

# Unsupervised Image-to-Video Adaptation via Category-aware Flow Memory Bank and Realistic Video Generation

Anonymous Author(s)

## ABSTRACT

Image-to-Video adaptation is proposed to train a model using labeled images and unlabeled videos to facilitate the classification of unlabeled videos. The latest work synthesizes videos using still images to mitigate the modality gap between images and videos. However, the synthesized videos are not realistic due to the camera movements are only simulated in 2D space. Therefore, we generate realistic videos by simulating arbitrary camera movements in 3D scenes, and then the model can be trained using the generated source videos. Unfortunately, the optical flows from the generated videos have unexpected negative impacts, resulting in suboptimal performance. To address this issue, we propose the Category-aware Flow Memory Bank, which replaces optical flows in source videos with real target flows, and the new composed videos are beneficial for training. In addition, we leverage the video pace prediction task to enhance the speed awareness of the model in order to solve the problem that the model performs poorly in handling some categories with similar appearances but significant speed differences. Our method achieves state-of-the-art performance and comparable performance on three Image-to-Video benchmarks.

## CCS CONCEPTS

• Computing methodologies → Activity recognition and understanding.

## KEYWORDS

Image-to-Video Adaptation, Category-aware Flow Memory Bank, Realistic Video Generation, Speed Awareness Enhancement, Action Recognition

## ACM Reference Format:

Anonymous Author(s). 2024. Unsupervised Image-to-Video Adaptation via Category-aware Flow Memory Bank and Realistic Video Generation. In *Proceedings of the 32nd ACM International Conference on Multimedia (MM '24)*, October 28–November 1, 2024, Melbourne, Australia. *Proceedings of the 32nd ACM International Conference on Multimedia (MM'24)*, October 28–November 1, 2024, Melbourne, Australia. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Video recognition is currently an active research direction in the field of multimedia due to its wide-ranging applications, such as

**Unpublished working draft. Not for distribution.**

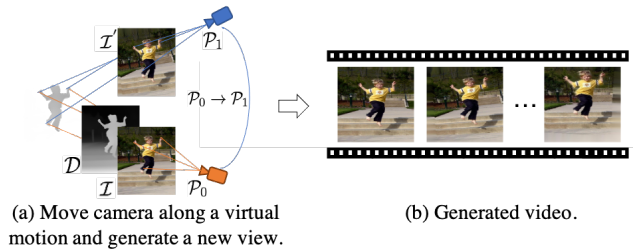
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted by ACM, provided that the copies are not made for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '24, October 28–November 1, 2024, Melbourne, Australia

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXX.XXXXXXX>

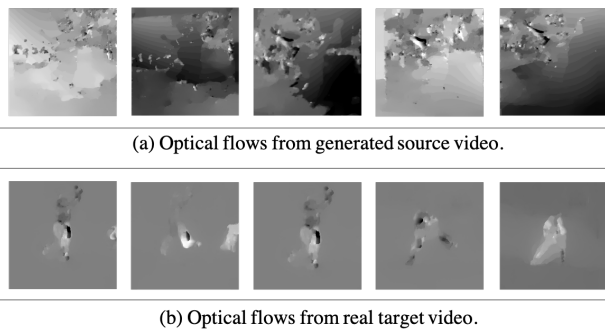


**Figure 1: Generating a video from a single source image is achieved through the Depthstillation [1] pipeline. Initially, we project the pixels in the input image  $I$  into 3D space, guided by the corresponding estimated depth map  $\mathcal{D}$ . Subsequently, we move the camera along a virtual motion path from  $P_0$  to  $P_1$ . Finally, this process yields a new view  $I'$ . By combining these synthesized views, we can construct a more realistic source video.**

video retrieval [11, 42], intelligent video surveillance [33, 46], and video captioning [28, 32]. However, training a high-performance video classifier requires collecting and annotating a large amount of video data, which is costly and time-consuming. As images are easier to annotate than videos, and there are numerous labeled image datasets accessible, image-to-video domain adaptation methods [4, 21, 22, 30, 51] that leverage the labeled images and unlabeled videos for training high performance video classifier appear as a challenging task and attract much attention.

The first challenge of image-to-video adaptation is the modality gap between images and videos. This gap refers to the fact that the temporal information in videos does not exist in source images. Bridging this modality gap is necessary for transferring knowledge from source domain to target domain effectively. Another challenge arises from domain discrepancy caused by variations in scenes, image styles and so on between source images and target video frames. Domain discrepancy is a key factor causing models trained in the source domain to perform poorly in the target domain. Overcoming these two key challenges of modality gap and domain discrepancy is crucial for achieving effective image-to-video adaptation.

Existing approaches [47, 48] predominantly employ a two-stage paradigm to address the challenges of domain discrepancy and modality gap. The first stage involves frame-level adaptation to reduce domain discrepancies between source images and target video frames. The second stage entails learning a spatio-temporal model to bridge the modality gap and incorporate temporal information. For example, Wei et al. [21] first employs DANN [9] for frame-level alignment and then leverages pseudo-labels from the first stage for self-supervised learning on the target videos. Recently, Zhuo et al. [51] proposes a single stage method ST-I2V by synthesizing source videos from static images with the help of



**Figure 2: We employ TV-L1 [49] to extract optical flows from both source and target videos for category ‘jump’, as shown in (a) and (b) respectively. It is evident that there is a significant discrepancy between the flows of the source and target videos. Specifically, the flows of the source video exhibit more interference and noise, while those of the target video appear cleaner.**

Grad-CAM [35], and solve the image-to-video adaptation problems with video-to-video domain adaptation methods.

Though being simple and effective, ST-I2V [51] randomly selects regions within image as intermediate frames to simulate camera displacement in 2D space which results in unrealistic synthesized video with improper temporal information. To address this issue, we rethink the imaging process in 3D space. As shown in Fig. 1, by recovering the position of camera, we can simulate the arbitrary movement of camera and generate realistic frame. Compared with ST-I2V [51], our method is simpler as it does not require training an additional classifier to locate major objects for an action. Besides, the complete original images are retained, avoiding any loss of crucial appearance information. The generated source video preserves static appearance and dynamic motion information, which is beneficial for training a discriminative spatio-temporal model.

Utilizing the generated source videos, we can train a simple but effective spatio-temporal baseline model through cross entropy loss with labels of source domain. However, we observe an unfavorable phenomenon that the optical flows extracted from generated videos are not helpful for training a discriminative model and even lead to suboptimal performance. As shown in Fig. 2, the optical flows from generated source videos exhibit more noise and interference compared with those of the target videos. Consequently, there is a noticeable distribution gap between the source and target videos, resulting in suboptimal performance.

