

Augmenting and Aligning Snippets for Few-Shot Video Domain Adaptation

Yuecong Xu^{1*} Jianfei Yang^{2*} Yunjiao Zhou² Zhenghua Chen^{1†} Min Wu¹ Xiaoli Li¹

¹Institute for Infocomm Research, A*STAR, Singapore ²Nanyang Technological University
{xuyu0014, yang0478, yunjiao001, chen0832}@e.ntu.edu.sg {wumin, xlli}@i2r.a-star.edu.sg

Abstract

For video models to be transferred and applied seamlessly across video tasks in varied environments, Video Unsupervised Domain Adaptation (VUDA) has been introduced to improve the robustness and transferability of video models. However, current VUDA methods rely on a vast amount of high-quality unlabeled target data, which may not be available in real-world cases. We thus consider a more realistic Few-Shot Video-based Domain Adaptation (FSVDA) scenario where we adapt video models with only a few target video samples. While a few methods have touched upon Few-Shot Domain Adaptation (FSDA) in images and in FSVDA, they rely primarily on spatial augmentation for target domain expansion with alignment performed statistically at the instance level. However, videos contain more knowledge in terms of rich temporal and semantic information, which should be fully considered while augmenting target domains and performing alignment in FSVDA. We propose a novel SSA²lign to address FSVDA at the snippet level, where the target domain is expanded through a simple snippet-level augmentation followed by the attentive alignment of snippets both semantically and statistically, where semantic alignment of snippets is conducted through multiple perspectives. Empirical results demonstrate state-of-the-art performance of SSA²lign across multiple cross-domain action recognition benchmarks. Code will be provided at: <https://github.com/xuyu0010/SSA2lign>.

1. Introduction

Video Unsupervised Domain Adaptation (VUDA) [4, 7, 51, 46, 53] aims to improve the generalizability and robustness of video models by transferring knowledge to new domains, and is widely applied in scenarios where massive labeled videos are unavailable. Current VUDA methods assume that sufficient target data are accessible which enables domain alignment by minimizing cross-domain distri-

bution discrepancies and obtaining domain invariant representations [4, 7, 54]. However, this assumption may not be feasible in real-world applications such as in smart hospitals and security surveillance where video models are leveraged for anomaly behavior recognition [35, 32], and are expected to be functional at all times even across different environments. It is more practical to obtain a few labeled videos during the early stage of model deployment to improve the transferred models' performances in the new (target) environment. A Few-Shot Video Domain Adaptation (FSVDA) task is hence formulated to enable knowledge learned from labeled source video to be transferred to the target video domain given only very limited labeled target videos.

With only several target domain samples, FSVDA is more challenging than VUDA, since aligning distributions with limited samples is harder. A few research have touched on the image-based Few-Shot Domain Adaptation (FSDA) [27, 45, 49, 11] by domain alignment, e.g., moment matching or adversarial training, between a spatial-augmented target domain and a filtered target-similar source domain. More recently, there have been a few early research on FSVDA [12, 13] which extends the above strategies to videos by viewing each video sample as a whole and obtaining frame-based video features.

However, there are two major shortcomings when the image-based FSDA is applied to video domains. Firstly, applying frame-level spatial augmentation towards individual video frames ignores and undermines temporal correlation across sequential frames, and we find that such augmentation would result in only minor or even negative effects on FSVDA performance. Secondly, the effectiveness of domain alignment methods is built upon sufficient source domain and target domain data that depicts the distribution discrepancy, which is not available in FSVDA. Even worse, statistical estimation of video data distribution is less accurate due to the complicated temporal structure of video data. In this paper, we aim to overcome these two challenges by designing more effective target domain augmentation and semantic alignment in the spatial-temporal domain.

To this end, we propose to address the FSVDA task by a Snippet-attentive Semantic-statistical **Alignment** with

*Equal contribution.

†Corresponding author.

Stochastic Sampling Augmentation (SSA²lign). Instead of aligning features of whole video samples at the video level or frame level [12, 13], we align source and target video features at the snippet level. Snippets are formed from a limited series of adjacent sequential frames, thus they contain both spatial and short-term temporal information. Leveraging snippet features for FSVDA brings two unique advantages: i) a larger amount of target domain samples could be obtained via spatial-temporal augmentations on snippets, obtaining more diverse features across the temporal dimension; ii) additional alignment of the diverse but highly correlated snippet features of each video could further improve the discriminability of the corresponding videos, which has been proven to benefit the effectiveness of video domain adaptation [6, 59, 21, 53]. SSA²lign is therefore proposed. It firstly augments the source and target domain data by a simple yet effective stochastic sampling process that makes full use of the abundance of snippet information and then performs semantic alignment from three perspectives: alignment based on semantic information within each snippet, cross-snippets of each video, and across snippet-level data distribution. Our method is demonstrated to be very effective for the FSVDA problem, surpassing state-of-the-art methods by large margins on two VUDA benchmarks.

In summary, our contributions are threefold. (i) We propose a novel SSA²lign to address FSVDA at the snippet level by both statistical and semantic alignments that are achieved from three perspectives. (ii) We propose to augment target domain data and the snippet-level alignments by a simple yet effective stochastic sampling of snippets for more robust video domain alignment. (iii) Extensive experiments show the efficacy of SSA²lign, achieving a remarkable average improvement of 13.1% and 4.2% over current state-of-the-art FSDA/FSVDA methods on two large-scale cross-domain action recognition benchmarks.

2. Related Work

(Video) Unsupervised Domain Adaptation ((V)UDA). Current UDA and VUDA methods aim to transfer knowledge from the source to the target domain given that both domains contain sufficient data, improving the transferability and robustness of models [50, 57]. They could be generally divided into four categories: a) reconstruction-based methods [14, 56], where domain-invariant features are obtained by encoders trained with data-reconstruction objectives; b) adversarial-based methods [4, 51, 7], where feature generators obtain domain-invariant features leveraging domain discriminators, trained jointly in an adversarial manner [17, 10]; c) semantic-based methods [58, 53], which exploit the shared semantics across domains such that domain-invariant features are obtained; and d) discrepancy-based methods [33, 62], which mitigate domain shifts by applying metric learning, minimizing metrics such as MMD [25]

and CORAL [36]. With the introduction of cross-domain video datasets such as Daily-DA [54] and Sports-DA [54], there has been a significant increase in research interest for VUDA [29, 5]. Despite the gain in video model robustness thanks to VUDA methods, they all assume that sufficient target data are accessible, which may not be feasible in real-world cases where a large amount of superior unlabeled target data are not available. A more related VUDA method concerns SAVA [7] which also utilizes the clips to design a self-supervised learning task (clip order prediction), but the adaptation is still performed with the video-level feature. We differ from SAVA [7] where our alignment is performed at the snippet level considering three different perspectives: within each snippet, cross-snippets of each video, and across snippet-level data distribution, therefore leading to better performances in FSVDA.

