

GLAD: Global-Local View Alignment and Background Debiasing for Unsupervised Video Domain Adaptation with Large Domain Gap

Hyogun Lee^{1*}, Kyungho Bae^{1*}, Seong Jong Ha², Yumin Ko³,
 Gyeong-Moon Park^{1†}, Jinwoo Choi^{1†},
¹Kyung Hee University, ²AI Center, CJ Corporation, ³NCSOFT
 {gunsbrother, kyungho.bae, gmpark, jinwoochoi}@khu.ac.kr
 oanchovy@cj.net, yuminko@ncsoft.com

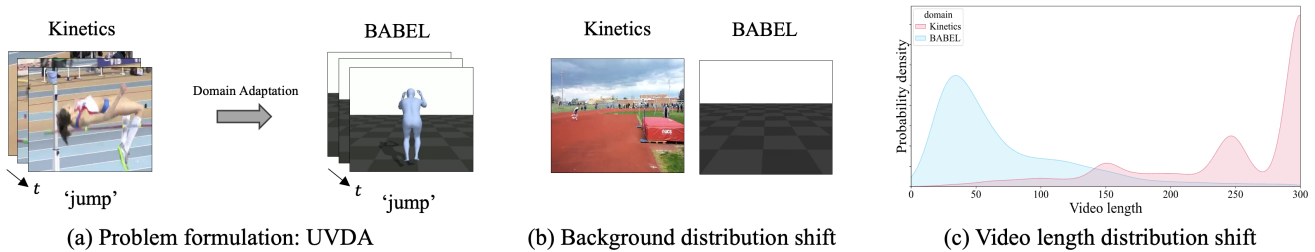


Figure 1. **Overview of Kinetics→BABEL dataset.** We introduce a challenging unsupervised domain adaptation (UDA) dataset, Kinetics→BABEL. (a) We formulate the problem of action recognition as UDA where we have labeled source dataset, e.g., Kinetics, and unlabeled target dataset, e.g., BABEL. The dataset presents two challenges: (b) Background distribution shift: The source dataset (Kinetics) exhibits diverse backgrounds, while the target dataset (BABEL) consistently features the same background across videos. (c) Video length distribution shift: Videos in the source dataset (Kinetics) tend to be longer, while videos in the target dataset (BABEL) are typically shorter. These challenges make the Kinetics→BABEL dataset a valuable benchmark for studying UDA for action recognition.

Abstract

In this work, we tackle the challenging problem of unsupervised video domain adaptation (UVDA) for action recognition. We specifically focus on scenarios with a substantial domain gap, in contrast to existing works primarily deal with small domain gaps between labeled source domains and unlabeled target domains. To establish a more realistic setting, we introduce a novel UVDA scenario, denoted as Kinetics→BABEL, with a more considerable domain gap in terms of both temporal dynamics and background shifts. To tackle the temporal shift, i.e., action duration difference between the source and target domains, we propose a global-local view alignment approach. To mitigate the background shift, we propose to learn temporal order sensitive representations by temporal order learning and background invariant representations by background augmentation. We empirically validate that the proposed method shows significant improvement over the existing methods on the Kinetics→BABEL dataset with a large domain gap.

*Equally contributed first authors.

†Corresponding authors.

The code is available at <https://github.com/KHU-VLL/GLAD>.

1. Introduction

Human action recognition in videos is an interesting problem in computer vision. There are immense practical applications of action recognition: video surveillance, retrieval, captioning, sports analysis, health care, and autonomous driving. Achieving accurate and robust action recognition performance enables improved security and efficient video analysis.

Recent advances in action recognition have witnessed remarkable progress, primarily attributed to the availability of extensive labeled datasets and the successful deployment of deep learning architectures, such as convolutional neural networks (CNNs) [4, 11, 24, 38] and transformers [1, 3, 25, 29]. However, collecting large-scale annotated video data remains a challenging and costly endeavor due to the additional temporal dimension compared to image annotation. Due to the high annotation cost, labeled video datasets do not scale sufficiently, resulting in poor generalization in unseen do-

main [6].

To address the aforementioned challenge of poor generalization, an effective approach is to formulate the action recognition task as an unsupervised domain adaptation (UDA) problem. In the UDA setting, we leverage a labeled source dataset to achieve good performance on an *unlabeled* target dataset. The recent works on unsupervised video domain adaptation (UVDA) for action recognition have shown impressive performance improvement [5, 7, 9, 26, 28, 33, 36, 43] on the standard UCF-HMDB [5] and EPIC-KITCHENS [26] datasets.

However, the impressive performance on the UCF-HMDB and EPIC-KITCHENS datasets may not necessarily reflect real-world scenarios. This discrepancy arises due to several reasons. Firstly, these datasets have a relatively small scale. The UCF-HMDB dataset consists of 3,209 videos from both the source and target domains, which is considerably smaller compared to the original UCF-101 [37] and HMDB-51 [21] datasets. This limited data can lead to overfitting issues as models struggle to effectively generalize. Secondly, the UCF-HMDB and EPIC-KITCHENS datasets do not exhibit significant domain gaps. As shown in Table 1, the accuracy gap between the model trained with target labels and the model trained with only the source data and labels is 11.4 points for UCF-HMDB and 26.2 points for EPIC-KITCHENS. However, real-world scenarios often involve more substantial domain gaps, such as the real-synthetic gap, day-night gap, sunny-snowy gap, and others. These domain gaps present additional challenges that need to be addressed for action recognition models to reliably perform in diverse and complex environments.

