

Multimodal Generation of Novel Action Appearances for Synthetic-to-Real Recognition of Activities of Daily Living

Zdravko Marinov*

David Schneider*

Alina Roitberg*

Rainer Stiefelhagen

Institute for Anthropomatics and Robotics
Karlsruhe Institute of Technology, Germany

zdravko.marinov@kit.edu

david.schneider@kit.edu

alina.roitberg@kit.edu

rainer.stiefelhagen@kit.edu

Abstract—Domain shifts, such as appearance changes, are a key challenge in real-world applications of activity recognition models, which range from assistive robotics and smart homes to driver observation in intelligent vehicles. For example, while simulations are an excellent way of economical data collection, a SYNTHETIC→REAL domain shift leads to > 60% drop in accuracy when recognizing Activities of Daily Living (ADLs).

We tackle this challenge and introduce an activity domain generation framework which creates novel ADL appearances (*novel* domains) from different existing activity modalities (*source* domains) inferred from video training data. Our framework computes human poses, heatmaps of body joints, and optical flow maps and uses them alongside the original RGB videos to learn the essence of source domains in order to generate completely new ADL domains. The model is optimized by *maximizing* the distance between the existing source appearances and the generated novel appearances while ensuring that the semantics of an activity is preserved through an additional classification loss. While source data multimodality is an important concept in this design, our setup does not rely on multi-sensor setups, (i.e., all source modalities are inferred from a single video only.) The newly created activity domains are then integrated in the training of the ADL classification networks, resulting in models far less susceptible to changes in data distributions. Extensive experiments on the SYNTHETIC→REAL benchmark Sims4Action demonstrate the potential of the domain generation paradigm for cross-domain ADL recognition, setting new state-of-the-art results. Our code is publicly available at https://github.com/Zrzz1997/syn2real_DG.

I. INTRODUCTION

When roboticists apply visual activity recognition models in practice, they will quickly discover the problem of domain shifts. In fact, a model is rarely deployed under conditions identical to the ones in the training set, as we face changes in illumination, camera type and -placement [1], [2]. One domain change vital in robotic ADL assistance is the transition from synthetic to real data, as simulations ease the burden of intrusive data collection and privacy concerns in domestic environments [3], [4], [2], [5]. Especially in the light of the ageing population, domain-invariant recognition of Activities of Daily Living (ADL) is an important element for household robot perception and human-tailored planning [6], [7].

Recently, the paradigm of novel domain generation has been proposed as a way of extensive data augmentation for

* denotes equal contribution

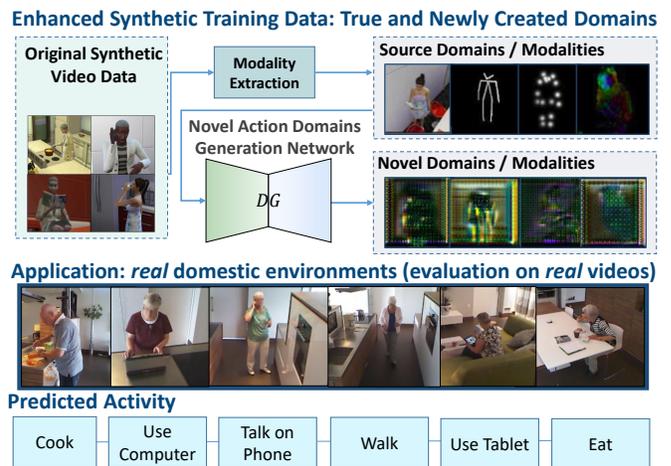


Fig. 1: An overview of the proposed SYNTHETIC→REAL activity recognition pipeline improved through neural network-based generation of novel activity domains. In addition to synthetic RGB videos, the training data is enriched with modalities explicitly extracted from videos (*source* modalities) and action representations from the *novel* modalities generated with our model. Such multimodal diversification of the training samples significantly mitigates adverse effects of the SYNTHETIC→REAL domain shift.

image recognition, leading to more domain-invariant models, e.g., for digit classification [8]. At the same time, activity analysis allows a plethora of true modalities based on human movement. Apart from raw RGB videos, modalities such as optical flow or body poses can be automatically extracted and used as *source* domains for learning to generate more diverse and activity-specific *novel* domains. Despite its high potential for mitigating domain shifts, such domain generation has been overlooked in activity analysis and is therefore the main motivation of our work.

We aim to make a step towards ADL recognition less susceptible to changes in data distribution and introduce a generative framework enriching the training data through generation of previously unseen activity domains. In our approach, multimodal action representations derived, e.g., from body pose, joint heatmaps and optical flow, are used to learn creating novel activity domains by *maximizing* the distance between the existing modalities (*source* domains) and the generated appearances (*novel* domains) while ensuring that the semantics of an activity are preserved through an additional *Classification Loss*. The newly generated novel modalities are then mixed with the initial source modalities

