

SADA: Semantic adversarial unsupervised domain adaptation for Temporal Action Localization

David Pujol-Perich, Albert Clapés, and Sergio Escalera,

Abstract—Temporal Action Localization (TAL) is a complex task that poses relevant challenges, particularly when attempting to generalize on new – unseen – domains in real-world applications. These scenarios, despite realistic, are often neglected in the literature, exposing these solutions to important performance degradation. In this work, we tackle this issue by introducing, for the first time, an approach for Unsupervised Domain Adaptation (UDA) in sparse TAL, which we refer to as Semantic Adversarial unsupervised Domain Adaptation (SADA). Our contributions are threefold: (1) we pioneer the development of a domain adaptation model that operates on realistic sparse action detection benchmarks; (2) we tackle the limitations of global-distribution alignment techniques by introducing a novel adversarial loss that is sensitive to local class distributions, ensuring finer-grained adaptation; and (3) we present a novel set of benchmarks based on EpicKitchens100 and CharadesEgo, that evaluate multiple domain shifts in a comprehensive manner. Our experiments indicate that SADA improves the adaptation across domains when compared to fully supervised state-of-the-art and alternative UDA methods, attaining a performance boost of up to 6.14% mAP.

Index Terms—Artificial Intelligence, Computer Vision, Video Understanding, Domain Adaptation

1 INTRODUCTION

RECENT advances in the field of video understanding have played a critical role in the surge of novel video-based applications –e.g., video indexing, summarization, recommendation, or surveillance. A critical task of this field is *Temporal Action Localization* (TAL), which involves identifying actions in a video consisting of both their time intervals and action categories. This is particularly difficult given the inherent variabilities of videos. These can be presented, among others, in the form of *appearance variability* –e.g., different kitchens and/or lighting conditions–, *acquisition variability* –e.g., different recording devices– or *viewpoint variability* –e.g., first- or third-person. All in all, these prompt a certain degree of confusion between similar actions to be discriminated.

Traditionally, fully supervised methods attempt to address this issue by leveraging enough training data to cover all the possible sources of variability. Unfortunately, this becomes virtually impossible when dealing with realistic scenarios. This compels these methods to operate under the influence of unseen data variations – i.e., *domain gaps* –, which exposes them to a considerable decline in performance. Overcoming this typically involves relabeling data from the new domain, so as to retrain and adapt the model. Unfortunately, this approach is impractical due to the considerable time and resource consumption involved, a

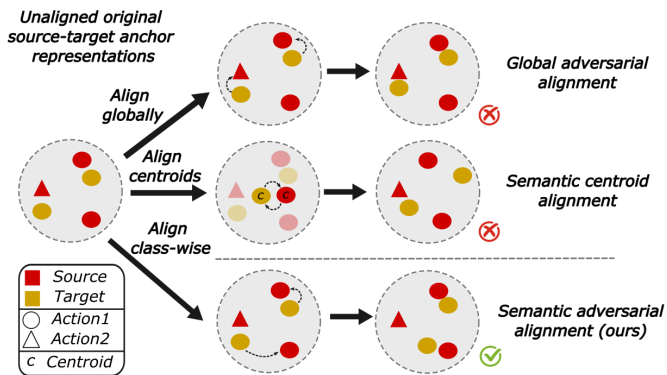


Fig. 1: Illustration of the differences between the two most similar domain-adaptation methods [1], [2], and our proposal, *SADA*. For this, we present a simple scenario with various anchor embeddings of different actions (identified by shapes) and domains (identified by colors). In this scenario, [1] (upper row) aligns embeddings in a class-agnostic manner, making it liable to aligning domain embeddings of unmatching action labels. [2] (middle row) computes class-wise mean centroids, and aligns them across domains, but as shown, minimizing their distance does not yield a proper adaptation. *SADA* (last row) improves [1] by aligning class-wise distributions, yielding the correct alignment by not aligning unmatching anchors.

challenge exacerbated when dealing with high-dimensional inputs like videos.

Unsupervised Domain Adaptation (UDA) has recently become a hot topic given its potential to leverage unlabelled data to mitigate this domain-induced degradation [3], [4], [5]. Unfortunately, despite its considerable success in image-based tasks [6], the application of UDA to video

- David Pujol-Perich is with the Computer Vision Center, Spain and with the Department of Mathematics and Informatics, University of Barcelona, Spain. E-mail: see <https://davidpujol.github.io>
- Albert Clapés is with the Computer Vision Center, Spain and with the Department of Mathematics and Informatics, University of Barcelona, Spain. E-mail: aclapes@ub.com
- Sergio Escalera is with the Computer Vision Center, Spain and with the Department of Mathematics and Informatics, University of Barcelona, Spain. E-mail: sescalera@ub.com

understanding remains underexplored. In fact, to the best of our knowledge, no prior work addresses UDA for TAL setups. The closest proposal is SSTDA [7], which approaches the problem of action segmentation. This focuses on making per-frame action predictions (i.e., *dense TAL*), tackling datasets where no concurrent actions take place [8], [9]. This approach inevitably requires additional smoothing techniques to preserve temporal coherence. This motivates the focus of this paper on *sparse TAL* –i.e., segment-level predictions– avoiding the need for additional losses, while intrinsically adapting to multi-label scenarios.

Consequently, in this paper, we propose the first Un-supervised Domain Adaptation method for sparse multi-label detection on TAL, which we name *Semantic Adversarial unsupervised Domain Adaptation* or *SADA* for short. Our goal is to minimize the discrepancy between feature representations of a labeled source domain and an unlabelled target domain. These features are extracted with a multi-resolution architecture [10] that we couple with a novel adversarial loss that improves the limitations of existing UDA works. Concretely, existing works normally align domain distributions globally [3]. We propose instead to use pseudo-labeling [11] to assign an action or background (no action) class to each feature representation. With this, our loss factorizes the global adaptation strategy into independent per-class and background alignments – i.e., aligning each action’s distribution across both domains. This results in a more sensitive alignment strategy compared to a global distribution approach, less prone to *feature misalignment* and, as we will show, better performing.

Assessing the effectiveness of UDA methods for video understanding is a challenging, still unresolved task. Existing proposals on action segmentation [7] follow a subject-based strategy where they aim to adapt a model to new unseen subjects. Here we refer to *subject* as a person appearing in a video. Nevertheless, little data is normally available from a single subject, which inevitably requires grouping several of them for training. This allows the model to generalize over the subject variability under study, making it unsuitable for domain adaptation. To address these limitations, we draw inspiration from the work of [12] on action recognition and investigate the impact of *viewpoint domain shifts* on sparse TAL in CharadesEgo [13]. However, we contend that a more comprehensive evaluation necessitates setups with more controllable shifts. For this, we also propose a suite of 6 new setups based on EpicKitchens100 [14] which study the effect *appearance* and *acquisition* domain shifts. These comprehensive benchmarks demonstrate that *SADA* mitigates the performance degradation, improving by a large margin the existing fully supervised (namely *source-only*) and UDA-based proposals. In short, our main contributions are:

- 1) We propose for the first time an unsupervised domain adaptation method suitable for sparse detection scenarios on TAL.
- 2) We introduce a novel adversarial loss that factorizes standard global alignment into independent class- and background-wise alignments (see Fig. 1).
- 3) We present new benchmarks to test sparse detection scenarios when facing 7 different domain shifts,

improving the state-of-the-art in all of them.

2 RELATED WORK

Temporal Action Localization. At the time of this writing, most of the literature on the task of *Temporal Action Localization* follows a traditional *source-only* approach. In other words, they restrict the models’ visibility solely to a training domain, while other domains seen during testing are not available. These works can be categorized as follows: **(1) Anchor-based methods** [10], [15], [16], [17], [18], [19], [20], [21] propose a two-stage pipeline, consisting of a proposal generation and classification. The first applies heuristic methods – e.g., uniform sampling [10], [15], [20], [21] or action boundaries’ grouping [22], [23] – to generate a dense set of proposals – i.e., temporal segments. In the second stage, they leverage a learnable classifier to predict the corresponding action class and localization offsets of every anchor. Our work falls into this category motivated by the recent success of these methods achieving state-of-the-art results in many TAL benchmarks [14], [24], [25] **(2) Anchor-free methods** [17], [26], [27], [28] avoid this two-stage approach making per-frame predictions of their corresponding action labels. These methods, however, often suffer from a tendency towards over-segmentation given the potential discrepancy between neighboring frames. Consequently, they require often complex smoothing techniques to improve the boundary predictions [7] **(3) Query-based methods** [29], [30], [31] recently emerged as an alternative paradigm that follows the principles presented by [32]. This approach exploits the use of a Transformer encoder-decoder architecture [33] to learn a fixed small set of queries given refined video features, each identifying one potential action segment. Intuitively, this results in a non-heuristic-based proposal generation. This comes with the limitation of an increased rigidity, as the number of proposals needs to be fixed beforehand.

Unsupervised Domain Adaptation. Domain Adaptation techniques emerge as an effective solution to bridge the gap between data collected from a source and a target distribution, respectively. A large suite of approaches has been proposed to perform this alignment between labeled and unlabeled domains – e.g., discrepancy minimization [34], [35], entropy minimization [36], [37] or contrastive learning [38]. The most popular approach is arguably the use of adversarial training methods [3], [4], [5], [39], [40], [41], [42], which incorporate a domain classifier trained to discern if samples come from the source or the target domain, in a min-max fashion. This results in the computation of domain-invariant embeddings [43]. Despite convenient, the simplicity of these methods might degrade the quality of the alignment [11], as they potentially align embeddings of source and target domain that represent different semantic information – e.g., different class labels. Few works have been proposed to do this alignment in a more sensitive way [11]. We highlight the importance of [2], proposing to minimize instead the distance of per-class centroids, computed as the mean of the set of predicted feature embeddings of a given class. Its effectiveness, however, relies on the assumption that the data is distributed somewhat homogeneously around the center, as otherwise, the centroids are not necessarily

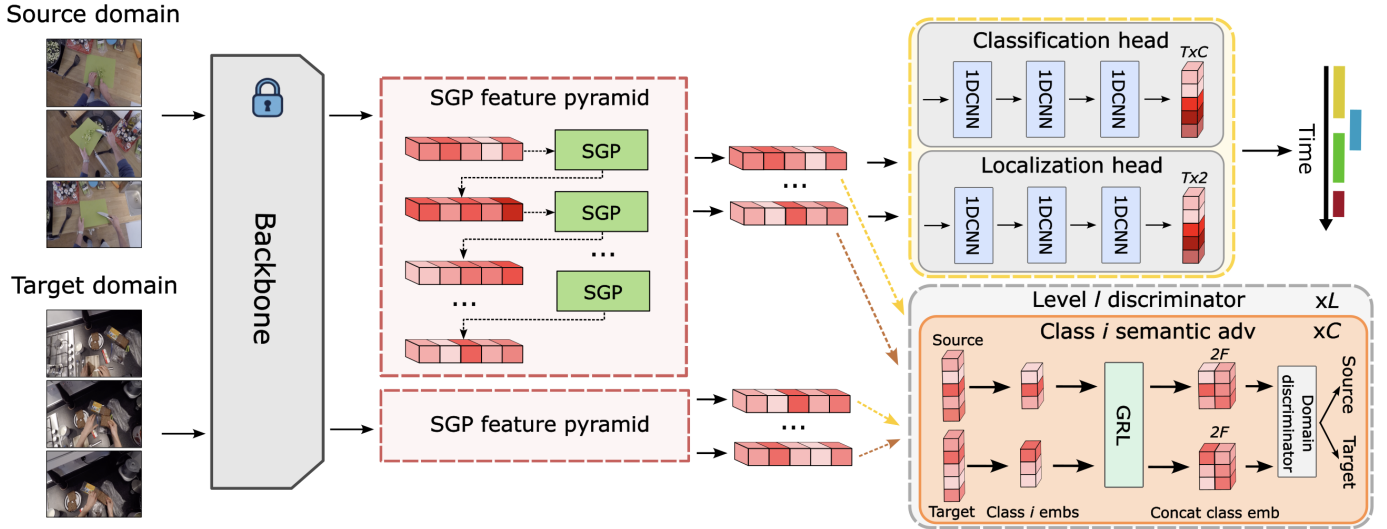


Fig. 2: Overview of the main model architecture of SADA. This takes as input videos from a Source and a Target domain, which are both fed to a shared multi-resolution feature extractor pyramid. The output embeddings of both of these domains are then aligned using the semantic alignment loss, *SADA*. This is done with a level and class-wise domain discriminator of the filtered embeddings, based on GT information and pseudo labels, for the source and target domains, respectively. Finally, the resulting domain invariant representations of the source domain are used to train a classification and localization head to learn the underlying task.

meaningful. In our work, we couple the advantages of both adversarial domain adaptation and semantic alignment and propose for the first time a pure adversarial semantic loss that yields domain invariant representations in a semantically meaningful way, without making explicit assumptions of the distributions (see Fig. 1).

Unsupervised Domain Adaptation for TAL: Despite the considerable success of UDA methods, their applicability has been mostly restricted to image-based scenarios [6] such as image classification [3], [4], [44], [45], [46], object detection [47], [48] or semantic segmentation [49], [50]. Much less attention has been dedicated to video-based applications such as action recognition [51], [52], [53] or spatio-temporal action segmentation [54], [55]. To the best of our knowledge, at the time of this writing, there is no direct comparison with our work focusing on UDA for *sparse* Temporal Action Localization. The closest work is SSTDA [7] that applies unsupervised domain adaptation for Action Segmentation. SSTDA proposes the use of two global-distribution-based auxiliary tasks to jointly align cross-domain feature spaces. Unlike our proposal, their work falls into the category of anchor-free, making per-frame action predictions. This restricts its applicability to action segmentation scenarios, where current datasets [8], [9] are designed to deal with frame-based single-action classification. In our work, we overcome this limitation by leveraging an anchor-based architecture that enables a natural adaptation to more realistic multi-label scenarios.

3 METHOD

3.1 Problem definition and notation

In this paper, we address the problem of unsupervised domain adaptation for TAL. For this, we define a source domain \mathcal{S} and a target domain \mathcal{T} . Domain \mathcal{S} consists

of N_S labeled input videos $\{(V_k^S, y_k^S)\}_{k=1}^{N_S}$, where each video V_k^S is a sequence of T frames $(X_{k,1}, \dots, X_{k,T})$ with $X_{k,t} \in \mathbb{R}^{H \times W \times C}$. Here $y_k = (b_k, e_k, c_k)$ contains the begin, end, and class actions of the ground-truth (GT) segments of the video, respectively. The target domain \mathcal{T} is similar to \mathcal{S} but lacks the GT information. Concretely, it consists of N_T unlabeled input videos $\{V_k\}_{k=1}^{N_T}$. Our goal is to train a model that can identify the action segments, including both segment coordinates and action labels, in videos from domain \mathcal{S} , while minimizing the performance degradation on the unlabeled domain \mathcal{T} .

