

Augmenting and Aligning Snippets for Few-Shot Video Domain Adaptation

Yuecong Xu, Jianfei Yang, Yunjiao Zhou

School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore
50 Nanyang Avenue, Singapore 639798

{xuyu0014, yang0478, yunjiao001}@e.ntu.edu.sg

Zhenghua Chen, Min Wu, Xiaoli Li

Institute for Infocomm Research (I²R), A*STAR, Singapore
1 Fusionopolis Way, #21-01, Connexis South, Singapore 138632

{chen.zhenghua, wumin}@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.

1. Introduction

Video Unsupervised Domain Adaptation (VUDA) [5, 8, 56, 49, 58] 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 distribution discrepancies and obtaining domain invariant representations [5, 8, 59]. 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, 31], 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 much more challenging than VUDA, since aligning distributions with limited samples is harder. A few research has touched on the image-based Few-Shot Domain Adaptation (FSDA) [26, 47, 54, 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 do-

main 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 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 [7, 64, 20, 58]. 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, outperforming the state-of-the-art methods by a large margin on two large-scale 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 [55, 62]. They could be generally divided into four categories: a) reconstruction-based

methods [14, 61], where domain-invariant features are obtained by encoders trained with data-reconstruction objectives; b) adversarial-based methods [5, 56, 8], where feature generators obtain domain-invariant features leveraging domain discriminators, trained jointly in an adversarial manner [16, 10]; c) semantic-based methods [63, 58], which exploit the shared semantics across domains such that domain-invariant features are obtained; and d) discrepancy-based methods [32, 67], which mitigate domain shifts by applying metric learning, minimizing metrics such as MMD [24] and CORAL [37]. With the introduction of cross-domain video datasets such as Daily-DA [59] and Sports-DA [59], there has been a significant increase in research interest for VUDA [8, 27, 6]. 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.

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 [26] is adversarial-based and augments the domain discriminator to classify 4 types of source-target pairs. d-SNE [54] learns a latent space through SNE [15] 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 [68, 40, 38, 65] 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 Adaptation task), thus $N_T = k \times \mathcal{C}$. \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 much more challenging than VUDA. Current VUDA methods [5, 56] 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 [46] and STPN [51]) 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 $N_T = k \times \mathcal{C}$ 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 much 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 [43]. 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 [23] and only 40% parameters of ViViT [1]. The feature of snippet s_b^a is $\mathbf{f}_b = Time(s_b^a)$ 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.

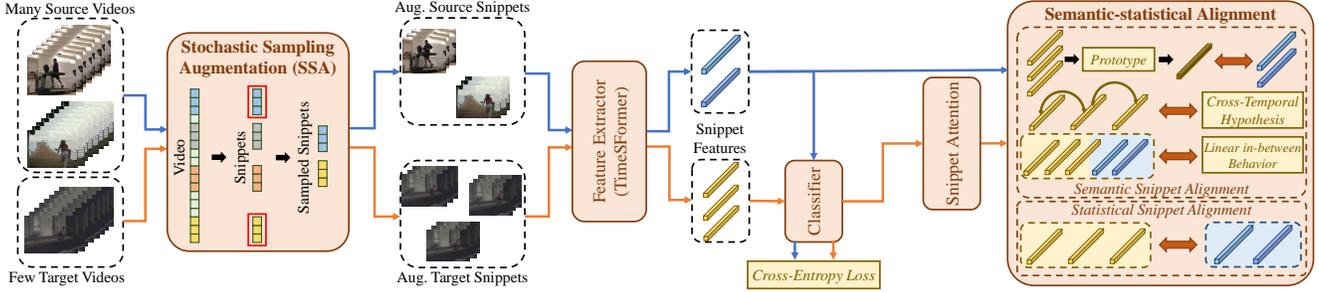


Figure 1. Pipeline of SSA^2lign . Source and target snippets are first obtained through the Stochastic Sampling Augmentation, whose features are obtained through the shared feature extractor. SSA^2lign then aligns the source and target domains at the snippet level with the Semantic-statistical Alignment, 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.

Following the above strategy, we propose the **Snippet-attentive Semantic-statistical Alignment (SSAlign)**, with the input obtained through SSA introduced in Sec. 3.1, forming the SSA^2lign . The overall pipeline of SSA^2lign 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 $f_{S,i}$, while the l -th target snippet ($l \in [1, r]$) from the j -th target video as $s_{T,jl}$ and its feature as $f_{T,jl}$. 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:

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

where $o_{S,i} = \sigma(H(f_{S,i}))$ and $o_{T,jl} = \sigma(H(f_{T,jl}))$ are the output predictions of snippet features $f_{S,i}$ and $f_{T,jl}$, 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 [33, 22] designed for few-shot learning [70, 45], 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:

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

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:

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

$n_{S,x}$ 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 [58] 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 [7, 64, 20, 58] 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 [7, 64, 20, 58]. We thus aim to drive the feature extractor towards obtaining snippet features that are distributed more discriminatively. Specifically, results in model robustness [66] 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 [44] 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:

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

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 [21, 60] 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

Algorithm 1 Training with SSA²lign for FSVDA

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

$(N_T \times r + N_S)$ 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 [24], CORAL [37], and MDD [67]. Compared to the adversarial-based adaptation strategy more commonly used in prior VUDA tasks [5, 56, 8], minimizing discrepancies does not require additional network structures (e.g., domain classifiers), thus is more stable. The MDD [67] 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 since it is intuitive that the importance of each target snippet is uneven. We thus propose a snippet attention to weigh the impact of different target snippets on the domain alignment dynamically. Intuitively, a snippet whose output prediction is the most accurate, i.e., whose classification is closest to its given ground truth, should be focused during alignment. A simple yet effective expression of how accurate the snippet’s output prediction is would be the inverse of the cross-entropy loss. The snippet attention weights are therefore built upon the inverse of the cross-entropy loss of the snippet, along with a residual connection for more stable opti-

