we incorporate various training techniques and multi-modal feature fusion to enhance the robustness and effectiveness of our AVGN. By combining these strategies, our AVGN sets new state-of-the-art performance in multiple video recognition benchmarks while achieving faster processing speed.

1. Introduction

The exponential growth of diverse video contents, particularly those involving human-related actions, has prompted the development of deep learning algorithms [2, 7, 35, 48] that can effectively process and understand such data. Video recognition can bring benefits to multiple fields, such as sports for performance analysis [4], military for situational awareness [9], transportation for traffic monitoring [1], security for threat detection [32] and surveillance for public safety [5]. However, the current state-of-the-art methods of-

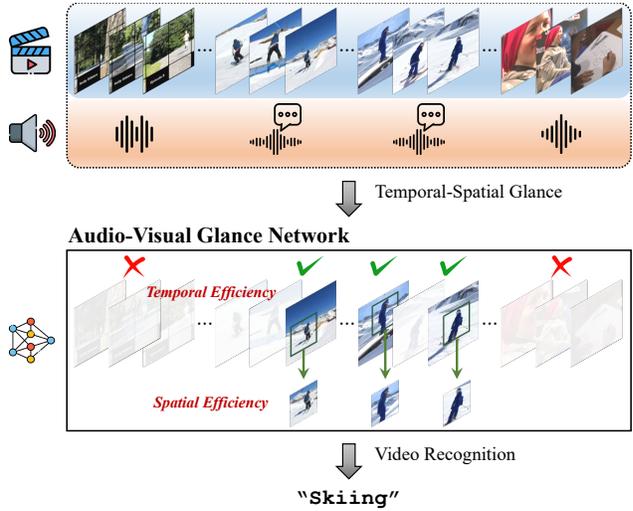


Figure 1. **Audio-Visual Glance Network (AVGN)** performs efficient video recognition by processing only a few highly salient frames based on the audio and visual information. Then, it determines and extracts the important spatial area in those frames to construct even more compact video representations.

ten require high computational costs, especially when dealing with lengthy and heavy videos, hindering their practical application in real-world scenarios. This has led to an increasing demand for efficient video understanding methods [34].

To address this issue, various approaches have been proposed, such as developing efficient and lightweight architecture [18], or adaptively selecting only the most informative subset of given videos [10, 46, 50], or training a policy network using policy gradient [25, 33]. On the other hand, some approaches manipulate the spatial resolution of input video frames to achieve efficiency, such as extracting only important spatial patches [39, 40] or adaptively changing the resolution frames based on the importance [25].

From all those approaches, we found that for efficient action recognition, we need to selectively locate and process only the important temporal and spatial location of the video. Intuitively, we need to know *when* and *where* to look at. For example in Fig. 1, to determine the action class “skiing”, we only need the frames that contain a person riding the ski, not the frames that only show the snowfield or the person

doing no action. This inspires us to create a network that can do efficient action recognition by *glancing* through the video to selectively find the important frames in a low-cost manner. Further, we aim to utilize the additional audio modality which is naturally available together with the visual modality in video recording, mimicking the human intuition of skimming through long sequences using both visual and audio cues to find important keyframes. Compared to the visual modality, the action-discriminative features of audio modality are easier to compute [34], and it also helps to distinguish actions that are visually similar [16].

We propose Audio-Visual Glance Network (AVGN), a comprehensive framework designed to enhance efficiency in both spatial and temporal dimensions. For temporal efficiency, AVGN aims to make the correct recognition of a video sequence with only a few important frames that actually contain distinctive cues. To achieve this, we construct an Audio-Visual Temporal Saliency Transformer (AV-TeST) that estimates the temporal saliency of a frame using coarse features generated by lightweight audio and visual backbones. Furthermore, to ensure spatial efficiency, we construct an Audio-Enhanced Spatial Patch Attention (AESPA) module that learns the relationship between audio feature sequences and visual features. This module generates audio-enhanced visual features that can be used by a patch extraction network to extract important spatial patches. For each frame, the patch contains only the important area of the image frame and has a lower pixel size compared to the original image. In addition to these modules, we also devise an appropriate feature fusion for classifier input and the training techniques to optimize these modules.

Our experimental results demonstrate that AVGN effectively incorporates audio and visual modalities for efficient action recognition. Our comparison with other state-of-the-art methods in Fig. 2 shows that AVGN achieves a higher mAP and lower FLOP cost on the ActivityNet dataset, indicating that it achieves pareto optimality. To conclude, our contributions are as follows:

- We show that incorporating audio modality into the video recognition process can lead to a pareto optimal solution, *i.e.*, improved accuracy without sacrificing efficiency.
- Our approach combines audio and visual information to improve temporal and spatial efficiency in a unified manner, and incorporates tailored training strategies that further optimize the performance of our AVGN.
- AVGN achieves state-of-the-art performance on multiple video recognition benchmarks as a result of the model building blocks and training techniques.

2. Related Work

Video recognition. Video recognition backbones are commonly used as part of the solution to the action recogni-

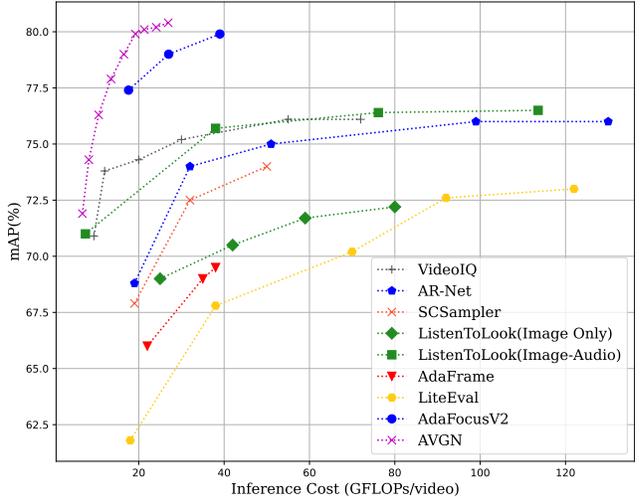


Figure 2. **Comparison of performance (mAP(%)) vs. cost (GFLOPs)** between AVGN and benchmark methods on the ActivityNet dataset [6]. AVGN achieves a pareto optimal by scoring the highest mAP, while also having lower GFLOPs.

tion task. Recent advances in deep-learning have led to the development of various techniques, such as C3D [35], I3D [2], and SlowFast [7], that can perform this task directly on trimmed video clips or snippets of untrimmed videos [37, 38]. These approaches have been extended to other related areas, such as action localization [30, 37], multi-modal action recognition [21, 34, 42], and efficient action recognition [17, 24, 43, 45]. In our work, we develop an efficient video recognition model that can achieve high accuracy on both trimmed and untrimmed videos while minimizing its computational cost.

Efficient video recognition. Temporal redundancy due to some frames being visually similar or containing irrelevant backgrounds leads to inefficient video recognition. To address this issue, various methods have been proposed [10, 11, 25, 33, 44, 46, 47, 49]. For example, temporal shift module [23] shifts feature maps along the temporal dimension to enable computationally-free temporal connections on top of 2D convolutions. AR-Net [25] consists of a policy network that decides which resolution to process a frame with, and multiple backbones with various resolutions. VideoIQ [33] contains lightweight policy network that can adjust the quantization precision of frames so that simpler frames are processed with lower precision, while Frame-Exit [10] uses an effective deterministic policy network and gating module to find the earliest exiting temporal point in a video sequence. MGSampler [51] and AdaFrame [46] use LSTMs to adaptively decide which frame to process next, while AdaFuse [26] dynamically chooses whether to reuse extracted features or fuse current and past features. The DSN [50] framework contains a sampling module and a policy maker to perform clip selection, forwarding only important frames to the classification module. In parallel direc-