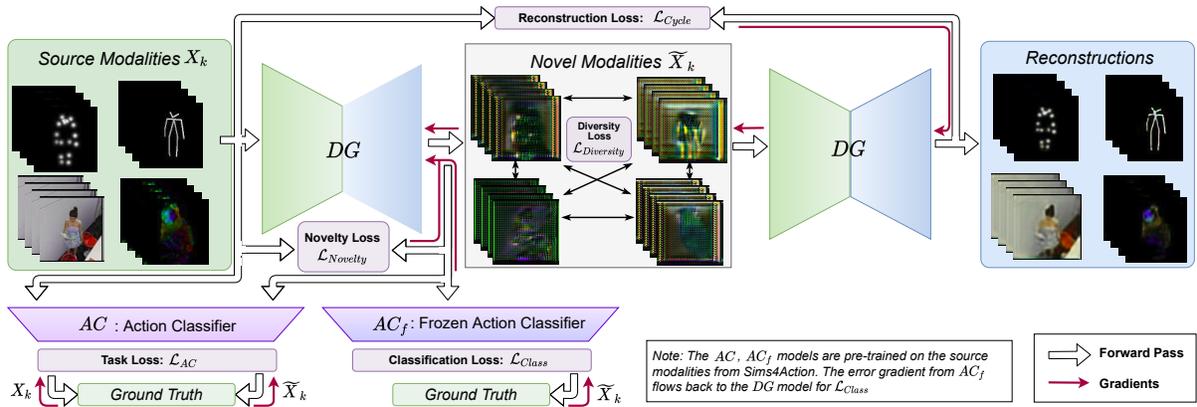


Fig. 2: Architecture of our multimodal framework for generation of novel ADL modalities. First, multiple modalities related to the body pose and movement are extracted and used alongside the original RGB videos as our source modalities X_k . Activity examples in these source modalities are then fed into the *Domain Generator (DG)* as single images in order to generate novel modalities. The distribution divergence between each source-novel domain pair is estimated via the *Novelty Loss*, whereas the *Diversity Loss* computes the divergence for each novel-novel domain pair. The source domains and the generated novel domains are both used as inputs to the *Action Classifier*, which makes separate predictions for both. The same generator is used to reconstruct the original input and the *Reconstruction Loss* is computed. The error signals for training our model via backpropagation are depicted as red arrows.

and constitute a diversified version of the training set (see overview in Figure 1). Multimodality is a central concept in our framework and we believe that the novel domain generation paradigm suits activity analysis especially well, as body pose and movement dynamics enable a wider range of existing source modalities, which in return encourages higher diversity of the generated novel domains. Nevertheless, our framework does not rely on multi-sensor setups since all source modalities are inferred from RGB videos only.

The problem of domain shift on the SYNTHETIC \rightarrow REAL benchmark is presented in Table I. The current state-of-the-art model is trained on the synthetic Sims4Action [4] dataset and experiences a large domain gap of $\approx 57\%$ when evaluated on real data. We conduct extensive experiments on the SYNTHETIC \rightarrow REAL ADL recognition benchmark [4] and make two observations. Firstly, multimodality itself is highly beneficial for cross-domain ADL recognition and using the body pose related data representations alone leads to a significant performance boost. Secondly, our idea of generating novel ADL appearances for training data enhancement consistently improves recognition results. Our framework outperforms state-of-the-art results, demonstrating the potential of multimodal domain generation for human activity analysis without using any pre-training on real data.

Training: Sims4Action [4]		Testing: Balanced Accuracy [%]		
Model	Pre-training	SYNTHETIC	REAL	Domain Gap
S3D [9]	Kinetics-400 [10]	84.61[†]	23.23	61.38
	None	56.52	12.40	44.12
I3D [11]	Kinetics-400 [10]	81.12	23.25[*]	57.87
	None	66.91	10.91	56.00

TABLE I: Current state-of-the-art results for the SYNTHETIC \rightarrow REAL^{*} and SYNTHETIC \rightarrow SYNTHETIC[†] benchmarks from *Let’s Play for Action* [4]. All four models are trained on Sims4Action. The SYNTHETIC and REAL test sets are Sims4Action [4] and Toyota Smarthome [12] respectively.

Note on our Terminology. Similarly to [13], [14], we define a *domain* as a joint distribution P_{XY} over a feature space \mathcal{X} and a label space \mathcal{Y} . In our work, we extract multiple *modalities* from synthetic data and learn to generate novel

modalities. Since each modality occupies a distinct feature subspace \mathcal{X}' and exhibits a unique appearance, we use the terms **modality** and **domain** interchangeably.

II. RELATED WORK

A. Recognizing Activities of Daily Living

To effectively interact with people, robots need to accurately perceive the current state of the human. Despite the impressive progress in general activity classification [15], [10], [16], [17], [9], [18], [19], [20] and a variety of frameworks introduced specifically for robotics applications [6], [21], [22], [23], [24], [25], this task remains very challenging in robotics, as agents often operate in a dynamic world where changes in concept-of-interest and data appearances may occur at any time [26]. In assistive robotics, recognizing Activities of Daily Living (ADL) is especially interesting and is often addressed by collecting and labelling new datasets tailored for the ADLs- and environments-of-interest [27], [28], [12], [21]. Creating such datasets which intend to realistically reflect real-world households requires larger efforts for sensory setups and data curation which results in datasets being smaller in comparison to general action classification benchmarks often created from web data [10], [29]. Methodologically, ADL recognition research is strongly influenced by architectures introduced in general video classification, with 3D Convolutional Neural Networks (CNNs) [11], [9], [30] being common backbone architectures [12], [31], but also more specialized approaches often derived from the body pose, have been introduced [32], [33], [24]. At the same time, recent research has raised alarming evidence, that deep learning-based ADL recognition approaches are very sensitive to changes in data distribution [4]. Mitigating this effect by exploring the domain generation paradigm [8] in the field of ADL recognition for the first time is the main contribution of our work.