3.2 Framework overview

In this work, we propose a model based on a feature pyramid and a classification and localization head (see Fig. 2). This architecture is coupled with a novel *semantic adversarial loss* that aligns the embeddings of the source and target domain in a semantically meaningful way. More in detail, the model takes as input two videos from domain \mathcal{S} and \mathcal{T} , respectively. The model first processes the raw input videos using a frozen pre-trained video backbone. The resulting embeddings of both domains are then passed through a shared SGP pyramid [10] that outputs a set of multi-resolution anchor embeddings for each of the pre-defined resolution levels. The main goal of our model is to make these anchor embeddings domain invariant. For this, we introduce a level-wise *semantic adversarial loss* that learns in an adversarial manner to align the embeddings of both domains belonging to a given action class i at a given resolution level l . Recall that GT information is only available for domain \mathcal{S} , therefore we rely on the use of pseudo labeling techniques [11] to infer the *probable* class labels of the data from domain \mathcal{T} . Finally, we use the domain invariant anchors of domain \mathcal{S} to train a classification

and localization head that learns the underlying tasks in a standard supervised fashion. In short, this permits to learn a classification and localization head that minimizes the decline of performance when applied to the unseen domain \mathcal{T} .

3.3 Backbone and SGP pyramid

Our model first takes two input videos V^S and V^T of both domains \mathcal{S} and \mathcal{T} . For simplicity, both videos V^S and V^T have length T , which we enforce using padding. The method then processes the two videos applying a frozen pre-trained backbone – e.g., I3D [56], Slowfast [57]. This permits to extract, in an effective way, temporal cues of the video into a set of refined video features. These embeddings are then fed to an SGP feature pyramid [10] which combines the use of SGP blocks and the progressive downsampling of the temporal length by a ratio of 2. This outputs a set of multi-resolution anchor embeddings $Z^S = \{Z_l^S\}_{l \in L}$ and $Z^T = \{Z_l^T\}_{l \in L}$, for the two domains, respectively. Here L denotes the set of predefined resolution levels and $Z_l \in \mathbb{R}^{T_i \times F}$ the anchor embeddings of level l of a given domain. Concretely, this permits to obtain embeddings for a set of uniformly sampled anchors at each of the $l \in L$ resolution levels. The use of a multi-resolution model is favorable to naturally adapt to different action lengths.

3.4 Classification and localization head

To learn the underlying TAL task, we train in a fully supervised manner a classification and localization module with the labeled source domain \mathcal{S} . Due to the anchor-based nature of our model, we first require a matching strategy between the set of candidate anchors to the actual GT segments. For this, we follow a center sampling strategy [10], [21], [58], [59]. In other words, for a given level l , we define an anchor as *positive* if the time instant t that it represents is near the center of an action. All the rest are marked as *negative*. We define B_l, \mathcal{E}_l and C_l as the begins, ends, and action classes of their matching GT segments. We identify negative anchors by assigning them to action class 0.

With this, we design a classification head $H_{cls} : \mathbb{R}^{T_i \times F} \rightarrow \mathbb{R}^{T_i \times C}$ that maps each of the anchor embeddings to their class distribution. More specifically, we model this as a sequence of 1D convolutions, and train it using a sigmoid focal loss (SFL) [60] as follows:

$$\mathcal{L}_{SFL}^l = SFL(H_{cls}(Z_l^S), C_l). \quad (1)$$

Similarly, we model a localization head $H_{loc} : \mathbb{R}^{T_i \times F} \rightarrow \mathbb{R}^{T_i \times 2}$ identically as H_{cls} , which predicts the begin-end offsets. We thus define the localization loss as a standard mean squared error (MSE) loss over the positive samples only:

$$\mathcal{L}_{loc}^l = MSE(H_{loc}(Z_{l_+}^S), (B_{l_+} \parallel \mathcal{E}_{l_+})), \quad (2)$$

where l_+ refers to the positive samples from the l -th level only and \parallel to the concatenation operation. This yields the final task loss defined as:

$$\mathcal{L}_{task} = \lambda_{cls} \sum_{l \in L} \mathcal{L}_{SFL}^l + \lambda_{loc} \sum_{l \in L} \mathcal{L}_{loc}^l. \quad (3)$$

Here λ_{cls} and λ_{loc} are two tunable hyperparameters.

3.5 Our proposal: Semantic adversarial multi-resolution alignment

One of the main contributions of this paper is the design of a novel adversarial-based loss that we name *SADA* loss, which attempts to overcome the limitations of the extensively used *global adversarial loss* [3]. Traditional adversarial domain adaptation relies on the idea of designing a domain classifier that learns to identify the domain that each of the embeddings belongs to. The rest of the model learns concurrently the opposite objective, which, as shown by [3], results in the learning of the domain invariant representations. While this approach has been shown to be effective in other fields – e.g., image classification, and object detection – we find that its performance in a more challenging setup like TAL still presents important challenges. One of the main issues, as argued by [11] is that this loss often suffers from *feature misalignment* which greatly declines the model’s performance. This refers to the cases where these methods align embeddings of non-matching class labels – e.g., aligning embeddings of domain \mathcal{S} of an action i with embeddings of domain \mathcal{T} of a class j .

Local adversarial alignment: To fix this *feature misalignment* in realistic scenarios like TAL, we propose an alternative adversarial loss formulation that provides a finer-grained alignment. This loss first attempts to perform a local class-aware alignment. This is, for every given resolution level l , we group the anchor embeddings of \mathcal{S} and \mathcal{T} corresponding to an action label i . This is straightforward for the source domain as we have the GT information. For the target domain, we use a hard-pseudo labeling strategy [11] that classifies a given embedding as class i if this is the highest-confidence score of the predicted class distribution, and this is above a threshold α . Formally we define the hard pseudo-label of a given anchor z as:

$$\hat{c}_z = \begin{cases} \operatorname{argmax}_i P_l[z, i] & \text{if } P_l[z, i] > \alpha \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

where $P_l = H_{cls}(Z_l^T) \in \mathbb{R}^{T_i \times C}$ are the predicted class probabilities of the anchors. We assign a class 0 to the *background anchors* – those not assigned to any class. From this, we obtain the new grouped embeddings of source and target domain of class i on level l as:

$$A_i^l = \{Z_l^S[z] : c_z = i\}_{z \in T_i}, \quad (5)$$

$$B_i^l = \{Z_l^T[z] : \hat{c}_z = i\}_{z \in T_i}, \quad (6)$$

for $A_i^l \in \mathbb{R}^{T_i, i \times F}$, $B_i^l \in \mathbb{R}^{T_i, i \times F}$. Also, c_z is the GT action label of anchor z and \hat{c}_z is its computed pseudo-label from Eq. 4. We then adversarially train a domain classifier $D : \mathbb{R}^{2F} \rightarrow \{0, 1\}$ to identify the domain of each of these embeddings using a binary cross entropy (BCE) loss:

$$\mathcal{L}_{local}^l = \sum_{i=1}^C (\mathcal{L}_{BCE}(D(A_i^l \parallel E_i), d_S)) + (\mathcal{L}_{BCE}(D(B_i^l \parallel E_i), d_T)), \quad (7)$$

where d_S and d_T are the labels of each of the domains. Following [3], we introduce a Reverse Gradient Layer (GRL) before the discriminator D to invert the sign of the gradient which creates the min-max game where the feature extractor learns to *confuse* the discriminator. We

condition the discriminator to class i using a learnable class embedding $e_i \in \mathbb{R}^F$ that we replicate for every selected anchor into an embedding E_i .

Local and global alignment (SADA): Eq. 7 aligns solely the anchor embeddings that are classified as one of the C classes, but *what happens with the background embeddings that fall below the threshold α ?* In this case, the loss ignores their influence, yielding a suboptimal partial alignment (see Table 4).

To overcome this issue, we propose our final SADA loss which attempts to combine the best of both *global alignment loss* [3] and Eq. 7. For this, we introduce a new loss term for the *background embeddings* as follows:

$$\mathcal{L}_{bkg}^l = \mathcal{L}_{\text{BCE}}(D(A_0^l \parallel E_0), d_S) + \mathcal{L}_{\text{BCE}}(D(B_0^l \parallel E_0), d_T), \quad (8)$$

where again A_0^l and B_0^l are the selected *background anchors*, and $E_0 \in \mathbb{R}^F$ is the learnable *background embedding*. Coupling Eq. 7 and Eq. 8 yields the final formulation of our proposed loss, combining class-wise alignment with the *background anchors alignment*. Formally:

$$\mathcal{L}_{sada} = \sum_{l \in L} \lambda_l (\mathcal{L}_{local}^l + \mathcal{L}_{bkg}^l), \quad (9)$$

where λ_l is a hyper-parameter that regulates the importance of level l on the final alignment loss.

3.6 Training

During training, we formulate the final loss as a min-max game where the main model architecture is optimized over the classification and localization loss while maximizing the adversarial loss. In parallel, the discriminator model D attempts to minimize the discriminator loss only. Formally,

$$\mathcal{L} = \lambda_{task} \mathcal{L}_{task} + \lambda_{sada} \mathcal{L}_{sada}. \quad (10)$$

Note again \mathcal{L}_{task} is optimized with domain \mathcal{S} while \mathcal{L}_{sada} promotes the alignment between both domains \mathcal{S} and \mathcal{T} . Moreover, λ_{task} and λ_{sada} are tunable parameters that control their influence on the loss.

4 DATASETS AND EXPERIMENTS

In this section, we present a novel comprehensive benchmark to evaluate the task of UDA for TAL. Our setup, for the first time, goes beyond action segmentation and evaluates the adaptation to different domain shifts in more realistic scenarios with sparse multi-label annotations. We then showcase the effectiveness of our model over the state-of-the-art methods together with several relevant ablations.

4.1 Benchmarks for UDA on sparse TAL

Evaluating domain adaptation-based methods in the context of video understanding is a challenging issue that requires the definition of a reasonable domain gap and identifying a sufficiently large set of intersecting action classes. Dividing existing datasets into different domains that comply with these conditions often restricts the amount of data to learn and adapt. SSTDA [7] approaches this problem on GTEA [8] and Breakfast [9] by defining a subject-based partitioning where they aim to adapt the model to new unseen subjects. However, as little data is available from a single subject in those datasets, they group several users for training. This

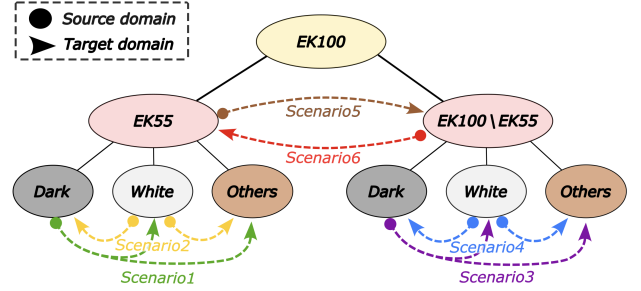


Fig. 3: Overview of the 6 proposed experimental setups for EpicKitchens100. Concretely, S1 and S2 evaluate the videos from the original EK55. They define the dark-counter and white-counter kitchens as Source, respectively, and the rest as Target. S3 and S4 are similar except that they consider only the *newest* videos from EK100. S5 and S6 use the *old* videos as Source and the *new* videos as Target, and vice versa.

allows the model to generalize over the subject variability under study, making it unsuitable to test for domain adaptation. Closely related to our work, [2], [61] propose several UDA scenarios for video classification based on EpicKitchens100 [14]. These define 3 domains based on the data of 3 different kitchens, thus performing cross-kitchen evaluation. This limits even further the amount of data in each domain –i.e., between 15 to 29 videos per domain.

EpicKitchens100: To overcome these limitations, we first propose a new set of 6 different scenarios ($S1, \dots, S6$) for sparse TAL based on the EpicKitchen100 [14] (see Fig. 3). EpicKitchen100 presents an ideal base for our tasks as it has become a gold standard to evaluate complex sparse detection scenarios on long egocentric videos (up to 45 minutes). We identify two domain gaps in this dataset: an *appearance domain shift* based on the different colors of the kitchen counters; and an *acquisition domain shift* that results in the differences of lighting and camera conditions when extending EpicKitchens55 [14] into its new version EpicKitchens100 [14]. This permits to define a rich set of benchmarks that provide a more fine-grained and comprehensive evaluation. Importantly, this strategy is also suitable for single-source settings which do not allow an easy generalization over the shift under study. For example, if a model is trained with dark-counter kitchens only, it cannot easily generalize to white-counter kitchens.

More in detail, we introduce 4 different benchmarks that derive from the aforementioned *appearance shift*. For this, we first define two different splits of $EK55$ and $\{EK100 \setminus EK55\}$ videos according to [14], put more plainly *old* versus *new* videos. We then further split each of them according to the color of the kitchen counter and define 4 different setups of single-source domain and multi-target domain –i.e., $dark \rightarrow white$ and $other$ kitchens; and $white \rightarrow dark$ and $other$ kitchens, for *old* and *new* videos, respectively. Similarly, given the two splits between $EK55$ and $\{EK100 \setminus EK55\}$ videos, we propose 2 additional setups of $EK55 \rightarrow \{EK100 \setminus EK55\}$ and $\{EK100 \setminus EK55\} \rightarrow EK55$. These measure the adaptability to different acquisition conditions –i.e., lighting and camera conditions. This results in domains that contain in the order of hundreds of videos. See Supp. for more details and their respective



Fig. 4: Comparison in two different scenarios between the egocentric perspective (left) and the third-person one (right).

statistics.

One important consideration is that EpicKitchens100 exhibits a very strong long-tail label distribution, where the vast majority of its 97 actions represent only a marginal percentage of the overall data (see Supp.). Adapting DA methods to these long tail distributions falls beyond the scope of this paper and most of the existing literature. For this reason, we follow a similar approach to that proposed by [?] and limit our evaluation to the 10 majority classes representing 80% of the original dataset. This still results in a very complex task proof of which is that the best-performing models to date attain less than a 30% mAP [10], [21].

CharadesEgo: One important aspect of the 6 previous setups is that all the domains share an egocentric viewpoint. For this reason, we evaluate the model degradation under extreme domain shifts caused by major changes in the perspective. For this, we use CharadesEgo [13] which extends the original third-person videos from Charades [62] by matching them with their corresponding egocentric videos. Following [12], we define the Source domain as the third-person videos, and the Target domain as the egocentric ones. Given the extreme shift that these present (see Fig. 4), we limit the task to predicting the verb of each of the action segments, and similarly to the previous case, we keep the 10 majority verbs only. This mitigates the effect of long-tail action distributions in our evaluation which falls beyond the scope of this paper. See the Supp. for more details.