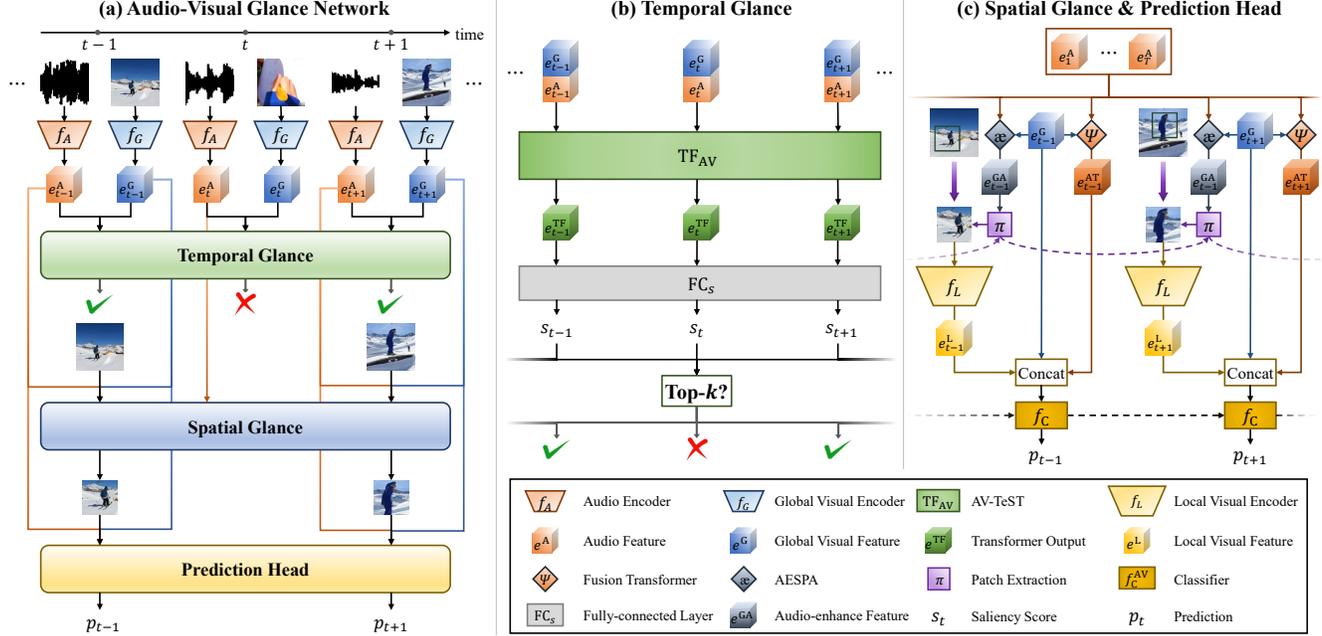


Figure 3. **Overview of AVGN.** (a) At each time step t , we extract the features e_t^A and e_t^G using the encoders f_A and f_G from the video. These features are then used for the temporal and spatial glance. (b) In temporal glance stage, Audio-Visual Temporal Saliency Transformer (AV-TeST) TF_{AV} generates features e_t^{TF} , which are then processed with FC_s to produce a temporal saliency score s_t . Only the top- k salient frames are passed through the spatial glance stage. (c) In spatial glance stage, Audio-Enhanced Spatial Patch Attention (AESPA) module \otimes (see Fig. 5 for more details) enhances the global visual features e_t^G with the sequence of audio features $\{e_1^A, \dots, e_T^A\}$, resulting in audio-enhanced visual features e_t^{GA} for the patch extraction network π to crop the most important areas of each frame. The patches are then fed to the local visual encoder f_L , which is heavier than f_G , to extract features e_t^L . Additionally, we use an audio fusion transformer ψ that selectively attends to the relevant parts of the audio feature sequence based on the global visual feature of the current time step e_t^G , producing \tilde{e}_t^A . Finally, we concatenate all the extracted features and feed them to the video classifier module f_c to output the prediction p_t .

tions, there are models that focus on reducing spatial redundancy. For example, AdaFocus [39, 40] uses a lightweight feature extractor and policy network to extract only the most important area or patch from an image, which is then fed into a heavier visual network with lower image resolution. AdaFocusV1 [39] uses reinforcement learning to train the policy network, while AdaFocusV2 [40] improves on this by utilizing a bilinear interpolation method to enable back-propagation along with their own training procedure. Our proposed network aims to reduce both temporal and spatial redundancy. For temporal redundancy reduction, we employ a multimodal transformer to assess and select only the most salient frames. Meanwhile, to address spatial redundancy we crop spatially only the most important areas for processing.

Audio in video understanding. In video understanding tasks, it has been demonstrated that incorporating additional modalities can improve performance beyond merely using visual modality alone [20, 28]. Audio is easy to acquire, requires low computational cost, and has distinctive characteristics. State-of-the-art video recognition models such as MM-ViT [3] have incorporated multiple modalities, utilizing cross-modal self-attention blocks, which can handle a large number of multimodal spatio-temporal tokens. Another approach, ListenToLook [8] efficiently integrates the audio

modality into an action recognition module by selectively sampling important frames and building a lightweight image-audio pair to mimic a more expensive clip-based model. SC-Sampler [19] combines both video and audio modalities by utilizing their saliency scores, which are obtained after processing with a lightweight clip sampler network. Our work also leverages the power of audio as an efficient modality for identifying important frames and areas in a video.

3. Audio-Visual Glance Network

The overview of AVGN is shown in Fig. 3. Given a video dataset $\mathcal{D} = \{(\mathbf{V}_n, \mathbf{A}_n, y_n)\}_{n=1}^N$ where each sample contains a sequence of video frames $\mathbf{V} = \{v_1, \dots, v_t, \dots, v_T\} \subset \mathbb{R}^{3 \times H_0 \times W_0}$ temporally paired with a sequence of audio spectrograms $\mathbf{A} = \{a_1, \dots, a_t, \dots, a_T\} \subset \mathbb{R}^{1 \times H_A \times W_A}$. Our goal is to correctly classify videos into their associated labels $y_n = \{0, 1\}^C \in \mathbb{R}^C$. Our model consists of an audio encoder f_A , a coarse global visual encoder f_G , a finer local level visual encoder f_L , an Audio-Visual Temporal Saliency Transformer (AV-TeST) module TF_{AV} , Audio-Enhanced Spatial Patch Attention (AESPA) module \otimes , audio fusion transformer ψ , and a spatial patch extraction network π .

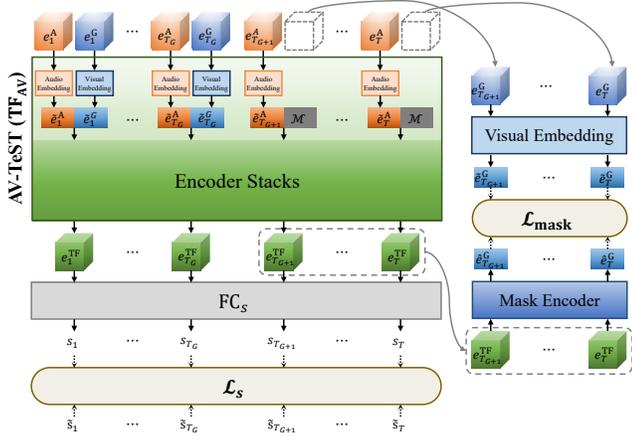


Figure 4. **AV-TeST components and auxiliary losses.** The saliency loss \mathcal{L}_s is used to train the model to predict the saliency score s_t that matches the confidence score \tilde{s}_t of each frame. The masked reconstruction loss $\mathcal{L}_{\text{mask}}$ is used to enhance the robustness of the AV-TeST by partially masking the visual features \hat{e}_t^G from frame index T_{G+1} up to the last frame T . We calculate the L2 loss between the embedded visual features \tilde{e}_t^G and the reconstructed visual features \hat{e}_t^G extracted by a mask encoder using the output features e_t^{TF} of AV-TeST.

3.1. Temporal Glance

The first stage of our network is to glance over the video sequence to find important temporal locations. The process starts with extracting audio and visual features using the lightweight encoders, f_A and f_G , respectively so that,

$$e_t^A = f_A(a_t), e_t^G = f_G(v_t), \quad (1)$$

where $e_t^G \in \mathbb{R}^{D_G \times H_G \times W_G}$, $e_t^A \in \mathbb{R}^{D_A \times H_A \times W_A}$. These encoders are designed to have low computational costs to prevent overburdening the glance process. We execute temporal glance using the AV-TeST which consists of modality embedding layers, and a stack of transformer encoders [36] (see Fig. 4). AV-TeST learns the temporal relationship among the audio-visual pairs, transforming them into audio-visual features ($\mathbf{e}_{1:T}^{\text{TF}}$). The process starts with averaging audio and visual features across their spatial dimension and then we use the embedding layers to produce audio tokens $\tilde{e}_t^A \in \mathbb{R}^{D_{AV}}$ and visual tokens $\tilde{e}_t^G \in \mathbb{R}^{D_{AV}}$. We then concatenate audio and global visual tokens from each time step and feed the resulting sequence into the Audio-Visual Temporal Saliency Transformer (AV-TeST) encoder module TF_{AV} .

$$\mathbf{e}_{1:T}^{\text{TF}} = \text{TF}_{AV}(\{[\tilde{e}_t^G, \tilde{e}_t^A], \dots, [\tilde{e}_t^G, \tilde{e}_t^A], \dots, [\tilde{e}_t^G, \tilde{e}_t^A]\}), \quad (2)$$

where $[\cdot, \cdot]$ denotes concatenation operation. Next, we pass each audio-visual feature at time t through a fully-connected layer FC_s to obtain a saliency score.

$$s_t = \text{FC}_s(e_t^{\text{TF}}). \quad (3)$$

Frames with high saliency scores are considered relevant for video recognition, while frames with low saliency scores are