Methods	Publication	Daily-DA													Sports-DA						
		H→A	M→A	KD→A	A→H	M→H	KD→H	H→M	A→M	KD→M	M→KD	H→KD	A→KD	Avg.	KS→U	S→U	U→S	KS→S	U→KS	S→KS	Avg.
TSF	-	37.859	32.584	31.110	44.583	57.083	45.833	36.500	30.000	34.500	61.656	58.897	75.724	45.527	91.657	91.069	76.368	77.737	87.768	85.118	84.953
TSF w/ T	-	39.565	39.488	39.410	61.667	62.917	62.500	41.500	38.750	36.500	77.793	80.276	83.586	55.329	92.480	93.420	78.947	79.052	88.021	87.003	86.487
TRX	CVPR-21	31.420	31.420	31.032	42.083	49.166	44.000	31.250	30.000	26.750	69.104	73.103	65.517	43.737	87.074	86.487	76.947	73.474	83.129	83.762	81.812
STRM	CVPR-22	33.825	32.351	32.894	43.333	50.833	44.417	30.750	29.500	28.250	72.138	74.620	68.965	45.156	91.539	90.012	78.579	75.158	86.901	84.628	84.470
HyRSM	CVPR-22	38.092	35.377	33.747	45.833	54.583	48.167	33.750	31.500	29.500	75.172	76.137	70.344	47.684	92.714	90.717	79.526	76.684	87.054	84.883	85.263
DANN	ICML-15	37.471	39.721	38.557	65.417	61.667	55.833	43.750	41.250	42.000	73.655	79.173	83.173	55.139	93.067	92.127	79.211	81.316	85.525	88.634	86.647
MK-MMD	ICML-15	35.299	42.746	35.609	64.167	63.333	56.667	44.000	41.750	36.500	76.690	81.931	79.862	54.879	92.597	93.420	80.737	77.842	84.760	88.124	86.247
MDD	ICML-19	42.514	42.281	42.901	64.583	64.167	57.917	45.000	39.500	37.750	75.173	81.517	84.276	56.465	93.184	93.067	78.474	79.105	86.697	87.716	86.374
SAVA	ECCV-20	39.178	41.660	41.738	63.333	63.333	60.000	42.750	41.500	39.250	77.517	81.242	80.690	56.016	93.302	91.540	79.263	80.474	86.378	87.512	86.617
ACAN	TNNLS-22	43.832	43.755	43.677	65.417	66.667	66.250	45.750	43.000	40.750	82.483	84.966	84.414	59.247	95.770	96.710	80.158	80.263	88.327	88.583	88.302
FADA	NeurIPS-17	39.100	42.126	32.351	46.250	58.750	47.500	37.250	30.750	35.250	77.241	81.103	77.517	50.432	93.655	93.655	76.947	78.316	88.736	86.086	86.233
d-SNE	CVPR-19	41.583	44.065	38.014	67.083	65.417	61.667	44.500	43.250	41.000	78.759	82.759	83.448	57.629	95.417	94.830	81.105	82.316	89.755	83.509	87.822
SSA ² lign	-	52.133	52.211	51.746	78.333	75.417	74.583	47.750	46.750	48.250	84.690	86.483	89.655	65.667	98.589	98.237	87.263	88.105	92.966	93.017	93.029

Table 1. Results for 10-shot ($k = 10$) FSVDA on Daily-DA and Sports-DA.

Methods	Publication	Daily-DA													Sports-DA						
		H→A	M→A	KD→A	A→H	M→H	KD→H	H→M	A→M	KD→M	M→KD	H→KD	A→KD	Avg.	KS→U	S→U	U→S	KS→S	U→KS	S→KS	Avg.
TSF	-	37.859	32.584	31.110	44.583	57.083	45.833	36.500	30.000	34.500	61.656	58.897	75.724	45.527	91.657	91.069	76.368	77.737	87.768	85.118	84.953
TSF w/ T	-	40.186	40.031	37.083	60.043	60.043	52.960	34.750	36.000	33.250	79.448	66.207	69.103	50.759	91.892	93.302	78.736	78.315	87.611	86.799	86.169
TRX	CVPR-21	32.794	30.260	29.887	39.425	47.446	40.349	29.000	27.750	24.750	69.545	63.032	55.882	40.768	86.531	86.955	76.342	70.946	81.827	80.411	80.502
STRM	CVPR-22	35.318	32.512	30.979	40.150	47.494	39.300	27.250	26.250	26.500	72.489	62.580	57.777	41.550	91.003	90.017	77.844	73.494	86.378	82.140	83.479
HyRSM	CVPR-22	39.065	35.088	31.061	42.736	52.098	43.740	31.000	29.250	28.000	75.646	63.889	58.598	44.181	92.166	90.769	79.122	75.272	86.251	81.713	84.215
DANN	ICML-15	40.496	38.789	36.385	60.833	58.750	52.917	41.750	38.500	39.750	74.345	66.759	69.655	51.577	92.245	93.067	78.421	75.631	85.117	82.161	84.440
MK-MMD	ICML-15	38.867	43.910	34.600	58.333	57.083	54.583	42.250	35.000	35.500	75.311	68.138	69.380	51.080	92.010	93.184	78.737	75.000	84.505	83.486	84.487
MDD	ICML-19	41.893	42.669	38.402	61.250	62.500	55.417	43.250	40.000	38.500	75.724	68.552	70.896	53.254	92.715	92.597	79.105	79.790	86.544	83.588	85.723
SAVA	ECCV-20	40.962	37.238	38.247	60.000	62.917	55.833	40.750	38.750	35.250	77.793	67.448	68.965	52.013	92.480	92.832	78.053	76.211	83.129	81.702	84.068
ACAN	TNNLS-22	44.453	44.298	41.350	63.333	63.333	56.250	39.000	40.250	37.500	80.552	70.897	73.793	54.584	95.182	96.592	79.947	79.526	88.175	88.379	87.967
FADA	NeurIPS-17	40.747	41.584	30.665	43.034	55.534	41.058	33.000	27.000	32.500	77.792	67.185	64.736	46.186	93.009	93.860	76.047	76.173	87.502	83.375	84.993
d-SNE	CVPR-19	41.994	44.162	35.738	64.927	63.365	56.936	41.250	41.750	39.750	79.448	72.530	73.296	54.596	94.992	94.734	80.983	81.768	89.371	81.749	87.266
SSA ² lign	-	52.366	51.978	47.401	76.667	72.917	70.417	47.000	46.250	47.500	86.759	79.310	81.793	63.363	97.062	97.885	84.053	86.211	91.182	90.214	91.101