To address the limitations of existing datasets, we introduce Kinetics→BABEL, a new and comprehensive dataset designed to present greater challenges for unsupervised video domain adaptation. The Kinetics→BABEL dataset significantly expands the scale, comprising a total of 18,946 videos. As depicted in Figure 1, the Kinetics→BABEL dataset exhibits substantial temporal and background distribution shifts between the source and target domains. In Figure 1 (c), it is evident that the videos from the Kinetics dataset tend to be longer compared to the videos from BABEL. Furthermore, the background distributions differ between the two datasets, with Kinetics displaying real but biased backgrounds for different actions, while BABEL features a consistent gray-scale checkerboard background across actions as shown in Figure 1 (b). In Figure 2, we compare the proposed Kinetics→BABEL dataset with existing datasets in terms of the scene distance (Δ_{bg}), temporal distance (Δ_{temp}), and scale. The Kinetics→BABEL dataset shows more substantial domain gaps between the source and target, and is much larger than the existing datasets. The proposed dataset is much more realistic and challenging compared to the existing datasets. Please refer to Section 3 for more details on the dataset.

To tackle the challenging UVDA with a large domain

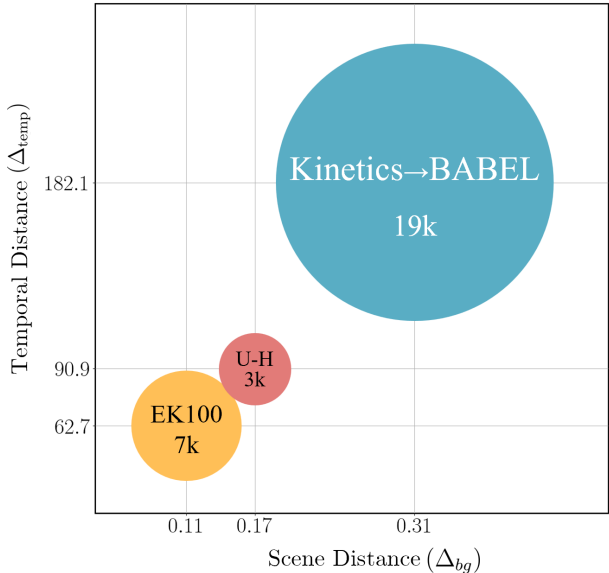


Figure 2. **Comparison between Kinetics→BABEL and the existing UVDA datasets.** We compare the scene distance (Δ_{bg}), the temporal distance (Δ_{temp}), and the scale of the UCF-HMDB, EPIC-KITCHENS, and the proposed Kinetics→BABEL datasets. The Kinetics→BABEL dataset shows more substantial domain gaps between the source and target, and is much larger than the existing datasets.

gap in Kinetics→BABEL, we propose i) **Global-Local view Alignment** and ii) **background Debiasing** for unsupervised video domain adaptation (**GLAD**). i) To address the temporal duration shift between the source and target domains, we propose a Global-Local temporal view Alignment approach, GLA. GLA aligns a set of source clips, sampled at diverse temporal sampling rates, with a set of target clips that also exhibit varying sampling rates. By considering global and local temporal perspectives, our approach facilitates the learning of domain-invariant representations, particularly effective in scenarios with large temporal shifts. ii) To address the background distribution shift between the source and target domains, we propose a background-invariant representation learning to debias background bias, inspired by prior works [7, 33]. The proposed debiasing method leverages both background augmentation via background mixing and temporal order learning. By incorporating these techniques, we mitigate the impact of background distribution shift between domains, thereby improving the performance on the target domain. To validate the efficacy of our proposed method, we conduct extensive empirical evaluations on the challenging Kinetics→BABEL dataset. Our experimental results demonstrate the superiority of GLAD in handling UVDA with a significant domain gap, showcasing its effectiveness in achieving robust action recognition performance in real-world scenarios. To facilitate further research, we

plan to publicly release the Kinetics→BABEL dataset and code upon acceptance of this paper.

In summary, our work makes the following key contributions:

- We introduce the novel Kinetics→BABEL dataset, specifically designed for unsupervised video domain adaptation with a substantial domain gap. The Kinetics→BABEL dataset exhibits significantly larger temporal and background distribution shifts compared to existing datasets, making it a more challenging and realistic benchmark.
- To tackle the temporal and background shifts between the source and target domains, we propose a novel approach called Global-Local view Alignment and background Debiasing (GLAD). GLAD incorporates global-local view alignment techniques to address temporal shifts and employs background debiasing methods to mitigate the background distribution shift.
- We empirically demonstrate the effectiveness of the proposed method via extensive experiments on the challenging Kinetics→BABEL dataset.

2. Related Work

Action Recognition. Deep neural networks have demonstrated remarkable progress in the field of action recognition. Various approaches have been explored to recognize actions from videos. One common approach is to utilize 2D CNNs [10, 18, 24, 27, 35, 46], which extract features from individual frames of a video and incorporate temporal modeling techniques. Another popular approach involves 3D CNNs [4, 11, 16, 38], which learn to capture spatio-temporal features from short video clips. Recently, transformers with spatio-temporal attention mechanisms have also demonstrated impressive performance in action recognition [1, 3, 29]. However, most of the existing action recognition methods heavily rely on large amounts of labeled data. In contrast, our work takes a different approach by formulating the action recognition problem as unsupervised domain adaptation. In this setting, we no longer require labeled data from the target domain, but instead leverage labeled data from the source domain.

Unsupervised Domain Adaptation. In recent years, substantial efforts have been dedicated to unsupervised domain adaptation (UDA) for both image domains [13, 34, 44, 45] and video domains (UVDA) [5, 7, 20, 33, 40, 42, 43]. To tackle the UVDA problem, adversarial-based methods [5, 26, 28], semantic-based methods [9, 33], and self-supervised methods [7, 26] have shown significant progress. However, the majority of existing UVDA works evaluate their performance on small-scale and less challenging datasets such as UCF-HMDB [5] or EPIC-KITCHENS [26]. This limitation hampers the comprehensive evaluation of UVDA

methods in more demanding scenarios. To address this gap, we introduce a novel and large-scale UVDA dataset called Kinetics→BABEL, which exhibits a significant domain gap. Our proposed method is specifically designed to tackle the challenges presented by this dataset. We anticipate that the Kinetics→BABEL dataset serves as a new standard benchmark for evaluating UVDA methods, facilitating further advancements in this field.