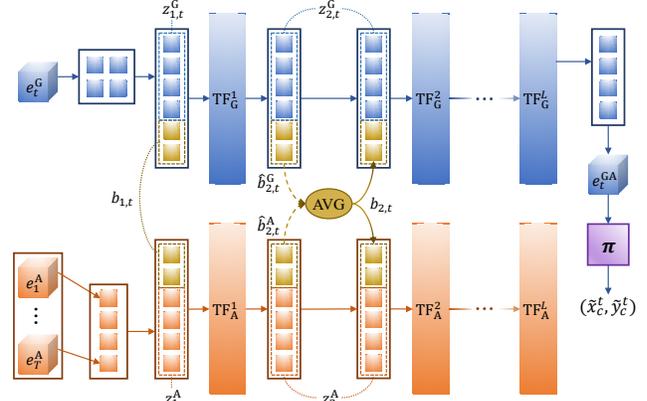


Figure 5. **The AESPA module** enhances visual features with a sequence of highly correlated audio features through a stack of transformers [36]. First, the global visual feature e_t^G is flattened, and a sequence of audio features $\{e_1^A, \dots, e_T^A\}$ is pooled, flattened, and stacked along the temporal axis. Bottleneck tokens $b_{1,t}$ are then appended. Next, these augmented features are passed through a series of transformer encoders TF. After each encoder, the bottleneck tokens are averaged and appended to the audio and visual tokens again. The output visual token of the last transformer layer is reshaped back to match its original shape, resulting in e_t^{GA} . Finally, the patch extraction network π utilizes e_t^{GA} to produce the visual patch center coordinates $(\tilde{x}_c^t, \tilde{y}_c^t)$.

considered non-relevant. We process only the k frames with the highest saliency scores for efficient inference.

3.2. Spatial Glance

In the second stage of AVGN, we aim to identify important spatial patches in the video frames. We employ a recurrent patch extraction network π inspired by [40], which utilizes a set of audio-enhanced visual tokens e_t^{GA} produced by our proposed Audio-Enhanced Spatial Patch Attention (AESPA) module (see Fig. 5). AESPA enhances the visual feature of each frame by the whole sequence of audio modality. It consists of two transformer [36] stacks: audio transformers TF_A and visual transformers TF_G .

Formally, we flatten the spatial dimensions of global visual feature e_t^G , resulting in a feature vector $z_{1,t}^G \in \mathbb{R}^{H_G \cdot W_G \times D_G}$. For a sequence of audio features $[e_1^A, \dots, e_T^A]$, we apply average pooling across the height and width dimension, then stack along the temporal axis to obtain the feature vector $z_1^A \in \mathbb{R}^{T \times D_A}$. These feature vectors serve as the initial inputs for the AESPA module. To efficiently exchange information between the two modalities we utilize a set of learnable shared bottleneck tokens $b_{l,t}$. The process inside l -th layer of AESPA can be formulated as follows:

$$\begin{aligned} z_{l+1,t}^A, \hat{b}_{l+1,t}^A &= \text{TF}_A^l([z_{l,t}^A, b_{l,t,t}]), \\ z_{l+1,t}^G, \hat{b}_{l+1,t}^G &= \text{TF}_G^l([z_{l,t}^G, b_{l,t,t}]), \\ b_{l+1,t} &= \text{AVG}(\hat{b}_{l+1,t}^A, \hat{b}_{l+1,t}^G), \end{aligned} \quad (4)$$

where TF denotes transformer encoder layer. The bottleneck

tokens are appended to the token sets of both modalities and processed together throughout L layers of transformer stacks. Finally, we reshape the output of visual transformer $z_{L+1,t}^G$ to obtain the audio-enhanced visual feature $e_t^{GA} \in \mathbb{R}^{D_A \times H_G \times W_G}$.

The enhanced visual feature e_t^{GA} is then fed into the patch extraction network π , which then produces the center coordinates $(\tilde{x}_c^t, \tilde{y}_c^t)$ of an important patch \tilde{v}_t on the image V . These center coordinates are continuous values and are obtained as:

$$(\tilde{x}_c^t, \tilde{y}_c^t) = \pi(\{e_1^{GA}, \dots, e_t^{GA}\}). \quad (5)$$

Then, we obtain the coordinates of each pixel in the patch $(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$ by adding a fixed offset o_{ij} to $(\tilde{x}_c^t, \tilde{y}_c^t)$.

$$(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t) = (\tilde{x}_c^t, \tilde{y}_c^t) + o_{ij}. \quad (6)$$

As the corresponding coordinates $(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$ have continuous values and need to be differentiable, we use bilinear interpolation [40] to obtain the patch pixel value from the four pixels surrounding the coordinate. Given the center coordinate of an important patch, we crop an original image to a patch of size $P \times P$. We process this image patch with a heavier local visual network f_L , which is larger in parameter size to extract fine-grained features of the patch.

$$e_t^L = f_L(\tilde{v}_t). \quad (7)$$

As the cropped patch is much smaller than the original image, the extraction process in network f_L requires significantly less cost than processing the original image.

3.3. Prediction Head

Lastly, we build a classifier module based on the fusion of the extracted features. Before the classifier, we use an audio fusion transformer ψ that has the same architecture as the standard transformer encoder [36]. It transforms the sequence of audio features $\{e_1^A, \dots, e_T^A\}$, which are used as key and value, using the global visual feature e_t^G as the query. The transformed audio feature at time t is denoted as e_t^{AT} . We concatenate the transformed audio feature with the global feature e^G and the local feature e^L . We then feed the resulting feature into our classifier module f_C^{AV} .

$$p_t = f_C^{AV}(\{[e_1^G, e_1^L, e_1^{AT}], \dots, [e_t^G, e_t^L, e_t^{AT}]\}). \quad (8)$$

At each time step t , the classifier f_C^{AV} generates a softmax prediction denoted as p_t . The classifier consists of fully-connected layers and aggregates features across the time steps using the *max* operation.

3.4. Training Techniques

Video classification loss. Firstly, we calculate the main loss of our network (\mathcal{L}_p), which is the cross-entropy loss of predicted output p_t for all values of t .

Auxiliary visual loss. This loss aims to better train the visual modality encoders. We pass the global visual feature e_t^G and the local visual feature e_t^L to separate fully-connected layers FC^G and FC^L , respectively. Additionally, we use a classifier f_C^V , which has the same structure as f_C^{AV} , and takes as input the concatenated global and local visual features $e_t^V = [e_t^G, e_t^L]$ extracted up until time step t . We calculate cross-entropy losses using these outputs.

$$\mathcal{L}_V = \frac{1}{T} \sum_{t=1}^T \left(\begin{array}{l} L_{CE}(FC^G(e_t^G), y) + \\ L_{CE}(FC^L(e_t^L), y) + \\ L_{CE}(f_C^V(\{e_1^V, \dots, e_t^V\}), y) \end{array} \right). \quad (9)$$

Auxiliary audio loss. Similarly, we also apply auxiliary loss on the audio modality. We use FC classifier FC^A and audio sequence classifier f_C^A then calculate the losses as,

$$\mathcal{L}_A = \frac{1}{T} \sum_{t=1}^T \left(\begin{array}{l} L_{CE}(FC^A(e_t^A), y) + \\ L_{CE}(f_C^A(\{e_1^A, \dots, e_t^A\}), y) \end{array} \right). \quad (10)$$

Masked visual token reconstruction. To enhance the robustness of the AV-TeST, we drop a subset of global visual tokens $\{\tilde{e}_{T_{G+1}}^G, \dots, \tilde{e}_T^G\}$, and train the model to recover the missing part using the remaining features. We first embed the remaining features, and concatenate audio embeddings $\{\tilde{e}_{T_{G+1}}^A, \dots, \tilde{e}_T^A\}$ with mask tokens \mathcal{M} to match the input dimension. These are then passed to the transformer encoders, producing $\{e_1^{TF}, \dots, e_T^{TF}\}$. Next, the partial outputs $\{e_{G+1}^{TF}, \dots, \tilde{e}_T^{TF}\}$ are fed to a mask encoder to reconstruct the missing global visual tokens, generating $\{\hat{e}_{T_{G+1}}^G, \dots, \hat{e}_T^G\}$. We calculate the L2 loss between the embedded visual tokens $\{\tilde{e}_{T_{G+1}}^G, \dots, \tilde{e}_T^G\}$ and the reconstructed visual tokens $\{\hat{e}_{T_{G+1}}^G, \dots, \hat{e}_T^G\}$ as $\mathcal{L}_{\text{mask}}$ (see Fig. 4). This loss enables AV-TeST to work robustly even with limited numbers of global visual tokens at the inference stage.