B. Synthetic Human Actions

Given the difficulty of collecting labeled datasets, learning from simulated data has been researched in many different

fields of computer vision and is also emerging in video-based learning tasks, for example for pose recognition [34], [35], [36] and more recently in the domain of human activity recognition [37], [38], [4], [39], [5]. The latter works focus on either augmenting existing training data by mixing it with generated data [37], [40], [5], learning action categories on synthetic data only [4] or on learning compositions of actions within virtual domains [38]. [41] make use of a hybrid approach and combine real videos with rendered synthetic humans shown from different viewpoints. While synthetic examples are an excellent alternative to intrusive and time-consuming ADL dataset creation, the transition from simulations to real data at test-time comes with a remarkable performance drop [4] (see Table I). In this work, we focus on the SYNTHETIC→REAL transition in ADL recognition. We introduce a multimodal framework which leads to more domain-invariant recognition models by learning to generate new appearance versions of the synthetic training samples and by using them to diversify the training dataset.

C. Domain Generalization and Adaptation

Unsupervised domain adaptation methods learn a task on a source domain and try to solve this task on a target domain by learning a mapping given unlabelled data from the target domain, a task which has seen significant development in recent years in the field of video-based learning [42], [43], [44], [45], [46], [47], [48], [49]. In contrast, domain generalization describes the ability to maintain performance on a target domain despite not having access to any training data from this domain. Recently, Zhou et al. [8] proposed the domain generation paradigm, which is fundamentally different from previous work, as it learns to map source data to *unseen, newly generated* domains. Our work extends the image-based technique of [8], for the first time exploring it in the scope of ADL recognition and video recognition in general, which opens many additional possibilities of ADL-related source modalities, such as body poses or optical flow.

III. MULTIMODAL GENERATION OF ACTIVITY DOMAINS

In this section, we introduce a multimodal framework for better domain generalization in activity recognition, aiming to lighten the impact of appearance changes when moving from synthetic ADL training data to real-world robotic applications. We follow the domain generation paradigm [8] recently proposed for image recognition, extending it to the scope of human activity analysis. Conceptually, our

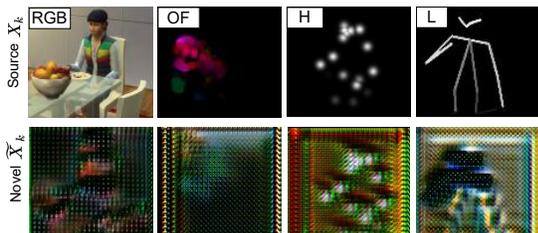


Fig. 3: Examples of the *source* X_k (top row) and *novel* modalities \tilde{X}_k (bottom row). H: Heatmaps, L: Limbs, OF: Optical Flow. We use all 8 modalities $X_k \cup \tilde{X}_k$ for training and evaluate our models on REAL data.

framework produces *novel* modalities \tilde{X}_k by learning to transform given *source* modalities X_k and comprises of four main building blocks: 1) A modality extraction module used to compute multiple source modalities from the original RGB video 2) a pre-trained domain classifier DC , 3) an action classifier AC , 4) and a domain generator DG . As we are using each extracted modality as a source *domain* we refer to the source and novel domains as source and novel *modalities*.

The task we address is SYNTHETIC→REAL domain generalization. We train an action classifier AC on samples from SYNTHETIC domains $x_s \in X_k \cup \tilde{X}_k$ with action labels $y \in Y$. In domain generalization, training and test data originate from distinct probability distributions, in our case $x_s \sim p_{synthetic}$ and $x_r \sim p_{real}$, and test samples x_r neither have labels, nor are seen during training. Our goal is to classify each instance x_r in the REAL target test domain X_r , which has a shared action label set Y with the training set. For this, we utilize the synthetic Sims4Action [4] dataset for training and the real Toyota Smarthome [12] and ETRI-Activity3D-LivingLab (ETRI) [21] as two separate test sets.

A. Extracting Source Modalities

The nature of activity recognition and video data in general allows us to leverage a wide range of modalities, such as body pose and movement dynamics, which would not be applicable in conventional image classification. We utilize four source modalities $X_k, k \in \{0, 1, 2, 3\}$ which are extracted directly from the training data (*i.e.*, RGB videos). The source modalities consist of 1) heatmaps of the body joint locations, 2) limbs connecting the joints as lines, 3) dense optical flow extracted between each two frames, and 4) raw RGB images (see top row of Figure 3). The heatmaps and limbs are extracted using the AlphaPose [50], [51], [52] pose detector, which infers 17 joint locations. The heatmaps modality $h(x, y)$ at pixel (x, y) is obtained by applying 2D gaussian maps, centered at each joint location (x_i, y_i) and weighted by its confidence c_i as seen in Equation 1.

$$h(x, y) = \exp\left(\frac{-((x - x_i)^2 + (y - y_i)^2)}{2\sigma^2}\right) \cdot c_i \quad (1)$$

The limbs domain is composed by connecting the joints with white lines and weighting each line by the smaller confidence of its endpoints. The optical flow is estimated using the Gunner-Farneback method [53]. We refer to these modalities as the four *source* modalities X_k (see top row in Figure 3).

B. Domain Classifier

The domain classifier DC is trained on the synthetic Sims4Action [4] to classify the source modalities X_k with

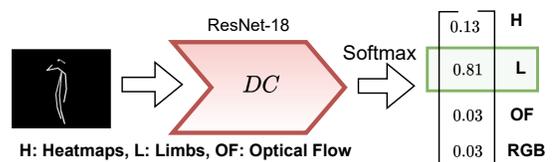


Fig. 4: Domain Classifier (DC): Overview of inference and training. The DC model is pre-trained end-to-end on Sims4Action.