To address the significant discrepancies between the flows in source and target videos, we construct a Category-aware Flow Memory Bank (CFB). The memory bank stores real flow data for each category within the real videos from target domain, where the category is determined by the pseudo label of target video. For a source video with ground-truth label  $\tilde{c}$ , we randomly select a flow of corresponding category  $c$  from the CFB. Then the selected flow is used to replace the original flow in the source video. As the new source video is more similar to the target video, the performance of the model is greatly improved.

Nevertheless, it is still difficult for the trained spatio-temporal model to distinguish categories with similar visual appearances but significant differences in speeds, such as ‘walk’ and ‘run’. So we leverage the video pace prediction task [41] to enhance the model’s perception of speed by altering video playback speeds. Specifically, we sample video clips at varying pace rates and treat the pace rates as labels. Subsequently, the spatio-temporal model is also trained with cross-entropy loss with video playback speed labels. Additionally, integrating video playback speed prediction task prevents the model overfitting on the source domain and enables the model to better generalize to target domain.

We validate our method on three widely used image-to-video adaptation benchmarks. The experimental results show that our method performs favorably against the current state-of-the-art approaches. We achieve the best-published results on the challenging E→H and B→U benchmarks, and competitive results on the S→U benchmark. Ablation studies are presented to verify the contribution of each key component in our approach.

In a nutshell, our contributions are as follows:

- To generate realistic source videos, we simulate camera movements in a 3D scene and capture new camera views that serve as source video frames, inspired by [1]. The generated source videos are very promising in training a discriminative spatio-temporal model.
- We propose a Category-aware Flow Memory Bank (CFB) to compensate the improper temporal information of the generated videos in source domain. By replacing the original flows of source videos with the retrieved flows from CFB for training, a remarkable improvement in the performance of model is achieved.
- We integrate video pace prediction task [41] to enhance the model’s perception of speed, which enables the model to distinguish categories with similar visual appearances but differences in speeds.
- Extensive experimental results show that our method achieves the best performance on the challenging E→H and B→U benchmarks and attains comparable results on the S→U benchmark.

## 2 RELATED WORK

**Image-to-video adaptation** methods focus on transferring knowledge from the image domain to the video domain. Existing unsupervised image-to-video adaptation tasks assume that only the labels of images are accessible, while the labels of the target videos are inaccessible. Mitigating the modality gap and reducing the distribution discrepancy are the primary objectives of image-to-video adaptation approaches [16, 20, 47]. For example, generative adversarial network (GAN) [12] is used to learn the mapping between image features and video features in HiGAN [48] and SymGAN [47]. The spatio-temporal causal graph [4] pursues similar goals. In order to mitigate the inherent modality gap between images and videos during domain adaptation, these methods leverage the strong generative modeling capabilities of GANs to transfer knowledge across modalities. CycDA [21] employs a four-stage method for adaptation. Class-agnostic alignment is performed in the first stage to derive pseudo-labels for training an independent spatio-temporal model in the second stage. The next two stages conduct iterative spatial alignment and spatio-temporal learning, with bidirectional knowledge transfer between the two components. Zhuo et al. [51] propose a new framework ST-I2V which synthesizes videos from

source static images, thereby converting the image-to-video adaptation task into video-to-video adaptation task. ST-I2V only simulates the transformation of the camera position in the 2D space, resulting in suboptimal performance. In contrast, we generate more realistic video frames by adjusting camera viewpoints at different positions in the 3D space. Additionally, we adopt both RGB and flow branches to construct our spatio-temporal model.

**Video-to-video adaptation.** Different from image-to-video adaptation, video-to-video adaptation methods are proposed to adapt labeled source videos to unlabeled target videos [25, 29], with their primary focus on addressing the challenges of domain alignment. Discrepancy-based methods are introduced to explicitly minimize domain discrepancies. For example, PTQ [3] minimizes the Maximum Mean Discrepancy (MMD) loss across both RGB and optical flow modalities to reduce the domain shift, resulting in improved performance. DVM [3] employs MixUp [5] to mitigate the domain-wise gap. This is achieved by progressively fusing the target videos with the source videos at the pixel-level, allowing for better alignment and adaptation between the domains. In our method, we propose a category-aware flow memory bank to replace the flow data in the generated videos in source domain, thereby reducing the domain gap.

**Video self-supervised learning** offers a promising annotation-free approach for representation learning in video domain. However, learning video representations is challenging due to temporal dynamics, motion, and other environmental factors. One key motivation behind defining pretext tasks is the idea that if a model can perform well on a complex task that requires a high-level understanding of video content, then it will learn more generalizable representations. For example, Jing et al. [15] and Wang et al. [4] design a task that applies appearance augmentations to video clips, and then the model is asked to classify the specific augmentation method that is applied. Fernando et al. [28] and Xu et al. [44] both design their approaches in a way that involves shuffling the order of video frames and having the model predict whether the video segment has been frame-shuffled. In order to enhance the model's awareness of speed and avoid overfitting problem in source domain, we introduce another simple video self-supervised learning method called video pace prediction [41].

### 3 METHODOLOGY

The goal of our method is to train a model that can achieve effective classification performance on target videos where the ground-truth annotations come from the labeled source image domain only. We train and evaluate our model in closed-set setting which means only the data from common categories in source and target domains will be used. It is supposed that there are a labeled source image domain  $\mathcal{I}_B = f^1 \mathcal{I}_B^0 \dots \mathcal{I}_B^{B-1}$  and an unlabeled target video domain  $\mathcal{I}_C = f^1 \mathcal{I}_C^0 \dots \mathcal{I}_C^{C-1}$ . Both domains contain the same classes.

Our overall framework is shown in Fig. 3. Since a large modality gap exists between images and videos, we first convert the image-to-video task into a video-to-video task. To simply construct a spatio-temporal model, we employ the I3D [1] network pretrained on the Kinetics dataset [7] as the backbone for both the RGB and flow branches. Following the instructions of I3D, we extract optical flows from generated source videos and real target videos. The

extracted flows are denoted by  $\mathcal{F}_B$  in source domain and  $\mathcal{F}_C$  in target domain. We address the issue of suboptimal performance caused by utilizing the original flows from the generated videos, with our proposed CFB. Furthermore, we enhance the model's ability to perceive speed by applying video pace prediction task [41].

#### 3.1 Source video generation

Although Zhuo et al. [51] have provided an effective method for source video generation, the generated videos may not be sufficiently realistic. This is because the approach only simulates the movement of the camera's viewpoint in 2D space, neglecting the fact that the body actions within the video should exist in 3D space.