Saliency loss. In order to train the AV-TeST (TF_{AV}) to produce the saliency score without any ground truth frame importance, we generate pseudo labels \tilde{s}_t . We first obtain softmax predictions p_t' using the classifier f_C^{AV} with only the features at time step t , without using features from previous time steps or the *max* operation for feature aggregation.

$$p_t' = f_C^{AV}([e_t^G, e_t^L, e_t^{AT}]) \quad (11)$$

We then obtain a confidence score of each frame by normalizing p_t' with the maximum prediction logit across the classes and across the time steps. We minimize the difference between the predicted saliency scores s_t and confidence score \tilde{s}_t with L1 loss (see Fig. 4).

$$\mathcal{L}_s = L_1(s_t, \tilde{s}_t), \quad \tilde{s}_t = \frac{\max_c p_{c,t}'}{\max_t \max_c p_{c,t}'}, \quad (12)$$

Methods	Published on	Backbones	ActivityNet		FCVID		Mini-Kinetics	
			mAP \uparrow	GFLOPs \downarrow	mAP \uparrow	GFLOPs \downarrow	Acc. \uparrow	GFLOPs \downarrow
LiteEval [47]	NeurIPS'19	MN2+RN	72.7%	95.1	80.0%	94.3	61.0%	99.0
SCSampler [19]	ICCV'19	MN2+RN	72.9%	42.0	81.0%	42.0	70.8%	42.0
ListenToLook [8]	CVPR'20	MN2+RN	72.3%	81.4	–	–	–	–
ListenToLook [8]	CVPR'20	IA	76.4%	76.1	–	–	–	–
AR-Net [25]	ECCV'20	MN2+RN	73.8%	33.5	81.3%	35.1	71.7%	32.0
AdaFrame [46]	T-PAMI'21	MN2+RN	71.5%	79.0	80.2%	75.1	–	–
VideoIQ [33]	ICCV'21	MN2+RN	74.8%	28.1	82.7%	27.0	72.3%	20.4
OCSampler [22]	CVPR'22	MN2 [†] +RN	76.9%	21.7	82.7%	26.8	72.9%	17.5
AdaFocusV2 [40]	CVPR'22	MN2+RN	79.0%	27.0	<u>85.0%</u>	27.0	<u>75.4%</u>	27.0
AdaFocusV2 ⁺ [40]	CVPR'22	MN2+RN	74.8%	9.9	82.7%	<u>10.1</u>	72.3%	<u>6.3</u>
AdaFocusV3 [41]	ECCV'22	MN2+RN	<u>79.5%</u>	21.7	85.9%	26.8	75.0%	17.5
AdaFocusV3 ⁺ [41]	ECCV'22	MN2+RN	76.9%	10.9	82.7%	7.8	72.9%	8.6
AVGN ($k=14$)	–	MN2+RN	80.2%	24.4	84.1%	24.4	77.9%	25.0
AVGN* ($T_G = 4$)	–	MN2+RN	74.7%	9.6	82.2%	10.4	72.5%	6.4
AVGN ⁺ ($T_G = 2$)	–	MN2+RN	75.8%	<u>9.7</u>	82.4%	10.2	73.1%	6.1

Table 1. **Comparison of AVGN and baselines on three benchmark datasets.** MN2,IA,RN denotes MobileNet-V2,ImgAud[8], and ResNet respectively. The best two results are **bold-faced** and underlined, respectively. GFLOPs refer to the average computational cost for processing every videos. For AVGN*, we use $k = 5$ for ActivityNet and FCVID, $k = 3$ for Mini-Kinetics. \dagger represents the addition of TSM [23] and $+$ represents the addition of early exit [10]. \uparrow : The higher the better. \downarrow : The lower the better.

where $\hat{p}_{c,t}$ denotes the prediction logit value of class c at time step t . The denominator normalizes the pseudo labels, ensuring that the maximum value across all time steps and classes is 1.

Ordered AV logits loss. We reorder e^{AV} based on the saliency scores (s), followed by simulating the limited frame inference by applying Gumbell-Softmax sampling [13] on the saliency scores, which selects only a few frames from the entire sequence. We then compute the cross-entropy loss for the video classification using the reordered sequence.

$$\mathcal{L}_{\text{ord}} = L_{\text{CE}}(f_C^{\text{AV}}(\{e_{i_1}^{\text{AV}}, \dots, e_{i_T}^{\text{AV}}\}), y), \quad (13)$$

where i_t represents the temporal index of the frame with the t -th highest saliency score that is selected after the sampling.

We train AVGN by summing all the loss functions ($\mathcal{L}_p, \mathcal{L}_V, \mathcal{L}_A, \mathcal{L}_s, \mathcal{L}_{\text{mask}}, \mathcal{L}_{\text{ord}}$). To stabilize the training, we stop the gradient from π and TF_{AV} to f_G and f_A , and from ψ to f_G .

3.5. Inference

At the initial temporal glance stage, we can achieve efficiency by limiting global visual feature extractions up to only T_g frames, and instead, put mask tokens \mathcal{M} to replace the unextracted features.

$$e_t^{\text{AV}} = \begin{cases} [e_t^G, e_t^A] & \text{if } t \leq T_g, \\ [\mathcal{M}, e_t^A] & \text{otherwise.} \end{cases} \quad (14)$$

Then, for the spatial glance and the subsequent process we can process only k frames based on their saliency scores.

4. Experiments

4.1. Setup

Datasets. We used three large-scale video recognition datasets, *i.e.*, ActivityNet [6], FCVID [14], and Mini-Kinetics [15], and adopted their official training-validation splits. Following the common practice [8, 25, 26, 38, 39, 40, 50], we evaluated the performance of different methods via Top-1 accuracy (Top-1 Acc.) for Mini-Kinetics and mean average precision (mAP) for the other datasets.

Implementation details. For every video, we temporally sampled 16 pairs of image frames and audio clips. In the training stage, the frames are firstly grouped into 16 bins, and then one sample is randomly taken out from each bin. Then, we set the order of the frames with strategy following [10]. For the audio clips, we convert them to 1-channel audio-spectrograms of size 96×64 (960 msecs. and 64 frequency bins). We followed [39, 40] for the visual data pre-processing and augmentation, and performed time masking and frequency masking on the audio. For inference, we resized all frames to 256×256 and performed center-crop with size of 224×224 . The value of T_G is set to 12 unless stated otherwise.

4.2. Comparison with state-of-the-art

In Table 1, we compared AVGN with several competitive action recognition baselines [8, 19, 25, 26, 33, 40, 41, 46, 47] on multiple datasets: ActivityNet [6], FCVID [14], and Mini Kinetics [15]. We evaluated the model based on the mAP and computational cost measured in terms of floating point operations per second (FLOPS) averaged over the entire videos during the inference stage. We report the performance of AVGN in three settings: (i) AVGN with the maximum

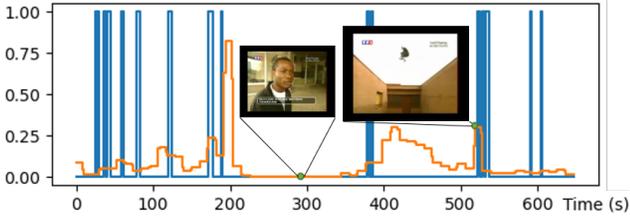


Figure 6. **Qualitative result of highlight detection on a long video.** We show **predicted** important and unimportant frames that correlate with the **ground truth** for the category `parkour`.

Method	VESD	LM-A	LM-S	MN	CHD	AVGN
Property	WS,ED	WS,ED	WS,ED	WS,ED	US	ZS, ED
mAP	0.423	0.524	0.563	0.698	0.527	0.557

Table 2. **Highlight detection results on the TVSum Dataset.** We report the top-5 mAP results along with their properties (WS = Weakly Spv., US = Unspv., ZS = Zero-shot, ED = External Data).

number of frames allowed for the prediction set to 14 ($k = 14$); (ii) AVGN with $T_G = 4$ and $k = 5$ for ActivityNet and FCVID, $k = 3$ for Mini-Kinetics; (iii) AVGN with an early exit [10], with the number of glances set to 2 ($T_G = 2$) for cost-effective inference.

We see that AVGN outperformed ListenToLook [8], which also utilizes the audio modality in its `ImgAud` backbones, with an mAP gain of 3.7% on ActivityNet using less than a third of GFLOPs. Additionally, AVGN performs competitively with AdaFocus variants and outperforms both in trimmed video recognition with better efficiency on Mini-Kinetics (73.1% mAP with 6.1 GFLOPs vs. 72.9% mAP with 8.6 GFLOPs on AdaFocusV3 [41]). AVGN achieved 1.2% mAP gain over AdaFocusV2 [40] while requiring 2.6 fewer GFLOPs on ActivityNet. On Mini-Kinetics, AVGN outperformed AdaFocusV2 by 2.5% accuracy with 2.6 fewer GFLOPs. Focusing only on the mAP metric, AVGN achieves the highest scores in ActivityNet and Mini-Kinetics. From the results of setting (ii) and (iii), we see the reliable performance of AVGN for low-cost inference, either by simple method of limiting the k or threshold method as [10]. AVGN⁺ also shows competitive results with significantly fewer GFLOPs than other methods. For example, it achieves comparable results with OCSampler [22] with half the GFLOPs on the FCVID dataset. It is worth noting that the slightly lower performance of AVGN on FCVID can be attributed to the fact that FCVID is a large-scale dataset that contains diverse video genres, whereas ActivityNet and Mini-Kinetics focus more on human action. Nonetheless, our experiments demonstrate that AVGN provides an efficient solution for video recognition tasks, surpassing existing state-of-the-art methods in terms of both accuracy and computational cost.