Few-Shot (Video) Domain Adaptation (FS(V)DA). It is more practical to obtain a few labeled target data to aid video models to adapt. There have been a few research that explores image-based FSDA. Among them, FADA [27] is adversarial-based and augments the domain discriminator to classify 4 types of source-target pairs. d-SNE [49] learns a latent space through SNE [16] with large-margin nearest neighborhood [9], and utilizes spatial augmentations to create sibling target samples. AcroFOD [11] explores FSDA for object detection by applying multi-level spatial augmentation and filtering target-irrelevant source data. There are also works as in [63, 38, 37, 60] that combine domain adaptation (DA) with few-shot learning (FSL), yet we differ them in the assumption of similar target and source classes and only limited target data accessible, which is more realistic. More recently, there have been a few early research on FSVDA, including PASTN [12] that constructs pairwise adversarial networks performed across source-target video pairs, while PTC [13] further leverages optical flow features. Both PASTN and PTC obtain video features from a frame-based video model. Despite some advances made in FS(V)DA, the above methods have not tackled FSVDA effectively by leveraging the rich temporal information as well as semantic information embedded within videos. We propose to engage in FSVDA by augmenting and attentively aligning snippet-level features which contain temporal information via both semantic and statistical alignments.

3. Proposed Method

For *Few-Shot Video Domain Adaptation*, we are given a labeled source domain $\mathcal{D}_S = \{(V_{S,i}, y_{S,i})\}_{i=1}^{N_S}$ with sufficient N_S i.i.d. source videos across \mathcal{C} classes, characterized by a probability distribution of p_S . We are also given a labeled target domain $\mathcal{D}_T = \{(V_{T,j}, y_{T,j})\}_{j=1}^{N_T}$ with a limited number of $N_T \ll N_S$ i.i.d. target videos across the same \mathcal{C} classes, where each video class only contains k target video samples (corresponding to the k -shot Video Domain Adap-

tation task), thus $N_T = k \times C$. \mathcal{D}_T is characterized by a probability distribution of p_T .

Owing to the absence of sufficient target data and the lack of target information, FSVDA is more challenging than VUDA. Current VUDA methods [4, 51] that are primarily moment matching-based are ineffective without target information for domain alignment. FSVDA should be tackled by leveraging the temporal information of videos fully for more temporally diverse features while aligning with the embedded semantic information to improve video discriminability for effective video domain adaptation. We propose SSA²lign, a novel method to transfer knowledge from the source domain to the target domain with only limited labeled target data by obtaining, augmenting, and aligning snippet features attentively. We start by introducing how snippet features are obtained and augmented through the Stochastic Sampling Augmentation (SSA), followed by a detailed illustration of the proposed SSA²lign.

3.1. Snippet Features with the Stochastic Sampling Augmentation

The key to effective target domain expansion and domain alignment in FSVDA is to obtain and augment features with temporal information such that the augmented features are diverse temporally. While various spatial augmentation methods (e.g., color jittering, flipping, cropping) have been adopted in supervised action recognition thanks to their capability in improving the robustness of video models, and in prior FSDA for expanding the target domain \mathcal{D}_T , they are performed at the frame-level across randomly selected individual frames. Meanwhile, the temporal information corresponds to the correlation of sequential frames and would be undermined by spatial augmentation since sequential frames may not be equally augmented. Augmentations for FSVDA must be performed above the frame level.

Snippets are formed from a limited series of adjacent sequential frames and have been utilized in multiple supervised action recognition methods (e.g., TSN [44] and STPN [48]) thanks to their ability in including both spatial and short-term temporal information. Therefore, we align source and target video features at the snippet level. Mathematically, given a target video $V = [f^1, f^2, \dots, f^n]$ that contains n frames, we denote the i -th frame as f^i . We denote the length of a snippet s to be m , then video V would contain $n - m + 1$ snippets in total. We define a snippet $s^j = [f^j, f^{j+1}, \dots, f^{j+m-1}]$ as the snippet starting from the j -th frame. While given only $N_T = k \times C$ target videos, there are $N_T \times (n - m + 1)$ target snippets, which can greatly expand the target domain for domain alignment while preserving essential temporal information.

While the target domain is largely expanded, utilizing all snippets for alignment is computationally inefficient (a 10-second 30-fps video contains more than 290 8-frame

snippets). Moreover, snippets that are obtained adjacently would differ over only ONE frame, resulting in high redundancy in temporal information. To ensure that diverse temporal information is utilized, we adopt a simple Stochastic Sampling Augmentation (SSA) over the snippets. Formally, during training we sample $r > 1$ snippets s_b^a stochastically per target video per mini-batch, where $a \in [1, n - m + 1]$ denotes the starting frame of the snippet and $b \in [1, r]$ denotes the b -th snippet sampled. SSA further ensures that the sampled snippets are diverse from two perspectives. Firstly, SSA samples snippets with a minimum of \hat{m} difference between the starting frame of any two snippets from the same target video, that is $\forall b_x \in [1, r], b_y \in [1, r]$ with $b_x \neq b_y$, we set $|a_x - a_y| \geq \hat{m}$. Secondly, since there are more source videos than target videos during training, it is likely that the same target video would be encountered across different mini-batches. SSA ensures that different snippets are sampled each time the same target video is included in a mini-batch across the same training epoch.

The SSA is also applied to the source videos to obtain source snippets. However, since there are sufficient source videos, it is more reasonable and efficient to exploit source knowledge with different source videos rather than the different snippets of a source video that would contain redundant source knowledge. Therefore, we only sample $r = 1$ snippet stochastically per source video via SSA.

Another crucial step towards transferring source knowledge to the target domain is to obtain rigorous snippet features that include both spatial and temporal information. We resort to the Transformer-based TimeSFormer [2] which extracts spatial and temporal features with separate space-time attention blocks based on self-attention [41]. While various Transformer-based video models achieve competitive performances on action recognition, TimeSFormer possesses the least amount of parameters, requiring only 60% parameters of Swin [24] and only 40% parameters of ViViT [1]. The feature of snippet s_b^a is $\mathbf{f}_b = \text{Time}(s_b^a)$ where Time denotes the TimeSFormer.

3.2. Snippet-attentive Semantic-statistical Alignment with SSA

With the absence of sufficient target data, conventional VUDA methods that are primarily moment matching-based would not be fully effective since target data distribution is unknown. Alternatively, we tackle FSVDA at the snippet level by aligning the embedded semantic information from three perspectives: aligning based on the semantic information within each snippet, cross-snippets of each video, and across snippet-level data distribution. Statistical alignment is also adopted for more stable domain alignment, while both alignments attend to the more impactful snippets.