Background bias. The research community has recognized background bias as a significant challenge in video action recognition [6, 22, 23]. When an action recognition model is biased toward the background, it relies on spurious correlations between actions and backgrounds rather than understanding the true semantics of the human actions. The background bias becomes even more detrimental in the context of UVDA, where the model needs to adapt to a target domain with different background distributions without action labels. Several approaches demonstrate the benefits of background debiasing in UVDA [7, 33]. In this work, we also address the significant background bias present in the source domain, Kinetics, aiming to achieve favorable performance on the target domain, BABEL, which exhibits entirely different background distributions. By mitigating the background bias, we encourage the action recognition model to focus on genuine action semantics and enhance its ability to adapt to diverse target domains with varying background characteristics.

3. Kinetics→BABEL Dataset

We introduce a new dataset called Kinetics→BABEL, designed to evaluate the performance of UVDA methods in a more realistic and challenging setting. In this work, we set Kinetics as the source domain and BABEL as the target domain. The Kinetics→BABEL dataset is constructed by re-organizing two existing datasets: Kinetics [19] and BABEL [31]. Kinetics→BABEL consists of 12 classes, specifically selected from the overlapping classes of Kinetics and BABEL: jump, run, throw, kick, bend, dance, clean something, squat, punch, crawl, clap, pick up. The dataset comprises 14,881 training and 650 test videos from the Kinetics dataset, and 2,963 training and 452 test videos from the BABEL dataset.

The proposed UVDA dataset encompasses both the real-world Kinetics dataset and the synthetic BABEL dataset. Leveraging synthetic datasets is cost-effective compared to real-world data collection, making their integration as source or target datasets a commonly adopted approach. As shown in the previous works [5, 15, 39], real-to-synthetic and synthetic-to-real domain adaptation problems are quite challenging which makes the proposed dataset interesting. In this work, we focus on the Kinetics→BABEL domain adaptation setting, leaving BABEL→Kinetics domain adaptation setting as a future work.

The Kinetics→BABEL domain adaptation presents two significant challenges: the appearance gap and the temporal gap between the source and target data. The BABEL dataset lacks background information, in contrast to the Kinetics dataset which consists of videos with realistic backgrounds. Moreover, while Kinetics videos exhibit similar durations, BABEL videos encompass a wider range of durations. Consequently, addressing both the background and temporal gaps in a comprehensive domain adaptation strategy becomes crucial to achieve a good performance on the Kinetics→BABEL dataset.

Notably, the proposed Kinetics→BABEL dataset exhibits a larger domain gap compared to existing UVDA datasets, such as UCF-HMDB [5] and EPIC-KITCHENS [26]. To quantify the background gap, denoted as Δ_{bg} , we calculate the average minimum scene feature distance between each source video and all target videos and vice versa as follows:

$$\Delta_{bg} = \frac{1}{2} \left[\frac{1}{L_S} \sum_{i=1}^{L_S} \min_j d(\mathbf{u}_i, \mathbf{v}_j) + \frac{1}{L_T} \sum_{j=1}^{L_T} \min_i d(\mathbf{u}_i, \mathbf{v}_j) \right]. \quad (1)$$

Here, \mathbf{u}_i represents the scene feature vector of the source domain with L_S videos, \mathbf{v}_j denotes the scene feature vector of the target domain with L_T videos, and $d(\mathbf{u}, \mathbf{v}) = 1 - \mathbf{u}^T \mathbf{v}$ is the cosine distance between them. We employ a ResNet-50 [14] model pre-trained on the Places365 dataset [47] to extract scene features.

Furthermore, Kinetics→BABEL also shows the huge domain gap in the temporal perspective. To assess the temporal gap, we leverage the earth mover’s distance (EMD) [17, 32]. The EMD quantifies the minimal cost required to transform one distribution into another, providing an intuitive measure of similarity between distributions. We compute the EMD between two video length distributions p, q as follows:

$$\Delta_{temp} = \text{EMD}(p, q) = \int |CDF_p(x) - CDF_q(x)| dx. \quad (2)$$

In Table 1, we show three domain gaps between the source and target data: the scene distance (Δ_{bg}), the temporal distance (Δ_{temp}), and the accuracy gap (Δ_{Acc}) for various UVDA datasets. It is evident that both the UCF-HMDB and EPIC-KITCHENS datasets exhibit relatively smaller scene distances of 0.17 and 0.11 respectively. In contrast, the proposed Kinetics→BABEL dataset demonstrates a significantly larger scene distance of 0.31, indicating a more pronounced background gap between the domains. Furthermore, Kinetics→BABEL shows a more realistic temporal gap for UVDA settings. The temporal distance of Kinetics→BABEL is 182.1 frames which is $2\times$ bigger than the temporal gap of the UCF-HMDB and $3\times$ bigger than the temporal gap of the EPIC-KITCHENS. To achieve good performance on the Kinetics→BABEL dataset, a model should be able to focus on the action instead of the background as well as learn to represent videos with various lengths.

Table 1. **UVDA dataset statistics.** We provide a quantitative evaluation of commonly used benchmarks in the field of UVDA. The table includes the number of shared classes (# classes), the total number of videos (# videos), the scene distance (Δ_{bg}) in frames calculated by (1) and the temporal distance (Δ_{temp}) in frames calculated by (2), and the accuracy gap (Δ_{Acc}) between “target only” and “source only” performances. The best quantities are in bold.

Dataset	# classes	# videos	Δ_{bg}	Δ_{temp}	Δ_{Acc}
UCF-HMDB [5]	<u>12</u>	3,209	0.17	<u>90.9</u>	11.4
EPIC-KITCHENS UDA [26]	8	*6,729	0.11	62.7	26.2
Mixamo→Kinetics [9]	14	36,195	0.24	66.7	†68.1
Kinetics→BABEL	<u>12</u>	<u>18,946</u>	0.31	182.1	<u>65.0</u>

*The average number of videos across 6 settings.