4.3. Long Videos Understanding

To demonstrate the effectiveness of AVGN on long videos, we use the ActivityNet-trained AVGN for the task of zero-

Number of glances (T_G)	$k=1$		$k=6$		$k=16$	
	mAP	GFLOPs	mAP	GFLOPs	mAP	GFLOPs
2	55.3%	2.7	75.6%	10.9	80.2%	27.6
4	56.4%	3.3	75.9%	11.1	80.1%	27.6
8	57.5%	4.7	76.8%	11.6	80.0%	27.6
12	57.7%	6.0	77.8%	12.8	80.2%	27.6
16	57.6%	7.3	78.3%	14.1	80.2%	27.6

Table 3. **Effect of the number of glances (T_G) on performance (mAP) and computation (FLOPs).** Experiments are conducted on ActivityNet [6] and k denotes the maximum number of image-audio pairs passed to the spatial glance stage.

shot Highlight Detection (HD). We use the TVSum[31] dataset, consisting of 2 to 10 minutes videos, with 4.2 minutes on average. The results in Table 2 and Figure 6 show AVGN performs competitively to models dedicatedly trained for HD, despite being trained for action recognition on another dataset. This highlights AVGN’s effectiveness in understanding long videos across datasets and tasks.

4.4. Ablation Studies

Effect of the number of glanced frames. To investigate the effect of the number of glances, we conducted experiments with varying numbers of glanced frames (T_G) and the maximum number of image-audio pairs (k) passed to the spatial glance stage. Table 3 presents the results. As shown in the table, increasing the number of glanced frames (T_G) leads to higher mAP due to more precise saliency score estimations. However, this comes at the expense of higher GFLOPs, indicating a trade-off between computational cost and accuracy. For instance, at $k = 6$, increasing T_G from 2 to 16 results in a mAP gain, from 75.6% to 78.3%, but at the same time, the GFLOPs value increases from 10.9 to 14.1. When all frames are passed through ($k = 16$) so that there is no temporal filtering, the difference in performance between different T_G values is minuscule. The results show that depending on the k , a moderate increase in T_G can result in significant gain in mAP, but beyond a certain point, the improvement in performance becomes marginal.

Contribution of AVGN components. Table 4 presents the results of ablation studies on the AVGN model. We first compared the performance of the model with and without AV-TeST (in Exp. 2 vs. Exp 1 and Exp. 3 vs. Exp. 4). The results show that AV-TeST improves the performance even when the number of frames for inference is limited and irrespective of whether is available or not. Exp. 1 vs. Exp. 3 demonstrates that incorporating audio modality can increase the mAP by around 2%. Exp. 4 vs. Exp. 5 shows that AESPA module boosts the performance about 1%. Exp. 5 vs. Exp. 6 indicates that ψ has a significant impact on the performance, especially when k is low, e.g., 4.6% mAP gain when $k = 1$. Comparing Exp. 3 to 6, the performance improves progressively with the addition of more components, especially when the audio modality is fully utilized. The

Exp.	Audio	AV-TeST	AESPA	ψ	$k=1$	$k=4$	$k=6$	$k=8$	$k=12$	$k=16$
1	✗	✗	✗	✗	46.6%	63.6%	69.0%	72.7%	75.6%	76.9%
2	✗	✓	✗	✗	50.1%	70.4%	73.0%	74.5%	76.0%	76.9%
3	✓	✗	✗	✗	48.6%	65.2%	70.9%	74.4%	77.1%	78.4%
4	✓	✓	✗	✗	51.5%	70.3%	73.9%	76.0%	78.1%	78.4%
5	✓	✓	✓	✗	52.1%	71.5%	75.1%	76.7%	79.0%	79.7%
6	✓	✓	✓	✓	57.7%	75.4%	77.8%	78.8%	79.9%	80.2%

Table 4. **Ablation study on AVGN key components.** We report the mAP results on ActivityNet [6]. We alternately add or remove audio modality, AV-TeST transformer (TF_{AV}), AESPA module (\ae), and the audio transformation network (ψ) at the classifier.

Temporal Sampling	$k = 4$		$k = 8$		$k = 12$	
	mAP	GFs	mAP	GFs	mAP	GFs
Uniform	67.3%	7.4	76.6%	14.2	78.8%	20.9
AV-TeST($T_G=12$)	75.4%	10.5	78.8%	16.0	79.9%	22.0

Table 5. **Ablation study on temporal glance.**

Spatial Sampling	$k = 1$	$k = 4$	$k = 8$	$k = 12$	$k = 16$
Center Crop	44.2%	59.3%	67.6%	70.0%	71.4%
Patch Network (π)	48.0%	62.2%	69.8%	71.9%	73.6%
AESPA + π	48.6%	62.9%	70.4%	72.8%	74.1%

Table 6. **Ablation study on spatial glance.**

\mathcal{L}_V	\mathcal{L}_U	$\mathcal{L}_{\text{mask}}$	\mathcal{L}_{ord}	$k = 1$	$k = 4$	$k = 8$	$k = 16$
✗	✓	✓	✓	56.9%	73.0%	76.5%	78.4%
✓	✗	✓	✓	54.6%	73.0%	77.2%	78.6%
✓	✓	✗	✓	55.0%	73.7%	77.6%	79.5%
✓	✓	✓	✗	57.0%	75.2%	78.8%	80.2%
✓	✓	✓	✓	57.7%	75.4%	78.8%	80.2%

Table 7. **Ablation study on training losses.** We report the mAP results on ActivityNet [6].

model achieves the highest accuracy of 80.2% in Exp. 6 when all components are used with $k = 16$. Comparing AV-TeST sampling with uniform temporal sampling in Table 5, AVGN achieves better mAP while not adding much GFLOPs (GFs). To assess the effectiveness of spatial glancing, we studied the predictive power of the extracted visual patch in the experiment presented in Table 6. We use local features (e_t^l) that have been aggregated through the temporal axis using a linear layer classifier for prediction. Improvements are achieved with the extraction patch network, especially when combined with our AESPA for audio.

Ablation study of losses. Table 7 shows the ablation experiments of training losses. First, we confirm the importance of the modality auxiliary losses, specifically the visual \mathcal{L}_V and audio \mathcal{L}_A losses. Removing either of these losses led to significant drop in mAP. The masking loss $\mathcal{L}_{\text{mask}}$ significantly improved the performance, especially when processing a lower number of k , e.g., 2.7% mAP gain when $k = 1$. Lastly, the ordered AV logits loss \mathcal{L}_{ord} contributed to a slight yet noteworthy improvement to the performance.

4.5. Qualitative result

We present the qualitative result of our model in the form of salient frames and their important spatial patches extracted from a video. In Fig. 7 (a), we show the first 8 frames directly sampled from the video, and in Fig. 7 (b), we display

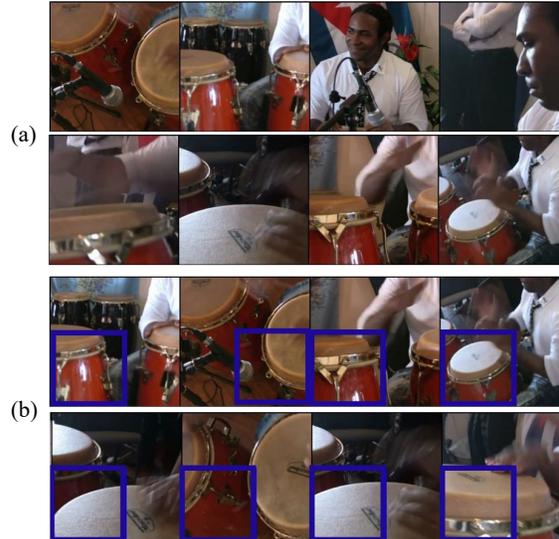


Figure 7. **Qualitative result.** Part (a) shows the original first 8 frames of the video and part (b) shows the Top-8 salient frames. The images are ordered from left to right and the bounding box on each frame highlights the important spatial area.

8 highest saliency score frames, with bounding boxes that highlight the important patches. The results demonstrate that our AVGN effectively identified the important part of the “playing conga” action class, which is the musical instrument. From a temporal perspective, the model prioritized frames where the conga is clearly visible.