Following the above strategy, we propose the **Snippet-attentive Semantic-statistical Alignment (SSAlign)**, with

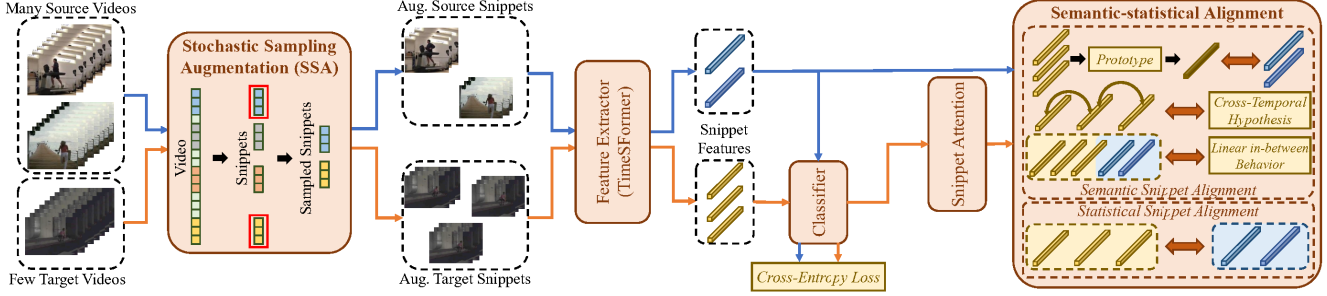


Figure 1. Pipeline of SSA²lign. Source and target snippets are first obtained through the Stochastic Sampling Augmentation, whose features are obtained through the shared feature extractor. SSA²lign then aligns the source and target domains at the snippet level with the Semantic-statistical Alignment, while weighing the impact of different target snippets through snippet attention, whose weight is built based on the output prediction of target snippets, obtained from a shared classifier with source snippets. The blue and orange lines imply the data flow for source and target videos respectively.

the input obtained through SSA introduced in Sec. 3.1, forming the **SSA²lign**. The overall pipeline of SSA²lign is presented in Fig. 1. We obtain the augmented source and target snippets through SSA whose features are extracted by applying TimeSFormer. We denote a source snippet from the i -th source video as $s_{S,i}$ and its feature as $\mathbf{f}_{S,i}$, while the l -th target snippet ($l \in [1, r]$) from the j -th target video as $s_{T,jl}$ and its feature as $\mathbf{f}_{T,jl}$. The superscript of the snippet expression is omitted for clarity. Domain alignment is achieved by performing both the Semantic Snippet Alignment and the Statistical Snippet Alignment. The snippet attention is applied to the augmented target snippets to weigh the snippets dynamically. The TimeSFormer feature extractor *Time* is shared across source and target domains while a shared classifier *H* outputs a prediction o for the source and target snippets, optimized through a cross-entropy loss:

$$\mathcal{L}_{pred} = \frac{1}{N_S} \sum_{i=1}^{N_S} l_{ce}(o_{S,i}, y_{S,i}) + \frac{1}{N_T \times r} \sum_{j=1}^{N_T} \sum_{l=1}^r l_{ce}(o_{T,jl}, y_{T,j}), \quad (1)$$

where $o_{S,i} = \sigma(H(\mathbf{f}_{S,i}))$ and $o_{T,jl} = \sigma(H(\mathbf{f}_{T,jl}))$ are the output predictions of snippet features $\mathbf{f}_{S,i}$ and $\mathbf{f}_{T,jl}$, while σ denotes the SoftMax function.

Semantic Snippet Alignment. The purpose of applying semantic alignment at the snippet level is to match the embedded semantic information (e.g., each individual feature or characteristic over a set of features) across source and target domains. Since both domains share the same TimeSFormer feature extractor, this implies that for each individual snippet feature, those of the same class should be close together across both domains. However, it is computationally expensive to compute the distances between each source and target snippet features given their large quantity. Inspired by the Prototypical Network [34, 23] designed for few-shot learning [65, 43], we resort to a more efficient solution where semantic alignment across each snippet is per-

formed by minimizing the distance between source snippet features and target prototypes. The target prototypes are obtained for each individual class C_x as the mean feature of all target snippet features classified as C_x , formulated as:

$$Pr_x = \frac{1}{n_{T,x}} \sum_{\forall s_{T,jl} \in C_x} \mathbf{f}_{T,jl}, \quad (2)$$

where $n_{T,x}$ is the number of target snippets classified as class C_x . For stable and effective alignment, the snippet features for computing the target prototypes are obtained after e training epochs. Target prototypes are subsequently updated per epoch by their exponential moving average as:

$$Pr_x \leftarrow \lambda_P Pr_x + (1 - \lambda_P) Pr'_x, \quad (3)$$

where Pr_x and Pr'_x denote the target prototype of class C_x computed at the current and previous epochs. Aligning source snippet features towards target prototypes is thus achieved by minimizing the Euclidean distances between them and denoted as the prototype alignment loss as:

$$\mathcal{L}_{proto} = \frac{1}{N_S} \sum_{x=1}^C \sum_{i=1}^{n_{S,x}} \sqrt{(\mathbf{f}_{S,i} - Pr_x)^2}. \quad (4)$$

$n_{S,x}$ is the number of source snippets classified as class C_x .

Besides the capability of obtaining temporally diverse features via SSA, leveraging snippet features for FSVDA is also more advantageous due to the inclusion of additional semantic information that exists across the diverse but highly correlated snippet features obtained from the same video, which should also be aligned. However, since we aim to exploit more source information with different source videos, the source cross-snippet semantic information cannot be directly obtained. Alternatively, the *cross-temporal hypothesis* introduced in [53] provides a thorough description of the cross-snippet semantic information for the source videos. Therefore, the equivalence of aligning the cross-snippet semantic information across source and

target domains is to align the cross-snippet semantic information of the target domain to the *cross-temporal hypothesis*, that is the snippet features over the snippets obtained from the same target video through SSA must be consistent. Meanwhile, aligning the *cross-temporal hypothesis* would also drive target videos to be discriminative, while previous studies [6, 59, 21, 53] have proven that improving discriminability can benefit the effectiveness of domain adaptation.

Formally, the cross-snippet consistency is achieved by minimizing the Kullback–Leibler (KL) divergence of the predictions of target snippets corresponding to the same target video. It is computed between each snippet against the key snippet of each target video, which is identified such that it is classified correctly and is certain in its prediction (i.e., low prediction entropy). In cases where no snippets are classified correctly, the snippet with the lowest prediction entropy is identified as the key snippet. The cross-snippet consistency loss is computed as:

$$\mathcal{L}_{cross} = \frac{1}{N_T(r-1)} \sum_{j=1}^{N_T} \sum_{l=1, l \neq k}^r KL(\log(o_{T,jy}) || \log(o_{T,jl})), \quad (5)$$

where $KL(p||q)$ denotes the KL-divergence while y denotes y -th snippet corresponding to the target video $V_{T,j}$ identified as the key snippet.