Table 2. Results for 5-shot ($k = 5$) FSVDA on Daily-DA and Sports-DA.

Methods	Publication	Daily-DA													Sports-DA						
		H→A	M→A	KD→A	A→H	M→H	KD→H	H→M	A→M	KD→M	M→KD	H→KD	A→KD	Avg.	KS→U	S→U	U→S	KS→S	U→KS	S→KS	Avg.
TSF	-	37.859	32.584	31.110	44.583	57.083	45.833	36.500	30.000	34.500	61.656	58.897	75.724	45.527	91.657	91.069	76.368	77.737	87.768	85.118	84.953
TSF w/ T	-	37.937	34.135	33.049	51.250	58.750	46.667	37.750	34.500	35.250	74.069	60.414	63.724	47.291	91.165	92.832	75.368	76.578	86.595	85.372	84.652
TRX	CVPR-21	24.679	23.331	25.059	32.052	43.341	30.229	28.000	27.500	23.250	66.187	52.086	49.378	35.424	84.898	86.300	71.852	69.974	80.815	79.784	78.937
STRM	CVPR-22	28.037	24.181	26.924	33.853	45.022	30.038	25.750	26.000	25.000	69.137	52.130	50.492	36.380	89.640	89.176	73.719	72.496	85.682	81.630	82.057
HyRSM	CVPR-22	30.939	27.142	26.970	35.670	48.420	34.360	30.000	28.750	26.500	71.710	52.927	50.582	38.664	90.679	90.241	74.732	73.273	84.706	81.047	82.446
DANN	ICML-15	30.566	28.627	34.057	53.750	51.667	42.083	39.750	37.250	33.000	73.655	52.966	64.276	45.137	91.637	91.422	72.895	76.895	84.709	82.545	83.350
MK-MMD	ICML-15	29.403	32.506	31.963	54.583	55.417	44.167	38.500	37.000	33.750	72.000	56.000	63.035	45.694	90.717	91.070	74.684	74.211	85.830	84.047	83.426
MDD	ICML-19	31.652	33.592	34.253	54.167	56.667	47.083	42.250	38.500	34.750	70.621	56.138	59.448	46.616	91.422	92.715	74.369	74.895	82.823	82.925	83.191
SAVA	ECCV-20	31.031	33.436	32.971	50.833	58.333	42.917	40.500	39.750	37.750	72.138	55.724	62.069	46.454	89.307	92.832	73.211	73.842	81.753	80.377	81.887
ACAN	TNNLS-22	38.635	35.609	35.299	55.000	61.667	46.667	38.250	38.750	35.750	76.276	59.311	62.621	48.653	93.750	96.240	75.368	77.052	85.658	87.105	85.862
FADA	NeurIPS-17	33.881	34.136	25.565	35.523	53.232	31.256	32.500	27.000	32.750	73.861	57.150	57.908	41.230	91.353	93.126	71.829	75.042	86.750	82.508	83.435
d-SNE	CVPR-19	36.263	37.859	32.131	59.195	60.914	50.006	40.500	41.000	38.500	76.293	64.298	62.934	49.991	93.896	94.492	76.440	77.111	87.857	81.049	85.141
SSA ² lign	-	44.831	46.780	45.306	68.750	70.833	62.083	46.750	46.000	45.000	79.724	68.793	71.586	57.828	96.592	97.415	80.053	80.947	88.940	89.755	88.950

Table 3. Results for 3-shot ($k = 3$) FSVDA on Daily-DA and Sports-DA.

snippet, expressed as:

$$w_{jl} = 1 + \frac{1}{l_{ce}(o_{T,jl}, y_{j,T})}. \quad (9)$$

The snippet attention weights are subsequently normalized across the r snippets corresponding to the same target video, expressed as $\bar{w}_{jl} = w_{jl} / \sum_{l'=1}^r w_{jl'}$. The normalized snippet attention weight \bar{w}_{jl} is then applied to the target snippet features, forming the weighted target snippet features by $\mathbf{f}'_{T,jl} = \bar{w}_{jl} \mathbf{f}_{T,jl}$, which are then aligned with the source domain through the semantic and statistical snippet alignments by replacing the features $\mathbf{f}_{T,jl}$ with $\mathbf{f}'_{T,jl}$.