4.2 Experimental results

In this section, we present the main experimental results evaluating *appearance* and *acquisition shifts*. All these experiments follow the standard transductive unsupervised DA protocol [63], [64], and report the mean average precision (mAP) at different intersections over union (IOU) thresholds (10% – 50%).

EpicKitchens100. In Table 1 we first present the results obtained in the 4 different scenarios that we designed to evaluate the performance of our method when facing different *appearance shifts* induced by changes in the background information. Concretely, we first evaluate the performance of the Actionformer [21] and TriDet [10], the two best-performing methods on EpicKitchens100 dataset [14], as well as our proposed architecture without our SADA loss. We then compare these 3 architectures with their extensions which include the SADA loss. Observe that our proposed loss improves the source-only version of all the architectures in all these 4 scenarios, improving by up to a 2.49% mAP when evaluating our architecture on S1. Additionally, we observe that our final architecture –i.e., Ours(SADA)– yields an improvement over the second-best performing source-only method of up to a 3.4% mAP for the *black-counter kitchens* scenarios and of 1.82% mAP for the *white-counter kitchens* ones.

Similarly, the last two scenarios of Table 1 depict the setups containing *acquisition* domain shifts. These scenarios present a similar behavior than before. In all but one case on the TriDet [10] architecture, the use of the SADA loss yields

| Scenario | Model | mAP10% | mAP20% | mAP30% | mAP40% | mAP50% | Avg |
|--|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| $B \rightarrow \{\text{all} \setminus B\}$ (EK100 \setminus EK55) | Actionformer [21] | 28.11 | 26.94 | 24.89 | 21.43 | 16.51 | 23.57 |
| | Tridet [10] | 29.47 | 28.32 | 25.50 | 21.99 | 16.34 | 24.32 |
| | Ours (source-only) | 30.03 | 28.70 | 26.62 | 23.03 | 17.79 | 25.23 |
| | Actionformer+SADA [21] | 30.54 | 29.31 | 27.36 | 23.75 | 18.80 | 25.95 |
| | Tridet+SADA [10] | 31.34 | 30.16 | 27.94 | 24.48 | 19.31 | 26.64 |
| | Ours (SADA) | 32.69 | 31.49 | 29.17 | 25.51 | 19.72 | 27.72 |
| $B \rightarrow \{\text{all} \setminus B\}$ (EK55) | Actionformer [21] | 30.21 | 28.73 | 26.39 | 22.60 | 17.09 | 25.00 |
| | Tridet [10] | 29.87 | 28.39 | 25.97 | 22.06 | 16.94 | 24.65 |
| | Ours (source-only) | 30.74 | 29.47 | 27.23 | 23.53 | 18.07 | 25.80 |
| | Actionformer+SADA [21] | 31.71 | 30.30 | 27.95 | 24.08 | 18.69 | 26.55 |
| | Tridet+SADA [10] | 29.55 | 28.27 | 26.16 | 22.89 | 17.63 | 24.90 |
| | Ours (SADA) | 31.60 | 30.29 | 28.22 | 24.47 | 18.98 | 26.72 |
| $W \rightarrow \{\text{all} \setminus W\}$ (EK100 \setminus EK55) | Actionformer [21] | 33.52 | 32.31 | 29.84 | 26.48 | 20.11 | 28.45 |
| | Tridet [10] | 34.01 | 32.52 | 30.07 | 26.40 | 19.78 | 28.55 |
| | Ours (source-only) | 34.41 | 33.38 | 30.58 | 26.99 | 21.00 | 29.27 |
| | Actionformer+SADA [21] | 34.11 | 32.87 | 30.60 | 27.05 | 20.93 | 29.11 |
| | Tridet+SADA [10] | 34.50 | 33.19 | 30.65 | 27.25 | 20.83 | 29.29 |
| | Ours (SADA) | 34.86 | 33.73 | 31.16 | 27.45 | 21.46 | 29.73 |
| $W \rightarrow \{\text{all} \setminus W\}$ (EK55) | Actionformer [21] | 27.46 | 26.54 | 24.61 | 21.85 | 17.19 | 23.53 |
| | Tridet [10] | 30.03 | 28.97 | 26.96 | 23.48 | 18.18 | 25.52 |
| | Ours (source-only) | 29.65 | 28.69 | 26.86 | 23.88 | 19.14 | 25.64 |
| | Actionformer+SADA [21] | 29.84 | 28.89 | 26.83 | 23.76 | 18.85 | 25.63 |
| | Tridet+SADA [10] | 30.21 | 29.33 | 27.58 | 24.35 | 19.44 | 26.18 |
| | Ours (SADA) | 31.54 | 30.68 | 28.77 | 25.52 | 20.22 | 27.34 |
| EK55 \rightarrow (EK100 \setminus EK55) | Actionformer [21] | 22.87 | 21.87 | 20.10 | 17.23 | 13.33 | 19.08 |
| | Tridet [10] | 24.77 | 22.93 | 21.49 | 19.09 | 15.15 | 20.48 |
| | Ours (source-only) | 25.58 | 24.79 | 23.08 | 19.56 | 15.15 | 21.63 |
| | Actionformer+SADA [21] | 24.85 | 23.99 | 22.21 | 19.11 | 15.28 | 21.09 |
| | Tridet+SADA [10] | 24.88 | 24.00 | 22.20 | 19.02 | 14.88 | 21.00 |
| | Ours (SADA) | 25.93 | 25.06 | 23.47 | 20.45 | 16.12 | 22.21 |
| (EK100 \setminus EK55) \rightarrow EK55 | Actionformer [21] | 22.16 | 21.22 | 19.71 | 17.44 | 14.08 | 18.92 |
| | Tridet [10] | 22.47 | 21.57 | 20.19 | 17.87 | 14.41 | 19.30 |
| | Ours (source-only) | 20.96 | 20.22 | 19.08 | 16.97 | 14.09 | 18.27 |
| | Actionformer+SADA [21] | 23.05 | 22.10 | 20.71 | 18.31 | 14.72 | 19.78 |
| | Tridet+SADA [10] | 21.00 | 20.08 | 18.64 | 16.39 | 13.13 | 17.85 |
| | Ours (SADA) | 23.94 | 22.95 | 21.47 | 19.16 | 15.24 | 20.55 |

TABLE 1: Comparison with the state-of-the-art for the 4 appearance-shift scenarios (1-4) and the 2 acquisition-shift scenarios (5-6) on EpicKitchens100.

a performance gain of up to 2.28% mAP. This is a relative 12.48% improvement. Our final model, moreover, improves the second-best performing source-only baselines by 1.73% mAP and 1.25% mAP, respectively.

CharadesEgo. In Table 2 we report a similar experimental comparison of our proposed loss when evaluated on CharadesEgo. Concretely, we showcase the performance boost that SADA reports on the three test architectures – i.e., Actionformer [21], Tridet [10] and Ours. In the case of Tridet [10], for instance, the improvement reaches a 1.22% mAP. In our proposed architecture, SADA yields an improvement over the source-only version of 0.75% mAP. We additionally highlight that our proposed architecture *Ours (SADA)* attains the state-of-the-art results over all the baseline methods.

4.3 Ablation studies

Comparing to other domain adaptation methods. In Sec. 4.2 we showed that *SADA* consistently improves the

performance of the three tested state-of-the-art source-only architectures when evaluated on our newly proposed setups. The question remains however of how well does *SADA* perform compared to existing domain adaptation methods. In this regard, following existing video-based domain adaptation works, we first evaluate various canonical UDA domain adaptation methods –i.e., ADDA [40], WDGRL [42], DANN [1] and MSTN [2]. These methods are integrated into our proposed underlying architecture resulting in a fair comparison with *SADA*. Additionally, we find that to the best of our knowledge, there is no existing UDA method for TAL in the literature that is directly comparable to this work. Nevertheless, to provide a richer comparison, we adapt the closest action segmentation proposal, SSTDA [7] and a state-of-the-art domain-adaptation method for video classification –i.e., TranSVAE [61]– to our proposed setup. We refer to the Supp. for all the details of these baselines.

Concretely, in Table 3 we compare *SADA* with these baselines on S1 and S2. These results empirically demon-

| Model | mAP10% | mAP20% | mAP30% | mAP40% | mAP50% | Avg |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Actionformer [21] | 31.22 | 28.51 | 23.82 | 18.91 | 13.94 | 23.28 |
| Tridet [10] | 30.18 | 27.06 | 22.58 | 17.81 | 13.10 | 22.15 |
| Ours (source-only) | 30.68 | 27.74 | 23.41 | 18.46 | 13.58 | 22.77 |
| Actionformer+SADA [21] | 31.46 | 28.57 | 24.17 | 19.04 | 13.87 | 23.42 |
| Tridet+SADA [10] | 31.53 | 28.41 | 24.04 | 18.88 | 13.98 | 23.37 |
| Ours (SADA) | 31.68 | 28.64 | 24.09 | 19.06 | 14.14 | 23.52 |

TABLE 2: Comparison with the state-of-the-art for the perspective-shift (third-person to egocentric) on CharadesEgo.

| $B \rightarrow \{\text{all} \setminus B\}$ (EK100-EK55) | | | | |
|---|-----------------|--------------|--------------|--------------|
| Model | mAP {10,30,50}% | | | Avg |
| ADDA [40] | 31.84 | 28.53 | 19.11 | 26.49 |
| WDGRL [42] | 25.05 | 22.08 | 13.91 | 20.35 |
| DANN [1] | 30.48 | 27.07 | 18.12 | 25.22 |
| MSTN [2] | 31.07 | 27.16 | 17.21 | 25.14 |
| SSTDA [7] | 31.17 | 28.01 | 18.98 | 26.05 |
| TranSVAE [65] | 29.91 | 26.12 | 16.16 | 24.06 |
| Ours (SADA) | 32.69 | 29.17 | 19.72 | 27.19 |

| $B \rightarrow \{\text{all} \setminus B\}$ (EK100-EK55) | | | | |
|---|-----------------|--------------|--------------|--------------|
| Model | mAP {10,30,50}% | | | Avg |
| ADDA [40] | 30.55 | 27.09 | 18.17 | 25.27 |
| WDGRL [42] | 24.85 | 21.87 | 13.68 | 20.13 |
| DANN [1] | 30.87 | 27.45 | 18.29 | 25.97 |
| MSTN [2] | 30.88 | 27.44 | 18.25 | 25.52 |
| SSTDA [7] | 29.92 | 26.49 | 17.55 | 24.65 |
| TranSVAE [65] | 25.88 | 21.74 | 13.12 | 20.25 |
| Ours (SADA) | 31.60 | 28.22 | 18.98 | 26.27 |

TABLE 3: Comparison with SOTA UDA methods on S1 and S3.

strate the effectiveness of our method by attaining the best results over all tested UDA methods. More in detail, *SADA* yields an improvement of up to 6.84% mAP and 6.13% mAP on each of the scenarios respectively, and a 0.7% mAP and 1.0% mAP over the second best-performing method.

Analysis of our loss. Next, we ablate over different variants of our proposed *SADA* loss (see Eq. 9). Concretely, in Table 4 we study the effect of class-wise distribution alignment (Eq. 7), global distribution alignment (equivalent to DANN [1]) and semantic background alignment (Eq. 8). In this regard, we highlight that the global adaptation consistently improves upon aligning only background embeddings. This is because the latter yields only a partial alignment of the embeddings, not considering *class anchors*. Therefore, a *rougher* yet complete adaptation might seem beneficial. We also observe that aligning local class-wise distributions has a considerable positive effect, especially in the first scenario. Its effect, however, considerably decreases when combined with a global alignment [1], as all the *class anchors* are then subject to the concurrent alignments of domain-level and class-wise distributions. This observation is also consistent with the performance decrease that we observe when combining global and background alignment, which again suggests that the concurrent alignment of background embeddings with two adaptation losses is harmful to performance. Finally, we observe that we consistently obtain the best results when the local alignment loss is coupled with the complementary (non-overlapping) background loss – i.e., *SADA* loss – indicating that this semantic fine-grained,

| $B \rightarrow \{\text{all} \setminus B\}$ (EK100-EK55) | | | | | | |
|---|--------|-----|-----------------|--------------|--------------|--------------|
| Local | Global | Bkg | mAP {10,30,50}% | | | Avg |
| | | | 30.03 | 26.62 | 17.79 | 24.81 |
| | ✓ | | 30.48 | 27.07 | 18.12 | 25.22 |
| | | ✓ | 30.34 | 26.66 | 17.04 | 24.68 |
| | ✓ | ✓ | 29.76 | 26.63 | 17.52 | 24.64 |
| ✓ | | | 31.36 | 28.01 | 18.67 | 26.01 |
| ✓ | ✓ | | 30.09 | 26.88 | 17.43 | 24.80 |
| ✓ | | ✓ | 32.69 | 29.17 | 19.72 | 27.19 |

| $B \rightarrow \{\text{all} \setminus B\}$ (EK55) | | | | | | |
|---|--------|-----|-----------------|--------------|--------------|--------------|
| Local | Global | Bkg | mAP {10,30,50}% | | | Avg |
| | | | 30.74 | 27.23 | 18.07 | 25.34 |
| | ✓ | | 30.87 | 27.45 | 18.29 | 25.97 |
| | | ✓ | 30.97 | 27.64 | 18.39 | 25.67 |
| | ✓ | ✓ | 30.36 | 26.82 | 18.01 | 25.06 |
| ✓ | | | 30.82 | 27.57 | 18.30 | 25.56 |
| ✓ | ✓ | | 30.30 | 26.98 | 17.85 | 25.04 |
| ✓ | | ✓ | 31.60 | 28.22 | 18.98 | 26.27 |

TABLE 4: Analysis of the effectiveness of the components of *SADA* on S1 and S3.

yet complete, alignment is the most desirable approach.

Ablation of the effect of the λ hyper parameters. As described in Eq. 9, we define our *SADA* loss using a set of hyperparameters $\{\lambda_l\}_{l \in L}$, where each of the parameters $\lambda_l \in (0, 1)$ controls the influence of a given resolution level in the overall loss. Given the relevance of these parameters, below we ablate over different variations to showcase the effect that they have in the final performance of the model.

Concretely, we define the three following scenarios:

- 1) **All levels to 1:** This scenario sets all the parameters $\lambda_l = 1$, keeping the contribution of each of the levels the same.
- 2) **First 3 levels 1:** This scenario sets the first half of the levels to 1 and the rest to 0. This is $\lambda_0 = \lambda_1 = \lambda_2 = 1$ and $\lambda_3 = \lambda_4 = \lambda_5 = 0$.
- 3) **Last 3 levels 1:** This scenario sets the first half of the levels to 0 and the rest to 1. This is $\lambda_0 = \lambda_1 = \lambda_2 = 0$ and $\lambda_3 = \lambda_4 = \lambda_5 = 1$.