Aligning semantically via matching the characteristics over differed snippet features could be further performed across the snippet-level data distribution. Since source snippets for training are obtained stochastically at each training epoch, semantic information embedded across the source snippet-level data distribution changes continuously, and would therefore be ineffective for the target snippet-level data distribution to be directly aligned. Alternatively, snippet features that are highly discriminative would imply effective domain adaptation since it has been proven that improving discriminability benefits domain adaptation [6, 59, 21, 53]. We thus aim to drive the feature extractor towards obtaining snippet features that are distributed more discriminatively. Specifically, results in model robustness [61] suggest that the discriminability of features can be improved if the feature extractor behaves linearly in-between training samples. The linear in-between behavior can be complied by employing the interpolation consistency training (ICT) technique [42] across both source and target snippets, which encourages the linearly interpolated features to produce a linearly interpolated prediction. Formally, given a pair of snippet features \mathbf{f}_* , $\mathbf{f}_{*'}$, and their corresponding output predictions o_* , $o_{*'}$, the ICT is conducted with the following process and optimization loss:

$$\begin{aligned} \tilde{\mathbf{f}} &= \lambda_v \mathbf{f}_* + (1 - \lambda_v) \mathbf{f}_{*'}. \\ \tilde{\mathbf{o}} &= \lambda_v \mathbf{o}_* + (1 - \lambda_v) \mathbf{o}_{*'}. \\ \mathcal{L}_{ICT}(*, *) &= l_{ce}(\sigma(H(\tilde{\mathbf{f}})), \tilde{\mathbf{o}}), \end{aligned} \quad (6)$$

Algorithm 1 Training with SSA²lign for FSVDA

Input: $\mathcal{D}_S = \{(V_{S,i}, y_{S,i})\}_{i=1}^{N_S}$, $\mathcal{D}_T = \{(V_{T,j}, y_{T,j})\}_{j=1}^{N_T}$, $N_T \ll N_S$.
1: while Training **do**
2: Obtain r target snippets $s_{T,jl}$ from $V_{T,j}$ and one source snippet $s_{S,i}$ from $V_{S,i}$ via SSA.
3: Obtain features $\mathbf{f}_{S,i}$, $\mathbf{f}_{T,jl}$, predictions $o_{S,i}$, $o_{T,jl}$.
4: Compute prediction loss as Eq. 1.
5: Obtain snippet attention as Eq. 9 and normalize. Update $\mathbf{f}_{T,jl}$ to $\mathbf{f}'_{T,jl}$.
6: if epoch $> e$ **then**
7: Obtain target prototypes Pr_x as Eq. 2-3.
8: Compute prototype alignment loss as Eq. 4.
9: end if
10: Compute cross-snippet consistency loss as Eq. 5.
11: Compute snippet distribution loss as Eq. 6-7.
12: Compute and optimize overall loss as Eq. 8.
13: end while
Output: Trained feature extractor *Time* and classifier *H*.

where $\lambda_v \in \text{Beta}(\alpha_v, \alpha_v)$ is the weight assigned to $\mathbf{f}_{T,j_1 l_1}$ sampled from a Beta distribution with α_v as the parameter. We refer to previous works [22, 55] and set $\alpha_v = 0.3$. $\tilde{\mathbf{f}}$ and $\tilde{\mathbf{o}}$ are the linearly interpolated features and the interpolated output predictions. In practice, we drive snippets to comply with the linear in-between behaviour by forming a single stochastic snippet pair for every snippet, forming $(N_T \times r + N_S)$ snippet pairs. Aligning the snippet-level data distribution with the linear in-between behavior is achieved by optimizing the snippet distribution loss as:

$$\mathcal{L}_{sn-dist} = \frac{1}{N_T \times r + N_S} \sum_{*,*' \in \{i \cup j\}} \mathcal{L}_{ICT}(*, *'). \quad (7)$$

It is possible that a snippet pair will include two snippets from the same target video. In such case, the corresponding \mathcal{L}_{ICT} across the snippet pair can be viewed as a low-ordered cross-snippet consistency loss. This implies that optimizing \mathcal{L}_{cross} and $\mathcal{L}_{sn-dist}$ share the common goal of improving feature discriminability for more effective video domain adaptation.

Statistical Snippet Alignment. To improve the stability of snippet-level alignment, we adopt a statistical alignment strategy apart from the aforementioned semantic alignment strategies. The statistical alignment is performed by minimizing the snippet-level distribution discrepancies $\mathcal{L}_{sn-stat}$ formulated as metrics such as MMD [25], CORAL [36], and MDD [62]. Compared to the adversarial-based adaptation strategy more commonly used in prior VUDA tasks [4, 51, 7], minimizing discrepancies does not require additional network structures (e.g., domain classifiers), thus is more stable. The MDD [62] metric is empirically selected. The overall optimization loss function for FSVDA is therefore:

$$\mathcal{L} = \mathcal{L}_{pred} + \lambda_{sem}(\mathcal{L}_{proto} + \mathcal{L}_{cross} + \mathcal{L}_{sn-dist}) + \lambda_{stat} \mathcal{L}_{sn-stat}, \quad (8)$$

where λ_{sem} and λ_{stat} are the tradeoff hyper-parameters for the semantic and statistical snippet alignment losses.

Snippet Attention. With multiple snippets leveraged per target video for both semantic and statistical snippet alignments, it is unreasonable to leverage each snippet equally

a wide range of cross-domain scenarios. We present superior results on both benchmarks. Further, ablation studies and analysis of SSA²lign are also presented to justify the design of SSA²lign. *Code is provided in the appendix.*

4.1. Experimental Settings

Daily-DA is a challenging dataset that has been leveraged in prior VUDA works [54, 53, 55]. It covers both normal and low-illumination videos and is constructed from four datasets: ARID (A) [52], HMDB51 (H) [20], Moments-in-Time (M) [26], and Kinetics-600 (KD) [3]. HMDB51, Moments-in-Time, and Kinetics-600 are widely used for action recognition benchmarking, while ARID is a recent dark dataset, with videos shot under adverse illumination. Daily-DA contains 18,949 videos from 8 classes, with 12 cross-domain action recognition tasks. **Sports-DA** is a large-scale cross-domain video dataset, built from UCF101 (U), Sports-1M (S) [18], and Kinetics-600 (KS), with 40,718 videos from 23 action classes, and includes 6 cross-domain action recognition tasks. Refer to prior FSDA/FSVDA works [12, 13, 11], we evaluate SSA²lign on both benchmarks with $k = (3, 5, 10)$ target videos per action class (i.e., 3-shot, 5-shot and 10-shot VDA tasks).

For a fair comparison, all methods examined and experiments conducted in this section adopt the TimeSFormer [64] as the feature extractor, pre-trained on Kinetics-400 [19]. All experiments are implemented with the PyTorch [30] library. We set the length of snippets and the number of snippets per target video via SSA empirically as $m = 8, r = 3$. Hyper-parameters λ_{sem} , λ_{stat} and λ_P are empirically set to 1.0, 1.0, and 0.6 and are fixed. *More specifications on benchmark details and network implementation are provided in the appendix.*

4.2. Overall Results and Comparisons

We compare SSA²lign with state-of-the-art FSDA approaches, and prevailing UDA/VUDA and few-shot action recognition (FSAR) approaches. These methods include: FADA [27], d-SNE [49] designed for image-based FSDA; DANN [10], MK-MMD [25], MDD [62], SAVA [7] and ACAN [51], designed for UDA/VUDA; and TRX [31], STRM [39], and HyRSM [47] proposed for FSAR. To adapt the FSAR approaches for FSVDA, the source domain is used for meta-training and the target domain is used for the meta-testing, while target labels are available for optimizing the cross-entropy loss to adapt UDA/VUDA approaches for FSVDA. We also report the results of the source-only model (denoted as TSF) by applying the model trained with only source data directly to the target data; and the source with few-shot target model (denoted as TSF w/ T) by optimizing only the prediction loss \mathcal{L}_{pred} for training. We report the top-1 accuracy on the target domains, averaged on 5 different settings of available target data randomly selected and

each with 5 runs (25 runs in total). Tables 1-3 show comparison of SSA²lign against the above methods.