SSA²lign. Finally, we sum up our proposed SSA²lign in Algorithm 1. The snippet features, SSA, and snippet attention are leveraged only during training. During testing, target video representations are obtained by uniform sampling across the target testing videos, while the video features and their output predictions are obtained by directly applying the trained feature extractor and classifier to the uniformly sampled target video representations.

4. Experiments

In this section, we evaluate our proposed SSA²lign across two challenging cross-domain action recognition benchmarks: Daily-DA and Sports-DA [59], which cover 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.

4.1. Experimental Settings

Daily-DA is a challenging dataset that has been leveraged in prior VUDA works [59, 58, 60]. It covers both normal and low

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 [69] as the feature extractor, pre-trained on Kinetics-400 [18]. All experiments are implemented with the PyTorch [29] 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 [26], d-SNE [54] designed for image-based FSDA; DANN [10], MK-MMD [24], MDD [67], SAVA [8] and ACAN [56], designed for UDA/VUDA; and TRX [30], STRM [41], and HyRSM [50] 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. No-

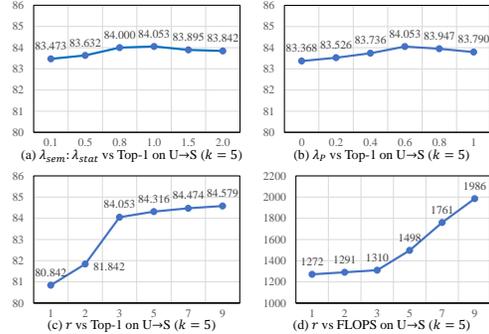


Figure 2. Sensitivity of hyper-parameters on U→S task.

tably, 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 and Analysis

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 [37] and MDD [67]; 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} ; SSAlign (w/ spatial aug.) augments target domain 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

Methods	Components						Daily-DA									Sports-DA						Avg.
	SSA	Sn-Attn	\mathcal{L}_{proto}	\mathcal{L}_{cross}	$\mathcal{L}_{sn-stat}$	$\mathcal{L}_{sn-stat}$	$k=10$			$k=5$			$k=3$			$k=10$		$k=5$		$k=3$		
							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	✓						41.660	41.971	42.048	42.824	42.281	38.247	38.790	36.385	35.221	81.053	81.210	80.789	80.421	75.894	77.210	55.734
SSA ² lign	✓	✓	✓	✓	✓	✓	45.616	46.315	45.695	46.470	45.686	41.738	39.168	40.962	39.488	84.316	85.473	81.053	83.579	77.264	77.368	58.679
	✓	✓	✓	✓	✓	✓	51.047	51.125	50.427	50.970	50.729	46.392	43.590	45.850	44.220	86.579	87.368	83.211	85.632	79.316	80.368	62.455
	✓	✓	✓	✓	✓	✓	49.883	49.806	49.729	50.349	49.255	45.074	42.193	44.530	42.901	85.579	86.737	82.737	84.737	78.685	79.473	61.445
	✓	✓	✓	✓	✓	✓	50.194	50.504	49.651	50.271	49.875	45.772	42.659	45.151	43.056	85.842	87.052	82.579	85.053	78.948	79.579	61.746
	✓	✓	✓	✓	✓	✓	48.176	48.565	46.858	47.556	48.091	43.289	40.642	42.513	41.970	84.684	85.947	81.948	83.737	77.737	78.631	60.023
	✓	✓	✓	✓	✓	✓	51.357	51.435	50.893	51.823	51.427	46.548	43.900	46.004	44.918	86.684	87.631	83.632	85.685	79.685	80.421	62.798
	✓	✓	✓	✓	✓	✓	52.133	52.211	51.746	52.366	51.978	47.401	44.831	46.780	45.306	87.263	88.105	84.053	86.211	80.053	80.947	63.425

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	51.900	51.978	51.513	51.978	51.582	47.091	44.521	46.392	45.073	87.052	87.842	83.895	86.000	79.895	80.684	63.160	-0.265	1302	-8
SSA ² lign-MMD	51.668	51.901	51.281	51.978	51.505	47.013	44.521	46.315	44.841	87.052	87.842	83.842	85.895	79.790	80.631	63.072	-0.353	1312	+2
SSA ² lign-FC	52.831	52.909	52.367	52.987	52.358	47.556	45.296	47.245	45.306	87.158	88.263	84.527	86.685	79.790	81.263	63.769	+0.344	1472	+162
SSA ² lign-SP	50.969	51.202	50.970	51.280	51.272	46.470	43.822	45.616	44.685	86.895	87.421	83.369	85.737	79.264	80.473	62.630	-0.795	1390	+80
SSA ² lign (w/ spatial aug.)	45.383	46.703	45.230	45.772	45.841	41.117	39.633	40.264	38.635	83.527	85.631	80.579	83.368	76.948	77.736	58.424	-5.001	1325	+15
SSA ² lign	52.133	52.211	51.746	52.366	51.978	47.401	44.831	46.780	45.306	87.263	88.105	84.053	86.211	80.053	80.947	63.425	-	1310	-

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

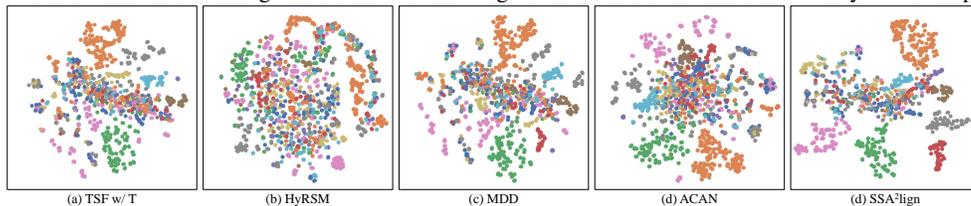


Figure 3. t-SNE visualizations of target features from (a) TSF w/T, (b) HyRSM, (c) MDD, (d) ACAN, (e) SSA²lign. Colors denote classes.