To address this issue, we first generate source videos from source images via a virtual camera motion engine module, inspired by [34]. For a given source image, we employ MiDaS [14] to estimate the depth map  $D$ . The estimated  $D$  is then utilized to project pixels in  $I$  into 3D space based on the inverse intrinsic matrix  $K^{-1}$  (the intrinsic matrix is used to transform 3D world coordinates into 2D image coordinates captured by a camera). Assuming that captured by the camera at 3D location  $P_0$ , we apply an arbitrary virtual motion to the camera, moving it to a new position  $P_1$ . Specifically, we generate a rotation matrix  $R$  and a translation vector  $T$  by sampling a random triplet of Euler angles and a random 3D vector, respectively. The transformation matrix is then defined as  $T_{P_0 \rightarrow P_1} = {}^1R_0 T_{01}$  which is corresponding to the virtual motion path  $P_0 \rightarrow P_1$ . For each pixel with coordinate  $\theta$  in  $I$ , the coordinate  $\theta^0$  of its corresponding pixel in  $I^0$  acquired from the new viewpoint  $P_1$  can be obtained by:

$$\theta^0 = MT_{01}^{-1} D^{-1} T_{01}^{-1} \theta \quad (1)$$

where  $D^{-1} \theta$  is the depth value with coordinate  $\theta$ . Finally, the new image  $I^0$  is obtained through forward warping.

Following these steps, we alternately generate the subsequent video frame by utilizing the previously generated frame, thereby generating a source video  $\mathcal{F}_B$  with continuous action. Subsequently, we utilize the TV-L1 algorithm [49] to compute optical flow  $\mathcal{F}_B^i$  from the  $i$ -th source video  $\mathcal{F}_B^i$  and optical flow  $\mathcal{F}_C^i$  from the  $i$ -th target video  $\mathcal{F}_C^i$ , respectively. Therefore, the source domain is denoted by  $\mathcal{I}_B = f^1 \mathcal{I}_B^0 \dots \mathcal{I}_B^{B-1}$ , and the target domain is denoted by  $\mathcal{I}_C = f^1 \mathcal{I}_C^0 \dots \mathcal{I}_C^{C-1}$ .

Given the source video  $\mathcal{F}_B, \mathcal{F}_B^B, \dots, \mathcal{F}_B^B$ , we use the cross-entropy loss as the classification loss at the frame-level and video-level referring to I3D [3] and ST-I2V [51]. The frame-level classification entropy loss named local loss  $L_{>2}$  is defined as below,

$$L_{>2} = \frac{1}{B} \sum_{B=1}^B \sum_{C=1}^C \text{CE}(\mathcal{F}_B^B, \mathcal{F}_C^B) \quad (2)$$

Here,  $\text{CE}(\cdot, \cdot)$  represents the cross-entropy loss, and  $\mathcal{F}_B^B$  corresponds to the network's logits over RGB frames from a generated video  $\mathcal{F}_B^B$ . The variable  $B$  denotes the batch size.

Following the instructions in I3D [1], we train the RGB branch and the flow branch individually. So the video-level classification entropy loss  $L_2$  is defined as below,

$$L_2 = \frac{1}{B} \sum_{B=1}^B \text{CE}(\mathcal{F}_B^B, \mathcal{F}_C^B) \quad (3)$$

Figure 3: The overall framework of our approach begins by generating videos from static images for source domain. Subsequently, we replace the flow input of the generated source video with the retrieved flow data from the proposed category-aware memory bank. We sample a new video segment from the target video, for instance, by applying a 2x speedup (i.e. VP=2) from the original video clip. Then, the RGB frames and the optical flows are fed input into the RGB branch and flow branch separately. Finally, the representations from source domain are used to compute the cross-entropy losses for classification (i.e. CE Loss). And the representations from target domain are used to calculate the cross-entropy loss for video pace prediction task (i.e. VC Loss).

where  $\overline{\logit_{F_2^0}} \in \mathbb{R}^K$  is an average over logits of RGB frames from video  $v_2^B$  and  $\overline{\logit_{F_2^0}} \in \mathbb{R}^K$  is an average over logits of flows from video  $v_2^B$ .

During the inference phase, the class logits  $\overline{\logit_{F_2^0}}$  and  $\overline{\logit_{F_2^0}}$  obtained from both RGB and flow branches are normalized via softmax activation function. The normalized results are then summed together, yielding probabilities for each action category. Finally, the action category with the highest probability is selected as the predicted category for the current video sample.

### 3.2 Category-aware flow memory bank

As shown in Fig. 2, the optical flow frames extracted from generated source videos contain more interference in comparison to the cleaner optical flow frames present in target videos. This phenomenon indicates the significant distribution gap between the source and target domains.

To address the negative impacts of the original flows from source videos on the performance of model, we propose a Category-aware Flow memory Bank (CFB). As shown in Fig. 3, we train the model using the original source videos  $v_2^B$  and  $v_2^B$  during the warm-up epochs. After the warm-up phase, we utilize the trained model to assign pseudo label for each target video sample. Subsequently, we sort the samples based on the confidences of these pseudo-labels and retain only the top  $\#$  samples for each pseudo-category. In this manner, we construct a memory bank of size  $\#$ , where  $\#$  represents the number of categories. We denote the flow of category  $c$  as  $F_2^c$ .

Given the  $\#$ th generated video sample  $v_2^B = \{F_2^B, F_2^B, \dots, F_2^B\}$  in source domain, we replace the original flows  $F_2^B$  in  $v_2^B$  with a randomly selected flow  $F_2^c$ , resulting in a new video sample  $v_2^B = \{F_2^c, F_2^c, \dots, F_2^c\}$  where  $\#$  is equal to  $\#$ . Then the Eqn. (3) is modified as below,

$$L_2^0 = \frac{1}{\#} \sum_{c=1}^{\#} \text{CE}(\overline{\logit_{F_2^0}}, \overline{\logit_{F_2^c}}), \text{CE}(\overline{\logit_{F_2^0}}, \overline{\logit_{F_2^c}}) \quad (4)$$

where  $\overline{\logit_{F_2^c}} \in \mathbb{R}^K$  is an average over logits of  $K$  flows of selected flow  $F_2^c$ .

### 3.3 Speed awareness enhancement

Distinguishing categories with similar visual appearances but significant differences in speeds, such as 'walk' and 'run', poses a challenge for the model. To address this issue, we leverage video pace prediction task [41] to empower the model with the capability to perceive speed by altering video playback speeds. When provided with a video in its natural pace containing frames, we sample video segments  $v_2^B$  by various video pace rates. These pace rates correspond to labels from a predefined pace label space  $\mathcal{P}$ . For example, we produce three pace rate candidates {normal, fast, superfast}, where the corresponding pace labels are {1, 2, 3}, respectively. We randomly choose the starting frame overframes for each target video and then loop over the video at a regular interval  $\Delta$  until we obtain the desired number of frames for training.