Results in Tables 1-3 show that the novel SSA²lign achieves the state-of-the-art results on all 18 cross-domain action recognition tasks across both cross-domain benchmarks, outperforming prior UDA/VUDA, FSDA or FSAR approaches by noticeable margins. Notably, SSA²lign outperforms all prior FSDA approaches originally designed for image-based FSDA (i.e., FADA and d-SNE) consistently on all tasks, by a relative average of 13% over the second-best performances on Daily-DA (across 3 k -shot settings and 12 tasks), and a relative average of 4.2% on Sports-DA (across 3 k -shot settings and 6 tasks). The consistent improvements justify empirically the effectiveness of augmenting and aligning both semantic information and statistical distribution at the snippet level for FSVDA.

It is also observed that prior FSDA and UDA/VUDA methods could not perform well on FSVDA tasks. Notably, even when $k = 10$ target videos are available per class, all but one of the evaluated FSDA and UDA/VUDA approaches result in performances inferior to that trained with only \mathcal{L}_{pred} without any adaptation (i.e., TSF w/ T). Prior FSDA approaches do not incorporate temporal features and their related semantic information, which are crucial for tackling FSVDA, while UDA/VUDA methods are not effective when target information is not fully available. Negative improvements are more severe when k decreases. It is also noted that at small k values (e.g., $k = 3$), the performance of TSF w/ T could be inferior to that trained without target data (i.e., TSF). This suggests that the few target data could be outliers of the target domain, whose distribution differs greatly from the other target data, resulting in a severe negative impact. Prior FSAR approaches could not tackle FSVDA as well, producing even poorer results than all UDA/VUDA approaches examined. This can be caused by domain shift that exists between data for the meta-training and meta-testing. Feature extractors trained via meta-training on the source domain could not be simply applied to the meta-testing phase on the target domain.

4.3. Ablation Studies, Analysis, and Discussion

To gain a comprehensive understanding of SSA²lign and justify its design, we perform extensive ablation studies as in Tables 4-5. The ablation studies explore the effects brought by its components, namely the semantic and statistical alignments, the SSA, and the snippet attention. It further validates the alignment details by assessing against 5 variants: SSA²lign-CORAL and SSA²lign-MMD formulate $\mathcal{L}_{sn-stat}$ as CORAL [36] and MDD [62]; SSA²lign-FC computes \mathcal{L}_{cross} over all $r \times (r - 1)$ snippet pairs for the same target video; SSA²lign-SP minimizes the distance between target snippet features and source class prototypes for \mathcal{L}_{proto} ; SSAalign (w/ spatial aug.) augments target domain

Methods	Components						Daily-DA									Sports-DA						Avg.
	SSA	Sn-Attn	\mathcal{L}_{proto}	\mathcal{L}_{cross}	$\mathcal{L}_{sn-dist}$	$\mathcal{L}_{sn-stat}$	$k=10$			$k=5$			$k=3$			$k=10$		$k=5$		$k=3$		
							H→A	M→A	KD→A	H→A	M→A	KD→A	H→A	M→A	KD→A	U→S	KS→S	U→S	KS→S	U→S	KS→S	
TSF w/ T	✓						39.57	39.49	39.41	40.19	40.03	37.08	37.94	34.14	33.05	78.95	79.05	78.74	78.32	75.37	76.58	53.86
SSA ² lign	✓	✓	✓	✓	✓	✓	45.62	46.32	45.70	46.47	45.69	41.74	39.17	40.96	39.49	84.32	85.47	81.05	83.58	77.26	77.37	58.68
	✓	✓	✓	✓	✓	✓	51.05	51.13	50.43	50.97	50.73	46.39	43.59	45.85	44.22	86.58	87.37	83.21	85.63	79.32	80.37	62.46
	✓	✓	✓	✓	✓	✓	49.88	49.81	49.73	50.35	49.26	45.07	42.19	44.53	42.90	85.58	86.74	82.74	84.74	78.69	79.47	61.45
	✓	✓	✓	✓	✓	✓	50.19	50.50	49.65	50.27	49.88	45.77	42.66	45.15	43.06	85.84	87.05	82.58	85.05	78.95	79.58	61.75
	✓	✓	✓	✓	✓	✓	48.18	48.57	46.86	47.56	48.09	43.29	40.64	42.51	41.97	84.68	85.95	81.95	83.74	77.74	78.63	60.03
	✓	✓	✓	✓	✓	✓	51.36	51.44	50.89	51.82	51.43	46.55	43.90	46.00	44.92	86.68	87.63	83.63	85.69	79.69	80.42	62.80
	✓	✓	✓	✓	✓	✓	52.13	52.21	51.75	52.37	51.98	47.40	44.83	46.78	45.31	87.26	88.11	84.05	86.21	80.05	80.95	63.43

Table 4. Ablation studies of the components of SSA²lign on 5 cross-domain tasks over Daily-DA and Sports-DA.

Methods	Daily-DA									Sports-DA						Avg.	Δ Avg.	GFLOPS	Δ GFLOPS
	$k=10$			$k=5$			$k=3$			$k=10$		$k=5$		$k=3$					
	H→A	M→A	KD→A	H→A	M→A	KD→A	H→A	M→A	KD→A	U→S	KS→S	U→S	KS→S	U→S	KS→S				
SSA ² lign-CORAL	51.90	51.98	51.51	51.98	51.58	47.09	44.52	46.39	45.07	87.05	87.84	83.90	86.00	79.90	80.68	63.16	-0.27	1302	-8
SSA ² lign-MMD	51.67	51.90	51.28	51.98	51.51	47.01	44.52	46.32	44.84	87.05	87.84	83.84	85.90	79.79	80.63	63.07	-0.36	1312	+2
SSA ² lign-FC	52.83	52.91	52.37	52.99	52.36	47.56	45.30	47.25	45.31	87.16	88.26	84.53	86.69	79.79	81.26	63.78	+0.35	1472	+162
SSA ² lign-SP	50.97	51.20	50.97	51.28	51.27	46.47	43.82	45.62	44.69	86.90	87.42	83.37	85.74	79.26	80.47	62.63	-0.80	1390	+80
SSAlign (w/ spatial aug.)	45.38	46.70	45.23	45.77	45.84	41.12	39.63	40.26	38.64	83.53	85.63	80.58	83.37	76.95	77.74	58.43	-5.00	1325	+15
SSA ² lign	52.13	52.21	51.75	52.37	51.98	47.40	44.83	46.78	45.31	87.26	88.11	84.05	86.21	80.05	80.95	63.43	-	1310	-

Table 5. Ablation studies of the alignment details of SSA²lign on 5 cross-domain tasks over Daily-DA and Sports-DA.