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 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.

Alignment Methods. Table 5 shows that while formulating $\mathcal{L}_{sn-stat}$ as MDD [67] 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 overhead (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 [42] 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.

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.

Appendix

This appendix presents more details of the proposed Snippet-attentive Semantic-statistical Alignment with Stochastic Sampling Augmentation (SSA²lign) and is organized as follows: first, we introduce the detailed implementation of SSA²lign with specific hyper-parameter settings, supported by additional results of hyper-parameter sensitivity analysis to show the robustness of SSA²lign. Subsequently, we present details of the cross-domain action recognition benchmarks for evaluating SSA²lign, including Daily-DA and Sports-DA; lastly, we compare in detail our SSA²lign with related but different FSDA and UDA/VUDA methods to highlight our novelty.

Implementation Details

Brief Review of SSA²lign. In this work, we propose the Snippet-attentive Semantic-statistical **Alignment** with **Stochastic Sampling Augmentation (SSA²lign)** to address *Few-Shot Video Domain Adaptation* (FSVDA) by augmenting the source and target domains and performing domain alignment at the snippet level. SSA²lign 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. To further improve the stability of snippet-level alignment, a statistical alignment strategy is additionally adopted, while snippet attention is proposed to weigh the impact of different target snippets on the domain alignment dynamically. In this section, we present the detailed implementation of SSA²lign, whose pipeline is demonstrated in Fig. 4.

TimeSFormer as Feature Extractor. To obtain features from snippets during training and videos during testing, we instantiate the Transformer-based TimeSFormer [2] as the feature extractor thanks to its capability in obtaining features that include both spatial and temporal infor-

mation. TimeSFormer extracts spatial and temporal features with separate space-time attention blocks based on self-attention [43] and obtains very competitive results on various action recognition benchmarks [2]. While other Transformer-based video models, such as Swin [23] and ViViT [1], also achieve competitive performances on action recognition, TimeSFormer possesses the least amount of parameters, requiring only 60% parameters of Swin and only 40% parameters of ViViT. The final classifier is implemented as a single fully connected layer. Both the feature extractor and the subsequent classifier are shared across source and target data.

Training Details and Hyper-parameters. For training, we initialize the TimeSFormer feature extractor from pre-trained weights obtained by pre-training on Kinetics-400 [18]. For more efficient training, we freeze the first 8 blocks of TimeSFormer, leaving the last 4 blocks to be fully trainable, with the learning rate set at 0.005. 11 new layers are trained from scratch, with their learning rates set to be 10 times that of the pretrained-loaded trainable layers (blocks). For the tasks constructed from the Daily-DA dataset [59], we train a total of 30 epochs, while we train a total of 50 epochs for tasks constructed from the Sports-DA dataset [59]. The stochastic gradient descent (SGD) algorithm [3] is used for optimization, with the weight decay set to 0.0001 and the momentum set to 0.9. During the training phase of SSA²lign, the batch size is set to 24 input snippets per GPU, with 12 source snippets from 12 source videos and 12 target snippets from 4 target videos ($r = 3$ by default). For a fair comparison, the batch size is set to 24 input videos per GPU when training all comparing methods. All experiments are implemented with the PyTorch [29] library and conducted on 2 NVIDIA A6000 GPUs. We set the length of snippets and the number of snippets per target video via SSA empirically as $m = 8, r = 3$. Hyper-parameters $\lambda_{sem} = 1.0, \lambda_{stat} = 1.0$ and $\lambda_P = 0.6$ are empirically set and are fixed. As shown in Section 4.3 and Fig. 2 of the paper, the performance of SSA²lign is robust to hyper-parameters $\lambda_{sem}, \lambda_{stat}$ and λ_P as well as r when $r \geq 3$, with minimal variations and maintains the best results with high computation efficiency with all the default hyper-parameter settings. To further illustrate the robustness of SSA²lign towards 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, we present the additional results of hyper-parameter sensitivity analysis under different experimental settings. Specifically, we present the results of the U→S task with $k = 10, k = 5$ (the same as presented in Fig. 2 of the paper), and the results of the KS→S task with $k = 10$, as shown in Fig. 5 of this appendix.