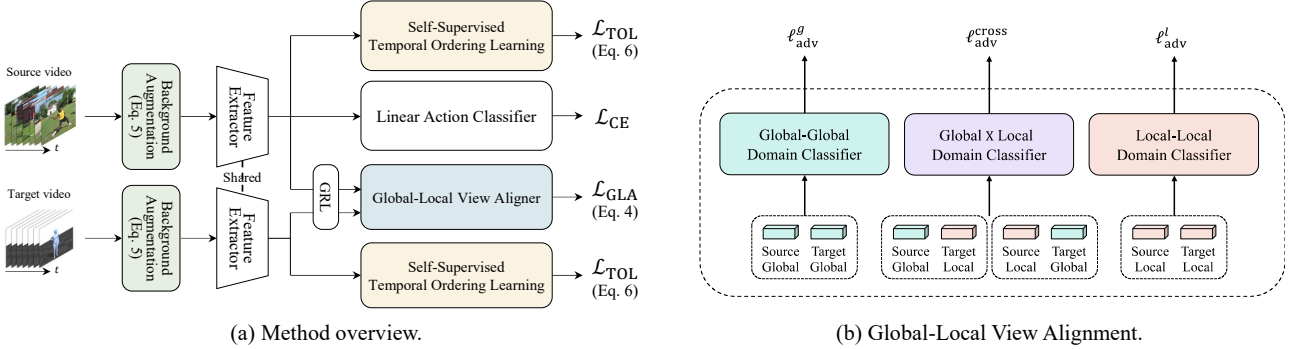
†The reported value from the paper.

Due to the presence of background and temporal gaps, the model performance decreases on the target domain, as indicated by Δ_{Acc} . Δ_{Acc} denotes the performance gap between the model trained with target labels and the model trained with the source data and labels only. The Kinetics→BABEL shows a significant gap of 65.0 points while UCF-HMDB shows 11.4 points and EPIC-KITCHENS shows 26.2 points respectively. Moreover, the Kinetics→BABEL dataset comprises 18,946 videos, making it substantially larger in scale compared to both UCF-HMDB (3,209 videos) and EPIC-KITCHENS (6,729 videos). These observations clearly demonstrate that the proposed Kinetics→BABEL dataset is a large-scale and challenging benchmark to properly evaluate the performance of unsupervised video domain adaptation methods.

Comparison with other synthetic-real datasets. There are a few synthetic-real datasets for the problem of domain-adaptive action recognition: Kinetics-Gameplay [5] and Mixamo-Kinetics [9]. Compared to these existing datasets, the proposed dataset has some advantages. As shown in Table 1, the proposed Kinetics→BABEL dataset offers the larger scene distance (Δ_{bg}) and temporal distance (Δ_{temp}) between domains, compared to the existing synthetic-real datasets. Also, note that the raw RGB data of the Kinetics-Gameplay dataset is not publicly available while we make the raw data of the Kinetics→BABEL dataset public.

4. Method

We formulate the video action recognition task as an unsupervised video domain adaptation (UVDA). In UVDA, we have a labeled source video dataset $\mathcal{D}^s = \{(x_i^s, y_i^s)\}$, where x_i^s represents the input video and y_i^s denotes the corresponding label, as well as an unlabeled target video dataset $\mathcal{D}^t = \{x_i^t\}$. The source and target datasets share the same label space \mathcal{K} between the source and target data. Our objective is to learn a model that performs well in the target domain. Simply applying a model trained solely on the source data to the target data leads to suboptimal perfor-



(a) Method overview.

(b) Global-Local View Alignment.

Figure 3. Overview of GLAD. (a) GLAD consists of several key components. Firstly, we mix a video with a different background from another video to mitigate background bias. Next, a feature extractor extracts spatio-temporal feature vectors from the augmented videos. Then we feed the source feature vectors into a linear classifier to learn action labels. We employ a global-local view alignment module following a gradient reversal layer to align the source and target features. To further address background bias, the model learns the temporal order of shuffled clips in a self-supervised manner. (b) To tackle the temporal shift between the source and target domains, GLAD utilizes three temporal view alignment methods: global-global, local-local, and global \times local. Each method employs dedicated domain classifiers to align the source and target features.

mance [12, 13]. Therefore, a UVDA method should effectively leverage not only the labeled source data but also the unlabeled target data to achieve superior performance in the target domain.

We show an overview of the proposed method, GLAD, in Figure 3 (a). Given a video, we mix it with a different background from another video for background debiasing (Section 4.2). Then we extract a spatio-temporal feature vector from the augmented background mixed. We feed a source video feature vector into a linear classifier to learn actions with the standard cross-entropy loss. To align the source and target domains, we feed both the source and target feature vectors into the global-local view alignment module following a gradient reversal layer [12, 13] (Section 4.1). To further mitigate the background shift between domains, we encourage the model to learn the temporal order of multiple clips in either a source or target video (Section 4.2). We provide more details on each component in the following subsections.

4.1. Global-Local View Alignment

We propose Global-Local view Alignment (GLA) to align features of different domains even if action durations are significantly different across domains. As illustrated in Figure 1 (c), we observe action duration shifts across different domains, such as in the Kinetics \rightarrow BABEL dataset. For example, the `jump` action in Kinetics spans a duration of 10 seconds, involving a sequence of a run-up, a jump, and a landing. In contrast, the `jump` action in BABEL lasts only 1 second, consisting of a brief jump. Due to these temporal shifts, simply aligning the source and target feature vectors of clips using the same sampling strategy across domains may lead to suboptimal performance in UVDA, particularly when a large temporal distribution shift exists, as in the case

of the Kinetics \rightarrow BABEL dataset.