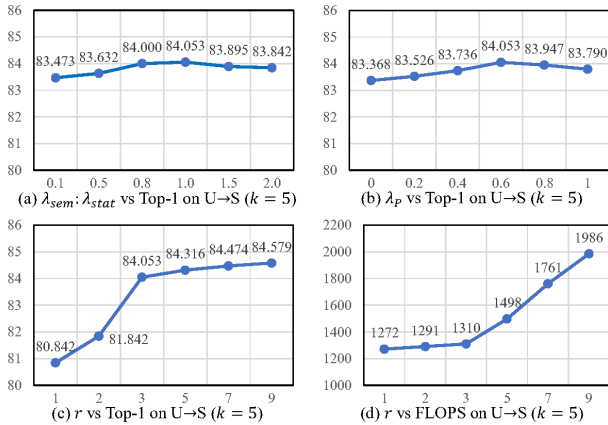


Figure 2. Sensitivity of hyper-parameters on U→S task.

through random spatial augmentation across the frames of r snippets. The ablation studies are conducted on 5 tasks over Daily-DA and Sports-DA. If SSA is not applied, we sample r snippets sequentially from the 1st frame of each target video and remain unchanged during training.

Semantic Alignment. As shown in Table 4, with only snippet-level semantic alignment (whether in full or any one of the three perspectives), the performance still surpasses all previous FSDA and UDA/VUDA methods compared. This conforms to our motivation that applying semantic alignment could tackle FSVDA more effectively. Moreover, statistical alignment and snippet attention further improve SSA²lign, but only by a marginal degree.

Superiority of SSA. Notably, a significant performance drop is observed when SSA is not applied, which proves the importance of expanding target domain data through SSA for subsequent alignment. The importance of SSA is further verified when we apply SSA for training with augmented

Methods	Sports-DA						Avg.	Δ Avg.	GFLOPS	Δ GFLOPS
	$k=10$		$k=5$		$k=3$					
	U→S	KS→S	U→S	KS→S	U→S	KS→S				
ACAN+MixUp	80.63	80.89	80.42	79.79	75.95	77.37	79.18			
ACAN+RandAugment	78.68	79.84	78.26	80.05	75.74	76.58	78.19			
ACAN+TrivialAugment	81.53	81.05	80.53	80.68	76.84	77.63	79.71			
ACAN+MixUp+TrivialAugment	82.16	81.95	81.21	81.58	78.05	78.32	80.55			
SSA ² lign	87.26	88.11	84.05	86.21	80.05	80.95	84.44			

Table 6. Compare with ACAN with up-to-date augmentations.

snippets but without adaptation which shows a noticeable gain compared to the original TSF w/ T. Further, the significantly inferior performance of SSAlign (w/ spatial aug.) as shown in Table 5 conforms with the motivation of SSA, which aims for more effective target video domain augmentation while spatial augmentation may undermine temporal correlation across sequential frames.

While the success of SSA²lign is built upon augmenting target videos with SSA, there have been more complex augmentation practices introduced for images, such as MixUp [61], RandAugment [8], and TrivialAugment [28]. To further prove the efficacy of SSA for FSVDA, we compare SSA²lign against a competitive VUDA method ACAN [51] with the aforementioned up-to-date augmentation practices applied to the target domain on Sports-DA, as shown in Table 6. Note that TimeSFormer [64] is leveraged in ACAN as the feature extractor for fair comparison. Results in Table 6 show that while there are improvements by applying augmentations to the target domain, SSA²lign with simple augmentation still outperforms multiple augmentations due to its snippet-level alignment, which further validates the superiority of SSA.

Alignment Methods. Table 5 shows that while formulating $\mathcal{L}_{sn-stat}$ as MDD [62] brings the best performance, selecting other metrics brings negligible impact. Further, computing \mathcal{L}_{cross} with all target snippet pairs only brings trivial performance gain at a cost of significant computation over-

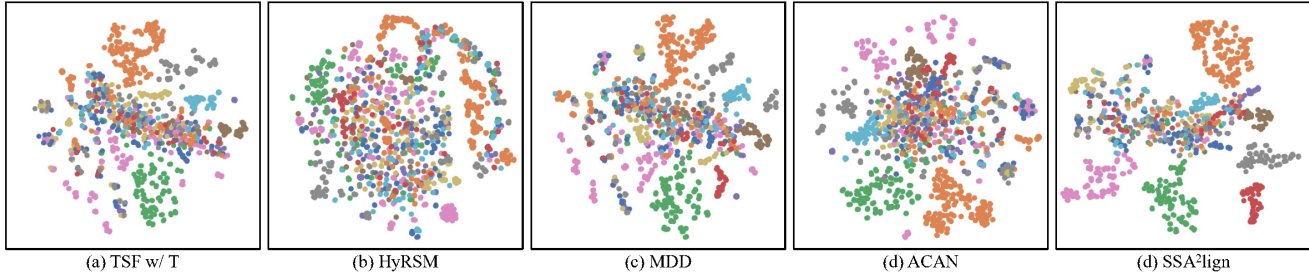


Figure 3. t-SNE visualizations of target features from (a) TSF w/T, (b) HyRSM, (c) MDD, (d) ACAN, (e) SSA²lign. Colors denote classes.

head (12% computation increase for 0.54% gain). Further, matching target snippet features to source class prototypes for \mathcal{L}_{proto} results in a performance drop with more computation. The inferior performance could be due to outliers in the source domain which could affect source class prototypes, bringing in source noise that should not be aligned.

Hyper-parameter Sensitivity. We focus on studying the sensitivity of λ_{sem} and λ_{stat} which control the strength of the semantic and statistical snippet alignment losses, λ_P which relates to the update of target prototypes and r the number of snippets per target video. Without loss of generality, we fix $\lambda_{stat} = 1.0$ and study the ratio $\lambda_{sem} : \lambda_{stat}$ in the range of 0.1 to 1.5. λ_P is in the range of 0 to 1 which corresponds to using only the initial prototypes or the updated prototypes, and r is in the range of 1 to 9. As shown in Fig. 2, SSA²lign is robust to ratio $\lambda_{sem} : \lambda_{stat}$ and λ_P , falling within a margin of 0.683%, with the best results obtained at the current default where $\lambda_{sem} : \lambda_{stat} = 1.0$ and $\lambda_P = 0.6$. SSA²lign is also robust to r when $r \geq 3$, i.e., when there are multiple snippets obtained via SSA per target video. $r = 3$ is selected as significant computation overhead would occur for $r > 3$ with marginal gain. Notably, SSA²lign cannot perform when $r < 3$, especially when $r = 1$ where the \mathcal{L}_{cross} does not work and the target domain is not expanded.