The additional results further justify that SSA²lign is ro-

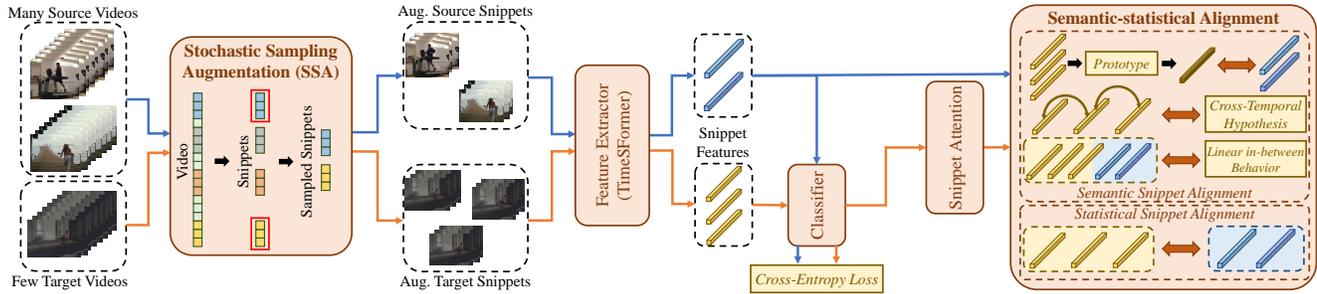


Figure 4. Pipeline of SSA²lign. Source and target snippets are first obtained through the Stochastic Sampling Augmentation, whose features are obtained through the shared feature extractor. SSA²lign then aligns the source and target domains at the snippet level with the Semantic-statistical Alignment, 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. *Best viewed in color.*

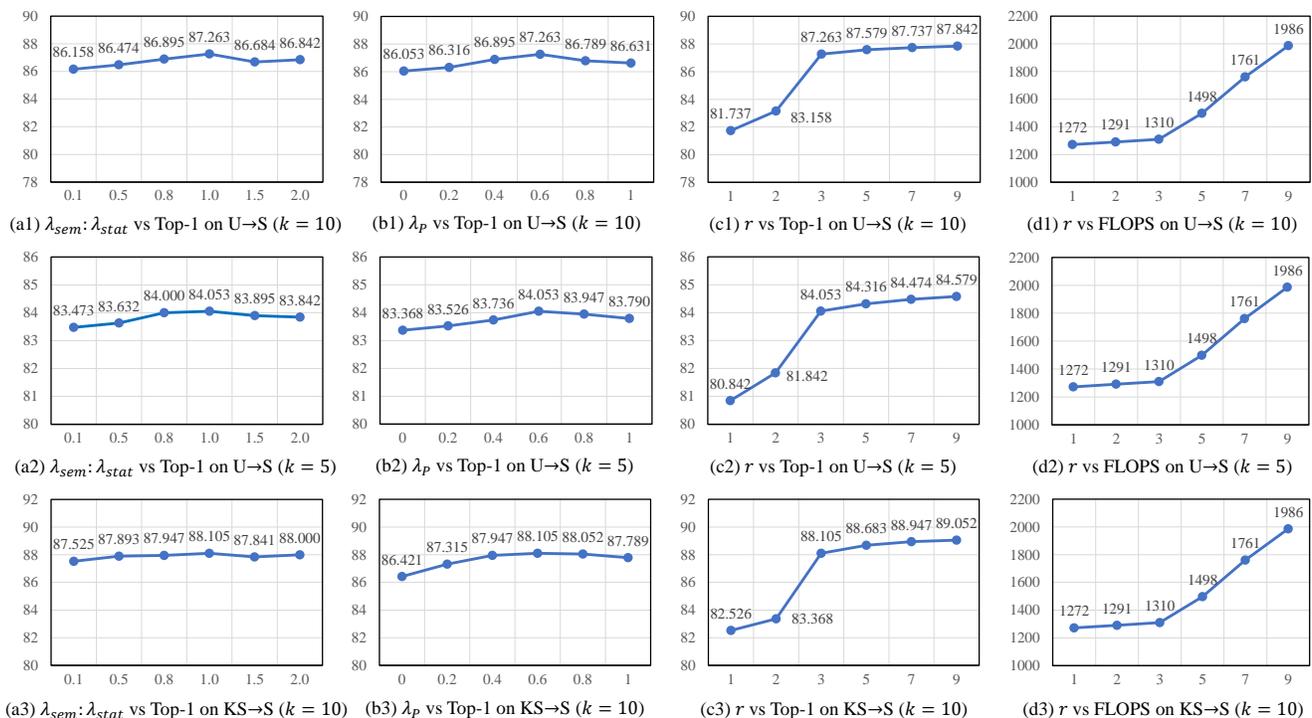


Figure 5. Hyper-parameter sensitivity on the U→S task with $k = 10$ (top), $k = 5$ (mid), and the KS→S task with $k = 10$ (bot).

bust to hyper-parameters λ_{sem} , λ_{stat} and λ_P as well as r when $r \geq 3$ under all examined experimental settings, while achieving the best results with high computation efficiency with the default hyper-parameter settings.

Cross-domain Action Recognition Benchmarks

In this paper, to evaluate our proposed SSA²lign, we utilized two cross-domain action recognition benchmarks: the Daily-DA and Sports-DA [59]. In this section, we provide more details on each benchmark.

Daily-DA

The Daily-DA dataset is a recently proposed cross-domain action recognition dataset for VUDA [59]. It is more comprehensive and challenging compared to prior benchmarks such as UCF-Olympic [36] and UCF-HMDB_{full} [5] which have resulted in saturated performance due to limited domains (only 2 domains in each dataset) and number of videos per domain. Daily-DA includes videos of daily actions from four domains and incorporates both normal videos and low-illumination videos. Specifically, Daily-DA is built from four datasets: the dark dataset ARID (A) [57],

Statistics	Daily-DA	Sports-DA
Video Classes #	8	23
Training Video #	A:2,776 / H:560 / M:4,000 / KD:8,959	U:2,145 / S:14,754 / KS:19,104
Testing Video #	A:1,289 / H:240 / M:400 / KD:725	U:851 / S:1,900 / KS:1,961

Table 6. Summary of cross-domain action recognition benchmarks statistics.