Global and local temporal views. We define a uniformly sampled clip as a global clip and a densely sampled clip as a local clip. For uniform sampling, we divide a video into equal-sized subsequences and randomly select one frame from each subsequence to construct a clip. On the other hand, for dense sampling, we select frames with regular intervals, starting from a randomly chosen point, to construct a clip.

Let ϕ_m^g, ϕ_n^l denote global and local clip feature vectors, respectively, extracted by the feature extractor. Then, we can define aggregated global/local feature vectors ψ as follows:

$$\psi^g = \frac{1}{M} \sum_{m=1}^M \phi_m^g, \quad \psi^l = \frac{1}{N} \sum_{n=1}^N \phi_n^l, \quad (3)$$

where M, N are the number of global and local clips sampled from a single video respectively.

Domain alignment. As shown in Figure 3 (b), we employ individual domain classifiers to align feature vectors with different temporal granularities from the source and target domains. Specifically, we align the global feature vectors from the source and target domains (global-global), the local feature vectors from the source and target domains (local-local), and a global feature vector from one domain with a local feature vector from another domain (global-local). For the global-global alignment, we use an MLP denoted as \mathcal{F}^g . Similarly, we employ another MLP \mathcal{F}^l for local-local alignment and yet another MLP $\mathcal{F}^{\text{cross}}$ for cross-scale (global-local) alignment. To introduce adversarial training, we insert a gradient reversal layer (GRL) [12] between the feature extractor and the domain classifiers. The GRL negates gradients during backpropagation, effectively making the domain classifier adversarial.

Then, we define an adversarial loss ℓ_{adv} for an arbitrary temporal view as follows:

$$\ell_{\text{adv}}(\mathcal{F}, \psi) = -\frac{1}{2B} \left[\sum_{i \leq B} \log \mathcal{F}(\psi_i) + \sum_{i > B} \log(1 - \mathcal{F}(\psi_i)) \right]. \quad (4)$$

Here, B represents the batch size, and we differentiate between the two domains using their batch indices: $1 \leq i \leq B$ for the source domain and $B < i \leq 2B$ for the target domain. We define the final global-local view alignment loss as follows:

$$\mathcal{L}_{GLA} = \ell_{\text{adv}}(\mathcal{F}^g, \psi^g) + \ell_{\text{adv}}(\mathcal{F}^l, \psi^l) + \ell_{\text{adv}}(\mathcal{F}^{\text{cross}}, \psi^{\text{cross}}), \quad (5)$$

where ψ^g, ψ^l , and $\psi^{\text{cross}} = (\psi^g, \psi^l)$ denote the feature vectors for global-global, local-local, and global-local alignments, respectively. With the loss function (5), we effectively align the feature vectors from different domains with significantly different action durations.

4.2. Background Debiasing

As depicted in Figure 1, the Kinetics→BABEL dataset exhibits a significant and realistic background distribution shift. To effectively address this background distribution shift, we incorporate two essential debiasing methods: i) background augmentation and ii) temporal order learning. These debiasing methods play a crucial role in enhancing the performance of UVDA, as demonstrated in the ablation study presented in Section 5.2. It is worth emphasizing that the careful selection and utilization of these debiasing methods contribute to achieving superior performance in UVDA tasks.

Background augmentation. To encourage a model to learn background-invariant representations, we employ a background augmentation technique. For each video in the dataset, we extract a background frame b using a temporal median filter (TMF) [30] and store these background frames for later use. The backgrounds obtained through TMF typically exhibit clear and appropriate backgrounds for the majority of videos.

During training, we randomly select a background from the stored background database and mix it with each frame of every video in a minibatch. We define the mixing process as follows:

$$\tilde{x}(t) = (1 - \lambda)x(t) + \lambda b, \quad t = 1, \dots, T. \quad (6)$$

Here, $x(t)$ represents the t -th frame of the input video $x \in \mathbb{R}^{T \times H \times W \times C}$, λ is a mix-up ratio uniformly sampled from the range $[0, 1]$. By providing action sequences against diverse backgrounds, we encourage the model to focus on the actions themselves rather than being overly influenced by the background context. This facilitates the learning of background-invariant representations that are essential for domain-adaptive action recognition [7, 33].

Temporal ordering learning. To account for significant background shifts across different datasets [7], we incorporate an additional learning objective, namely temporal order learning, to further regularize the model training in conjunction with background augmentation. We adopt the temporal clip order prediction [7, 41] as a pre-text task for this purpose.

In the clip order prediction task, the model tries to solve a puzzle of predicting the true order of N shuffled clips. By solving this clip order prediction task, the model is encouraged to focus more on the action itself rather than being influenced by the static background. As illustrated in Figure 3 (a), we feed both the source and target videos into the temporal order learning (TOL) module.

The TOL module shuffles the order of N clip features $\phi = (\phi_n)_{n=1}^N$ for each video. Consequently, we obtain $\tilde{\phi} = (\phi_{\sigma(n)})_{n=1}^N$, where σ denotes a permutation randomly chosen from the set of all possible permutations \mathcal{S}_N . We pass the shuffled clip features $\tilde{\phi}$ through a simple MLP, denoted as \mathcal{F}_Ω , followed by a softmax operation to predict the correct order $\hat{\omega}_i \in [0, 1]^{N!}$, where $\sum_j \hat{\omega}_{i,j} = 1$. We define the TOL loss as follows:

$$\mathcal{L}_{\text{TOL}} = -\frac{1}{2B \cdot N!} \sum_{i=1}^{2B} \sum_{j=1}^{N!} \omega_{i,j} \log \hat{\omega}_{i,j}. \quad (7)$$

A background-biased model is likely to struggle in predicting the correct order of clips, as its focus remains on the static background. Conversely, a model that focuses on the actions is more likely to predict the correct order. By incorporating the TOL loss, we encourage the model to learn background-invariant representations.