With the video segments  $v_2^B$  which is sampled by pace rate  $p$ , the objective of pace prediction task is to understand the content of

the video segment and predict the correct pace rate. Subsequently, we utilize the video pace labels to train our model using the cross-entropy loss  $L_{E2}$  which is defined as follows:

$$L_{E2} = -\frac{1}{\#} \sum_{s=1}^{\#} \text{CE}(\mathbf{A}_s^{\text{pred}}, \mathbf{A}_s^{\text{gt}}) \quad (5)$$

where  $\mathbf{A}_s$  denotes the video pace label of the  $s$ th video,  $\mathbf{A}_s^{\text{pred}}$  is the predicted pace logits and  $\#$  denotes the size of pace label space.

Overall, all loss functions mentioned above form the complete objective:

$$L = L_2 + \lambda_1 L_{>2} + \lambda_2 L_{E2} \quad (6)$$

where  $\lambda_1$  and  $\lambda_2$  are trade-off parameters.

## 4 EXPERIMENTS

### 4.1 Datasets and setup

We evaluate our method through experiments on three standard image-to-video adaptation benchmarks: E H, B U and S U. In the case of the E H, we utilize the EADs [5] dataset, which comprises Stanford40 [5] and the HII dataset [37], as the source image domain, and HMDB51 [9] as the target video domain. There are 13 common classes between EADs and HMDB51 for image-to-video adaptation. The labeled source images and the unlabeled target videos are used to train a model. Regarding B, we employ the BU101 [26] dataset as the labeled source image domain and UCF101 [36] as the unlabeled target video domain. We use a total of 101 classes for the image-to-video adaptation task, as the classes in BU101 completely correspond to those in UCF101. For the S U benchmark, we substitute the source image domain from the B U benchmark with the Stanford40 [5] dataset. To perform the image-to-video adaptation task, the 12 common classes between Stanford40 and UCF101 are selected for training and evaluation.

### 4.2 Implementation details

To generate videos for source domain, we utilize MiD39 [1] whose backbone is pretrained BEiT-Large-5 [21] to extract depth maps from still images in source domain. Subsequently, we utilize Depth-stillation [1] to generate 16 video frames using the extracted depth maps and still images. The coefficient of translation vector is set to 0.01. Some generated frames can be found in Supplementary material. We use all 16 frames of the generated source videos during training following ST-I2V [51] for fair comparison.

For constructing a category-aware memory bank with high-quality pseudo-labels, we train the model for 10 epochs as warm-up phase. We then employ the model to assign pseudo-labels for target videos in each subsequent training epoch. We select the top 60 samples with the highest confidence for each category based on the confidences of the pseudo-labels and store their raw data in the memory bank. The influence of number of selected samples can be found in Parameter sensitivity analysis of subsection 4.4.

For building a spatio-temporal model, we use I3D model [6] with both RGB and flow branches pretrained on the Kinetics dataset [6]. We replace the last classifier layer with a fully connected layer that includes neurons. We freeze the first three Unit3D blocks following ST-I2V [51] to accelerate the training process.

We train the model with mini-batch stochastic gradient descent optimizer where the momentum and weight decay are set to 0.9 and 0.0001, respectively. The initial learning rates, batch sizes and total epochs are set to (0.05, 0.1, 0.015), (16, 32, 32) and (60, 30, 20) for E H, B U and S U benchmarks, respectively. We also adopt multistep decaying learning rate with a 0.1 decay rate where the milestones are half of the total epochs and the 2/3 of the total epochs. After the warm-up training phase, CFB is applied when the pseudo labels of target videos are more accurate. Following ST-I2V [51], the values of hyper-parameter  $\eta_1$  are set to (1, 20, 1) for E H, B U and S U benchmarks, respectively.

We set the size of video pace label space to 5 for all benchmarks which includes five video pace labels {1, 2, 3, 4, 5}. We randomly select a beginning frame for each target video and loop over the video at the generated video pace rate until the training video clip contains 16 frames. The generated video pace rate is treated as the pace label. For example, if there are 30 frames in the video, and loop over it starting from the 20th frame at 2speed, the indices of the selected frames are {20, 22, 24, 26, 28, 30, 2, 4, 6, 8, 10, 14, 16, 18, 20}. And the pace is regarded as pace label. During inference, we follow the approach of ST-I2V [51] and extract 32 frames uniformly from each target video for fair comparison. The values of hyper-parameter  $\eta_2$  are set to (0.2, 0.001, 0.01) for E, B U and S U benchmarks, respectively.

### 4.3 Competitors and results

In our experiments, we conduct comparative evaluations against several prevailing approaches: The DANN [6] pioneers domain adversarial training for classical image-level adaptation. JADA [52] reduces the image-level domain shift by aligning the joint distributions of multiple domain-specific layers. DAL27 [54] is another image-level domain adaptation method that introduces a novel domain adaptation layer to align source and target distributions with a reference distribution. MEDA39 [55] minimizes structural risk to train a domain-agnostic classifier on the Grassmann manifold and dynamically aligns the distributions of multiple domains while evaluating the significance of marginal and conditional distributions. HiGAN [48] and SymGAN [47] attempt at bridging the modal gap by mapping image embeddings to video space using GAN [56]. DANN+I3D baseline leverages DANN adapted image features to train an I3D architecture with pseudo-labels which is implemented by Lin et al. [21]. CycDA [21] is a four-stage method that reduces domain discrepancies by using both class-agnostic and class-aware domain alignment techniques. It also utilizes pseudo labels to train a I3D model, effectively bridging the modality gap. ST-I2V [51] is the recent state-of-the-art approach, which employs Grad-CAM [56] and an additional classifier to generate source videos. It transforms the image-to-video domain adaptation task into a video-to-video domain adaptation task. Additionally, for reference, we include the lower bound (SO (Img), where SO stands for Source-only.) and the upper bound (ground truth supervised target) from work [51].