Class ID	ARID Class	HMDB51 Class	Moments-in-Time Class	Kinetics-600 Class
0	Drink	drink	drinking	drinking shots
1	Jump	jump	jumping	jumping bicycle, jumping into pool, jumping jacks
2	Pick	pick	picking	picking fruit
3	Pour	pour	pouring	pouring beer
4	Push	push	pushing	pushing car, pushing cart, pushing wheelbarrow, pushing wheelchair
5	Run	run	running	running on treadmill
6	Walk	walk	walking	walking the dog, walking through snow
7	Wave	wave	waving	waving hand

Table 7. List of action classes for Daily-DA.

Class ID	UCF101 Class	Sports-1M Class	Kinetics-600 Class
0	Archery	archery	archery
1	Baseball Pitch	baseball	catching or throwing baseball, hitting baseball
2	Basketball Shooting	basketball	playing basketball, shooting basketball
3	Biking	bicycle	riding a bike
4	Bowling	bowling	bowling
5	Breaststroke	breaststroke	swimming breast stroke
6	Diving	diving	springboard diving
7	Fencing	fencing	fencing (sport)
8	Field Hockey Penalty	field hockey	playing field hockey
9	Floor Gymnastics	floor (gymnastics)	gymnastics tumbling
10	Golf Swing	golf	golf chipping, golf driving, golf putting
11	Horse Race	horse racing	riding or walking with horse
12	Kayaking	kayaking	canoeing or kayaking
13	Rock Climbing Indoor	rock climbing	rock climbing
14	Rope Climbing	rope climbing	climbing a rope
15	Skate Boarding	skateboarding	skateboarding
16	Skiing	skiing	skiing crosscountry, skiing mono
17	Sumo Wrestling	sumo	wrestling
18	Surfing	surfing	surfing water
19	Tai Chi	t'ai chi ch'uan	tai chi
20	Tennis Swing	tennis	playing tennis
21	Trampoline Jumping	trampolining	bouncing on trampoline
22	Volleyball Spiking	volleyball	playing volleyball

Table 8. List of action classes for Sports-DA.

as well as HMDB51 (H), Moments-in-Time (M) [25], and Kinetics-600 (KD) [4], which are video datasets widely used for action recognition benchmarking [28]. Compared with other action recognition datasets such as Moments-in-Time and Kinetics, ARID is comprised of videos shot under adverse illumination conditions, characterized by low brightness and low contrast. Statistically, the RGB mean and standard deviation values (std) of videos in ARID are much lower among datasets leveraged in Daily-DA [56], which strongly suggests a larger domain shift between ARID and the other action recognition datasets. The Daily-DA includes a total of 16,295 training videos and 2,654 testing videos from 8 categories as listed in Table 6, with each category corresponding to one or more categories in

the original datasets as demonstrated in Table 7.

Sports-DA

To further demonstrate the efficacy of our proposed SSA²lign on large-scale cross-domain datasets, we further adopt the Sports-DA dataset as another cross-domain action recognition benchmark. Comparatively, Sports-DA contains almost double the amount of training and testing videos of Daily-DA. Specifically, it includes a total of 36,003 training videos and 4,721 testing videos from 23 categories of sports actions, collected from three large-scale datasets: UCF101 (U) [34], Sports-1M (S) [17], and Kinetics-600 (KS) [4], as shown in Table 6. Similar to

Method	Publication	Task	Techniques
d-SNE [54]	CVPR-19	Few-Shot Domain Adaptation (FSDA): source image data available with labels, a few (very limited) target image data available with labels, image-based.	(a) d-SNE learns a latent domain-agnostic space through SNE [15] with large-margin nearest neighborhood [9]; (b) d-SNE conducts FSDA in a min-max formulation with a modified-Hausdorff distance; (c) d-SNE creates sibling target samples with spatial augmentations, and trains feature extractor with the Mean-Teacher technique [39].
PASTN [12]	TIP-20	Few-Shot Video Domain Adaptation (FSVDA): source video data available with labels, a few (very limited) target video data available with labels, video-based.	(a) PASTN obtains video features from a frame-based video model; (b) PASTN forms source-target video pairs to address insufficient target video data; (c) PASTN constructs pairwise adversarial networks performed across source-target video pairs optimized by a pairwise margin discrimination loss [52].
DM-ADA [53]	AAAI-20	Unsupervised Domain Adaptation (UDA): source image data available with labels, sufficient target images available without labels, image-based.	(a) DM-ADA augments the target domain with the source domain by domain mixup [66]; (b) DM-ADA improves the feature extractor by leveraging soft domain labels; (c) DM-ADA jointly trains a domain discriminator which judges the samples' differences relative to the two domains with refined scores.
ACAN [56]	TNNLS-22	Video Unsupervised Domain Adaptation (VUDA): source video data available with labels, sufficient target videos available without labels, video-based.	(a) ACAN applies adversarial-based domain adaptation across spatio-temporal video features; (b) ACAN additionally aligns video correlation features in the form of long-range spatiotemporal dependencies [48]; (c) ACAN further aligns the joint distribution of correlation information of different domains by minimizing pixel correlation discrepancy (PCD).
SSA ² lign (Ours)	-	Few-Shot Video Domain Adaptation (FSVDA): source video data available with labels, a few (very limited) target video data available with labels, video-based.	(a) SSA ² lign addresses FSVDA at the snippet level instead of the frame or video-levels; (b) SSA ² lign augments target domain data and the snippet-level alignments by a simple yet effective stochastic sampling of snippets; (c) SSA ² lign performs both semantic and statistical alignments attentively, with the semantic alignments achieved by alignment based on the semantic information within each snippet, cross-snippets of each video, and across snippet-level data distribution.