5. Conclusion

We have proposed an efficient video recognition network called Audio-Visual Glance Network (AVGN) that selectively processes spatiotemporally important parts of a video. To improve network efficiency, we incorporated a cost-effective audio modality in addition to the visual modality. AVGN utilizes AV-TeST, a multimodal transformer that estimates salient frame saliency to achieve temporal efficiency. For spatial efficiency, we combine a recurrent patch extraction network and Audio Enhanced Spatial Patch Attention (AESPA) module to find important spatial patch using audio-enhanced coarse visual features. Our experiments have shown that AVGN performs competitively among state-of-the-art methods. Moreover, our ablation studies indicate that all of our model building blocks work collaboratively to contribute to model achievement.

6. Acknowledgement

This work was conducted by Center for Applied Research in Artificial Intelligence (CARAI) grant funded by DAPA and ADD (UD230017TD).

References

- [1] Mikhail Bortnikov, Adil Khan, Asad Masood Khattak, and Muhammad Ahmad. Accident recognition via 3d cnns for automated traffic monitoring in smart cities. In *Advances in Computer Vision: Proceedings of the 2019 Computer Vision Conference (CVC), Volume 2 1*, pages 256–264. Springer, 2020. 1
- [2] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733, Honolulu, HI, 2017. IEEE. 1, 2
- [3] Jiawei Chen and Chiu Man Ho. MM-ViT: Multi-Modal video transformer for compressed video action recognition. In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. IEEE, Jan. 2022. 3
- [4] Jun Chen, R Dinesh Jackson Samuel, and Parthasarathy Poovendran. LSTM with bio inspired algorithm for action recognition in sports videos. *Image Vis. Comput.*, 112:104214, Aug. 2021. 1
- [5] Ishan Dave, Zacchaeus Scheffer, Akash Kumar, Sarah Shiraz, Yogesh Singh Rawat, and Mubarak Shah. Gabriellav2: Towards better generalization in surveillance videos for action detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 122–132, 2022. 1
- [6] Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015. 2, 6, 7, 8, 12
- [7] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. SlowFast networks for video recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6201–6210, Seoul, Korea (South), 2019. IEEE. 1, 2
- [8] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. Listen to look: Action recognition by previewing audio. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) 2020*, pages 10457–10467, 2020. 3, 6, 7
- [9] Rúben Geraldes, Artur Gonçalves, Tin Lai, Mathias Villerebel, Wenlong Deng, Ana Salta, Kotaro Nakayama, Yutaka Matsuo, and Helmut Prendinger. UAV-Based situational awareness system using deep learning. *IEEE Access*, 7:122583–122594, 2019. 1
- [10] Amir Ghodrati, Babak Ehteshami Bejnordi, and Amirhossein Habibian. FrameExit: Conditional early exiting for efficient video recognition. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2021. 1, 2, 6, 7, 12
- [11] Shreyank N Gowda, Marcus Rohrbach, and Laura Sevilla-Lara. SMART frame selection for action recognition. *AAAI*, 35(2):1451–1459, May 2021. 2
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 12
- [13] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. 6
- [14] Yu-Gang Jiang, Zuxuan Wu, Jun Wang, Xiangyang Xue, and Shih-Fu Chang. Exploiting feature and class relationships in video categorization with regularized deep neural networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(2):352–364, Feb. 2018. 6, 12
- [15] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 6, 12
- [16] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. EPIC-Fusion: Audio-Visual temporal binding for egocentric action recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5491–5500, Seoul, Korea (South), 2019. IEEE. 2
- [17] Hanul Kim, Mihir Jain, Jun-Tae Lee, Sungrack Yun, and Fatih Porikli. Efficient action recognition via dynamic knowledge propagation. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2021. 2
- [18] Dan Kondratyuk, Liangzhe Yuan, Yandong Li, Li Zhang, Mingxing Tan, Matthew Brown, and Boqing Gong. MoViNets: Mobile video networks for efficient video recognition. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2021. 1
- [19] Bruno Korbar, Du Tran, and Lorenzo Torresani. SCSampler: Sampling salient clips from video for efficient action recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2019. 3, 6
- [20] Jun-Tae Lee, Mihir Jain, Hyoungwoo Park, and Sungrack Yun. Cross-Attentional Audio-Visual fusion for Weakly-Supervised action localization. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 3
- [21] Sumin Lee, Sangmin Woo, Yeonju Park, Muhammad Adi Nugroho, and Changick Kim. Modality mixer for multi-modal action recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3298–3307, 2023. 2
- [22] Jintao Lin, Haodong Duan, Kai Chen, Dahua Lin, and Limin Wang. OCSampler: Compressing videos to one clip with single-step sampling. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2022. 6, 7
- [23] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7083–7093. openaccess.thecvf.com, 2019. 2, 6
- [24] Zhaoyang Liu, Donghao Luo, Yabiao Wang, Limin Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Tong Lu. Teinet: Towards an efficient architecture for video recog-

- niton. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 11669–11676. AAAI Press, 2020. [2](#)
- [25] Yue Meng, Chung-Ching Lin, Rameswar Panda, Prasanna Sattigeri, Leonid Karlinsky, Aude Oliva, Kate Saenko, and Rogerio Feris. AR-net: Adaptive frame resolution for efficient action recognition. In *ECCV, Lecture notes in computer science*, pages 86–104. Springer International Publishing, Cham, 2020. [1](#), [2](#), [6](#), [12](#)
- [26] Yue Meng, Rameswar Panda, Chung-Ching Lin, Prasanna Sattigeri, Leonid Karlinsky, Kate Saenko, Aude Oliva, and Rogério Feris. AdaFuse: Adaptive temporal fusion network for efficient action recognition. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. [2](#), [6](#), [12](#)
- [27] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. [12](#)
- [28] Chiara Plizzari, Emanuele Alberti, and Barbara Caputo. Domain generalization through Audio-Visual relative norm alignment in first person action recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1807–1818, 2022. [3](#)
- [29] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4510–4520. IEEE, June 2018. [12](#)
- [30] Zheng Shou, Hang Gao, Lei Zhang, Kazuyuki Miyazawa, and Shih-Fu Chang. AutoLoc: Weakly-supervised temporal action localization in untrimmed videos. In *Computer Vision – ECCV 2018, Lecture notes in computer science*, pages 162–179. Springer International Publishing, Cham, 2018. [2](#)
- [31] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimés. Tvsum: Summarizing web videos using titles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5179–5187, 2015. [7](#)
- [32] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6479–6488, 2018. [1](#)
- [33] Ximeng Sun, Rameswar Panda, Chun-Fu Richard Chen, Aude Oliva, Rogerio Feris, and Kate Saenko. Dynamic network quantization for efficient video inference. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2021. [1](#), [2](#), [6](#)
- [34] Zehua Sun, Qihong Ke, Hossein Rahmani, Mohammed Benamoun, Gang Wang, and Jun Liu. Human action recognition from various data modalities: A review. *IEEE Trans. Pattern Anal. Mach. Intell.*, PP, June 2022. [1](#), [2](#)
- [35] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6450–6459, June 2018. [1](#), [2](#)
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv:1706.03762 [cs]*, Dec. 2017. [4](#), [5](#), [12](#)
- [37] Limin Wang, Yuanjun Xiong, Dahua Lin, and Luc Van Gool. UntrimmedNets for weakly supervised action recognition and detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, July 2017. [2](#)
- [38] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks for action recognition in videos. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(11):2740–2755, Nov. 2019. [2](#), [6](#)
- [39] Yulin Wang, Zhaoxi Chen, Haojun Jiang, Shiji Song, Yizeng Han, and Gao Huang. Adaptive focus for efficient video recognition. In *ICCV*. IEEE, Oct. 2021. [1](#), [3](#), [6](#)
- [40] Yulin Wang, Yang Yue, Yuanze Lin, Haojun Jiang, Zihang Lai, Victor Kulikov, Nikita Orlov, Humphrey Shi, and Gao Huang. AdaFocus v2: End-to-End training of spatial dynamic networks for video recognition. In *CVPR*, pages 20062–20072, 2022. [1](#), [3](#), [4](#), [5](#), [6](#), [7](#), [12](#)
- [41] Yulin Wang, Yang Yue, Xinhong Xu, Ali Hassani, Victor Kulikov, Nikita Orlov, Shiji Song, Humphrey Shi, and Gao Huang. AdaFocusV3: On unified Spatial-Temporal dynamic video recognition. In *Computer Vision – ECCV 2022*, pages 226–243. Springer Nature Switzerland, 2022. [6](#), [7](#)
- [42] Sangmin Woo, Sumin Lee, Yeonju Park, Muhammad Adi Nugroho, and Changick Kim. Towards good practices for missing modality robust action recognition. *AAAI*, 37(3):2776–2784, June 2023. [2](#)
- [43] Wenhao Wu, Dongliang He, Tianwei Lin, Fu Li, Chuang Gan, and Errui Ding. MVFNet: Multi-View fusion network for efficient video recognition. *AAAI*, 35(4):2943–2951, May 2021. [2](#)
- [44] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6222–6231. IEEE, Oct. 2019. [2](#)
- [45] Wenhao Wu, Yuxiang Zhao, Yanwu Xu, Xiao Tan, Dongliang He, Zhikang Zou, Jin Ye, Yingying Li, Mingde Yao, Zichao Dong, and Yifeng Shi. DSANet: Dynamic segment aggregation network for Video-Level representation learning. In *Proceedings of the 29th ACM International Conference on Multimedia, MM ’21*, pages 1903–1911, New York, NY, USA, Oct. 2021. Association for Computing Machinery. [2](#)
- [46] Zuxuan Wu, Hengduo Li, Caiming Xiong, Yu-Gang Jiang, and Larry S Davis. A dynamic frame selection framework for fast video recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(4):1699–1711, Apr. 2022. [1](#), [2](#), [6](#)
- [47] Zuxuan Wu, Caiming Xiong, Yu-Gang Jiang, and Larry S Davis. LiteEval: A Coarse-to-Fine framework for resource efficient video recognition. In Hanna M Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché Buc, Emily B Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7778–7787, 2019. [2](#), [6](#), [12](#)
- [48] Fanyi Xiao, Yong Jae Lee, Kristen Grauman, Jitendra Malik, and Christoph Feichtenhofer. Audiovisual SlowFast networks