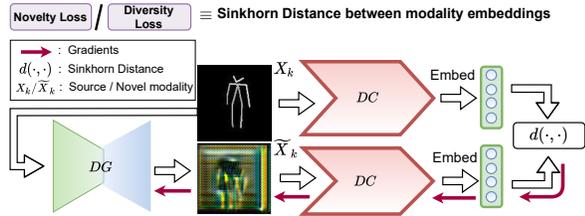


Fig. 5: Domain Classifier: Computation of the *Novelty Loss* for all source-novel modality pairs (X_k, \tilde{X}_k) . Note that the *Diversity Loss* is computed the same way with all novel-novel pairs $(\tilde{X}_k, \tilde{X}_l)_{k \neq l}$.

labels $k \in \{0, 1, 2, 3\}$ as shown in Figure 4. Afterwards, its weights are frozen and it is utilized for training the domain generator DG . The frozen DC is used to obtain embeddings from the source X_k and novel modalities \tilde{X}_k as seen in Figure 5. The Sinkhorn distance [54] between the embeddings is utilized as a distribution divergence metric and it is used to compute the *Novelty Loss* for each source-novel modality pair (X_k, \tilde{X}_k) and the *Diversity Loss* for all novel-novel pairs $(\tilde{X}_k, \tilde{X}_l)_{k \neq l}$ (see Equations 2, 3). The error gradient is propagated back to the domain generator DG and conditions it to produce novel modalities, which are both diverse and different from the source modalities.

C. Domain Generator

The goal of the domain generator DG is to extend and diversify the synthetic training data from Sims4Action [4]. The DG model is trained on the four source modalities X_k to generate four novel modalities $DG(X_k) = \tilde{X}_k$, which should be as diverse as possible, while remaining semantically and structurally consistent. The diversity of the new modalities is enforced by the *Novelty* and *Diversity* loss terms $\mathcal{L}_{Novelty}$ and $\mathcal{L}_{Diversity}$. The *Novelty Loss* maximizes the distribution divergence between the source and novel modalities, while the *Diversity Loss* maximizes the distribution divergence between each pair of generated novel modalities. The divergence measure we use is the Sinkhorn distance [54] $d(\cdot, \cdot)$ between the embeddings obtained by the domain classifier and is computed as illustrated in Figure 5. To ensure that the new modalities are dissimilar to the source modalities and are also dissimilar to each other these two loss terms are maximized w.r.t. DG as shown in Equations 2 and 3.

$$\mathcal{L}_{Novelty} = \max_{DG} \sum_{k=0}^3 d(\tilde{X}_k, X_k) \quad (2)$$

$$\mathcal{L}_{Diversity} = \max_{DG} \sum_{k=0}^3 \sum_{l=0}^3 d(\tilde{X}_k, \tilde{X}_l) \quad (3)$$

where $k, l \in \{0, 1, 2, 3\}$, $k \neq l$, and $\tilde{X}_k = DG(X_k)$.

Furthermore, the domain generator is conditioned to preserve the semantic and structural consistency of the actions in the novel modalities by minimizing the action *Classification Loss* \mathcal{L}_{Class} and the *Reconstruction Loss* \mathcal{L}_{Cycle} .

$$\mathcal{L}_{Class} = \min_{DG} \sum_{k=0}^3 \mathcal{L}_{CE}(AC_f(\tilde{X}_k), Y_k) \quad (4)$$

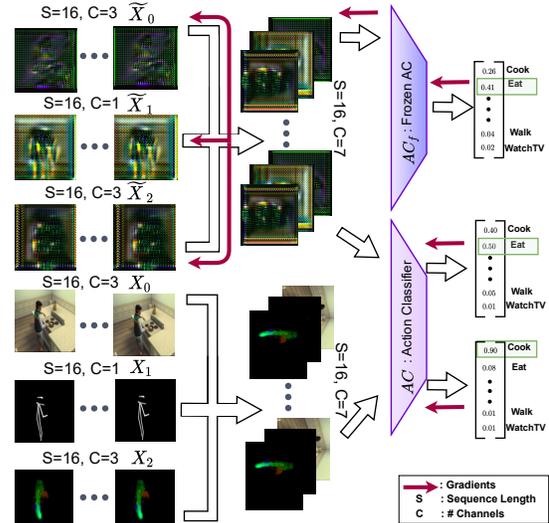


Fig. 6: Example of AC training with three modalities ($k = 2$). The input sequences are concatenated along the channel dimension C . The action classifier AC is trained on both the source X_k and novel \tilde{X}_k modalities for action recognition and makes separate predictions for both. Note that the input to the AC model is sequence-based ($S = 16$) as opposed to the image-based DC and DG ($S = 1$). Its frozen copy AC_f is pre-trained on the source modalities from Sims4Action and is only used for the computation of \mathcal{L}_{Class} in Equation 4. The error gradient from AC_f flows back to the novel modalities as its prediction is used to train DG with \mathcal{L}_{Class} (Eq. 4).

$$\mathcal{L}_{Cycle} = \min_{DG} \sum_{k=0}^3 \|DG(DG(X_k)), X_k\|_1 \quad (5)$$

where Y_k is the ground-truth action label, \mathcal{L}_{CE} is the cross-entropy loss, and AC_f is the frozen action classifier pre-trained on Sims4Action’s [4] source modalities X_k (explained in Section III-F). Note that we train on mini-batches. For this reason, X_k, \tilde{X}_k, Y_k refer to individual batches of size B in all equations, rather than the whole dataset.