Table 9. Detailed comparison of SSA²lign with related but different FS(V)DA and (V)UDA methods.

Daily-DA, each action class corresponds to one or more categories in the original datasets as presented in Table 8. With more than 40,000 training and testing videos, the Sports-DA benchmark is one of the largest cross-domain action recognition benchmarks introduced.

Detailed Comparison with Related FS(V)DA and (V)UDA Methods

In this paper, we proposed SSA²lign to address the more realistic and challenging FSVDA task, which achieves state-of-the-art performances with outstanding improvements on both cross-domain action recognition benchmarks (average 13.1% on Daily-DA tasks and average 4.2% on Sports-DA tasks). To further highlight the novelty of SSA²lign, we compare our proposed SSA²lign with prior FSDA/FSVDA and UDA/VUDA methods. Specifically, we compare with d-SNE [54] proposed for FSDA, PASTN [12] designed for FSVDA, ACAN [56] introduced for ACAN, and DM-ADA [53] which is an image-based UDA method that leverages MixUp [66]. These methods are all compared

from two perspectives: the tasks they tackle and the techniques leveraged, as displayed in Table 9.

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, 9
- [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, 9
- [3] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010. 9
- [4] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note about kinetics-600. *arXiv preprint arXiv:1808.01340*, 2018. 6, 11
- [5] Min-Hung Chen, Zsolt Kira, Ghassan AlRegib, Jaekwon Yoo, Ruxin Chen, and Jian Zheng. Temporal attentive alignment for large-scale video domain adaptation. In *Proceed-*

- ings of the *IEEE International Conference on Computer Vision*, pages 6321–6330, 2019. [1](#), [2](#), [3](#), [5](#), [10](#)
- [6] 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](#)
- [7] 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](#)
- [8] 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](#)
- [9] Carlotta Domeniconi, Dimitrios Gunopulos, and Jing Peng. Large margin nearest neighbor classifiers. *IEEE transactions on neural networks*, 16(4):899–909, 2005. [2](#), [12](#)
- [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](#), [12](#)
- [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] Geoffrey E Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15, 2002. [2](#), [12](#)
- [16] 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](#)
- [17] 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. [6](#), [11](#)
- [18] 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](#), [9](#)
- [19] 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. [6](#)
- [20] 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](#)
- [21] 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](#)
- [22] 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](#)
- [23] 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](#), [9](#)
- [24] 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](#)
- [25] 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. [6](#), [11](#)
- [26] 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](#)
- [27] Boxiao Pan, Zhangjie Cao, Ehsan Adeli, and Juan Carlos Nieves. Adversarial cross-domain action recognition with co-attention. In *AAAI*, pages 11815–11822, 2020. [3](#), [9](#)
- [28] Preksha Pareek and Ankit Thakkar. A survey on video-based human action recognition: recent updates, datasets, challenges, and applications. *Artificial Intelligence Review*, 54:2259–2322, 2021. [11](#)
- [29] 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](#), [9](#)
- [30] Toby Perrett, Alessandro Masullo, Tilo Burghardt, Majid Mirmehdi, and Dima Damen. Temporal-relational crosstransformers for few-shot action recognition. In *Pro-*

- ceedings of the *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 475–484, 2021. 7
- [31] Abdel Mlak Said, Aymen Yahyaoui, and Takoua Abdellatif. Efficient anomaly detection for smart hospital iot systems. *Sensors*, 21(4):1026, 2021. 1
- [32] 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
- [33] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 4
- [34] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 11
- [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] Waqas Sultani and Imran Saleemi. Human action recognition across datasets by foreground-weighted histogram decomposition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 764–771, 2014. 10
- [37] 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
- [38] 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
- [39] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. 12
- [40] 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
- [41] 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
- [42] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 8
- [43] 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, 9
- [44] 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
- [45] 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
- [46] 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
- [47] 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
- [48] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 12
- [49] 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
- [50] 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
- [51] 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
- [52] Chao-Yuan Wu, R Manmatha, Alexander J Smola, and Philipp Krahenbuhl. Sampling matters in deep embedding learning. In *Proceedings of the IEEE international conference on computer vision*, pages 2840–2848, 2017. 12
- [53] Minghao Xu, Jian Zhang, Bingbing Ni, Teng Li, Chengjie Wang, Qi Tian, and Wenjun Zhang. Adversarial domain adaptation with domain mixup. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6502–6509, 2020. 12
- [54] 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, 12
- [55] 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
- [56] 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, 11, 12
- [57] Yuecong Xu, Jianfei Yang, Haozhi Cao, Kezhi Mao, Jianxiong Yin, and Simon See. Arid: A new dataset for rec-

- ognizing action in the dark. In *International Workshop on Deep Learning for Human Activity Recognition*, pages 70–84. Springer, 2021. [6](#), [10](#)
- [58] 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](#), [6](#)
- [59] 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](#), [9](#), [10](#)
- [60] 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](#), [6](#)
- [61] 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](#)
- [62] 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](#)
- [63] 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](#)
- [64] 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](#)
- [65] 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](#)
- [66] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *International Conference on Learning Representations*, 2018. [5](#), [12](#)
- [67] 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](#)
- [68] 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](#)
- [69] 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](#)
- [70] 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](#)