4.3. Training

We define the final optimization objective as follows:

$$\mathcal{L} := \mathcal{L}_{\text{CE}}(\theta_f, \theta_c) + \mathcal{L}_{\text{TOL}}(\theta_f, \theta_\sigma) - \mathcal{L}_{GLA}(\theta_f, \theta_d),$$

$$(\theta_f^*, \theta_c^*, \theta_\sigma^*) = \underset{\theta_f, \theta_c, \theta_\sigma}{\text{argmin}} \mathcal{L}(\theta_d^*), \quad \theta_d^* = \underset{\theta_d}{\text{argmax}} \mathcal{L}(\theta_f^*, \theta_c^*, \theta_\sigma^*), \quad (8)$$

where $\theta_f, \theta_c, \theta_\sigma$, and θ_d denote the parameters of the feature extractor, action classifier, an MLP of TOL, and domain classifiers of GLA, respectively.

4.4. Inference

During the inference stage, we remove all auxiliary components, including TOL and GLA, and retain only the feature extractor and linear action classifier. We do not utilize background augmentation.

Given an input video during inference, we extract one global feature vector and two local feature vectors. These features capture both global and local temporal information.

Table 2. **Ablation study.** To validate the effect of each component, we show experimental results on the Kinetics→BABEL dataset. We conduct all experiments using the TSM [24] backbone. We report the mean class accuracy (MCA) with the corresponding standard deviation. The best performance is in **bold** and the second best is underscored.

(a) Effect of various temporal alignments.				(b) Effect of different temporal views.			(c) Effect of various debiasing methods.			(d) Effect of combining GLA and background debiasing.		
Temporal Alignment Strategies			MCA	Temporal View		MCA	Debias		MCA	Method		MCA
Global-Global	Local-Local	Cross	K→B	Global	Local	K→B	Bg. Aug.	TOL	K→B	Debiasing	GLA	K→B
Source-only baseline			18.5 ± 1.5	1	0	18.5 ± 1.5			26.9 ± 3.2			26.4 ± 2.4
✓			25.5 ± 4.9	0	1	21.5 ± 5.1			29.6 ± 1.7			36.7 ± 3.6
	✓		27.6 ± 3.7	3	0	28.3 ± 1.8	✓		<u>33.4</u> ± 1.7	✓		26.9 ± 3.2
		✓	26.7 ± 1.6	2	1	30.1 ± 4.4		✓	37.7 ± 2.5		✓	37.7 ± 2.5
✓	✓		28.4 ± 6.1	1	2	29.6 ± 1.7	✓	✓		✓	✓	
✓	✓	✓	29.6 ± 1.7	0	3	25.1 ± 2.4						

We then average these feature vectors to obtain a single consensus feature vector that effectively represents the entire video. Finally, we feed the consensus feature vector into a linear classifier to predict the corresponding action label.

5. Experiments

We conduct all the experiments on the Kinetics→BABEL dataset. We use mean-class accuracy as an evaluation metric.

5.1. Implementation details

We implement the proposed method using PyTorch and the mmaction library [8]. We choose I3D [4] as the feature extractor for benchmarking against state-of-the-art methods, and TSM [24] for conducting ablation studies. The feature extractors are initialized with Kinetics400 pre-trained weights. In the GLA module, we employ a 4-layer MLP for each domain classifier. To stabilize the training process, we employ curriculum learning [2]. We first pre-train the model with \mathcal{L}_{TOL} for 500 epochs using 3 local clips to warm up the model. Then, we train the model with the final training objective (8) for 50 epochs. We use SGD as the optimizer with a momentum of 0.9, a weight decay of $1e-4$, and an initial learning rate of $2e-3$. The learning rate is reduced by a factor of 10 at the 5th and 10th epochs. During warm-up, the batch size is set to 384 per GPU, while during the main training, it is set to 24 per GPU for both the source and target domains. Background augmentation is applied only to the source domain clips, with a probability of 25% and a fixed λ value of 0.75. To better capture the temporal context in videos, we adopt two different sampling strategies: uniform sampling for global clips and dense sampling for local clips, maintaining a frame interval of 2 in both domains. All experiments are conducted using 8 NVIDIA RTX 3090 GPUs.

5.2. Ablation Study

We conduct an extensive ablation study to verify the effectiveness of each component and show the results in Table 2.

Effect of various temporal alignment methods in GLA.

In Table 2 (a), we show experimental results demonstrating

the impact of different temporal alignment methods in the GLA module. The *Global-Global* refers to employing a domain classifier with global clips from source and target domains. *Local-Local* refers to employing a domain classifier with local clips. *Cross* refers to employing a domain classifier with a global clip from one domain and a local clip from another domain. As shown in the table, incorporating both global and local alignments leads to superior performance (28.4%) compared to focusing on either one alone (25.5%, 27.6%). Notably, we achieve the highest performance of 29.6% when we employ all three alignments together. Furthermore, the collaborative operation of these alignment methods results in a relatively more stable performance, as indicated by the lower standard deviation value.

Effect of the number of global and local views. From the results presented in Table 2 (b), aligning only local clips surpasses aligning only global clips (21.5% versus 18.5%). However, combining both global and local alignments leads to even higher accuracy. Specifically, employing a combination of two global and one local view per video for alignment achieves the highest accuracy of 30.1% with a standard deviation of 4.4. Notably, using one global and two local views per video for alignment demonstrates comparable accuracy of 29.6%, with a lower standard deviation of 1.7. Based on these findings, we utilize one global and two local views per video for alignment in the subsequent experiments.

Effect of background debiasing. As shown in Table 2 (c), both background augmentation and TOL demonstrate performance improvements of 2.7 and 6.5 points, respectively, compared to the baseline without debiasing. Furthermore, when we combine both debiasing methods, we observe a substantial gain of 10.8 points. These results highlight the complementary nature of the two debiasing methods, emphasizing the importance of employing them together. Please note that GLA is enabled for all experiments conducted.