D. Action Classifier

The action classifier AC is firstly pre-trained only on the source modalities X_k and then trained further on both X_k and \tilde{X}_k to assign the correct activity label. This encourages the model to learn representations in both the source and novel modalities and results in a larger and a more heterogeneous training set. The novel modalities increase the diversity of the training samples as their distribution differences are maximized by $\mathcal{L}_{Novelty}$ and $\mathcal{L}_{Diversity}$, which leads to a more versatile training dataset. The generation of \tilde{X}_k can hence be viewed as a method for data augmentation.

E. Model Architectures

As our action classifier AC , we utilize a 3D-CNN model to exploit the temporal features in the videos, which are important for action recognition [55], [30], [56]. 3D-CNN models such as the Inflated 3D-ConvNet (I3D) [11] have shown remarkable performance in recognition for ADL [12], [11]. Our action classification 3D-CNN model is implemented with the Separable 3D architecture (S3D) [9]. The S3D model leverages separable convolutions and replaces most of the 3D convolutions of I3D with cheaper 2D convolutions to reduce the complexity while boosting the performance

of I3D. For the domain classifier DC we employ ResNet18 [57] to learn to classify all source modalities (see Figure 4). Lastly, the domain generator DG , which also operates on single images, is modelled after Zhou et al. [8] and consists of two down-sampling conv-layers, two residual blocks [57] with instance normalization [58] and two transposed conv-layers to up-sample back to the input’s size.

F. Training Procedure

All of our models are trained *only* on Sims4Action [4]. The AC model is evaluated in the end on the real Toyota Smarthome [12] and ETRI [21] datasets to test its capability for domain generalization. The rest of the models are used only for the computation of loss terms - AC_f, DC , or for enhancing AC ’s training data - DG . The domain classifier DC is trained on Sims4Action to distinguish the source modalities X_k as seen in Figure 4. Then the DC model is frozen and used solely for the computation of $\mathcal{L}_{Novelty}$ and $\mathcal{L}_{Diversity}$ as in Figure 5 for the rest of the training.

Before we jointly train the other two models AC and DG , we pre-train an action classifier on all source modalities X_k . We then produce two copies of the pre-trained action classifier - the first copy AC_f is frozen and used only to enforce the semantic consistency of the domain generator in \mathcal{L}_{Class} (see Equation 4). The second copy AC is further fine-tuned on both the source X_k and novel modalities \tilde{X}_k with the cross-entropy loss. Only the AC model is evaluated in the end on the real Toyota Smarthome and ETRI datasets.

Loss Computation. In each iteration, a sequence of images is sampled from random chunks of the training videos from all source modalities X_k . The modality images from the sequence are concatenated along the channel dimension and are used to train the AC model. Each individual image from the sequence is transformed by DG into the novel modalities \tilde{X}_k , which are reshaped into image sequences (see Figure 6 top-left). The novel modality sequences are again concatenated along the channels and fed to AC and AC_f . The weights of AC_f are not updated, but its error is propagated back to the domain generator to compute \mathcal{L}_{Class} as seen in the red arrows in Figure 6. The *Novelty* and *Diversity* loss terms are computed with the help of the frozen DC model and the Sinkhorn distance [54] (see Figure 5). Finally the *Reconstruction Loss* is computed by iterating over all source and novel modality images and applying Equation 5. The final loss function \mathcal{L}_{DG} for the DG model is:

$$\mathcal{L}_{DG} = \min_{DG} \lambda_c \mathcal{L}_{Class} + \lambda_r \mathcal{L}_{Cycle} - \lambda_d (\mathcal{L}_{Novelty} + \mathcal{L}_{Diversity}) \quad (6)$$

where $\lambda_c, \lambda_r, \lambda_d$ are balancing parameters for the loss terms.

The AC model is trained on both the source and novel modality sequences. Hence, the *Task Loss* \mathcal{L}_{AC} is:

$$\mathcal{L}_{AC} = \min_{AC} \sum_{k=0}^3 \alpha \mathcal{L}_{CE}(AC(X_k), Y_k) + (1 - \alpha) \mathcal{L}_{CE}(AC(\tilde{X}_k), Y_k) \quad (7)$$

where α balances training on source X_k and novel modalities \tilde{X}_k and the other terms follow the notation of Equation 4.

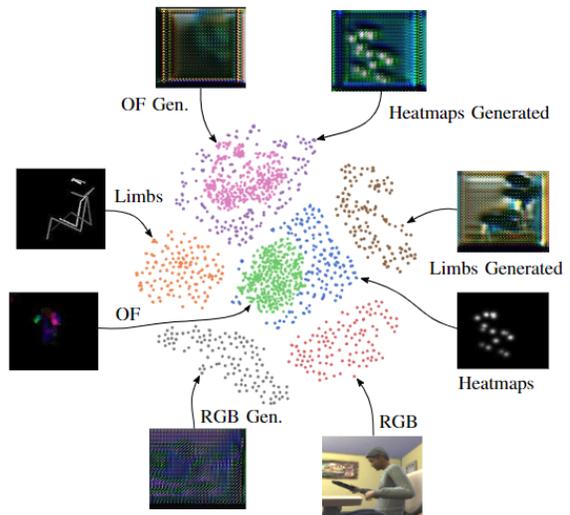


Fig. 7: Visualizations using t-SNE [59] embeddings from DG in both, *real* and *generated* ADL domains. All domains are marked with different colors.

G. Evaluation Procedure

Firstly, we evaluate the performance of S3D-based AC models without our domain generalization method. To investigate the effect of multimodality, we train 15 AC models on all $\sum_{i=1}^4 \binom{i}{4} = 15$ source modality combinations from Sims4Action. We use early fusion via channel concatenation for each of these combinations. Then, we evaluate all AC models on the real Toyota Smarthome [12] and ETRI [21] datasets to estimate their domain generalization capabilities.