In Tab. 1, we present comparison results. Our approach achieves new state-of-the-art performances on the E H and B U benchmarks and demonstrates comparable result on the S U benchmark. Specifically, our method outperforms ST-I2V [51] by 6.1% and 2.2% on the E H and B U benchmarks, respectively. It's important to

465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522

523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580

Table 1: Results on E! H, B! U and S! U, averaged over 3 random trials.

method	E! H	B! U	S! U
SO (Img)	37.2	54.8	76.8
DANN [9]	39.6	55.3	80.3
JAN [24]	40.9	-	91.4
HiGAN [48]	44.6	-	95.4
DAL [27]	45.5	-	97.6
MEDA [39]	43.1	-	94.3
SymGAN [47]	55.0	-	97.7
DANN+I3D	53.8	68.3	97.9
CycDA [21]	62.0	72.6	99.1
ST-I2V [51]	71.3	78.9	98.6
Ours	77.4	81.1	97.3
supervised target	83.2	93.1	99.3

Table 2: Ablation study results on E! H, B! U, and S! U, averaged over 3 random trials.

method	E! H	B! U	S! U
SO ([51])	59.0	60.2	96.3
SO (RGB)	59.8	76.4	96.6
SO (RGB+ ow)	60.7	74.6	96.3
SO (RGB+ ow) + CFB	74.1	80.2	97.2
Full Model	77.4	81.1	97.3

note that the E! H benchmark is considerably more challenging than B! U and S! U, given the difficulties in distinguishing categories within the HMDB51 dataset. Nevertheless, the performance of our model on the S! U benchmark still lags behind the current state-of-the-art method. This gap may arise from the fact that our approach primarily focuses on enhancing temporal information, while S! U benchmark relies less on temporal information which can be verified from the superior performance of SO (RGB) in Tab. 2. On the other hand, we can adopt some existing domain adaptation techniques like BNM [7] and MCC [14] to further improve our model on S! U benchmark, achieving new state-of-the-art result, as shown in Tab. 4.

The experimental results indicate the effectiveness of our source video generation method, the proposed category-aware flow memory bank and the speed awareness enhancement approach, in the context of image-to-video domain adaptation learning.

#### 4.4 Ablation study

To study the contribution of each component in our approach towards the overall performance, we conduct the ablation study of our proposed approach. We evaluate the following variants of model: (1) SO (RGB) which denotes the model that contains RGB branch only and is trained with the labeled generated source videos. (2) SO (RGB+ ow), the model is trained simultaneously using both the RGB branch and the ow branch with source data. (3) (RGB+ ow) + CFB, is the model trained by replacing the original

ow data in source generated video samples with the retrieved ow data from CFB after the warm-up phase. (4) Full Model, that is trained with incorporating VPT (video pace prediction task) into the baseline model SO (RGB+ ow) + CFB. Additionally, we include the performance of SO ([51]) for reference, which is trained with synthesized videos generated by Zhuo et al. [51].

The ablation study results are shown in Tab. 2. Comparing with SO ([51]), we can observe that our model SO (RGB), when using only the RGB branch (the same as SO ([51])), has brought improvements of 0.8%, 16.2% and 0.3% on E! H, B! U and S! U, respectively. This means that our video generation approach is more effective to learn a discriminative spatio-temporal model against ST-I2V. When training the ow branch using the original ow data from source videos, it still provides some improvements on the E! H benchmark but has negative effects on B! U and S! U benchmarks. This phenomenon demonstrates the improper ow from source video is one of the key factors contributing to the poor performance of spatio-temporal models.

Our proposed Category-aware Flow Memory Bank (CFB) aims to address this issue. The results of SO (RGB+ ow) + CFB show that our approach brings performance improvements across three benchmarks, is specially with gains of 13.4% and 5.6% in E! H and B! U benchmarks, respectively. After enhancing the model's perception of speed by introducing VPT, the performance of Full Model is further improved, validating the effectiveness of our approach. Parameter sensitivity analysis. To evaluate the parameter sensitivity, we conduct a series of experiments on the E! H benchmark. Fig. 4 reports the results of parameter sensitivity analysis, and more results can be found in Supplementary material.

The weight<sub>1</sub> is one of the key factors that influences the performance of the model. We evaluate the impact of different values of weight<sub>1</sub> using the baseline SO (RGB) and present the results in Fig. 4 (a). As observed, a small weight<sub>1</sub> leads to the model ignoring appearance characteristics of video frames, while a large weight<sub>1</sub> results in the model excessively focusing on appearance characteristics at the expense of temporal features. This further hinders the model's ability to comprehend the content of the video.

An appropriate dimension of CFB is a critical factor in determining the performance of model. So we investigate the impact on performance by setting different numbers of ow samples for each category in baseline SO (RGB+ ow) + CFB and the results are shown in Fig. 4 (b). It is evident that the value of n can significantly influence the performance of model. An excessively large n may introduce noisy ow samples with low-quality pseudo labels, which leads to inferior performance. Conversely, a small n may lead to insufficient model generalization.

To investigate the proper settings for VPT, we explore different video pace prediction loss weights<sub>2</sub> and sizes of video pace label space E<sub>2</sub> based on baseline SO (RGB+ ow) and the results are shown in Fig. 4 (c) and Fig. 4 (d) respectively. From the above experimental results, it is observed that employing video pace prediction task on SO (RGB+ ow) can achieve competitive results under an appropriate range of weight<sub>2</sub> values. Additionally, the small size of the video pace label space restricts the model's ability to perceive speed. On the other hand, an over-sized video pace label space increases the difficulty of the speed prediction task, making it challenging for the model to learn meaningful semantic representations.

(a) Acc. of SO (RGB) w.r.t.  $\lambda$ . (b) Acc. of SO (RGB+ ow) + CFB w.r.t.  $\beta$ . (c) Acc. of SO (RGB+ ow) + VPT w.r.t.  $\alpha$ . (d) Acc. of SO (RGB+ ow) + VPT w.r.t.  $\gamma$ .

Figure 4: The plots of parameter sensitivity analysis. We obtain the results on E! H benchmark.

(a) SO (RGB+ ow) (b) SO (RGB+ ow) + CFB

Figure 5: t-SNE visualizations of video representations (colored w.r.t. ground truth) from source and target domain in E! H benchmarks. We plot the representations of SO (RGB+ ow) (a) and the representations of SO (RGB+ ow) + CFB (b). We use 13 different colors to represent each category. We use 'o' and '+' to represent source representations and target representations respectively.