All in all, this study attempts to clarify how sensitive is our model to this design decision. In this regard, these results (see Table 5) indicate a mild sensitivity to this choice. Specifically, we observe that the second best-performing method attains an absolute difference over the optimal choice of only -0.08% mAP and -0.57% mAP for the two considered scenarios S1 and S3, respectively. Moreover, we observe that the worst performing naive strategy still outperforms in both scenarios 4 out of the 6 tested UDA baselines.

| $B \rightarrow \{\text{all} \setminus B\}$ (EK100-EK55) | | | | |
|---|-----------------|--------------|--------------|--------------|
| Strategy | mAP {10,30,50}% | | | Avg |
| All levs 1 | 32.42 | 29.59 | 19.34 | 27.11 |
| First 3 levs 1 | 30.96 | 27.48 | 18.29 | 25.58 |
| Last 3 levs 1 | 30.54 | 27.19 | 18.28 | 25.34 |
| “Optimal choice” | 32.69 | 29.17 | 19.72 | 27.19 |
| $B \rightarrow \{\text{all} \setminus B\}$ (EK55) | | | | |
| Strategy | mAP {10,30,50}% | | | Avg |
| All levs 1 | 30.60 | 27.36 | 18.41 | 25.45 |
| First 3 levs 1 | 31.05 | 27.64 | 18.30 | 25.66 |
| Last 3 levs 1 | 30.95 | 27.59 | 18.56 | 25.70 |
| “Optimal choice” | 31.60 | 28.22 | 18.98 | 26.27 |

TABLE 5: Ablation of the effect of the λ parameters, which regulate the influence of each of the resolution levels on the overall SADA loss. The “optimal choice” refers to the hyperparameter choice resulting from a Bayesian optimization process that leverages the labeled source domain to identify the best-performing set of parameters. Moreover, for efficiency purposes, we limit the search space to a uniform sampling in the (0,1] interval –i.e., $\{0.1, 0.2, \dots, 0.9, 1.0\}$.

Study of the class embedding. One of the main contributions of our work is to adversarially align distributions in a class-wise fashion. Intuitively, this requires that our level-wise domain discriminator *knows* the class distribution that it is aligning. In this regard, in Sec. 3.5 we propose to concatenate to every anchor a learnable embeddings $e_i \in \mathbb{R}^F$ of its corresponding class i (see Eq. 7). In this ablation, we empirically justify our choice. For this, we compare our approach against several non-learnable alternatives to encode a given class i : a naive one-hot encoding, a random class-wise dense initialization, and the sinusoidal encoding originally proposed by [33].

In Table 6 we show the experimental results of each of the variants tested in S1 and S3. These indicate that a naive one-hot encoding of a class is the best-performing non-learnable strategy, obtaining considerable improvements over the other two tested non-learnable baselines. This improvement is especially prominent in the first scenario, boosting the performance over the other non-learnable strategies by 1.05% mAP and 2.32% mAP, respectively. Moreover, to our surprise, the popular Sinusoidal encoding [33] proves to be the worst-performing method, being consistently outperformed by even the random dense encoding. Finally, we highlight that in both scenarios the use of a learnable class embedding yields the best results, improving by 0.44% mAP and 0.36% mAP the performance of the one-hot encoding strategy, which justifies our choice.

Ablation of the effect of the *negative anchors*. In this work we focus on applying domain adaptation over anchor-based localization methods, the current state-of-the-art architectures on datasets like EpicKitchens100 [66]. The success of these architectures is normally explained by the generation of a very considerable number of anchors, largely surpassing the number of actual ground-truth labels to predict. As described in Sec. 3 these are the excess anchors –i.e., *negative anchors*– that are not directly matched with a GT segment, consequently training the model to predict class 0 instead. In other words, in these cases we hope for the model to

| $B \rightarrow \{\text{all} \setminus B\}$ (EK100-EK55) | | | | | | |
|---|-----------------------|--------------|--------------|--------------|--------------|--------------|
| Strategy | mAP {10,20,30,40,50}% | | | | | Avg |
| One-hot | 32.29 | 31.09 | 28.58 | 24.88 | 19.56 | 27.28 |
| Random emb | 30.71 | 29.50 | 27.43 | 23.78 | 18.52 | 25.99 |
| Sinusoidal [33] | 29.93 | 28.57 | 26.44 | 22.87 | 17.01 | 24.96 |
| Learnable | 32.69 | 31.49 | 29.17 | 25.51 | 19.72 | 27.72 |
| $B \rightarrow \{\text{all} \setminus B\}$ (EK55) | | | | | | |
| Strategy | mAP {10,20,30,40,50}% | | | | | Avg |
| One-hot | 31.31 | 29.97 | 27.84 | 24.28 | 18.81 | 26.44 |
| Random emb | 31.38 | 30.01 | 27.94 | 24.13 | 18.70 | 26.43 |
| Sinusoidal [33] | 31.10 | 29.78 | 27.71 | 23.99 | 18.56 | 26.23 |
| Learnable | 31.68 | 30.32 | 28.37 | 24.55 | 19.09 | 26.80 |

TABLE 6: Ablation of the effect of the use of a learnable class embedding over other static strategies.

| $B \rightarrow \{\text{all} \setminus B\}$ (EK100-EK55) | | | | | | |
|---|-----------------------|--------------|--------------|--------------|--------------|--------------|
| Filtered % | mAP {10,20,30,40,50}% | | | | | Avg % |
| 0% | 32.69 | 31.49 | 29.17 | 25.51 | 19.72 | 27.72 |
| 25% | 32.54 | 31.87 | 30.06 | 26.99 | 21.91 | 28.67 |
| 50% | 34.46 | 33.55 | 31.29 | 27.59 | 21.18 | 29.61 |
| 75% | 35.09 | 34.53 | 32.97 | 30.05 | 23.97 | 31.32 |
| 100% | 35.57 | 35.24 | 34.16 | 30.72 | 24.50 | 32.04 |
| $B \rightarrow \{\text{all} \setminus B\}$ (EK55) | | | | | | |
| Filtered % | mAP {10,20,30,40,50}% | | | | | Avg % |
| 0% | 31.60 | 30.29 | 28.22 | 24.47 | 18.98 | 26.71 |
| 25% | 32.76 | 31.60 | 29.81 | 26.14 | 20.24 | 28.11 |
| 50% | 35.25 | 34.25 | 32.68 | 29.25 | 23.01 | 30.89 |
| 75% | 38.39 | 37.77 | 36.40 | 33.19 | 26.33 | 34.42 |
| 100% | 43.93 | 43.59 | 42.61 | 39.26 | 31.85 | 40.25 |

TABLE 7: Ablation of the effect of masking out different percentages of *negative anchors* during inference only.

learn to identify these anchor embeddings as *background*, not predicting any action class during inference. This, however, is a cumbersome issue that most of the anchor-based methods face, which is exacerbated when doing domain adaptation in this complicated setup. In this case, it is not clear the amount of noise that these embeddings induce in the alignment strategy. Even though we leave a more thorough analysis of this issue as future work, this ablation aims to provide minimal proof of the potential negative effect that this has on the model’s performance.

Concretely, we propose an experiment where we progressively mask out the set of *negative anchors*. For this, in Table 7 we randomly mask out these only during inference. Put differently, we leave all the training routines unchanged – thus considering all the anchors – but ensure that a percentage of the *negative anchors* are not considered during inference. This measures the predicting potential of a model that is able to perfectly deal with these embeddings. As observed in Table 7, in both considered scenarios we observe a linear performance increase, reaching an absolute improvement of 4.32% mAP and 13.54% mAP of the version with 100% filtered negative anchors over the 0% version. This is, the variant where only the *positive* anchors are considered over that containing both *positive* and all the *negative* ones.

In Table 8 we propose an alternative analysis. In this case, we do not only mask the *negative anchors* during

| B → {all \ B} (EK100-EK55) | | | | | | |
|----------------------------|-----------------------|--------------|--------------|--------------|--------------|--------------|
| Filtered % | mAP {10,20,30,40,50}% | | | | | Avg |
| 0% | 32.69 | 31.49 | 29.17 | 25.51 | 19.72 | 27.72 |
| 25% | 31.48 | 30.34 | 27.97 | 24.07 | 18.93 | 26.56 |
| 50% | 32.11 | 31.11 | 29.25 | 25.85 | 20.72 | 27.81 |
| 75% | 35.99 | 35.28 | 33.48 | 29.90 | 23.64 | 31.66 |
| 100% | 37.06 | 35.73 | 32.34 | 26.45 | 33.82 | 47.08 |
| B → {all \ B} (EK55) | | | | | | |
| Filtered % | mAP {10,20,30,40,50}% | | | | | Avg |
| 0% | 31.60 | 30.29 | 28.22 | 24.47 | 18.98 | 26.71 |
| 25% | 32.25 | 30.46 | 28.04 | 24.54 | 19.03 | 26.86 |
| 50% | 33.43 | 32.06 | 29.98 | 26.44 | 20.55 | 28.49 |
| 75% | 38.26 | 37.33 | 35.72 | 32.53 | 25.61 | 33.89 |
| 100% | 47.08 | 46.69 | 45.65 | 42.31 | 34.42 | 43.23 |

TABLE 8: Ablation of the effect of masking out different percentages of *negative anchors* during training and inference.

inference but also during training. In this case, for instance, masking out 50% of the *negative anchors* additionally implies not training any of the task losses on them, nor using them on our alignment loss. As observed in the results, masking them out entirely – i.e., 100% – has an even bigger impact, reaching an absolute improvement of 6.1% mAP and 16.52% mAP, respectively, when comparing it to the 0% filtering version. Nevertheless, in this case, we observe that the improvement is no longer linear, and even worsens the performance for small percentages of masking. We hypothesize this is induced by a reduction of the training data of *background class*. Thus, while masking some *negative embeddings* has a positive effect, we must also consider the negative effect that having less training data has on the ones that we do consider. Hence, the improvement in performance remains little to non-existent until this trade-off is dominated by the positive effect of masking out the majority of *negative anchors*. In conclusion, these experiments emphasize the need for devising carefully designed models that can properly cope with this intrinsic limitations of anchor-based methods.

Per-class analysis. Finally, we provide a more detailed class-wise analysis of the performance of SADA. Concretely, in Table 9 we show the respective class-wise mAP scores of the 10 considered classes on both S1 and S3. For the analysis, we also include the results obtained with the source-only variant of our model as well as DANN [3] baseline.

Observe that in S1 (see Table 9 bottom) our model attains the best class-wise performance in 5 of the classes, while DANN [3] does so on 3, and the source-only model on 2. In contrast, in S3 (Table 9 top) our method obtains a much clearer improvement over the chosen baselines, yielding the best results in 8 of the 10 classes. Overall, we can observe that our method performs very well in the 3 majority classes of both scenarios, where we highlight the absolute improvement of SADA over DANN [3] of 6.8% mAP for the class *put* in S3. Moreover, our method fails to improve the performance of action *pour* in both scenarios, which we attribute to the lack of sufficient data to operate on our proposed methodology. We also highlight that S1 presents a more challenging setup, degrading the performance in other actions such as *open*, *close*, *insert* or *cut*. This indicates the existence of intrinsic qualitative aspects that harden the

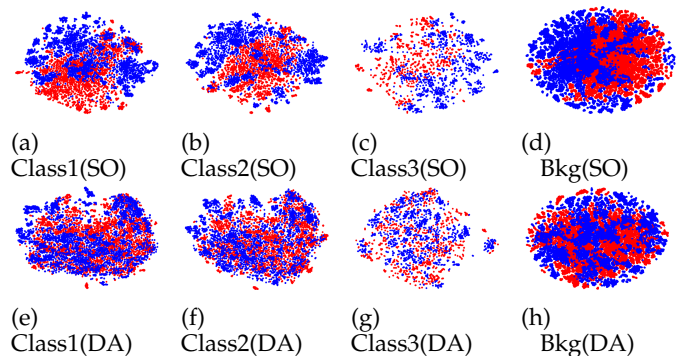


Fig. 5: TSNE plots of the source-only variation of our model (top row) and our proposed model (bottom row). Find in the first 3 columns the TSNE plots of action classes 1 to 3 of the source (red) and target (blue) domain anchors. The last column shows the plot of the background anchors, so those not assigned to any GT label.

adaptation when dealing with *old* videos. This is not the case in scenario S3, as aside from the aforementioned *pour* segments, it only fails to achieve the best results for the *turn-on* action.

4.4 Qualitative analysis

TSNE alignment study. In this study we first perform a qualitative per-class analysis of the effect of our proposed alignment loss SADA, visually depicting the level of alignment between Source and Target domain embeddings. To this end, in Fig. 5 we show the TSNE plots of the domain-invariant embeddings of the 3 majority classes (first 3 columns) and background embeddings (last column) when trained with SADA and when learning source-only. Observe that source-only (top row) yields clearly unaligned distributions with scarce to no overlap in the projected space. Our method, in contrast, presents a considerable distribution mix improving the alignment of class-wise distributions across domains. Given the anchor-based nature of our method, we have numerous background anchors –i.e., not assigned to any true action label. As observed in the last column, these are also aligned by our method (resulting from Eq. 8) effectively aligning the entire data distributions but in a semantically sensitive way.

In Fig. 5 we only include the visualizations of the embeddings that correspond to the level 0 of our multi-resolution architecture. It remains unclear, however, whether this alignment behavior remains consistent across the different resolution levels. To showcase this behavior we include in Fig. 6 - Fig. 9 the TSNE plots of the 3 majority classes and of background embeddings, respectively. Each of these figures includes the different plots at each of the first 3 resolution levels, where level 0 is the shallowest, and level 3 is the deepest studied level. We highlight that deeper levels are not meaningful to plot given the scarce number of resulting domain-invariant embeddings, caused by the downsampling that each of the levels of the architecture performs.

Observe that the influence of the alignment loss follows a similar pattern in all the resolutions. Concretely, in the three

| | | $\mathbf{B} \rightarrow \{\text{all} \setminus \mathbf{B}\} (\text{EK100-EK55})$ | | | | | | | | | |
|--------------------|--|--|--------------|--------------|--------------|--------------|---------------|----------------|--------------|-----------------|--------------|
| | | <i>take</i> | <i>put</i> | <i>wash</i> | <i>open</i> | <i>close</i> | <i>insert</i> | <i>turn-on</i> | <i>cut</i> | <i>turn-off</i> | <i>pour</i> |
| Ours (source-only) | | 24.11 | 25.69 | 29.77 | 31.96 | 27.8 | 6.69 | 25.52 | 34.16 | 17.48 | 29.18 |
| DANN [1] | | 24.78 | 24.38 | 28.75 | 32.73 | 29.32 | 5.38 | 27.52 | 36.73 | 17.70 | 29.05 |
| Ours (SADA) | | 28.32 | 31.25 | 31.18 | 34.52 | 31.82 | 7.29 | 26.66 | 39.12 | 18.68 | 28.34 |

| | | $\mathbf{B} \rightarrow \{\text{all} \setminus \mathbf{B}\} (\text{EK55})$ | | | | | | | | | |
|--------------------|--|--|--------------|--------------|--------------|--------------|---------------|----------------|--------------|-----------------|--------------|
| | | <i>take</i> | <i>put</i> | <i>wash</i> | <i>open</i> | <i>close</i> | <i>insert</i> | <i>turn-on</i> | <i>cut</i> | <i>turn-off</i> | <i>pour</i> |
| Ours (source-only) | | 24.20 | 29.30 | 37.30 | 32.69 | 23.89 | 10.98 | 37.39 | 18.18 | 26.24 | 18.55 |
| DANN [1] | | 24.85 | 30.51 | 37.09 | 35.08 | 22.93 | 11.68 | 35.80 | 17.23 | 24.19 | 20.35 |
| Ours (SADA) | | 26.31 | 30.70 | 39.14 | 34.68 | 23.42 | 10.58 | 39.86 | 17.05 | 26.58 | 18.78 |

TABLE 9: Per class mAPs (in percentage) obtained by the source only variation of our model, DANN [45] and our proposal SADA. These results correspond to S1 and S3.