Feature Visualization. We further understand the characteristics of SSA²lign by plotting the t-SNE embeddings [40] of target features with class information from the model trained without adaptation (TSF w/T), HyRSM, MDD, ACAN and SSA²lign for U→S with $k = 10$ in Fig. 3. It is observed that target features from SSA²lign are more clustered and discriminable, corresponding to better performance. Such observation intuitively proves that video domain adaptation can be improved when feature extractors possess stronger discriminability. However, SSA²lign is not designed to deal explicitly with classes that could be similar spatially or temporally, thus certain features observe lower discriminability, which denotes future work.

Limitations in Choice of Datasets. The datasets Sports-DA and Daily-DA [54] are leveraged as our benchmarks as they have been commonly used in the VUDA community in

prior VUDA works [53, 54, 55]. However, it is noted that benchmarks that have been used for action recognition with stronger temporal reasoning assessment and fine-grained action classes (such as Something-Something V2 [15]) are not included in any current cross-domain video datasets as there is few datasets that offer overlapped fine-grained action classes. The current results show that the FSVDA is still a challenging task in our proposed benchmarks, we believe that exploring how to adapt models from coarse-category datasets to fine-grained datasets (SSv2) denotes future exploration.

5. Conclusion

In this work, we propose a novel SSA²lign to tackle the challenging yet realistic Few-Shot Video Domain Adaptation (FSVDA), where only limited labeled target data are available. Without sufficient target data, SSA²lign tackles FSVDA at the snippet level via a simple SSA augmentation and performing the semantic and statistical alignments attentively, where the semantic alignment is further achieved from three perspectives based on semantic information within and across snippets. Extensive experiments and detailed ablation studies across cross-domain action recognition benchmarks validate the superiority of SSA²lign in addressing FSVDA.

Acknowledgements

This research is jointly supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG2-RP-2021-027), and A*STAR Singapore under its Career Development Award (Grant No. C210112046), and Nanyang Technological University, Singapore, under its NTU Presidential Postdoctoral Fellowship, “Adaptive Multimodal Learning for Robust Sensing and Recognition in Smart Cities” project fund.

References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video

- vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021. [3](#)
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, volume 2, page 4, 2021. [3](#)
- [3] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note about kinetics-600. *arXiv preprint arXiv:1808.01340*, 2018. [7](#)
- [4] Min-Hung Chen, Zsolt Kira, Ghassan AlRegib, Jaekwon Yoo, Ruxin Chen, and Jian Zheng. Temporal attentive alignment for large-scale video domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6321–6330, 2019. [1](#), [2](#), [3](#), [5](#)
- [5] Min-Hung Chen, Baopu Li, Yingze Bao, Ghassan Al-Regib, and Zsolt Kira. Action segmentation with joint self-supervised temporal domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9454–9463, 2020. [2](#)
- [6] Xinyang Chen, Sinan Wang, Mingsheng Long, and Jianmin Wang. Transferability vs. discriminability: Batch spectral penalization for adversarial domain adaptation. In *International conference on machine learning*, pages 1081–1090. PMLR, 2019. [2](#), [5](#)
- [7] Jinwoo Choi, Gaurav Sharma, Samuel Schulter, and Jia-Bin Huang. Shuffle and attend: Video domain adaptation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 678–695. Springer, 2020. [1](#), [2](#), [5](#), [7](#)
- [8] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020. [8](#)
- [9] Carlotta Domeniconi, Dimitrios Gunopulos, and Jing Peng. Large margin nearest neighbor classifiers. *IEEE transactions on neural networks*, 16(4):899–909, 2005. [2](#)
- [10] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. [2](#), [7](#)
- [11] Yipeng Gao, Lingxiao Yang, Yunmu Huang, Song Xie, Shiyong Li, and Wei-Shi Zheng. AcrofoD: An adaptive method for cross-domain few-shot object detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 673–690. Springer, 2022. [1](#), [2](#), [7](#)
- [12] Zan Gao, Leming Guo, Weili Guan, An-An Liu, Tongwei Ren, and Shengyong Chen. A pairwise attentive adversarial spatiotemporal network for cross-domain few-shot action recognition-r2. *IEEE Transactions on Image Processing*, 30:767–782, 2020. [1](#), [2](#), [7](#)
- [13] Zan Gao, Leming Guo, Tongwei Ren, An-An Liu, Zhi-Yong Cheng, and Shengyong Chen. Pairwise two-stream convnets for cross-domain action recognition with small data. *IEEE Transactions on Neural Networks and Learning Systems*, 33(3):1147–1161, 2020. [1](#), [2](#), [7](#)
- [14] Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, David Balduzzi, and Wen Li. Deep reconstruction-classification networks for unsupervised domain adaptation. In *European conference on computer vision*, pages 597–613. Springer, 2016. [2](#)
- [15] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The” something something” video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017. [9](#)
- [16] Geoffrey E Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15, 2002. [2](#)
- [17] Ling Huang, Anthony D Joseph, Blaine Nelson, Benjamin IP Rubinstein, and J Doug Tygar. Adversarial machine learning. In *Proceedings of the 4th ACM workshop on Security and artificial intelligence*, pages 43–58, 2011. [2](#)
- [18] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. [7](#)
- [19] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset, 2017. [7](#)
- [20] Hildegard Kuehne, Hueihan Jhuang, Estibaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International Conference on Computer Vision*, pages 2556–2563. IEEE, 2011. [7](#)
- [21] Jogendra Nath Kundu, Akshay R Kulkarni, Suvaansh Bhambri, Deepesh Mehta, Shreyas Anand Kulkarni, Varun Jampani, and Venkatesh Babu Radhakrishnan. Balancing discriminability and transferability for source-free domain adaptation. In *International Conference on Machine Learning*, pages 11710–11728. PMLR, 2022. [2](#), [5](#)
- [22] Jian Liang, Dapeng Hu, Jiashi Feng, and Ran He. Dine: Domain adaptation from single and multiple black-box predictors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. [5](#)
- [23] Jinlu Liu, Liang Song, and Yongqiang Qin. Prototype rectification for few-shot learning. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 741–756. Springer, 2020. [4](#)
- [24] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3202–3211, 2022. [3](#)
- [25] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. [2](#), [5](#), [7](#)
- [26] Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Adel Bargal, Tom Yan, Lisa Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, et al. Moments