Afterwards, we apply our domain generalization approach to each of these AC models and create two copies of each - AC and AC_f . For each modality combination, we initialize a new domain generator DG and train it alongside its corresponding AC and AC_f models on Sims4Action as described in Section III-F. A pre-trained domain classifier DC is re-used in all 15 training sessions. In the end, we evaluate all the AC models on the real Toyota Smarthome [12] and ETRI [21] datasets. The motivation for evaluating all modality combinations is to explore which modalities synergize well and to show that the novel domains improve the domain generalization for the vast majority of the 15 modality combinations.

H. Implementation Details

We use the same hyperparameters for all of our training sessions and experiments. The action classifiers are pre-trained on Sims4Action’s [4] source modalities for 200 epochs. Then the domain generator and action classifier are trained jointly on Sims4Action for 50 epochs as described in Section III-E. The input frame size is 112×112 and the sequence length for the S3D action classifier is $S = 16$ frames. The input videos are divided into chunks of 90 frames each following the protocol of [4]. During training we randomly sample 16 consecutive frames from the chunks and feed them to the domain generator and action classifier. We utilize the Adam optimizer [60] with $\eta = 10^{-4}$, $\beta_1 = 0.5$, and $\beta_2 = 0.999$. For regularization, we employ a weight

decay of $\lambda = 5 \cdot 10^{-5}$ for all models. We set $\lambda_c = \lambda_d = 1$, $\lambda_r = 10$, and $\alpha = 0.5$ for \mathcal{L}_{DG} and \mathcal{L}_{AC} . The embedding size of DC is 512 and we set $\sigma = 6$ for Equation 1.

IV. EXPERIMENTS

A. Benchmark Details

We focus on domain generalization between simulated and real data using Sims4Action [4] as a training dataset. For testing, we use real data, which originates from the Toyota Smarthome [12] and ETRI [21] datasets. Sims4Action provides ten hours of synthetic video material recorded from the computer game Sims 4, covering 10 basic human actions of daily living which have direct correspondences in the two real datasets. Toyota Smarthome [12] contains video material of 18 subjects performing 31 unscripted activities of daily living within a single apartment, and ETRI [21] is composed of 50 subjects performing 55 actions recorded from various perspectives that robots can be located in a home environment. However, we use only the 10 action correspondences to Sims4Action from both real datasets for our evaluation. We use the official cross-subject test split from Toyota Smarthome [31], [12], in which 7 of the subjects are reserved for testing. The same protocol is adopted in the state-of-the-art results [4] on the SYNTHETIC→REAL benchmark. For ETRI, we use the unsupervised domain adaptation protocol from [61] and use the whole dataset for testing. Note that during test-time we extract the pose and optical flow modalities from the real test samples so that they match the input format of the AC network.

As our evaluation data was collected in real households [12], [21], the number of samples per activity class is unbalanced and follows a Zipf-like distribution rather than a uniform one. We therefore focus on the *balanced accuracy* (mean per-class accuracy) as our main evaluation metric and additionally report the *unbalanced accuracy* (correct prediction rate among the complete test set). Note, that while we also evaluate the unbalanced accuracy to be consistent with [4], the metric is highly biased towards overrepresented categories in the test set and should be taken with caution especially in real-life datasets, where the categories are not evenly distributed. We therefore consider the balanced accuracy as a much more reliable and less biased metric, as per-class-averaged metrics are used in most of the unbalanced activity recognition datasets [12], [29], [62].

B. Qualitative Analysis

First, we inspect the domain generation results produced by the domain generator DG provided in Figure 3. The model has learned to produce novel appearances, which might seem unusual at first sight, but the activity semantics are preserved through the additional classification loss \mathcal{L}_{Class} . Note, that the movement typical for the activity is more obvious in video results, than in still images.

To better understand the learned representations, we visualize the output of the domain generator’s bottleneck layer with the t-distributed stochastic neighbor embeddings (tSNE) [59] in Figure 7. The embeddings are produced by sampling

10 random frames from each video from the Sims4Action dataset for each of the source modalities. Each image sample is fed to the domain generator and its embedding is used for the t-SNE visualization. The DG model is then applied again to the novel modalities \tilde{X}_k to obtain their embeddings. The plot in Figure 7 shows that each modality forms an individual cluster of points. This confirms that the generator has learned features which help differentiate between all 8 modalities $X_k \cup \tilde{X}_k$. Furthermore, the modalities span a wider and non-overlapping distribution in the latent space, i.e. the training dataset has been diversified by adding the novel modalities \tilde{X}_k , which leads to better domain generalization.

C. Quantitative Results

In Table III, we compare the ACs trained only on the source modalities X_k to the ones which are trained including the generated novel modalities \tilde{X}_k . All modalities $X_k \cup \tilde{X}_k$ are produced from Sims4Action [4] only. Additionally, we consider all 15 modality combinations and evaluate all models on the real Toyota Smarthome [12] and ETRI [21] datasets via the balanced and unbalanced accuracy metrics.

For most combinations, the extension of the training dataset with novel modalities increases the balanced accuracy significantly, in some cases by up to 13% points (ETRI: L). This is especially prevalent in multi-modal settings joined over early fusion, and we hypothesize that the generation of additional data alleviates overfitting problems which arise when training a larger, multi-modal model. Optical flow (OF) profits the least from this technique when combined with other modalities, in some cases performing worse than the baseline. This might be explained by OF’s low accuracy on the SYNTHETIC→SYNTHETIC benchmark (44%), i.e. OF is a weak modality on Sims4Action. However, performance losses in these cases are small in comparison to the significant gains which are achieved by the other combinations. We do list unbalanced accuracy as well, which provides similar results. Note that the improvement depends on the modality combination and test dataset. However, the average improvement for the balanced accuracy for Toyota Smarthome/ETRI is +3.1%/+2.8% and +4.7%/+8.0% for unbalanced accuracy.