- for video recognition. Technical Report arXiv:2001.08740, arXiv, Mar. 2020. [1](#)
- [49] Yeung, Russakovsky, Mori, and Fei-Fei. End-to-End learning of action detection from frame glimpses in videos. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 0, pages 2678–2687, June 2016. [2](#)
- [50] Yin-Dong Zheng, Zhaoyang Liu, Tong Lu, and Limin Wang. Dynamic sampling networks for efficient action recognition in videos. *IEEE Trans. Image Process.*, 29:7970–7983, 2020. [1](#), [2](#), [6](#)
- [51] Yuan Zhi, Zhan Tong, Limin Wang, and Gangshan Wu. MGSampler: An explainable sampling strategy for video action recognition. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2021. [2](#)

A. Notation List

For the convenience of the reader, we listed the Table of Notation containing frequently used notations along with their definition in Table 8.

B. Datasets

ActivityNet-1.3 [6] contains 10,024 training videos and 4,926 validation videos sorted into 200 human action categories. The average duration is 117 seconds.

FCVID [14] contains 45,611 videos for training and 45,612 videos for validation, which are annotated into 239 classes. The average duration is 167 seconds.

Mini-Kinetics is a subset of the Kinetics [15] dataset. We establish it following [10, 25, 26, 47]. The dataset include 200 classes of videos, 121k for training and 10k for validation. The average duration is around 10 seconds [15].

C. Implementation Detail

C.1. Network Architecture

Encoders. For audio encoder f_A and f_G we use MobileNetV2 [29] and for local visual encoder f_L we use ResNet-50 [12]. We use a patch size of 128×128 for the input to f_L , thus the size of the patch extracted by the patch extraction network is also the same. To encode a single image, the f_G requires 0.33 GFLOPs and f_L requires 1.35 GFLOPs, meanwhile to encode the whole audio sequence f_A requires 0.68 GFLOPs.

AV-TeST. In our implementation we construct TF_{AV} using a multi-head attention transformer [36] with 256 encoder dimension size, 2 stacks, and 4 heads. As the input to the transformer is concatenated audio-visual feature, for each modality we embed them to $128 - d$ vectors with separate linear embedding layers. To reconstruct the visual token, we utilize a transformer with the same architecture and append a linear embedding layer at the end.

AESPA. In our implementation of AESPA module, we use the same transformer architecture for both audio and visual modality. To minimize the computational burden, we reduce the incoming channel of both audio and visual modality to 256. Then we use the reduced feature maps as input to the bottleneck fusion transformers. Each modality transformer consists of 4 stack of encoder with 4 heads. We use 4 bottleneck tokens to be appended to the modality tokens.

Training Details. To train the network, we use an SGD optimizer with cosine learning rate annealing and a momentum of 0.9. The L2 regularization co-efficient is set to $1e-4$. The two encoders f_G and f_L are initialized using the ImageNet

pre-trained models¹, while the rest of the network is trained from random initialization. The size of the mini-batch is set to 24. The initial learning rates of f_G , f_A , f_L , f_C , π , ϖ , and TF_{AV} are set to 0.001, 0.001, 0.002, 0.01, $2e-4$, $2e-4$, and 0.01. We use a masking ratio of 0.75 for L_{mask} , and for Gumbell-Softmax we use 5 as the temperature value

C.2. Patch Extraction Network.

We explain in detail the process inside the spatial patch extraction network. To enable end-to-end training, we adopt the differentiable solution proposed in [40] to obtain \tilde{v}_t . Suppose that the size of the original frame v_t and the patch \tilde{v}_t is $H \times W$ and $P \times P$ ($P < H, W$), respectively². We assume that π outputs the continuous centre coordinates $(\tilde{x}_c^t, \tilde{y}_c^t)$ of \tilde{v}_t using audio-enhanced global visual feature up to t^{th} ($\{e_1^{\text{GA}}, \dots, e_t^{\text{GA}}\}$),

$$\begin{aligned} (\tilde{x}_c^t, \tilde{y}_c^t) &= \pi(\{e_1^{\text{GA}}, \dots, e_t^{\text{GA}}\}), \\ \tilde{x}_c^t \in [\frac{P}{2}, W - \frac{P}{2}], \quad \tilde{y}_c^t \in [\frac{P}{2}, H - \frac{P}{2}], \end{aligned} \quad (15)$$

We refer to the coordinates of the top-left corner of the frame as $(0, 0)$, and Eq. (15) ensures that \tilde{v}_t will never go outside of v_t .

The feed-forward process involves the bilinear interpolation method to enable backpropagation through $(\tilde{x}_c^t, \tilde{y}_c^t)$. As mentioned in the paper, the coordinates of a pixel in the patch \tilde{v}_t can be expressed as the addition of $(\tilde{x}_c^t, \tilde{y}_c^t)$ and a fixed offset:

$$\begin{aligned} (\tilde{x}_{ij}^t, \tilde{y}_{ij}^t) &= (\tilde{x}_c^t, \tilde{y}_c^t) + o_{ij}, \\ o_{ij} \in \left\{ -\frac{P}{2}, -\frac{P}{2} + 1, \dots, \frac{P}{2} \right\}. \end{aligned} \quad (16)$$

$(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$ denotes the corresponding horizontal and vertical coordinates in the original frame v_t to the i^{th} row and j^{th} column of \tilde{v}_t , while the offset o_{ij} is the vector from the patch center $(\tilde{x}_c^t, \tilde{y}_c^t)$ to this pixel. Given a fixed patch size, o_{ij} is a constant conditioned only on i, j , regardless of t or the inputs of π .

Since the values of $(\tilde{x}_c^t, \tilde{y}_c^t)$ are continuous, there does not exist a pixel of v_t exactly located at $(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$ to directly get the pixel value. Hence, we utilize the four adjacent pixels of $(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$ to obtain the pixel value using bilinear interpolation. We denote the four surrounding coordinates as $(\lfloor \tilde{x}_{ij}^t \rfloor, \lfloor \tilde{y}_{ij}^t \rfloor)$, $(\lfloor \tilde{x}_{ij}^t \rfloor + 1, \lfloor \tilde{y}_{ij}^t \rfloor)$, $(\lfloor \tilde{x}_{ij}^t \rfloor, \lfloor \tilde{y}_{ij}^t \rfloor + 1)$ and $(\lfloor \tilde{x}_{ij}^t \rfloor + 1, \lfloor \tilde{y}_{ij}^t \rfloor + 1)$, respectively, where $\lfloor \cdot \rfloor$ denotes the rounding-down operation. By assuming that the corresponding pixel values of these four pixels are $(m_{ij}^t)_{00}, (m_{ij}^t)_{01},$

¹In most cases, we use the 224x224 ImageNet pre-trained models provided by PyTorch [27].

²In our implementation, the height/width/coordinates are correspondingly normalized using the linear projection $[0, H] \rightarrow [0, 1]$ and $[0, W] \rightarrow [0, 1]$. Here we use the original values for the ease of understanding.