- in time dataset: one million videos for event understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):502–508, 2019. 7
- [27] Saeid Motiian, Quinn Jones, Seyed Iranmanesh, and Gianfranco Doretto. Few-shot adversarial domain adaptation. *Advances in neural information processing systems*, 30, 2017. 1, 2, 7
- [28] Samuel G Müller and Frank Hutter. Trivialaugment: Tuning-free yet state-of-the-art data augmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 774–782, 2021. 8
- [29] Boxiao Pan, Zhangjie Cao, Ehsan Adeli, and Juan Carlos Niebles. Adversarial cross-domain action recognition with co-attention. In *AAAI*, pages 11815–11822, 2020. 2
- [30] 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. In *Advances in neural information processing systems*, pages 8026–8037, 2019. 7
- [31] Toby Perrett, Alessandro Masullo, Tilo Burghardt, Majid Mirmehdi, and Dima Damen. Temporal-relational cross-transformers for few-shot action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 475–484, 2021. 7
- [32] Abdel Mlak Said, Aymen Yahyaoui, and Takoua Abdellatif. Efficient anomaly detection for smart hospital iot systems. *Sensors*, 21(4):1026, 2021. 1
- [33] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2018. 2
- [34] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 4
- [35] 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
- [36] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016. 2, 5, 7
- [37] Guangyu Sun, Zhang Liu, Lianggong Wen, Jing Shi, and Chenliang Xu. Anomaly crossing: New horizons for video anomaly detection as cross-domain few-shot learning. *arXiv e-prints*, pages arXiv–2112, 2021. 2
- [38] Takeshi Teshima, Issei Sato, and Masashi Sugiyama. Few-shot domain adaptation by causal mechanism transfer. In *International Conference on Machine Learning*, pages 9458–9469. PMLR, 2020. 2
- [39] Anirudh Thatipelli, Sanath Narayan, Salman Khan, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Bernard Ghanem. Spatio-temporal relation modeling for few-shot action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19958–19967, 2022. 7
- [40] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 9
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [42] Vikas Verma, Alex Lamb, Juho Kannala, Yoshua Bengio, and David Lopez-Paz. Interpolation consistency training for semi-supervised learning. In *International Joint Conference on Artificial Intelligence*, pages 3635–3641, 2019. 5
- [43] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *proceedings of the IEEE/CVF international conference on computer vision*, pages 9197–9206, 2019. 4
- [44] 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 transactions on pattern analysis and machine intelligence*, 41(11):2740–2755, 2018. 3
- [45] Tao Wang, Xiaopeng Zhang, Li Yuan, and Jiashi Feng. Few-shot adaptive faster r-cnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7173–7182, 2019. 1
- [46] Xiyu Wang, Yuecong Xu, Jianfei Yang, and Kezhi Mao. Calibrating class weights with multi-modal information for partial video domain adaptation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 3945–3954, 2022. 1
- [47] Xiang Wang, Shiwei Zhang, Zhiwu Qing, Mingqian Tang, Zhengrong Zuo, Changxin Gao, Rong Jin, and Nong Sang. Hybrid relation guided set matching for few-shot action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19948–19957, 2022. 7
- [48] Yunbo Wang, Mingsheng Long, Jianmin Wang, and Philip S Yu. Spatiotemporal pyramid network for video action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1529–1538, 2017. 3
- [49] Xiang Xu, Xiong Zhou, Ragav Venkatesan, Gurumurthy Swaminathan, and Orchid Majumder. d-sne: Domain adaptation using stochastic neighborhood embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2497–2506, 2019. 1, 2, 7
- [50] Yuecong Xu, Haozhi Cao, Zhenghua Chen, Xiaoli Li, Lihua Xie, and Jianfei Yan. Video unsupervised domain adaptation with deep learning: A comprehensive survey. *arXiv preprint arXiv:2211.10412*, 2022. 2
- [51] Yuecong Xu, Haozhi Cao, Kezhi Mao, Zhenghua Chen, Lihua Xie, and Jianfei Yang. Aligning correlation information for domain adaptation in action recognition. *IEEE Transactions on Neural Networks and Learning Systems*, 2022. 1, 2, 3, 5, 7, 8
- [52] Yuecong Xu, Jianfei Yang, Haozhi Cao, Kezhi Mao, Jianxiong Yin, and Simon See. Arid: A new dataset for recognizing action in the dark. In *International Workshop on*

- Deep Learning for Human Activity Recognition*, pages 70–84. Springer, 2021. [7](#)
- [53] Yuecong Xu, Jianfei Yang, Haozhi Cao, Keyu Wu, Min Wu, and Zhenghua Chen. Source-free video domain adaptation by learning temporal consistency for action recognition. In *European Conference on Computer Vision*, pages 147–164. Springer, 2022. [1](#), [2](#), [4](#), [5](#), [7](#), [9](#)
- [54] Yuecong Xu, Jianfei Yang, Haozhi Cao, Keyu Wu, Min Wu, Zhengguo Li, and Zhenghua Chen. Multi-source video domain adaptation with temporal attentive moment alignment network. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023. [1](#), [2](#), [6](#), [7](#), [9](#)
- [55] Yuecong Xu, Jianfei Yang, Min Wu, Xiaoli Li, Lihua Xie, and Zhenghua Chen. Extern: Leveraging endo-temporal regularization for black-box video domain adaptation. *arXiv preprint arXiv:2208.05187*, 2022. [5](#), [7](#), [9](#)
- [56] Jinyu Yang, Weizhi An, Sheng Wang, Xinliang Zhu, Chaochao Yan, and Junzhou Huang. Label-driven reconstruction for domain adaptation in semantic segmentation. In *European Conference on Computer Vision*, pages 480–498. Springer, 2020. [2](#)
- [57] Jianfei Yang, Yuecong Xu, Haozhi Cao, Han Zou, and Lihua Xie. Deep learning and transfer learning for device-free human activity recognition: A survey. *Journal of Automation and Intelligence*, 1(1):100007, 2022. [2](#)
- [58] Jianfei Yang, Jiangang Yang, Shizheng Wang, Shuxin Cao, Han Zou, and Lihua Xie. Advancing imbalanced domain adaptation: Cluster-level discrepancy minimization with a comprehensive benchmark. *IEEE Transactions on Cybernetics*, 0:1–12, 2021. [2](#)
- [59] Jianfei Yang, Han Zou, Yuxun Zhou, Zhaoyang Zeng, and Lihua Xie. Mind the discriminability: Asymmetric adversarial domain adaptation. In *European Conference on Computer Vision*, pages 589–606. Springer, 2020. [2](#), [5](#)
- [60] Xiangyu Yue, Zangwei Zheng, Hari Prasanna Das, Kurt Keutzer, and Alberto Sangiovanni Vincentelli. Multi-source few-shot domain adaptation. *arXiv preprint arXiv:2109.12391*, 2021. [2](#)
- [61] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *International Conference on Learning Representations*, 2018. [5](#), [8](#)
- [62] Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael Jordan. Bridging theory and algorithm for domain adaptation. In *International Conference on Machine Learning*, pages 7404–7413. PMLR, 2019. [2](#), [5](#), [7](#), [8](#)
- [63] An Zhao, Mingyu Ding, Zhiwu Lu, Tao Xiang, Yulei Niu, Jiechao Guan, and Ji-Rong Wen. Domain-adaptive few-shot learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1390–1399, 2021. [2](#)
- [64] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 803–818, 2018. [7](#), [8](#)
- [65] Xiatian Zhu, Antoine Toisoul, Juan-Manuel Perez-Rua, Li Zhang, Brais Martinez, and Tao Xiang. Few-shot action recognition with prototype-centered attentive learning. *arXiv preprint arXiv:2101.08085*, 2021. [4](#)