Approach	Pre-training	Toyota SmartHome		Etri Activity	
		Balanced Accuracy	Unbalanced Accuracy	Balanced Accuracy	Unbalanced Accuracy
Random Choice	None	10.00	10.00	10.00	10.00
Currently reported state-of-the-art results on the SYNTHETIC→REAL benchmark					
S3D [11], [4]	None	12.40	19.95	11.71	13.86
S3D [11], [4]	Kinetics-400 [10]	23.25	22.75	23.45	28.57
Other domain generalization methods on the SYNTHETIC→REAL benchmark					
TA ³ N [43]	ImageNet [63]	14.19	14.44	25.11	35.12
APN [64]	ImageNet [63]	22.09	17.19	27.97	37.67
VideoDG [64]	ImageNet [63]	25.71	21.12	25.55	41.09
Our model with domain generation (best modality combination)					
Ours	None	27.73	32.20	29.05	41.88

TABLE II: Comparison of our model to state-of-the-art results on the SYNTHETIC→REAL benchmark [4] and to three additional domain generalization methods - TA³N [43], APN [64], and VideoDG [64].

D. Comparison to State-of-the-art and Other Approaches

In Table II we compare the best results achieved by our approach (i.e., the best modality variants from Table III

	Balanced Accuracy		Unbalanced Accuracy		Balanced Accuracy		Unbalanced Accuracy		Bal. Acc.	Unbal. Acc.	
	SYNTHETIC→REAL										SYNTHETIC→SYNTHETIC
Test Set	Toyota Smarthome [12]				ETRI [21]				Sims4Action [4]		
Train Set	Only X_k	Ours: $X_k \cup \tilde{X}_k$	Only X_k	Ours: $X_k \cup \tilde{X}_k$	Only X_k	Ours: $X_k \cup \tilde{X}_k$	Only X_k	Ours: $X_k \cup \tilde{X}_k$	Only X_k		
Input											
Individual											
RGB	13.7	13.3 (-0.4)	18.5	17.6 (-0.9)	11.7	15.0 (+3.3)	13.9	15.6 (+1.7)	61.8	59.4	
Heatmaps (H)	20.2	25.8 (+5.6)	20.0	23.7 (+3.7)	15.2	17.9 (+2.7)	22.2	40.9 (+18.7)	71.4	70.4	
Limbs (L)	22.0	27.7 (+5.7)	21.9	21.7 (-0.2)	16.1	29.1 (+13.0)	16.2	38.7 (+22.5)	75.1	74.4	
Optical Flow (OF)	21.3	22.6 (+1.3)	31.5	32.2 (+0.7)	11.6	14.7 (+3.1)	13.5	17.5 (+4.0)	44.5	43.7	
Early Fusion by Channel Concatenation											
RGB + H	10.3	20.7 (+10.4)	6.1	22.8 (+16.7)	7.3	12.9 (+5.6)	7.6	14.8 (+7.2)	76.8	77.2	
RGB + L	13.5	18.5 (+5.0)	8.2	18.6 (+10.4)	13.4	15.7 (+2.3)	15.9	23.5 (+7.6)	81.5	81.0	
RGB + OF	19.4	17.0 (-2.4)	31.0	30.4 (-0.6)	14.6	14.1 (-0.5)	20.9	41.9 (+21.0)	70.4	72.3	
H + L	15.7	25.3 (+9.6)	15.7	20.2 (+4.5)	15.0	17.8 (+2.8)	16.5	17.9 (+1.4)	57.8	50.1	
H + OF	16.7	23.7 (+7.0)	20.0	24.6 (+4.6)	12.6	14.2 (+1.6)	15.7	17.9 (+2.2)	71.1	64.7	
L + OF	19.7	19.8 (+0.1)	24.8	22.1 (-2.7)	11.4	12.5 (+1.1)	13.4	16.7 (+3.3)	72.7	71.2	
RGB + H + L	15.3	13.9 (-1.4)	10.0	23.2 (+13.2)	5.8	12.2 (+6.4)	5.8	13.8 (+8.0)	73.7	73.1	
RGB + H + OF	11.6	13.4 (+1.8)	13.0	29.5 (+16.5)	12.8	14.2 (+1.4)	12.5	22.3 (+9.8)	70.7	73.1	
RGB + L + OF	12.0	14.1 (+2.1)	15.0	13.5 (-1.5)	11.5	10.1 (-1.4)	15.9	17.8 (+1.9)	83.0	79.8	
H + L + OF	25.1	20.7 (-4.4)	36.8	22.9 (-13.9)	15.2	12.5 (-2.7)	13.9	15.6 (+1.7)	77.1	72.9	
RGB + H + L + OF	12.1	18.6 (+6.5)	7.0	27.4 (+20.4)	10.3	13.0 (+2.7)	9.1	18.9 (+9.8)	66.8	65.1	

TABLE III: Evaluation results for all 15 modality combinations on the SYNTHETIC→REAL benchmark for the real Toyota Smarthome [12] and ETRI [21] datasets and on the SYNTHETIC→SYNTHETIC benchmark for Sims4Action [4]. Our method of generating novel modalities \tilde{X}_k for the training data leads to a significant improvement of the vast majority of the models on the SYNTHETIC→REAL benchmark. The performance on Sims4Action is listed to illustrate the domain gap. The performance boost by adding novel modalities \tilde{X}_k to the training data is indicated in brackets for each modality combination.

marked in bold), and compare them with current state-of-the-art results on the Sims4Action SYNTHETIC→REAL benchmark [4]. We also compare our method to three other domain generalization methods: TA³N [43], APN [64], and VideoDG [64]. To this end, we train all three models end-to-end on Sims4Action [4] and follow the same evaluation protocols for Toyota Smarthome [12] and ETRI [21], which we used for our approach in Section IV-C.