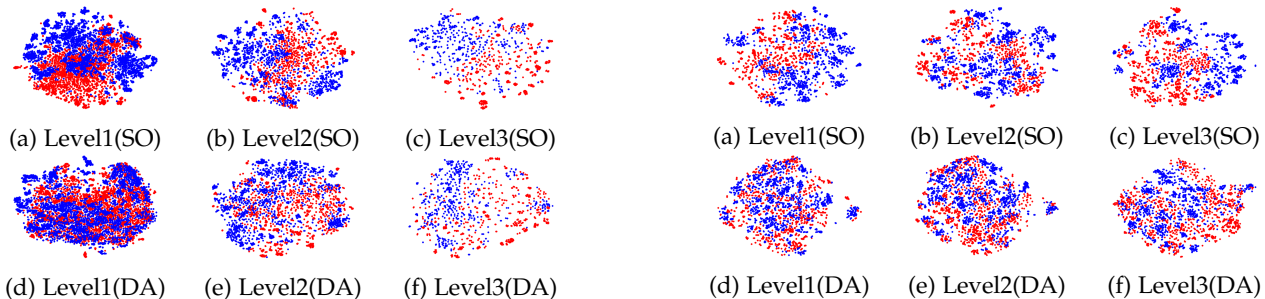


Fig. 6: TSNE plots of class 1 of the source-only variation of our model (top row) and our proposed model (bottom row). Concretely, find in the 3 columns the TSNE plots of class 1 on the first 3 resolution levels of the source (red) and target (blue) domain anchors.

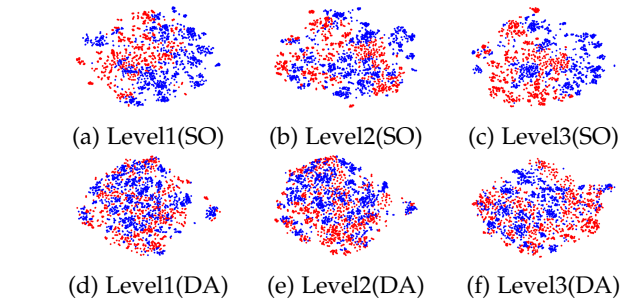


Fig. 8: TSNE plots of class 3 of the source-only variation of our model (top row) and our proposed model (bottom row). Concretely, find in the 3 columns the TSNE plots of class 3 on the first 3 resolution levels of the source (red) and target (blue) domain anchors.

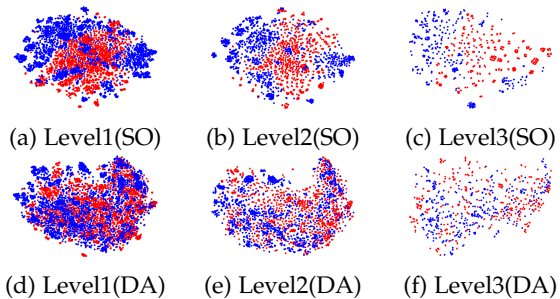


Fig. 7: TSNE plots of class 2 of the source-only variation of our model (top row) and our proposed model (bottom row). Concretely, find in the 3 columns the TSNE plots of class 2 on the first 3 resolution levels of the source (red) and target (blue) domain anchors.

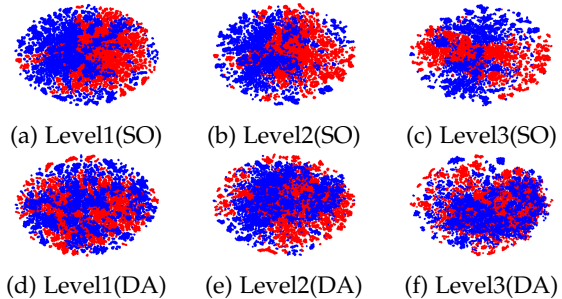


Fig. 9: TSNE plots of *negative class* of the source-only variation of our model (top row) and our proposed model (bottom row). Concretely, find in the 3 columns the TSNE plots of *negative class* on the first 3 resolution levels of the source (red) and target (blue) domain anchors.

presented, we observe a considerable mixing between the source and target distributions (red and blue, respectively). This contrasts with the very clear disentanglement of the representations of both domains in the source-only version, where these present little to no overlap. We highlight what seems to be the only exception which is resolution level 3 of class 1. This one shows little improvement in the mixing of the distributions compared to the source-only variant. We also find that our overall improvement of the mixing of the distributions is consistent when analyzing the *background*

class embeddings. Observe in Fig. 9 that in all 3 studied resolution levels, SADA improves very considerably the alignment, pushing the feature space of both domains to be domain invariant. This emphasizes the positive influence of the background alignment term of our loss (see Eq. 8).

Segment visualization. We complement our qualitative study by depicting in Fig. 10 a segment visualization of S3 (left) and one of our proposed setup for CharadesEgo (right). More specifically, in the segment visualization of S3 we can observe that the Actionformer [21] misses many of

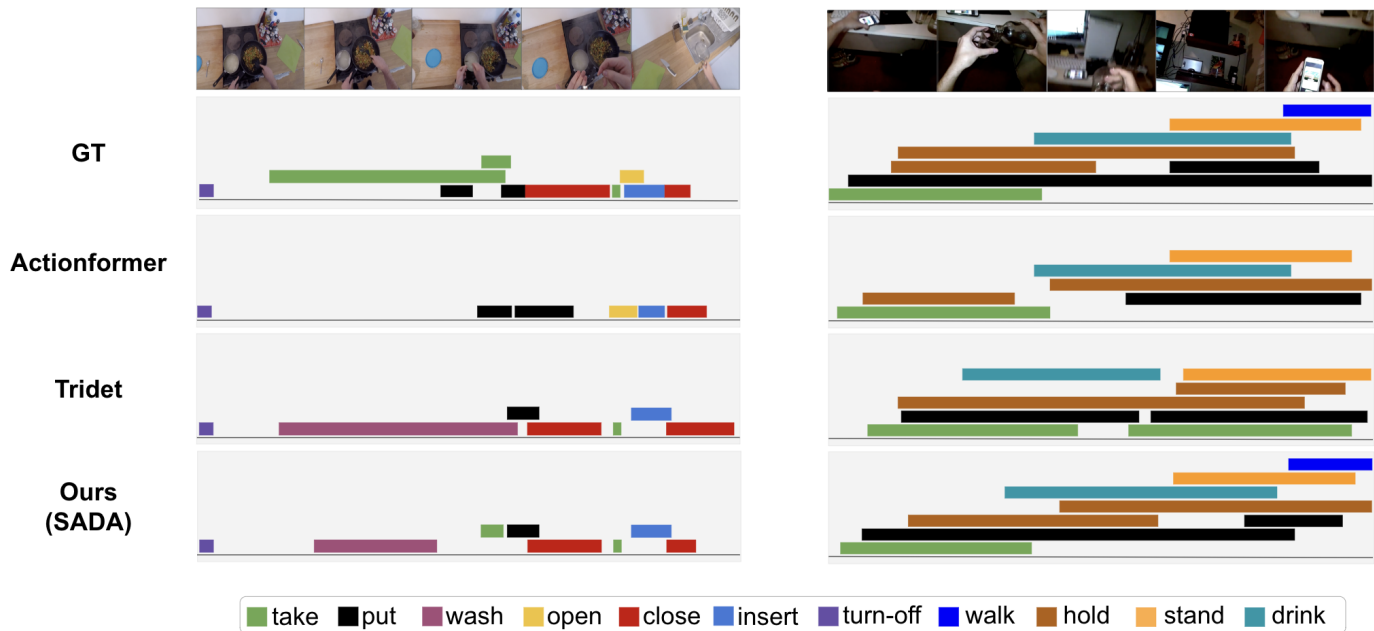


Fig. 10: Visualization of the predicted segments of our method and the chosen set of source-only. We include on top the ground-truth (GT) segments as a reference. Find on the left a segment visualization of S3 and on the right a segment visualization for the CharadesEgo benchmark.

the segments in the shown video clip ignoring all the *take* actions and mistaking a *close* action for a *take*. Tridet [10] performs better but misses all the *take* actions in the first half of the video while worsening the boundary prediction of the last *close*. SADA improves Tridet by predicting the second *take* action and considerably improving the boundary of the last *close* action.

Similarly, note that the segment visualization on CharadesEgo (right) shows a much more complex setup, where there are many overlapping actions. For the depicted example, we must point out that the worst-performing model is the ActionFormer [21]. This model misses the action *walk* and one of the actions *put*. It also has important limitations in terms of localization, especially for the action labels *drink* and *hold*. Then there is Tridet [10], which improves ActionFormer considerably but still misses the *walk* action and is unable to localize the longest *put* action. It also predicts an action *take* that does not correspond to any ground truth. Instead, our method produces much more accurate results, does not miss any of the ground truth actions, and in most cases accurately localizes most of the segments. Yet it confuses the two *hold* actions, which is a consistent failure across all tested baselines.

5 CONCLUSIONS AND LIMITATIONS

In this work, we deal for the first time with Unsupervised Domain Adaptation on realistic Temporal Action Localization (TAL) scenarios. We propose a novel semantic adversarial loss that enables a more fine-grained distribution alignment compared to existing global-distribution-based approaches. Given the lack of suitable evaluation setups for this scenario, we propose a suite of 7 different benchmarks that provide a comprehensive assessment of the model performance across various

domain shifts. These experiments indicate that our model yields a considerable improvement over state-of-the-art methods, which we support with extensive quantitative and qualitative results.

Limitations and future work. Despite the virtues of our proposed benchmarking scenarios, we leave as future work to complement this evaluation with long-tail scenarios, analyzing the effect of domain adaptation optimization in the presence of highly unbalanced data in terms of labels. We also leave as future work the analysis and possible benefits of the adaptation of our semantic adversarial loss to other domain adaptation problems, even to those not necessarily related to video tasks.

REFERENCES

- [1] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [2] S. Xie, Z. Zheng, L. Chen, and C. Chen, "Learning semantic representations for unsupervised domain adaptation," in *International conference on machine learning*. PMLR, 2018, pp. 5423–5432.
- [3] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *International conference on machine learning*. PMLR, 2015, pp. 1180–1189.
- [4] S. Motiian, M. Piccirilli, D. A. Adjeroh, and G. Doretto, "Unified deep supervised domain adaptation and generalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5715–5725.
- [5] M. HassanPour Zonoozi and V. Seydi, "A survey on adversarial domain adaptation," *Neural Processing Letters*, vol. 55, no. 3, pp. 2429–2469, 2023.
- [6] G. Csuska, "Domain adaptation for visual applications: A comprehensive survey," *arXiv preprint arXiv:1702.05374*, 2017.
- [7] M.-H. Chen, B. Li, Y. Bao, G. AlRegib, and Z. Kira, "Action segmentation with joint self-supervised temporal domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9454–9463.