Variables		Functions	
Symbol	Definition	Symbol	Definition
t	Frame or time index	f_A	Audio encoder
a_t	Audio spectrogram clip at time	f_G	Global visual encoder
v_t	Input image frame at time t	f_L	Local visual encoder
y	label class	TF_{AV}	AV-TeST Transformer Network
e_t^A	Audio feature at time t	FC_s	Saliency score prediction head
e_t^G	Coarse/Global visual feature at time t	æ	Audio Enhanced Spatial Patch Attention (AESPA) module
$z_{l,t}^A$	AESPA audio vector at layer l at time t	TF_A^l	AESPA audio transformer at layer l
$z_{l,t}^G$	AESPA visual vector at layer l at time t	TF_G^l	AESPA visual transformer at layer l
e_t^{GA}	Enhanced Coarse/Global visual feature at time t	π	Spatial patch extraction network
e_t^L	Fine/Local visual feature at time t	ψ	Fusion transformer
e_t^{TF}	Audio-visual feature for AV-TeST input t	f_C^{AV}	Audio-visual classifier
s_t	Frame saliency score at time t	f_C^A	Auxiliary audio prediction head
\tilde{e}_t^A	Transformed audio feature at time t	FC^G	Auxiliary frame-wise global visual prediction head
$(\tilde{x}_c^t, \tilde{y}_c^t)$	Center coordinates t	FC^L	Auxiliary frame-wise local visual prediction head
\tilde{v}_t	Visual patch at time t	FC^A	Auxiliary frame-wise audio prediction head
$(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$	Coordinates of pixel patch t	f_C^V	Auxiliary visual prediction head
o_{ij}	Fixed offset for coordinate (i, j)	Hyperparameters	
\tilde{e}_t^G	AV-TeST embedded visual token (i, j)	Symbol	Definition
\hat{e}_t^G	Reconstructed AV-TeST embedded visual token	T_G	Visual temporal glance limit
p_t	Softmax prediction of f_C^{AV} with feature only at time t	k	Number of selected frames for prediction
\tilde{s}_t	Pseudo-label saliency score	P	Patch size
p_t	Class prediction		

Table 8. Table of Notation

$(m_{ij}^t)_{10}$, and $(m_{ij}^t)_{11}$, the pixel value at $(\tilde{x}_{ij}^t, \tilde{y}_{ij}^t)$ (referred to as \tilde{m}_{ij}^t) can be obtained via differentiable bilinear interpolation:

$$\begin{aligned} \tilde{m}_{ij}^t = & (m_{ij}^t)_{00}(\lfloor \tilde{x}_{ij}^t \rfloor - \tilde{x}_{ij}^t + 1)(\lfloor \tilde{y}_{ij}^t \rfloor - \tilde{y}_{ij}^t + 1) \\ & + (m_{ij}^t)_{01}(\tilde{x}_{ij}^t - \lfloor \tilde{x}_{ij}^t \rfloor)(\lfloor \tilde{y}_{ij}^t \rfloor - \tilde{y}_{ij}^t + 1) \\ & + (m_{ij}^t)_{10}(\lfloor \tilde{x}_{ij}^t \rfloor - \tilde{x}_{ij}^t + 1)(\tilde{y}_{ij}^t - \lfloor \tilde{y}_{ij}^t \rfloor) \\ & + (m_{ij}^t)_{11}(\tilde{x}_{ij}^t - \lfloor \tilde{x}_{ij}^t \rfloor)(\tilde{y}_{ij}^t - \lfloor \tilde{y}_{ij}^t \rfloor). \end{aligned} \quad (17)$$

Consequently, we can obtain the image patch \tilde{v}_t by traversing all possible i, j in Eq. (17).

Assume we have the training loss \mathcal{L} , we can compute the gradient $\partial \mathcal{L} / \partial \tilde{m}_{ij}^t$ with standard back-propagation. Following the chain rule, we have

$$\frac{\partial \mathcal{L}}{\partial \tilde{x}_c^t} = \sum_{i,j} \frac{\partial \mathcal{L}}{\partial \tilde{m}_{ij}^t} \frac{\partial \tilde{m}_{ij}^t}{\partial \tilde{x}_c^t}, \quad \frac{\partial \mathcal{L}}{\partial \tilde{y}_c^t} = \sum_{i,j} \frac{\partial \mathcal{L}}{\partial \tilde{m}_{ij}^t} \frac{\partial \tilde{m}_{ij}^t}{\partial \tilde{y}_c^t}. \quad (18)$$

Combining Eq. (16) and Eq. (18), we can further derive

$$\frac{\partial \tilde{m}_{ij}^t}{\partial \tilde{x}_c^t} = \frac{\partial \tilde{m}_{ij}^t}{\partial \tilde{x}_{ij}^t}, \quad \frac{\partial \tilde{m}_{ij}^t}{\partial \tilde{y}_c^t} = \frac{\partial \tilde{m}_{ij}^t}{\partial \tilde{y}_{ij}^t}. \quad (19)$$

Given that \tilde{x}_c^t and \tilde{y}_c^t are the outputs of the network π , the back-propagation process is able to proceed in an ordinary way.

D. Qualitative Results

We present more qualitative results in image format in Fig. 8 and in video format. Our qualitative results show how the model is able to estimate the salient frames and prioritize them over the non-relevant ones, *e.g.* in (c) salient frames are the ones containing ice hockey-related actions and in (c) and (a) frames with only text are non-salient. From the examples in video format, we observe how strong audio cues are present in the salient frames. For example, in “playing ten pins” class sample, the sound of the ball crashing the pins provide strong cues to estimate saliency.

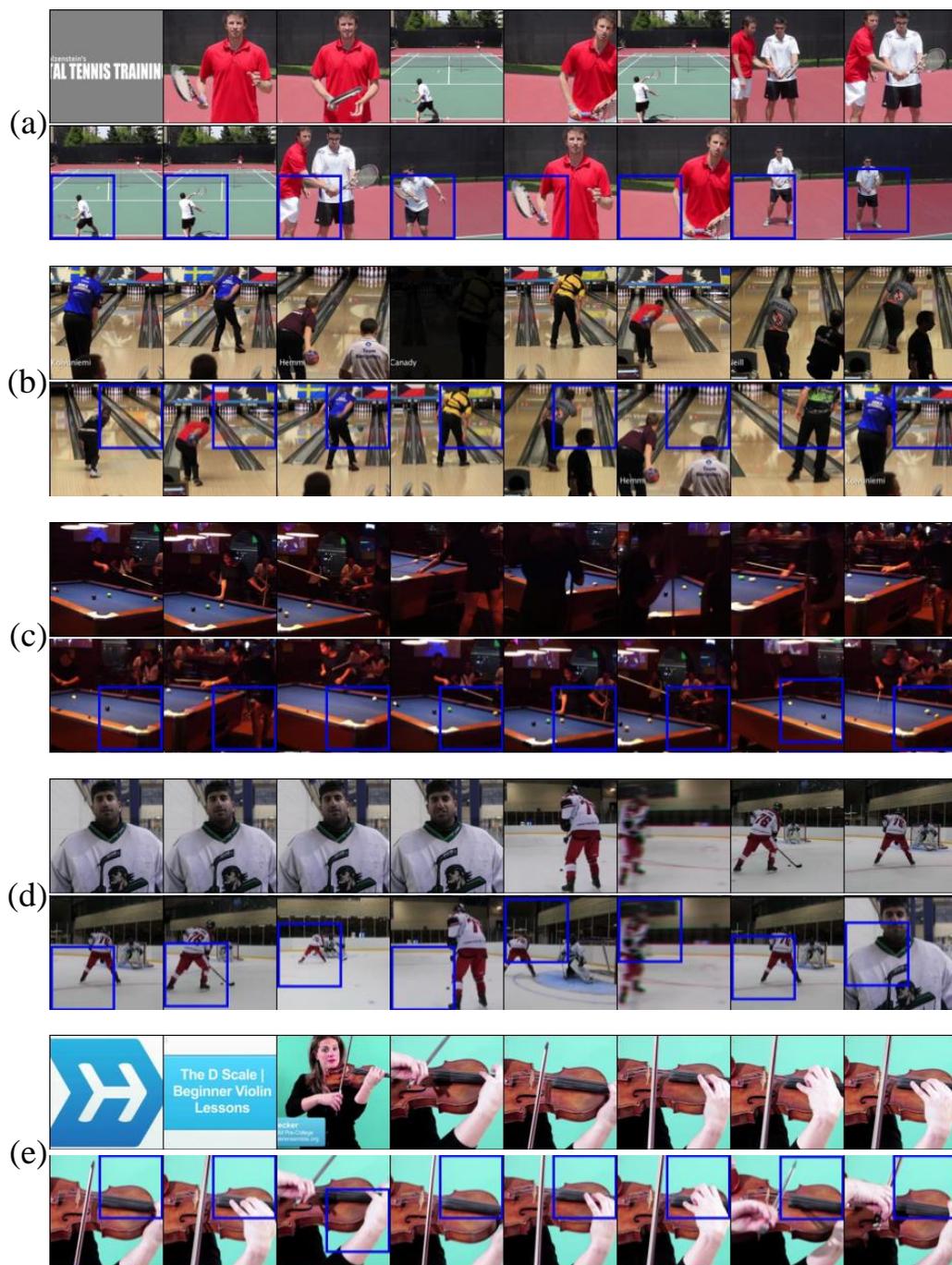


Figure 8. **Extended qualitative result** shows pair of the first 8 frames in original sequence and the Top-8 salient frames from classes (a) “tennis serve”, (b) “playing ten pins”, (c) “playing pool”, (d) “playing ice hockey”, and (e) “playing violin”. We also provide qualitative results in video format to better comprehend the effect of the audio.