Table 3: Accuracies of different aggregation methods on E! H, averaged over 3 random trials.

method	E! H
SO (RGB+ ow)	60.7
Mean 60	59.0
SimW Sum60	59.6
SimW Mean top5	63.8
SimW Sum top5	64.0
SimW top1	72.5
SimW Random 1/top5	73.2
Random 1/60	74.1

Different aggregation methods for retrieved ows. Based on SO (RGB+ ow) + CFB, we conduct a study to investigate the impact of different retrieved ow aggregation methods on the performance of model. We conduct these experiments on E! H benchmark. The results are presented in Tab. 3. After constructing a CFB which stores top 60 ows with the highest confidence for each category,

we first investigate the impact of aggregating all retrieved ows using mean pooling (denoted by Mean 60) on the performance of model. Next, we design various aggregation methods that automate the selection of retrieved ows through the cosine similarity of source and target RGB features. The cosine similarity is considered as the weight of each retrieved ow data. And we represent these methods using 'SimW' as the prefix. We evaluate the following methods: (1) SimW Sum60, means that we perform a weighted summation of all retrieved ows by using the weights assigned to each ow. (2) SimW Mean top5, we perform mean pooling on the top 5 retrieved ows with the highest weights. And then replace the source ow with the pooled ow. (3) SimW Sum top5, which denotes that we perform weighted summation of top 5 retrieved ows with the highest weights. (4) SimW top1, which means that we only choose the retrieved ow which own the highest weight. (5) SimW Random 1/top5, we randomly select one ow from the top 5 retrieved ows with the highest weights. At last, Random 1/60 represents the same aggregation setting as SO (RGB+ ow) + CFB in Tab. 2, which is the one we used in the manuscript. For reference, we also include the result of SO (RGB+ ow).

Table 4: Accuracies on E! H and S! U after combining with DA Method, averaged over 3 random trials.

method	E! H	S! U
SO (RGB+ ow) + CFB	74.1	97.2
SO (RGB+ ow) + CFB + DAN	74.7	97.3
SO (RGB+ ow) + CFB + MCC	75.0	99.2
SO (RGB+ ow) + CFB + BNM	74.5	99.3

The results indicate that regardless of using mean pooling or weighted summation to aggregate retrieved ows, the constructed ows lead to suboptimal results or even have significant negative impacts on the performance of model. This could be due to that the aggregated ows lose too much information and dissimilar to real ow data. With less (top5) ows for aggregation, the performance is improved. With only top1 ow, the performance is further improved. So we only use 1 ow without aggregation. We use Random 1/60 as we think that randomly choosing 1 ow may boost the robustness of the model and the experimental result verifies our conjecture. So we use Random 1/60 in our manuscript.

#### 4.5 Further remarks

**Integrating domain adaptation techniques.** We employ several typical domain adaptation techniques into the constructed baseline SO (RGB+ ow) + CFB, including DAN [3], MCC [14] and BNM [7]. Specifically, we use  $L_0 = L_2$ ,  $L_{1>2}$ ,  $L_{C5}$  to train the model.  $L_{C5}$  is the transfer loss like MMD [23] loss, and BNM [7] loss. The values of hyper-parameters are set to (0.05, 0.3) for EH and S U, respectively. The results on EH and S U are shown in Tab. 4 and we also report the results of SO (RGB+ ow) + CFB for better demonstration.

It is observed that the performance of our model can still be greatly improved by applying typical domain adaptation methods on our constructed video-to-video domain adaptation baseline which involves generating source videos and utilizing CFB. With recent state-of-the-art domain adaptation techniques MCC [14] and BNM [7], our method outperforms CycDA [1] on S U benchmark, achieving new state-of-the-art performance. In addition, the results also demonstrate that the combination of our method with DA techniques shows good adaptability on EH. **Compatibility with CLIP.** We replace the network of RGB branch in sec. 4.2 with visual encoder of CLIP [1] whose backbone architecture is ViT-B/32. Different from the experimental hyperparameter settings described in sec. 4.2, we set total epochs to (40, 20) and the number of warm-up training phase to (10, 5) for EH and S U benchmarks respectively. And setting the batch size to 8, the learning rate of CLIP to  $5e-5$  and the prompt for text encoder to `a video of a person {}.` for both benchmarks. {} in prompt indicates category names like `climb`, `run` and so on.

For reference, we include the results of our approach employing the I3D network in RGB backbone. The results shown in Tab. 5 indicates that the performance improvements can be attributed to the robust representation learning capability and the powerful knowledge base of CLIP itself. By integrating CLIP into our method,

Table 5: Accuracies on E! H and S! U after integrating with CLIP, averaged over 3 random trials.

RGB backbone	method	E! H	S! U
I3D	SO (RGB)	59.8	96.6
	Full Model	77.4	97.3
CLIP	SO (RGB)	66.8	97.7
	Full Model	78.0	98.3

we can leverage its knowledge and capabilities to enhance the performance. Additionally, the compatibility between our method and CLIP is crucial for achieving further performance gains. The ability of our method to effectively incorporate CLIP's features and merge them with the existing framework allows for a synergistic effect, especially resulting in improved performance on EH benchmark. Our approach's compatibility with prevailing large multimodal models like CLIP showcases its strength and demonstrates its ability to achieve better results.

**T-SNE visualization.** We visualize the representations of the baseline models, SO (RGB+ ow) and SO (RGB+ ow) + CFB, in Fig. 5 using t-SNE [38]. We project the source and target videos of EH benchmark into 2-dimensional representations. Intuitively, there are more confusions among the representations from SO (RGB+ ow) (Fig. 5 (a)), which can be improved by incorporating CFB, as shown in Fig. 5 (b). Specifically, for categories like `pour`, `kick`, and `push`, the incorporation of CFB enables the model to learn more discriminative representations, enabling better differentiation from others.

Moreover, for categories of `talk`, `smoke`, `climb` and `wave`, the representations from both the source and target domains become closer after replacing the ows retrieved from CFB, as depicted in Fig. 5 (b). It is suggested that our proposed CFB further reduces the distribution discrepancies between source and target domains, which is beneficial for training a model with great generalization.

## 5 CONCLUSION

We overcome the challenges of image-to-video domain adaptation task, aiming to enhance the spatio-temporal model's discriminative ability for unlabeled video classification in the target domain. To mitigate the modality gap between labeled source images and unlabeled target videos, we generate realistic source videos by simulating diverse camera movements in 3D scenes and the new perspectives are served as frames. To further mitigate the negative influences of the ows extracted from the generated source videos, we propose the category-aware ow memory bank (CFB). By replacing the optical ow in a generated source video with real target ow which is retrieved from CFB, we create a new video sample that closely resembles the target video. Additionally, we leverage the video pace prediction task to enhance the model's perception of speed. Our proposed method demonstrates promising results compared with the current state-of-the-art approaches. In our current method, the domain discrepancy is not fully concerned, which could be further improved with video-to-video domain adaptation methods in the future work.