- [8] A. Fathi, X. Ren, and J. M. Rehg, "Learning to recognize objects in egocentric activities," in *CVPR 2011*. IEEE, 2011, pp. 3281–3288.
- [9] H. Kuehne, A. B. Arslan, and T. Serre, "The language of actions: Recovering the syntax and semantics of goal-directed human activities," in *Proceedings of Computer Vision and Pattern Recognition Conference (CVPR)*, 2014.
- [10] D. Shi, Y. Zhong, Q. Cao, L. Ma, J. Li, and D. Tao, "Tridet: Temporal action detection with relative boundary modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18857–18866.
- [11] Y. Li, L. Guo, and Y. Ge, "Pseudo labels for unsupervised domain adaptation: A review," *Electronics*, vol. 12, no. 15, p. 3325, 2023.
- [12] Y. Zhang, H. Doughty, L. Shao, and C. G. Snoek, "Audio-adaptive activity recognition across video domains," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13791–13800.
- [13] G. A. Sigurdsson, A. Gupta, C. Schmid, A. Farhadi, and K. Alahari, "Charades-ego: A large-scale dataset of paired third and first person videos," *arXiv preprint arXiv:1804.09626*, 2018.
- [14] D. Damen, H. Doughty, G. M. Farinella, A. Furnari, J. Ma, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray, "Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100," *International Journal of Computer Vision (IJCV)*, vol. 130, p. 33–55, 2022. [Online]. Available: <https://doi.org/10.1007/s11263-021-01531-2>
- [15] Y.-W. Chao, S. Vijayanarasimhan, B. Seybold, D. A. Ross, J. Deng, and R. Sukthankar, "Rethinking the faster r-cnn architecture for temporal action localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1130–1139.
- [16] Z. Li and L. Yao, "Three birds with one stone: Multi-task temporal action detection via recycling temporal annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4751–4760.
- [17] C. Lin, C. Xu, D. Luo, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, and Y. Fu, "Learning salient boundary feature for anchor-free temporal action localization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3320–3329.
- [18] Z. Qing, H. Su, W. Gan, D. Wang, W. Wu, X. Wang, Y. Qiao, J. Yan, C. Gao, and N. Sang, "Temporal context aggregation network for temporal action proposal refinement," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 485–494.
- [19] D. Sridhar, N. Quader, S. Muralidharan, Y. Li, P. Dai, and J. Lu, "Class semantics-based attention for action detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13739–13748.
- [20] S. Buch, V. Escorcia, C. Shen, B. Ghanem, and J. Carlos Niebles, "Sst: Single-stream temporal action proposals," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 2911–2920.
- [21] C.-L. Zhang, J. Wu, and Y. Li, "Actionformer: Localizing moments of actions with transformers," in *European Conference on Computer Vision*. Springer, 2022, pp. 492–510.
- [22] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, X. Tang, and D. Lin, "Temporal action detection with structured segment networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2914–2923.
- [23] P. Zhao, L. Xie, C. Ju, Y. Zhang, Y. Wang, and Q. Tian, "Bottom-up temporal action localization with mutual regularization," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*. Springer, 2020, pp. 539–555.
- [24] H. Idrees, A. R. Zamir, Y.-G. Jiang, A. Gorban, I. Laptev, R. Sukthankar, and M. Shah, "The thumos challenge on action recognition for videos "in the wild"," *Computer Vision and Image Understanding*, vol. 155, pp. 1–23, 2017.
- [25] H. Zhao, A. Torralba, L. Torresani, and Z. Yan, "Hacs: Human action clips and segments dataset for recognition and temporal localization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8668–8678.
- [26] Z. Shou, J. Chan, A. Zareian, K. Miyazawa, and S.-F. Chang, "Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5734–5743.
- [27] T. Lin, X. Zhao, and Z. Shou, "Single shot temporal action detection," in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 988–996.
- [28] Z. Yuan, J. C. Stroud, T. Lu, and J. Deng, "Temporal action localization by structured maximal sums," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3684–3692.
- [29] J. Tan, J. Tang, L. Wang, and G. Wu, "Relaxed transformer decoders for direct action proposal generation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 13526–13535.
- [30] X. Liu, Q. Wang, Y. Hu, X. Tang, S. Bai, and X. Bai, "End-to-end temporal action detection with transformer," *arXiv preprint arXiv:2106.10271*, 2021.
- [31] X. Liu, S. Bai, and X. Bai, "An empirical study of end-to-end temporal action detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 20010–20019.
- [32] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 2020, pp. 213–229.
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [34] P. Koniusz, Y. Tas, and F. Porikli, "Domain adaptation by mixture of alignments of second-or higher-order scatter tensors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4478–4487.
- [35] Y. Xu, H. Cao, K. Mao, Z. Chen, L. Xie, and J. Yang, "Aligning correlation information for domain adaptation in action recognition," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [36] Y. Grandvalet and Y. Bengio, "Semi-supervised learning by entropy minimization," *Advances in neural information processing systems*, vol. 17, 2004.
- [37] Y. Xu, J. Yang, H. Cao, Z. Chen, Q. Li, and K. Mao, "Partial video domain adaptation with partial adversarial temporal attentive network," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9332–9341.
- [38] G. Kang, L. Jiang, Y. Yang, and A. G. Hauptmann, "Contrastive adaptation network for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4893–4902.
- [39] J. Hoffman, E. Tzeng, T. Darrell, and K. Saenko, "Simultaneous deep transfer across domains and tasks," *Domain Adaptation in Computer Vision Applications*, pp. 173–187, 2017.
- [40] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7167–7176.
- [41] M. Long, Z. Cao, J. Wang, and M. I. Jordan, "Domain adaptation with randomized multilinear adversarial networks," *arXiv preprint arXiv:1705.10667*, 2017.
- [42] J. Shen, Y. Qu, W. Zhang, and Y. Yu, "Wasserstein distance guided representation learning for domain adaptation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [43] S. Cui, S. Wang, J. Zhuo, C. Su, Q. Huang, and Q. Tian, "Gradually vanishing bridge for adversarial domain adaptation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 12455–12464.
- [44] P. Haeusser, T. Frerix, A. Mordvintsev, and D. Cremers, "Associative domain adaptation," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2765–2773.
- [45] M. Long, Y. Cao, J. Wang, and M. Jordan, "Learning transferable features with deep adaptation networks," in *International conference on machine learning*. PMLR, 2015, pp. 97–105.
- [46] O. Sener, H. O. Song, A. Saxena, and S. Savarese, "Learning transferable representations for unsupervised domain adaptation," *Advances in neural information processing systems*, vol. 29, 2016.
- [47] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Domain adaptive faster r-cnn for object detection in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3339–3348.
- [48] P. Oza, V. A. Sindagi, V. V. Sharmine, and V. M. Patel, "Unsupervised domain adaptation of object detectors: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [49] S. Sankaranarayanan, Y. Balaji, A. Jain, S. N. Lim, and R. Chellappa, "Learning from synthetic data: Addressing domain shift

- for semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3752–3761.
- [50] S. Hu, F. Bonardi, S. Bouchafa, and D. Sidibé, “Multi-modal unsupervised domain adaptation for semantic image segmentation,” *Pattern Recognition*, vol. 137, p. 109299, 2023.
- [51] M.-H. Chen, Z. Kira, G. AlRegib, J. Yoo, R. Chen, and J. Zheng, “Temporal attentive alignment for large-scale video domain adaptation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6321–6330.
- [52] A. Jamal, V. P. Namboodiri, D. Deodhare, and K. Venkatesh, “Deep domain adaptation in action space.” in *BMVC*, vol. 2, no. 3, 2018, p. 5.
- [53] B. Pan, Z. Cao, E. Adeli, and J. C. Niebles, “Adversarial cross-domain action recognition with co-attention,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 815–11 822.
- [54] N. Agarwal, Y.-T. Chen, B. Dariush, and M.-H. Yang, “Unsupervised domain adaptation for spatio-temporal action localization,” *arXiv preprint arXiv:2010.09211*, 2020.
- [55] Y. Lu, G. Singh, S. Saha, and L. Van Gool, “Exploiting instance-based mixed sampling via auxiliary source domain supervision for domain-adaptive action detection,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 4145–4156.
- [56] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [57] C. Feichtenhofer, H. Fan, J. Malik, and K. He, “Slowfast networks for video recognition,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6202–6211.
- [58] Z. Tian, C. Shen, H. Chen, and T. He, “Fcos: Fully convolutional one-stage object detection,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9627–9636.
- [59] S. Zhang, C. Chi, Y. Yao, Z. Lei, and S. Z. Li, “Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9759–9768.
- [60] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [61] P. Wei, L. Kong, X. Qu, Y. Ren, J. Jiang, X. Yin *et al.*, “Unsupervised video domain adaptation for action recognition: A disentanglement perspective,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [62] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta, “Hollywood in homes: Crowdsourcing data collection for activity understanding,” in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. Springer, 2016, pp. 510–526.
- [63] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.
- [64] G. Csurka, “A comprehensive survey on domain adaptation for visual applications,” *Domain adaptation in computer vision applications*, pp. 1–35, 2017.
- [65] P. Wei, L. Kong, X. Qu, Y. Ren, Z. Xu, J. Jiang, and X. Yin, “Unsupervised video domain adaptation for action recognition: A disentanglement perspective,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [66] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray, “Scaling egocentric vision: The epic-kitchens dataset,” in *European Conference on Computer Vision (ECCV)*, 2018.
- [67] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [68] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, “The kinetics human action video dataset,” *arXiv preprint arXiv:1705.06950*, 2017.
- [69] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [70] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis, “Soft-nms—improving object detection with one line of code,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5561–5569.
- [71] L. Rüschemdorf, “The wasserstein distance and approximation theorems,” *Probability Theory and Related Fields*, vol. 70, no. 1, pp. 117–129, 1985.

Supplementary Material

IN this supplementary material, we first provide the implementation details needed to reproduce our work (see Sec. A). We then extend the description of Sec. IV.A of the main paper regarding our 7 proposed benchmarking scenarios (see Sec. B). In Sec. C we describe in detail our chosen UDA baselines. Finally we extend some of the provided ablation studies from Sec. IV.C (see Sec. D.1 and Sec. D.2).

APPENDIX A IMPLEMENTATION DETAILS

In this section we include all the relevant implementation details that ensure the proper reproducibility of SADA in the different setups that we present in this paper. Note that we report only the hyperparameters of the best model only, which were obtained through a Bayesian optimization process to minimize the downstream loss –i.e., not including the domain adaptation losses– of the validation split of the Source domain. We follow this approach since by definition we do not have any label information on the Target domain. Nevertheless, we hypothesize that this is a suitable strategy as the goal of our adaptation loss is to *bring both domains closer*. Thus, a set of hyperparameters that performs well on the Source domain should similarly do so on the Target domain. We mitigate the presence of overfitting by using early stopping and a weight decay with a value of 0.05.

All the models that we present are implemented using PyTorch-2.0, CUDA 12.4 and trained for 50 epochs using AdamW [67] optimizer with a learning rate of $1e^{-4}$ with cosine decay and a warm-up phase of 5 epochs on one single NVIDIA GeForce RTX 3090. The training process in all our experiments took between 12 and 24 hours.

Architecturally, all our models define the main feature extractor as a 6-level SGP feature pyramid [10] that follows a max pooling strategy and uses a downsampling rate of 2 and an internal feature dimensionality of 1024. We model the classification and localization heads as a sequence of three 1D-CNNs with a kernel size of 3. These heads are shared across embeddings of different resolution levels. Finally, we introduce level-wise domain discriminators designed as a multi-layer perceptron (MLP) of depth 2 of width 512. Find in Table S1 the other relevant scenario-dependent hyperparameters related to our main contribution, the SADA loss. The rest of the hyperparameters follow the original implementation of Tridet [10].

During training, SADA first extracts video features using a frozen video backbone. For all the scenarios based on EpicKitchens100 this relies on a Slowfast [56] video backbone pre-trained on Kinetics [68]. This backbone is fed with raw videos from EpicKitchens100 [14] with a rate of 30 FPS, a feature stride of 16, and a maximum length of 2304 ensured by either padding – for shorter videos – or random cropping. Alternatively, in the case of CharadesEgo [13] we rely on a frozen I3D video backbone pre-trained on ImageNet [69]. Our model then shuffles the pre-extracted features of the source and target domain videos, respectively. It then feeds a batch of 2 videos of each of the domains, resulting in a batch of 4 videos per training iteration. Note that depending on the scenario under study, one of the domains may have more data than the other. For this reason, in each epoch we repeat the smaller domain until the other finishes, ensuring that all the data of each of the domains is leveraged at least once.

During inference, we follow a similar approach to [10], [21] and define an exponential moving average of the trained model, which we update every iteration with an exponential decay of 0.999. Moreover, given the excess of final predictions, we use the standard SoftNMS [70] with an IOU threshold of 0.1, a minimum score threshold of 0.001, and a sigma value of 0.4.

APPENDIX B EXTENDED BENCHMARK DESCRIPTION

In this section, we describe in more depth the 7 different benchmarking scenarios that we propose in Sec. IV.A of the main text. Concretely, in Sec. B.1 we expand on the 6 different setups that build on the EpicKitchens100 dataset. For this, in Sec. B.1 we first briefly describe EpicKitchens100 [14], the base of all our benchmarks. We then illustrate in Sec. B.1 and Sec. B.1 the 2 different domain shifts that we identify, which we complement with a more detailed description of the specific scenarios derived from each of them. In Sec. B.1 we also include a detailed set of statistics for each of the proposed setups. We proceed similarly for the scenario based on Charades-Ego (see Sec. B.2). This includes a brief description of the dataset and a detailed explanation of the domain adaptation benchmark that we propose in this work based on this dataset.

B.1 EpicKitchens100 scenarios

Description: EpicKitchens100 [14] is a large-scale ego-centric dataset that has recently gained a lot of attraction

| Dataset | Scenario | Backbone | λ_1 | λ_2 | λ_3 | λ_4 | λ_5 | λ_6 | α | λ_{reg} | λ_{cls} | λ_{sada} |
|-------------|----------|----------|-------------|-------------|-------------|-------------|-------------|-------------|----------|-----------------|-----------------|------------------|
| EK100 | S1 | SF | 0.7 | 0.3 | 0.9 | 0.8 | 0.7 | 0.0 | 0.3 | 1.0 | 1.0 | 3.0 |
| | S2 | SF | 0.7 | 0.6 | 1.0 | 0.7 | 0.2 | 0.2 | 0.3 | 1.0 | 1.0 | 3.0 |
| | S3 | SF | 0.4 | 0.8 | 0.7 | 0.7 | 0.9 | 0.6 | 0.6 | 1.0 | 1.0 | 1.0 |
| | S4 | SF | 0.4 | 0.9 | 0.1 | 0.4 | 0.1 | 0.6 | 0.5 | 1.0 | 1.0 | 0.5 |
| | S5 | SF | 1.0 | 0.4 | 0.8 | 0.7 | 0.9 | 0.6 | 0.5 | 1.0 | 1.0 | 2.0 |
| | S6 | SF | 0.6 | 0.1 | 1.0 | 0.6 | 0.2 | 0.4 | 0.3 | 1.0 | 1.0 | 2.0 |
| CharadesEgo | - | I3D | 0.4 | 0.8 | 0.7 | 0.7 | 0.9 | 0.6 | 0.6 | 1.0 | 1.0 | 1.0 |

TABLE S1: Summary of the main hyperparameters detailed for each of the presented benchmarks. Here I3D refers to [56] and SF to Slowfast [57].

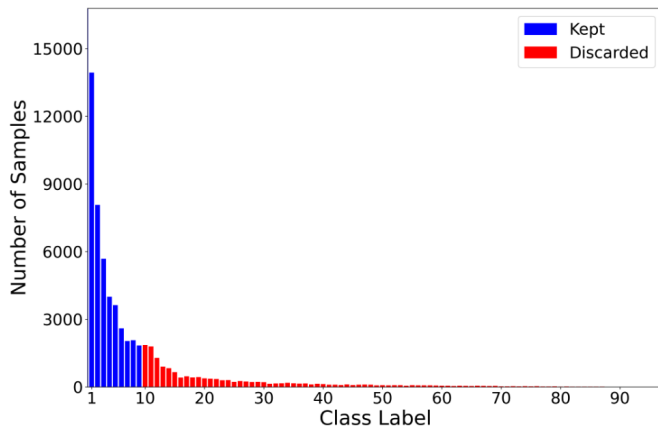


Fig. S1: Histogram of the number of GT action segments of every class. We also depict in blue the 10 majority classes –which we keep in our setups– while leaving as red the remaining class that we discard to avoid the long-tail action distribution problem.

in the Temporal Action Localization community to test for challenging detection scenarios. This dataset extends the previous version of the dataset, namely EpicKitchens55 [66]. This dataset contains videos of 32 subjects performing daily cooking activities – e.g., wash, put, open, close – in different kitchen environments. This results in 100 hours of non-scripted videos containing 20M frames annotated with 90K action segments.

One of the key challenges of this dataset is the presence of 97 different actions, forming a very long tail action distribution. This inherent characteristic is however an issue that falls beyond the scope of this paper, given the important additional challenges that this poses in existing domain adaptation-based methods. For this reason, in all our experiments we consider only the 10 majority classes – i.e., *take*, *put*, *wash*, *open*, *close*, *insert*, *turn-on*, *cut*, *turn-off* and *pour*. As shown in Fig. S1 this pruning keeps most of the labeled action segments of the original dataset, concretely keeping up to 80% of it (marked in blue in the histogram). Moreover, we find that this design decision does not diminish in a relevant way the challenging nature of this task. Proof of this is that the currently best-performing methods on this dataset – i.e., Actionformer [21] and Tridet [10]– still attained mAP metrics inferior to 30% in all our proposed scenarios.