Complementary nature of GLA and background debiasing.

Table 2 (d) demonstrates the complementary nature of background debiasing and GLA. When applying background debiasing without GLA, we observe a substantial improvement of 10.3 points compared to the baseline. Similarly, applying GLA without debiasing results in a modest

Table 3. **Comparison with state-of-the-art on the Kinetics→BABEL dataset.** We show the mean class accuracy (MCA) For a fair comparison, we indicate the number of clips N_c and the number of frames per clip N_f . All methods employ I3D [4] as the backbone. The best performance is in **bold** and the second best is underscored.

Method	$N_c \times N_f$	Kinetics→BABEL
Source only	$3 \times 8 = 24$	11.7 ± 0.7
DANN [12]	$3 \times 8 = 24$	<u>29.3 ± 1.5</u>
CoMix [33]	$16 \times 8 = 128$	21.4 ± 0.3
CO2A [9]	$4 \times 16 = 64$	24.7 ± 0.8
GLAD (Ours)	$3 \times 8 = 24$	33.7 ± 1.8
Supervised target	$3 \times 8 = 24$	76.7 ± 2.1

improvement of 0.5 points. However, when we employ both debiasing and GLA together, we achieve a remarkable improvement of 11.3 points compared to the baseline, with a lower standard deviation (3.6 vs. 2.5). These results clearly indicate that the two methods are complementary to each other, generating a synergistic effect that enhances the overall performance of the UVDA model.

5.3. Comparison with state-of-the-arts

In this section, we compare the proposed method with state-of-the-art UVDA methods. We show the results in Table 3. ‘‘Source only’’ refers to the baseline method of training on labeled source data and testing on target data, which sets the lower bound for UVDA. ‘‘Supervised target’’ is an upper bound performance: a model trained with target data with labels. DANN [12] is an image-based domain adaptation method extended to UVDA. CoMix [33] and CO2A [9] are state-of-the-art UVDA methods. Surprisingly, we observe that the simple DANN method outperforms CoMix and CO2A on the challenging Kinetics→BABEL dataset. Our proposed method, GLAD, achieves the highest performance of 33.7%, surpassing DANN by 4.4 points. Notably, we achieve superior results with significantly fewer clips and frames compared to CoMix and CO2A, which highlights the high efficiency and accuracy of the proposed method.

5.4. Qualitative evaluation

In Figure 4, we show some qualitative examples from the Kinetics→BABEL to validate the effectiveness of GLAD. We compare the predictions of the baseline (DANN [13]) and GLAD on the BABEL dataset. The ground-truths for the four example videos are dance, clean_something, crawl and pick_up with durations of 27.0, 10.0, 2.7 and 1.9 seconds, respectively. In the example shown in Figure 4 (a) with dance action, the baseline fails to understand a long video of 27.0 seconds. The prediction bend implies the model focuses only on the bending motion which lasts for only 3 seconds in the video. The result might imply that the baseline tries to focus on a few key frames or local motions instead of focusing on the global temporal context when the

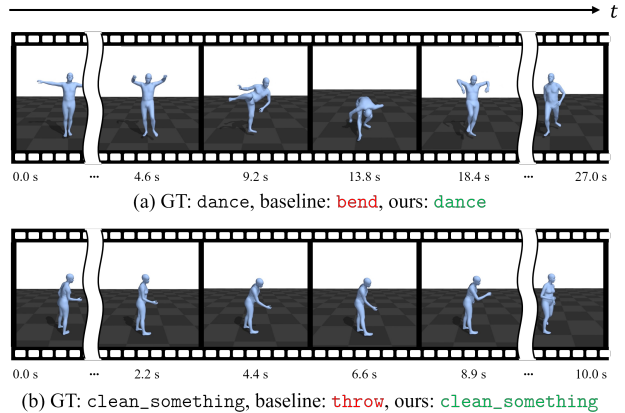


Figure 4. **Qualitative examples from Kinetics→BABEL.** We compare predictions of ours (GLAD) with predictions of a baseline (DANN [13]). GT denotes ground-truth, correct predictions are in green, and incorrect predictions are in red. We observe the baseline fails to predict correct actions due to the challenging temporal gap while GLAD consistently predicts correct actions.

action duration differs from the source data. Furthermore, for the example shown in Figure 4 (b), the baseline fails to distinguish clean_something from throw which involves understanding different speeds. In contrast, GLAD correctly predicts clean_something.

6. Conclusions

In this paper, we have addressed the challenging problem of unsupervised video domain adaptation for action recognition, specifically focusing on scenarios with a significant domain gap between the source and target domains. To overcome the limitations of existing datasets that are small in scale and lack significant domain gaps, we have introduced the Kinetics→BABEL dataset, which provides a more challenging, realistic, and large-scale benchmark. Our proposed method, GLAD, incorporates global-local view alignment to tackle temporal distribution shifts and background debiasing to address background distribution shifts. We have demonstrated the effectiveness of our proposed method through extensive experiments. Despite using fewer clips and frames compared to existing methods, GLAD has achieved favorable performance. The promising results highlight the efficacy and efficiency of our proposed method, paving the way for further advancements in unsupervised video domain adaptation for action recognition.

Acknowledgment. This work is supported by NCSOFT; by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea Government (MSIT) (Artificial Intelligence Innovation Hub) under Grant 2021-0-02068; by the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT) (No. 2022R1F1A1070997).