## REFERENCES

- [1] Filippo Aleotti, Matteo Poggi, and Stefano Mattoccia. 2021. Learning optical flow from still images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. 2021. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254* (2021).
- [3] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* 2017, 929–938.
- [4] Jin Chen, Xinxiao Wu, Yao Hu, and Jiebo Luo. 2021. Spatial-temporal causal inference for partial image-to-video adaptation. *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 35, 1027–1035.
- [5] Jin Chen, Xinxiao Wu, Yao Hu, and Jiebo Luo. 2021. Spatial-temporal causal inference for partial image-to-video adaptation. *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 35, 1027–1035.
- [6] Min-Hung Chen, Zsolt Kira, Ghassan AlRegib, Jaekwon Yoo, Ruxin Chen, and Jian Zheng. 2019. Temporal attentive alignment for large-scale video domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* 6321–6330.
- [7] Shuhao Cui, Shuhui Wang, Junbao Zhuo, Liang Li, Qingming Huang, and Qi Tian. 2020. Towards discriminability and diversity: Batch nuclear-norm maximization under label insu cient situations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 3941–3950.
- [8] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. 2017. Self-supervised video representation learning with odd-one-out networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* 3636–3645.
- [9] Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning* PMLR, 1180–1189.
- [10] Zan Gao, Leming Guo, Tongwei Ren, An-An Liu, Zhi-Yong Cheng, and Shengyong Chen. 2020. Pairwise two-stream convnets for cross-domain action recognition with small data. *IEEE Transactions on Neural Networks and Learning Systems* 3 (2020), 1147–1161.
- [11] Yuying Ge, Yixiao Ge, Xihui Liu, Dian Li, Ying Shan, Xiaohu Qie, and Ping Luo. 2022. Bridging video-text retrieval with multiple choice questions. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 16176.
- [12] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (Eds.), 2672–2680. <https://proceedings.neurips.cc/paper/2014/hash/5ca3e9b122f61f8f06494c97b1afccf3-Abstract.html>
- [13] Jiayuan Huang, Arthur Gretton, Karsten Borgwardt, Bernhard Schölkopf, and Alex Smola. 2006. Correcting sample selection bias by unlabeled data. *Advances in neural information processing systems* 19 (2006).
- [14] Ying Jin, Ximei Wang, Mingsheng Long, and Jianmin Wang. 2020. Minimum class confusion for versatile domain adaptation. *Computer Vision ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII* Springer, 464–480.
- [15] Longlong Jing, Xiaodong Yang, Jingen Liu, and Yingli Tian. 2018. Self-supervised spatiotemporal feature learning via video rotation prediction. *arXiv preprint arXiv:1811.11387* (2018).
- [16] Andrew Kae and Yale Song. 2020. Image to video domain adaptation using web supervision. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* 567–575.
- [17] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, and Andrew Senior. 2017. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06962* (2017).
- [18] Donghyun Kim, Yi-Hsuan Tsai, Bingbing Zhuang, Xiang Yu, Stan Sclaro, Kate Saenko, and Manmohan Chandraker. 2021. Learning cross-modal contrastive features for video domain adaptation. *Proceedings of the IEEE/CVF International Conference on Computer Vision* 1836–1845.
- [19] Hildegard Kuehne, Hueihan Jhuang, Estibaliz Garrote, Tomaso Poggio, and Thomas Serre. 2011. HMDB: a large video database for human motion recognition. In *2011 International conference on computer vision*, 355–364.
- [20] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. 2017. Attention transfer from web images for video recognition. *Proceedings of the 25th ACM international conference on multimedia* 109–118.
- [21] Wei Lin, Anna Kukleva, Kunyong Sun, Horst Possegger, Hilde Kuehne, and Horst Bischof. 2022. CycDA: Unsupervised Cycle Domain Adaptation to Learn from Image to Video. In *Computer Vision ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III* Springer, 698–715.
- [22] Yang Liu, Zhaoyang Lu, Jing Li, Tao Yang, and Chao Yao. 2019. Deep image-to-video adaptation and fusion networks for action recognition. *IEEE Transactions on Image Processing* 28 (2019), 3168–3182.
- [23] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. 2015. Learning Transferable Features with Deep Adaptation Networks. *Proceedings of the 32nd International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 37)*, Francis Bach and David Blei (Eds.), PMLR, Lille, France, 97–105. <https://proceedings.mlr.press/v37/long15.html>
- [24] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. 2017. Deep transfer learning with joint adaptation networks. In *International conference on machine learning* PMLR, 2208–2217.
- [25] Yadan Luo, Zi Huang, Zijian Wang, Zheng Zhang, and Mahsa Baktashmotlagh. 2020. Adversarial bipartite graph learning for video domain adaptation. In *Proceedings of the 28th ACM International Conference on Multimedia* 1072–1081.
- [26] Shugao Ma, Sarah Adel Bargal, Jianming Zhang, Leonid Sigal, and Stan Sclaro. 2017. Do less and achieve more: Training cnns for action recognition utilizing action images from the web. *Pattern Recognition* 68 (2017), 334–345.
- [27] Fabio Maria Carlucci, Lorenzo Porzi, Barbara Caputo, Elisa Ricci, and Samuel Rota Buló. 2017. Autodial: Automatic domain alignment layers. *Proceedings of the IEEE international conference on computer vision* 5067–5075.
- [28] Daniela Motezuma, Tania Ramírez-delReal, Guillermo Ruiz, and Othón González-Chávez. 2023. Video captioning: A comparative review of where we are and which could be the route. *Computer Vision and Image Understanding* 233 (2023), 103671.
- [29] Boxiao Pan, Zhangjie Cao, Ehsan Adeli, and Juan Carlos Niebles. 2020. Adversarial cross-domain action recognition with co-attention. *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 34, 11815–11822.
- [30] Jungin Park, Jiyoung Lee, and Kwanghoon Sohn. 2023. Dual-path Adaptation from Image to Video Transformers. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 2203–2213.
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning* PMLR, 8748–8763.
- [32] Ghazala Ra q, Muhammad Ra q, and Gyu Sang Choi. 2023. Video description: A comprehensive survey of deep learning approaches. *Artificial Intelligence Review* (2023), 1–80.
- [33] Rohit Raja, Prakash Chandra Sharma, Md Rashid Mahmood, and Dinesh Kumar Saini. 2023. Analysis of anomaly detection in surveillance video: recent trends and future vision. *Multimedia Tools and Applications* 82, 8 (2023), 12635–12651.
- [34] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. 2020. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence* 44, 3 (2020), 1623–1637.
- [35] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE international conference on computer vision* 618–626.
- [36] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0415* (2012).
- [37] Gokhan Tanisik, Cemil Zalluhoglu, and Nazli Ikizler-Cinbis. 2016. Facial descriptors for human interaction recognition in still images. *Pattern Recognition Letters* 73 (2016), 44–51.
- [38] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9 (2008), 2579–2601.
- [39] Jindong Wang, Wenjie Feng, Yiqiang Chen, Han Yu, Meiyu Huang, and Philip S Yu. 2018. Visual domain adaptation with manifold embedded distribution alignment. In *Proceedings of the 26th ACM international conference on Multimedia* 1011–1020.
- [40] Jiangliu Wang, Jianbo Jiao, Linchao Bao, Shengfeng He, Yunhui Liu, and Wei Liu. 2019. Self-supervised spatio-temporal representation learning for videos by predicting motion and appearance statistics. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 4100–4105.
- [41] Jiangliu Wang, Jianbo Jiao, and Yun-Hui Liu. 2020. Self-supervised video representation learning by pace prediction. *Computer Vision ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII* Springer, 504–521.
- [42] Zeyu Wang, Yu Wu, Karthik Narasimhan, and Olga Russakovsky. 2022. Multi-query Video Retrieval. [arXiv:2201.03639 \[cs.CV\]](https://arxiv.org/abs/2201.03639)
- [43] Han Wu, Chunfeng Song, Shaolong Yue, Zhenyu Wang, Jun Xiao, and Yanyang Liu. 2022. Dynamic video mix-up for cross-domain action recognition. *Neuro-computing* 471 (2022), 358–368.
- [44] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. 2019. Self-supervised spatiotemporal learning via video clip order prediction. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 10343–10352.
- [45] Bangpeng Yao, Xiaoye Jiang, Aditya Khosla, Andy Lai Lin, Leonidas Guibas, and Li Fei-Fei. 2011. Human action recognition by learning bases of action attributes