Studying acquisition shifts: The first domain shift that we identify in the EpicKitchens100 [14] is what we call the *acquisition shift*. This shift refers to the changes induced by differences in the acquisition conditions of the videos. In the case of EpicKitchens100 [14] this results from the extension of the original dataset EpicKitchens55 [66]. The original dataset, concretely, was formed by 55 hours of non-scripted videos, and nearly 40K action segments. Hence, we find that this presents a suitable setup to define one domain as the *old* videos – recorded in the original EpicKitchens55 [66] – and the other formed by the *new* videos – recorded for the extended version. As argued by [14] this results in several important domain gaps that are captured by these data splits:

- **Changes in the acquisition devices:** The *new* videos

that were recorded during the extension relied on a newer camera device. Importantly, this camera incorporates camera stabilization techniques. Fig. S2 visually depicts the improvement in the stability of *new* over *old* videos.

- **Lighting conditions:** Given the changes in the hour of the recording, these domains also present differences in the lighting conditions of the videos.

Studying appearance shifts: Another important domain shift that we are interested in capturing is one induced by changes in the background. In this regard, there are several possibilities for this, but we find that the vast majority allow only a vague intuition of the domain shift that is under study. For instance, [14] argues that in the extension, several kitchens might have some changes in the furniture or the order of the main kitchen utensils. But, *which kitchens really contain these changes?* Not understanding this issue in sufficient depth results in an obscure evaluation, which we try to avoid in this work.

For this reason, we identify a clear, understandable domain gap. This is the color of the kitchen counters. This is an essential part of the background information, and thus, we hypothesize that adapting to this factor is also critical to obtain high-performing models. In this regard, we split the data into three different domains: *dark*, *white*, and *other* types of kitchen counters. This provides a clear domain gap that we can visualize and thus, understand (see Fig. S3). We refer interested readers to the file *domain shifts* from the supplementary which shows video content that further depicts the differences between domains in all the proposed scenarios.

Intuitively, this permits to define the following 2 scenarios. Firstly, we can define a scenario for the black-counter kitchens as the source domain, while the rest are the target domain. Similarly, we define another scenario where the white-counter kitchens are the source domain, while the rest are the target domain. Notice, however, that this mixes acquisition conditions as we do not differentiate *old* and *new* videos. To ensure that this does not impede a clear understanding of the scenario, we consider *old* and *new* videos, independently. This results in our 4 final scenarios to measure appearance shifts.

All in all, we stress that all our proposed scenarios provide sufficiently even splits to design single-domain setups where we can ensure that there is no easy generalization of the target domain by means of the source training data (as argued in Sec. IV.A of the main text).

Study of the benchmark statistics: In Sec. IV.A of the main text we propose a set of 6 new scenarios to perform UDA on egocentric videos based on EpicKitchens100 [14]. To do so we rely on two main factors to split the original data from EpicKitchens100 [14] into domains: color of the kitchen counter and year of acquisition –i.e., if the data was already part of EpicKitchens55 [66]. In this regard, in Table S2 we detail the relevant statistics according to these splits where we present the number of videos, the number of segments, the average length of these segments, and finally the class-wise number of segments to depict the overall distribution of actions.

Given this data splits, we define 6 different scenarios

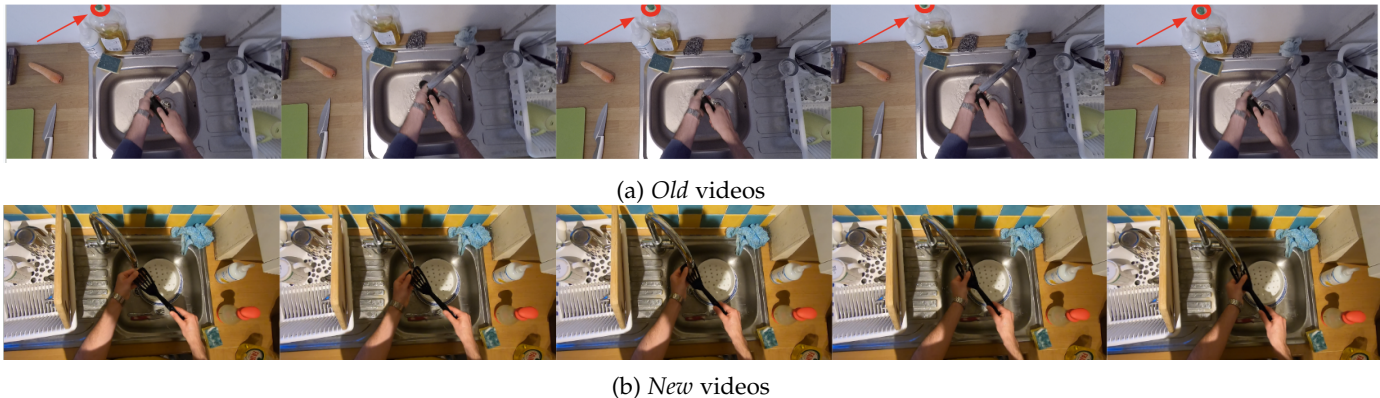


Fig. S2: Visualization of the differences in camera stability between *old* and *new* videos. These images were obtained with differences of 3 frames, which correspond to a period of 0.1 seconds. We highlight with a red circle in S2a a reference in the images to visualize the high instability of the video.



Fig. S3: Visualization of the 3 different appearance-based domain shifts resulting from the split of dark-counter kitchens, white-counter kitchens, and finally all the other kitchens.

according to the description of Sec. IV.B. In this respect, in Table S3 and Table S4 we present the aforementioned statistics for the source and target domain of each scenario, respectively. Observe that thanks to our careful experimental setup, we are able to define scenarios with considerably large source and target domains. This contrasts with other experimental setups like that proposed by [61] which define the source and target domain of a scenario as the videos of a single kitchen, respectively. We find that this, in consequence, yields domains with only 15-28 videos, which we deem insufficient to train state-of-the-art train large architectures.

B.2 CharadesEgo scenario

Description: Charades-Ego [13] is a dataset that was originally proposed in order to enable the transfer of knowledge from third-view videos to their corresponding first-person videos –a.k.a., egocentric videos. Concretely, they present a set of 4000 paired videos –i.e., 8000 videos in total–involving 112 people with an average length of 31.2 seconds. Importantly, this dataset follows the methodology proposed by [62] and asks each of the users to record two videos: 1) Recording the action from a third-person perspective based on a pre-defined script 2) Perform the same action while recording it with a camera fixed in the forehead. Overall, this dataset encompasses a set of 157 complex actions –e.g., *holding some clothes*, *Taking a book from somewhere*.

Third-to-egocentric setup: Based on this dataset, we are the first to present a benchmark on domain adaptation

for Temporal Action Localization. Our presented setup, concretely, explores the adaptation of third-person videos to egocentric ones. This involves an extreme domain shift. For this, following the work of [12], we define the third-person videos as Source domain and the egocentric ones as Target domain. This setup presents an important limitation that we find is critical to make this benchmark viable in this complex task. Concretely, observe that class actions are mostly defined as NLP descriptions, which creates very concrete actions, making it cumbersome to adapt classes across domains. We argue that one of the main reasons is the lack of sufficient data for each of the actions. Consequently, we rely on the originally extended annotations from [13] and propose to use instead the *verbs* of each of the corresponding action labels. This is, for a given action *holding some clothes*, we consider this action with its verb *hold*. One final consideration is that this still results in underrepresented *verb* classes which we fix following a similar approach to that presented by [61] and consider only the 10 majority classes being *take*, *put*, *sit*, *walk*, *hold*, *stand*, *open*, *drink*, *smile* and *laugh*. This still results in a very challenging setup as observed in Table II of the main text, which indicates that none of the state-of-the-art methods tested attain more than 25% mAP. Moreover, this setup retains 70% of the original annotations.

APPENDIX C UNSUPERVISED DOMAIN ADAPTATION BASELINES

In this section, we describe in greater detail the different baselines that we use in Sec. IV.C of the main text, to estab-

| Color | Year | Split | # vids | # segs | Avg len (s) | c1 | c2 | c3 | c4 | c5 | c6 | c7 | c8 | c9 | c10 |
|-------|------|-------|--------|--------|-------------|------|------|------|-----|-----|-----|-----|-----|-----|-----|
| Dark | 55 | Train | 126 | 8336 | 3.11 | 1982 | 1611 | 1208 | 868 | 561 | 768 | 315 | 492 | 243 | 288 |
| | | Val | 35 | 2944 | 3.05 | 699 | 603 | 446 | 273 | 194 | 271 | 131 | 124 | 92 | 111 |
| Dark | 100 | Train | 45 | 5487 | 2.34 | 1576 | 1209 | 672 | 429 | 314 | 453 | 262 | 186 | 226 | 160 |
| | | Val | 13 | 1504 | 1.95 | 430 | 320 | 177 | 127 | 80 | 113 | 98 | 11 | 69 | 79 |
| White | 55 | Train | 124 | 9995 | 2.95 | 2453 | 2398 | 1739 | 963 | 658 | 439 | 477 | 321 | 352 | 195 |
| | | Val | 31 | 1595 | 3.75 | 385 | 313 | 184 | 189 | 111 | 101 | 78 | 120 | 56 | 58 |
| White | 100 | Train | 49 | 6895 | 2.26 | 2236 | 1889 | 818 | 579 | 417 | 87 | 313 | 109 | 245 | 202 |
| | | Val | 19 | 2328 | 2.79 | 705 | 553 | 377 | 184 | 135 | 67 | 107 | 73 | 83 | 44 |
| Brown | 55 | Train | 30 | 3013 | 3.13 | 852 | 723 | 492 | 277 | 229 | 141 | 87 | 98 | 63 | 141 |
| | | Val | 7 | 640 | 2.95 | 177 | 135 | 55 | 78 | 67 | 43 | 18 | 21 | 16 | 30 |
| Brown | 100 | Train | 10 | 1383 | 2.12 | 422 | 330 | 228 | 108 | 92 | 47 | 55 | 18 | 52 | 31 |
| | | Val | 2 | 663 | 2.13 | 235 | 172 | 54 | 43 | 37 | 35 | 19 | 41 | 16 | 11 |
| Other | 55 | Train | 62 | 4200 | 3.53 | 1162 | 797 | 389 | 609 | 314 | 360 | 134 | 198 | 102 | 135 |
| | | Val | 17 | 331 | 3.81 | 91 | 59 | 30 | 47 | 28 | 36 | 5 | 19 | 6 | 10 |
| Other | 100 | Train | 46 | 6747 | 1.72 | 2113 | 1593 | 497 | 599 | 496 | 586 | 300 | 146 | 265 | 152 |
| | | Val | 10 | 2031 | 1.97 | 594 | 528 | 204 | 186 | 143 | 76 | 82 | 51 | 73 | 94 |

TABLE S2: Analysis of the different proposed domains, with their corresponding data splits.

| Scenario | Domain | Split | # vids | # segs | Avg len (s) | c1 | c2 | c3 | c4 | c5 | c6 | c7 | c8 | c9 | c10 |
|----------|------------|-------|--------|--------|-------------|------|------|------|------|------|------|------|-----|------|------|
| S1 | Dark(55) | Train | 126 | 8336 | 3.11 | 1982 | 1611 | 1208 | 868 | 561 | 768 | 315 | 492 | 243 | 288 |
| | | Val | 35 | 2944 | 3.05 | 699 | 603 | 446 | 273 | 194 | 271 | 131 | 124 | 92 | 111 |
| S2 | White(55) | Train | 124 | 9995 | 2.95 | 2453 | 2398 | 1739 | 963 | 658 | 439 | 477 | 321 | 352 | 195 |
| | | Val | 31 | 1595 | 3.75 | 385 | 313 | 184 | 189 | 111 | 101 | 78 | 120 | 56 | 58 |
| S3 | Dark(100) | Train | 45 | 5487 | 2.34 | 1576 | 1209 | 672 | 429 | 314 | 453 | 262 | 186 | 226 | 160 |
| | | Val | 13 | 1504 | 1.95 | 430 | 320 | 177 | 127 | 80 | 113 | 98 | 11 | 69 | 79 |
| S4 | White(100) | Train | 49 | 6895 | 2.26 | 2236 | 1889 | 818 | 579 | 417 | 87 | 313 | 109 | 245 | 202 |
| | | Val | 19 | 2328 | 2.79 | 705 | 553 | 377 | 184 | 135 | 67 | 107 | 73 | 83 | 44 |
| S5 | All(55) | Train | 342 | 25634 | 3.12 | 5529 | 6449 | 1762 | 2717 | 1708 | 1013 | 3828 | 760 | 1109 | 759 |
| | | Val | 90 | 5510 | 3.29 | 1352 | 1110 | 451 | 587 | 400 | 170 | 715 | 209 | 284 | 232 |
| S6 | All(100) | Train | 150 | 20512 | 2.10 | 930 | 6347 | 1319 | 788 | 1173 | 5021 | 1715 | 545 | 459 | 2215 |
| | | Val | 44 | 6526 | 2.27 | 1573 | 540 | 395 | 306 | 1964 | 241 | 228 | 291 | 812 | 176 |

TABLE S3: Analysis of the Source domains of each of the proposed scenarios.

lish a fair comparison with our work, *SADA*. Concretely, we first describe the adaptations to our architecture of multiple UDA methods like DANN [3], ADDA [40], WDGR [42] or MSTN [2]. We then describe the adaptation of the closest related work to *SADA*, SSTDA [7] and for comparison purposes, the adaptation of one of the latest video action recognition methods, TranSVAE [65] to the task of Temporal Action Localization.

DANN [3]: This work pioneered the use of adversarial losses to devise methods that learn domain invariant representations. They do so by proposing a domain classifier $D : \mathbb{R}^F \rightarrow \{0, 1\}$ that learns to identify the domain that each of the feature embeddings belongs to. To integrate this proposal into a multi-resolution architecture we define a level-wise domain classifier $D_l : \mathbb{R}^F \rightarrow \{0, 1\}$. Then, at a given level l , we apply the following adaptation loss on the feature representations of source domain $Z_l^S \in \mathbb{R}^{T_l \times F}$ and the target domain $Z_l^T \in \mathbb{R}^{T_l \times F}$:

$$\mathcal{L}_{DANN}^l = \mathcal{L}_{BCE}(D_l(Z_l^T), d_S) + \mathcal{L}_{BCE}(D_l(Z_l^T), d_T). \quad (11)$$

With this, we define the final loss for the multi-resolution DANN:

$$\mathcal{L}_{DANN} = \sum_{l \in L} \lambda_l \mathcal{L}_{DANN}^l, \quad (12)$$

where λ_l is a hyper-parameter that controls the influence of each of the resolution levels on the overall adaptation loss.