References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. ViViT: A Video Vision Transformer. In *ICCV*, 2021. 1, 3
- [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, 2009. 7
- [3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is Space-Time Attention All You Need for Video Understanding? In *ICML*, 2021. 1, 3
- [4] Joao Carreira and Andrew Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *CVPR*, 2018. 1, 3, 7, 8
- [5] Min-Hung Chen, Zsolt Kira, Ghassan AlRegib, Jaekwon Yoo, Ruxin Chen, and Jian Zheng. Temporal Attentive Alignment for Large-Scale Video Domain Adaptation. In *ICCV*, 2019. 2, 3, 4
- [6] Jinwoo Choi, Chen Gao, Joseph C. E. Messou, and Jia-Bin Huang. Why Can't I Dance in the Mall? Learning to Mitigate Scene Bias in Action Recognition. In *NeurIPS*, 2019. 2, 3
- [7] Jinwoo Choi, Gaurav Sharma, Samuel Schulter, and Jia-Bin Huang. Shuffle and attend: Video domain adaptation. In *ECCV*, 2020. 2, 3, 6
- [8] MMAction2 Contributors. Openmmlab's next generation video understanding toolbox and benchmark. <https://github.com/open-mmlab/mmaaction2>, 2020. 7
- [9] Victor G Turrissi da Costa, Giacomo Zara, Paolo Rota, Thiago Oliveira-Santos, Nicu Sebe, Vittorio Murino, and Elisa Ricci. Dual-head contrastive domain adaptation for video action recognition. In *WACV*, 2022. 2, 3, 4, 8
- [10] Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. Long-term Recurrent Convolutional Networks for Visual Recognition and Description. In *CVPR*, 2015. 3
- [11] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019. 1, 3
- [12] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. 5, 8
- [13] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 2016. 3, 5, 8
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4
- [15] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 3
- [16] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3D convolutional neural networks for human action recognition. *TPAMI*, 2013. 3
- [17] Leonid V Kantorovich. On the translocation of masses. *C. R. (Doklady) Acad. Sci. URSS (N. S.)*, 37:199–201, 1942. 4
- [18] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014. 3
- [19] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 3
- [20] Donghyun Kim, Yi-Hsuan Tsai, Bingbing Zhuang, Xiang Yu, Stan Sclaroff, Kate Saenko, and Manmohan Chandraker. Learning Cross-modal Contrastive Features for Video Domain Adaptation. In *ICCV*, 2021. 3
- [21] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *ICCV*, 2011. 2
- [22] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In *ECCV*, 2018. 3
- [23] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In *CVPR*, 2019. 3
- [24] Ji Lin, Chuang Gan, and Song Han. TSM: Temporal Shift Module for Efficient Video Understanding. In *ICCV*, 2019. 1, 3, 7
- [25] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *CVPR*, 2022. 1
- [26] Jonathan Munro and Dima Damen. Multi-Modal Domain Adaptation for Fine-Grained Action Recognition. In *CVPR*, 2020. 2, 3, 4
- [27] Joe Yue-Hei Ng, Matthew Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici. Beyond Short Snippets: Deep Networks for Video Classification. In *CVPR*, 2015. 3
- [28] Boxiao Pan, Zhangjie Cao, Ehsan Adeli, and Juan Carlos Nieves. Adversarial Cross-Domain Action Recognition with Co-Attention. In *AAAI*, 2020. 2, 3
- [29] Mandela Patrick, Dylan Campbell, Yuki M. Asano, Ishan Misra, Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi, and João F. Henriques. Keeping Your Eye on the Ball: Trajectory Attention in Video Transformers. In *NeurIPS*, 2021. 1, 3
- [30] Massimo Piccardi. Background subtraction techniques: a review. In *2004 IEEE international conference on systems, man and cybernetics (IEEE Cat. No. 04CH37583)*, 2004. 6
- [31] Abhinanda R. Punnakkal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra Quiros-Ramirez, and Michael J. Black. BABEL: Bodies, Action and Behavior with English Labels. In *CVPR*, 2021. 3
- [32] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. The earth mover's distance as a metric for image retrieval. *International Journal of Computer Vision*, 40(2):99–121, 2000. 4
- [33] Aadarsh Sahoo, Rutav Shah, Rameswar Panda, Kate Saenko, and Abir Das. Contrast and mix: Temporal contrastive video domain adaptation with background mixing. In *NeurIPS*, volume 34, 2021. 2, 3, 6, 8
- [34] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 2018. 3
- [35] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NeurIPS*, 2014. 3

- [36] Xiaolin Song, Sicheng Zhao, Jingyu Yang, Huanjing Yue, Pengfei Xu, Runbo Hu, and Hua Chai. Spatio-temporal contrastive domain adaptation for action recognition. In *CVPR*, 2021. 2
- [37] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 2
- [38] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015. 1, 3
- [39] Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schuster, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. Learning to adapt structured output space for semantic segmentation. In *CVPR*, 2018. 3
- [40] Pengfei Wei, Lingdong Kong, Xinghua Qu, Xiang Yin, Zhiqiang Xu, Jing Jiang, and Zejun Ma. Unsupervised video domain adaptation: A disentanglement perspective. *arXiv preprint arXiv:2208.07365*, 2022. 3
- [41] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. Self-supervised spatiotemporal learning via video clip order prediction. In *CVPR*, 2019. 6
- [42] Yuecong Xu, Jianfei Yang, Haozhi Cao, Keyu Wu, Wu Min, and Zhenghua Chen. Source-free Video Domain Adaptation by Learning Temporal Consistency for Action Recognition. In *ECCV*, 2022. 3
- [43] Lijin Yang, Yifei Huang, Yusuke Sugano, and Yoichi Sato. Interact before align: Leveraging cross-modal knowledge for domain adaptive action recognition. In *CVPR*, 2022. 2, 3
- [44] Youshan Zhang. A survey of unsupervised domain adaptation for visual recognition. *arXiv preprint arXiv:2112.06745*, 2021. 3
- [45] Yabin Zhang, Hui Tang, Kui Jia, and Mingkui Tan. Domain-symmetric networks for adversarial domain adaptation. In *CVPR*, 2019. 3
- [46] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal Relational Reasoning in Videos. In *ECCV*, 2018. 3
- [47] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *TPAMI*, 2017. 4