1045	and parts. In <i>2011 International conference on computer vision</i> . IEEE, 1331–1338.	1103
1046	[46] Huiling Yao and Xing Hu. 2023. A survey of video violence detection. <i>Cyber-Physical Systems</i> 9, 1 (2023), 1–24.	1104
1047	[47] Feiwu Yu, Xinxiao Wu, Jialu Chen, and Lixin Duan. 2019. Exploiting images for video recognition: heterogeneous feature augmentation via symmetric adversarial learning. <i>IEEE Transactions on Image Processing</i> 28, 11 (2019), 5308–5321.	1105
1048	[48] Feiwu Yu, Xinxiao Wu, Yuchao Sun, and Lixin Duan. 2018. Exploiting images for video recognition with hierarchical generative adversarial networks. In <i>Proceedings of the 27th International Joint Conference on Artificial Intelligence</i> . 1107–1113.	1106
1049	[49] Christopher Zach, Thomas Pock, and Horst Bischof. 2007. A duality based approach for realtime tv-l 1 optical flow. In <i>Pattern Recognition: 29th DAGM Symposium, Heidelberg, Germany, September 12-14, 2007. Proceedings 29</i> . Springer, 214–223.	1107
1050	[50] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. [n. d.]. mixup: Beyond Empirical Risk Minimization. In <i>International Conference on Learning Representations</i> .	1108
1051	[51] Junbao Zhuo, Xingyu Zhao, Shuhui Wang, Huimin Ma, and Qingming Huang. 2023. Synthesizing Videos from Images for Image-to-Video Adaptation. In <i>Proceedings of the 31st ACM International Conference on Multimedia</i> . 8294–8303.	1109
1052		1110
1053		1111
1054		1112
1055		1113
1056		1114
1057		1115
1058		1116
1059		1117
1060		1118
1061		1119
1062		1120
1063		1121
1064		1122
1065		1123
1066		1124
1067		1125
1068		1126
1069		1127
1070		1128
1071		1129
1072		1130
1073		1131
1074		1132
1075		1133
1076		1134
1077		1135
1078		1136
1079		1137
1080		1138
1081		1139
1082		1140
1083		1141
1084		1142
1085		1143
1086		1144
1087		1145
1088		1146
1089		1147
1090		1148
1091		1149
1092		1150
1093		1151
1094		1152
1095		1153
1096		1154
1097		1155
1098		1156
1099		1157
1100		1158
1101		1159
1102		1160

## A OVERVIEW

In this document, we conduct ‘Integrating domain adaptation technique’ experiments with Full Model on  $E \rightarrow H$  and  $S \rightarrow U$ . Moreover, we present the computational information of our method. Finally, we visualize the generated video frames. The source codes are also provided.

## B ADDITIONAL EXPERIMENTS

**Further evaluations of integrating domain adaptation (DA) techniques.** We further investigate the performance of combining our Full Model with DA methods on  $E \rightarrow H$  and  $S \rightarrow U$  benchmarks and the results are presented in Tab. 6. The results demonstrate that the integration with typical DA methods, significantly enhances the model’s recognition capability. The performance not only achieves new state-of-the-art results on both benchmarks but also attains comparable performance to supervised learning on  $S \rightarrow U$  benchmark. The improvement is attributed to the effective mitigation of distribution discrepancies by employing classical DA methods. Our experimental results further emphasize the high compatibility between our approach and DA methods.

**Table 6: Accuracies on  $E \rightarrow H$  and  $S \rightarrow U$ , averaged over 3 random trials.**

method	$E \rightarrow H$	$S \rightarrow U$
Full Model	77.4	97.3
Full Model + DAN	77.4	98.0
Full Model + MCC	78.4	99.2
Full Model + BNM	<b>79.2</b>	<b>99.3</b>

**Analysis of computational costs.** As shown in Table. 7, we present the computational information of video generation and our Full Model. Although the video generation requires more computational resources, it only needs to run once, and the generated frames can be used throughout the training process.

**Table 7: Computational information**

component	Params(M)	GFLOPs	Train/Inference time(s)
Video generation	266.13	371.85	- / 0.17
Full Model	24.59	65.62	0.16 / 0.13

## C VISUALIZATIONS OF GENERATED FRAMES

In this section, we visualize generated video frames for certain categories including ‘Hource Race’, ‘Running’, ‘Climb’ and ‘Ski-jet’. The visualization results are shown in Fig. 6. We exhibit two groups of generated frames for each category. The first column represents the input still images from source domain, and the subsequent columns show the sampled second, sixth, tenth and sixteenth frames, respectively.

Through visualization, we can observe that by simulating the arbitrary movements of the camera in 3D space, we can obtain more realistic video frames. These generated video frames are beneficial for learning a high-performance spatial-temporal model.

## D SOURCE CODES

We provide the source codes for training the spatial-temporal model. Details can be referred to readme.md in codes/.