Importantly, this domain classifier is trained through the use of a Gradient Reverse Layer (GRL) that is placed right before the domain classifier D . This allows to create a min-max game where the backbone optimizes the opposite goal of D , which can be seen as learning to *confuse* the domain discriminator.

ADDA [40]: Following a similar approach to the previously presented extension of DANN, we adapt the subsequent work ADDA [40] to our proposed architecture. In our adaptation, this mainly affects the training loss that is employed to optimize Eq. 12. For this, in this case, we do not include the GRL module before each of the domain classifiers but rather employ a GAN loss. This splits the original optimization into two independent objectives for the generator and the discriminator. As seen in the literature, this yields the same min-max loss but with stronger gradients, which is often beneficial. Note that we train these two objectives in an alternating manner, such that at each of the iterations we back-propagate in one or the other GAN losses.

WDGR [42]: In order to adapt WDGR to our architecture, we again follow similar principles to those presented in the adaptation of DANN. In this case, however, we substantially modify the underlying goal following the original work [42]. Concretely, in this case, the goal of D is not the

| Scenario | Domain | Split | # vids | # segs | Avg len (s) | c1 | c2 | c3 | c4 | c5 | c6 | c7 | c8 | c9 | c10 |
|----------|----------------|-------|--------|--------|-------------|------|------|------|------|------|------|------|-----|------|------|
| S1 | All-Dark(55) | Train | 216 | 17298 | 3.12 | 3918 | 4467 | 1201 | 1849 | 940 | 698 | 2620 | 517 | 617 | 471 |
| | | Val | 55 | 2566 | 3.56 | 653 | 507 | 180 | 314 | 206 | 78 | 269 | 98 | 160 | 101 |
| S2 | All-White(55) | Train | 218 | 15639 | 3.23 | 3131 | 3996 | 1104 | 1754 | 1269 | 536 | 2089 | 408 | 788 | 564 |
| | | Val | 59 | 3915 | 3.10 | 967 | 797 | 350 | 398 | 289 | 114 | 531 | 151 | 164 | 154 |
| S3 | All-Dark(100) | Train | 105 | 15025 | 2.02 | 668 | 4771 | 1005 | 562 | 720 | 3812 | 1286 | 385 | 273 | 1543 |
| | | Val | 31 | 5022 | 2.38 | 1253 | 413 | 315 | 208 | 1534 | 172 | 149 | 178 | 635 | 165 |
| S4 | All-White(100) | Train | 101 | 13617 | 2.01 | 617 | 4111 | 902 | 543 | 1086 | 3132 | 1136 | 343 | 350 | 1397 |
| | | Val | 25 | 4198 | 1.99 | 1020 | 356 | 260 | 199 | 1259 | 158 | 184 | 224 | 435 | 103 |
| S5 | All(100) | Train | 150 | 20512 | 2.10 | 930 | 6347 | 1319 | 788 | 1173 | 5021 | 1715 | 545 | 459 | 2215 |
| | | Val | 44 | 6526 | 2.27 | 1573 | 540 | 395 | 306 | 1964 | 241 | 228 | 291 | 812 | 176 |
| S6 | All(55) | Train | 342 | 25634 | 3.12 | 5529 | 6449 | 1762 | 2717 | 1708 | 1013 | 3828 | 760 | 1109 | 759 |
| | | Val | 90 | 5510 | 3.29 | 1352 | 1110 | 451 | 587 | 400 | 170 | 715 | 209 | 284 | 232 |

TABLE S4: Analysis of the Target domains of each of the proposed scenarios.

discriminate between the Source and Target domain, but rather to estimate the Wasserstein distance [71]. This is, we define a level-specific discriminator $D_l : \mathbb{R}^F \rightarrow \mathbb{R}$. Integrating this module into our multi-resolution architecture yields the following loss for a given level l :

$$\mathcal{L}_{WDGRL}^l = \mathcal{L}_{Wasserstein}(D_l(Z_l^S), D_l(Z_l^T)). \quad (13)$$

This yields the final loss for the multi-resolution WDGRL:

$$\mathcal{L}_{WDGRL} = \sum_{l \in L} \lambda_l \mathcal{L}_{DANN}^l, \quad (14)$$

Importantly, following [42] we enforce all domain discriminators to be 1-Lipschitz. This allows us to approximate the empirical Wasserstein distance as:

$$\mathcal{L}_{Wasserstein}(X^S, X^T) = \frac{1}{|X^S|} \sum_{x^S \in X^S} x^S - \frac{1}{|X^T|} \sum_{x^T \in X^T} x^T \quad (15)$$

where $X_l^S, X_l^T \in \mathbb{R}^{T_l}$ are the estimated Wasserstein distances of each of the T_l anchor embeddings of a given level l .

MSTN [2]: In the original work from [2], they propose to compute an online centroid estimate of each of the class embeddings. This keeps important resemblances with our approach as they attempt to reduce the discrepancy between embeddings across domains that belong to a given action class. For this, they also require classifying each of the anchor embeddings to their corresponding class labels. In the context of UDA, this is straightforward in the case of the source domain as we have their corresponding labels. Doing so on the target domain is more challenging as this domain is unlabelled. Consequently, we follow the same strategy presented in Eq. 11 and Eq. 12, classifying the target domain embeddings according to their pseudo-label. Hence, classifying them according to the highest confidence predicted class, if this is above a given threshold α . This results in the source embeddings A_i^l of class i and resolution level l , and the corresponding embeddings of the target domain B_i^l .

We then use these embeddings to update an exponential moving average (EMA) centroid of the given class i and level l as follows:

$$\hat{C}_{i,l}^S = \frac{1}{|A_i^l|} \sum_{z \in A_i^l} z, \quad (16)$$

$$C_{i,l}^S = \Theta C_{i,l}^S + (1 - \Theta) \hat{C}_{i,l}^S, \quad (17)$$

where Θ is the EMA decay factor. We define the centroids

of the target domain analogously:

$$\hat{C}_{i,l}^T = \frac{1}{|B_i^l|} \sum_{z \in B_i^l} z \quad (18)$$

$$C_{i,l}^T = \Theta C_{i,l}^T + (1 - \Theta) \hat{C}_{i,l}^T. \quad (19)$$

Finally, as proposed in the original paper, we apply an MSE loss to reduce the centroid discrepancy and computed the mean loss over every level and class:

$$\mathcal{L}_{cent} = \frac{1}{C} \frac{1}{|L|} \sum_{i=1}^C \sum_{l \in L} MSE(C_{i,l}^S, C_{i,l}^T) \quad (20)$$

SSTDA [7]: In the original work, SSTDA [7] proposes to couple what they called *Local SSTDA* with *Global SSTDA*. *Local SSTDA* refers to the standard frame-level alignment loss. To adapt this to our multi-resolution setup, we follow the same strategy as in our adaptation of DANN, applying adversarially training a classifier that learns to discriminate the domain that each of the frames/embeddings comes from.

Global SSTDA, in contrast, refers to a sequential domain prediction loss. This splits the original source and target domain into 2 clips of sequential frames. They then apply a *domain attentive temporal pooling* [7] to first weight the frames according to their entropy, and then pool them into a single compact clip embedding. The source and target clip embeddings are then randomly shuffled, and finally passed through a domain classifier that learns to predict the permutation that these clips form (see the original work [7] for more details). To adapt this methodology to our multi-resolution setup, we apply this idea for every level independently. Concretely, we split into two segments the domain-invariant feature representations of domains \mathcal{S} and \mathcal{T} for a given level l , yielding $Z_l^S = \{Z_{l,1}^S, Z_{l,2}^S\}$ and $Z_l^T = \{Z_{l,1}^T, Z_{l,2}^T\}$. Then, following the original proposal, we apply *domain attentive temporal pooling* to compute segment-level representations $V_l^S = \{V_{l,1}^S, V_{l,2}^S\}$ and $V_l^T = \{V_{l,1}^T, V_{l,2}^T\}$. Finally, we apply the sequential domain prediction on these segment embeddings of a given resolution level l . For this, we shuffle them randomly into – e.g., $\{V_{l,1}^S, V_{l,2}^T, V_{l,2}^S, V_{l,1}^T\}$ – and train a level-wise domain classifier $D_L : \mathbb{R}^{4F} \rightarrow \{0, 1\}^{\#perm}$ to predict which of the possible permutations was used.

TransVAE [65]: This work is one of the latest state-of-the-art methods on domain adaptation for video action recog-

nition. This is especially interesting given the innovative disentanglement-based approach. This approach, concretely, leverages the use of a Variational Autoencoder that learns two types of representations: 1) Frame (clip) level representations with dynamic information 2) Static video level representation that includes domain-specific information. The main goal of this model is to enforce domain invariance on the frame-level representations. For this, they leverage several domain adversarial losses (see the original [65] for more details). Finally, it leverages a classification head to aggregate the domain invariant frame-level representations to perform a single video-level prediction.

Adapting this work to our target task, Temporal Action Localization, is however nontrivial. Notice that we deal with sparse detection setups, so it is not enough to transform the goal of TransVAE into a frame-level classification problem. We find this Action Segmentation to be ineffective for our chosen benchmarks. Consequently, we propose instead to substitute their original classification head for the SGP pyramid and classification/localization head that we describe in Sec. III.C and Sec. III.D. Hence, we utilize the domain invariant representations computed by TransVAE as input of the multi-resolution pyramid and later prediction head. This makes this model architecture adaptable to our proposed experimental setup, and comparable to *SADA*.

APPENDIX D

EXTENDED RESULTS AND ABLATIONS

In this section, we extend the ablation studies presented in Sec. IV.C of the main text by analyzing two additional scenarios. Concretely, we compare the performance of *SADA* with the adaptations of different UDA methods, and with different variations of our loss, when evaluated in scenarios S2 and S4.

D.1 Domain adaptation methods

In Table S5 we present the comparison of our proposed method to the other considered UDA baselines. Concretely, the first scenario (left) refers to the *white new* videos, while the second scenario (right) refers to the *white old* videos. Similarly to Sec. IV.C of the main text, we observe that *SADA* improves the results of all the considered baselines for both scenarios. More specifically, *SADA* attains a performance boost of up to 2.04% mAP in the first scenario and 6.33% in the second. Moreover, our method improves by 0.72% mAP and 0.19% mAP the performance of the second best-performing model in both scenarios, respectively. Interestingly, in both cases, the second best-performing method is DANN [1] which considerably improves other more complex methods like SSTDA [7]. Another interesting note is that WDGRL [42] presents a considerable variability in the performance. While this is the third best-performing method in the first scenario, its performance drops drastically in the second attaining only a 20.51% mAP. This contrasts with the performance of similar methods like ADDA [40] or DANN [1].

| W → {all \ W} (EK100-EK55) | | | | |
|----------------------------|-----------------|--------------|--------------|--------------|
| Model | mAP {10,30,50}% | | | Avg |
| ADDA [40] | 33.51 | 29.82 | 20.05 | 27.79 |
| WDGRL [42] | 33.47 | 29.75 | 20.67 | 27.96 |
| DANN [1] | 34.11 | 30.48 | 20.75 | 28.44 |
| MSTN [2] | 32.95 | 29.22 | 19.20 | 27.12 |
| SSTDA [7] | 33.17 | 29.38 | 18.92 | 27.15 |
| TransVAE [65] | 32.9 | 28.67 | 17.44 | 26.34 |
| Ours (SADA) | 34.86 | 31.16 | 21.46 | 29.16 |
| W → {all \ W} (EK100-EK55) | | | | |
| Model | mAP {10,30,50}% | | | Avg |
| ADDA [40] | 30.16 | 27.58 | 19.72 | 25.82 |
| WDGRL [42] | 24.77 | 21.96 | 14.81 | 20.51 |
| DANN [1] | 31.32 | 28.55 | 20.10 | 26.65 |
| MSTN [2] | 30.17 | 27.32 | 18.55 | 25.34 |
| SSTDA [7] | 29.58 | 26.81 | 18.36 | 24.92 |
| TransVAE [65] | 28.31 | 25.67 | 16.72 | 23.57 |
| Ours (SADA) | 31.54 | 28.77 | 20.22 | 26.84 |

TABLE S5: Ablation study on S2 and S4, comparing the performance of our proposal *SADA* with the chosen domain adaptation methods coupled in our proposed architecture.

| W → {all \ W} (EK100-EK55) | | | | | | |
|----------------------------|--------|-----|-----------------|--------------|--------------|--------------|
| Local | Global | Bkg | mAP {10,30,50}% | | | Avg |
| | | | 34.41 | 30.58 | 21.00 | 28.66 |
| | ✓ | | 34.11 | 30.48 | 20.75 | 28.44 |
| | | ✓ | 34.46 | 30.85 | 20.90 | 28.73 |
| | ✓ | ✓ | 34.63 | 31.08 | 21.26 | 28.99 |
| ✓ | | | 33.47 | 29.91 | 20.17 | 27.85 |
| ✓ | ✓ | | 32.84 | 29.06 | 19.41 | 27.10 |
| ✓ | | ✓ | 34.86 | 31.16 | 21.46 | 29.16 |
| W → {all \ W} (EK55) | | | | | | |
| Local | Global | Bkg | mAP {10,30,50}% | | | Avg |
| | | | 29.65 | 26.86 | 19.14 | 25.21 |
| | ✓ | | 31.32 | 28.55 | 20.10 | 26.65 |
| | | ✓ | 29.94 | 27.39 | 19.40 | 25.57 |
| | ✓ | ✓ | 30.86 | 28.02 | 19.84 | 26.57 |
| ✓ | | | 30.84 | 27.98 | 19.37 | 26.06 |
| ✓ | ✓ | | 29.77 | 27.10 | 18.70 | 25.19 |
| ✓ | | ✓ | 31.54 | 28.77 | 20.22 | 26.84 |

TABLE S6: Ablation study of the effect of several variations of our proposed *SADA* loss on S2 and S4.

D.2 Ablation over variants of the *SADA* loss

Closely related to Sec. IV.C of the main text, in Table S6 we ablate over different variants of our loss on scenario S2 and S4. These scenarios define the source domain as *white old* and *new* counter kitchens and the rest as the target, respectively. Observe in this table that our main insights from Sec. IV.C of the main text are mostly confirmed. Concretely, one of our main observations was that aligning the complete *global* distribution was generally preferable over aligning the partial *background* distribution only. This observation is satisfied in most of the scenarios that we study –with the exception of scenario S3. Similarly, in Sec. IV.C of the main text we argued that combining alignment loss that presented overlaps yielded also a decrease in performance. More in detail, this happens when we couple either the *local* or the *background* loss with a *global* alignment approach. In these cases, the embeddings have a double (simultaneous)

alignment objective, which we argue is not desirable. As observed in Table S6, this observation is also mostly satisfied in these scenarios as the combination of these losses yields a performance degradation in 7 of the 8 cases studied in this paper.