Despite not making use of pre-training, we significantly improve the state-of-the-art on both accuracy metrics and on both real test datasets. Our method also consistently outperforms the other domain generalization methods TA³N [43], APN [64], and VideoDG [64]. APN and VideoDG also outperform the previous state-of-the-art on some metrics, which supports the claim that a domain generalization strategy for the SYNTHETIC→REAL benchmark is beneficial.

V. CONCLUSION AND LIMITATIONS

In this paper, we explored the paradigm of novel domain generation for recognizing Activities of Daily Living (ADL). Our work is motivated by the idea, that such synthesis of novel action appearances diversifies the training data and therefore mitigates the problem of domain shift. We specifically aim for SYNTHETIC→REAL ADL recognition and introduce a multimodal framework which leverages RGB, body pose, joint heatmaps and optical flow to learn generating novel modalities. Our experiments confirm that complementing training data with novel modalities leads to significant improvements in domain generalization, outperforming previous state-of-the-art results on the SYNTHETIC→REAL benchmark and other domain generalization approaches on both real test datasets Toyota Smarthome and ETRI.

While our method strongly improves the performance in case the data appearance has changed at test-time, it is not without limitations. First, we observe, that the model complexity increases with the number of source modalities, leading to a difficult optimization and longer training

times. Secondly, while we achieve state-of-the-art results on the SYNTHETIC→REAL benchmark, we acknowledge, that this comes with additional computational cost, as the body pose needs to be estimated first. Nevertheless, our work makes a step towards real-life utilization of synthetic datasets, which would enable far less intrusive data collection for raising action recognition capabilities of ADL robotic applications.

Acknowledgements. This work was supported by the JuBot project sponsored by the Carl Zeiss Stiftung and Competence Center Karlsruhe for AI Systems Engineering (CC-KING) sponsored by the Ministry of Economic Affairs, Labour and Housing Baden-Württemberg.

REFERENCES

- [1] M. Wulfmeier, A. Bewley, and I. Posner, "Addressing appearance change in outdoor robotics with adversarial domain adaptation," in *IROS*. IEEE, 2017, pp. 1551–1558.
- [2] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, "Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks," in *CVPR*, 2019, pp. 12627–12637.
- [3] C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. Vainio, C. Gokmen, G. Dharan, T. Jain, *et al.*, "Igbson 2.0: Object-centric simulation for robot learning of everyday household tasks," *arXiv preprint arXiv:2108.03272*, 2021.
- [4] A. Roitberg, D. Schneider, A. Djamal, C. Seibold, S. Reiß, and R. Stiefelhagen, "Let's play for action: Recognizing activities of daily living by learning from life simulation video games," in *IROS*, 2021.
- [5] H. Hwang, C. Jang, G. Park, J. Cho, and I.-J. Kim, "Eldersim: A synthetic data generation platform for human action recognition in eldercare applications," *IEEE Access*, pp. 1–1, 2021.
- [6] E. G. Christoforou, A. S. Panayides, S. Avgousti, P. Masouras, and C. S. Pattichis, "An overview of assistive robotics and technologies for elderly care," in *Mediterranean Conference on Medical and Biological Engineering and Computing*. Springer, 2019, pp. 971–976.
- [7] L. Marco and G. M. Farinella, *Computer vision for assistive health-care*. Academic Press, 2018.
- [8] K. Zhou, Y. Yang, T. Hospedales, and T. Xiang, "Learning to generate novel domains for domain generalization," in *ECCV*, 2020.
- [9] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy, "Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification," in *ECCV*, 2018, pp. 305–321.

- [10] 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.
- [11] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *CVPR*, 2017, pp. 6299–6308.
- [12] S. Das, R. Dai, M. Koperski, L. Minciullo, L. Garattoni, F. Bremond, and G. Francesca, “Toyota smarhome: Real-world activities of daily living,” in *ICCV*, 2019.
- [13] K. Zhou, Z. Liu, Y. Qiao, T. Xiang, and C. C. Loy, “Domain generalization: A survey,” 2021.
- [14] J. Wang, C. Lan, C. Liu, Y. Ouyang, T. Qin, W. Lu, Y. Chen, W. Zeng, and P. Yu, “Generalizing to unseen domains: A survey on domain generalization,” *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- [15] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, “Ntu rgb+ d: A large scale dataset for 3d human activity analysis,” in *CVPR*, 2016.
- [16] L. Smaira, J. Carreira, E. Noland, E. Clancy, A. Wu, and A. Zisserman, “A short note on the kinetics-700-2020 human action dataset,” 2020.
- [17] L. Wang, P. Koniusz, and D. Huynh, “Hallucinating IDT descriptors and I3D optical flow features for action recognition with CNNs,” *ICCV*, vol. 2019-October, pp. 8697–8707, 2019.
- [18] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, “A closer look at spatiotemporal convolutions for action recognition,” in *CVPR*, 2018, pp. 6450–6459.
- [19] J. Liu, A. Shahroudy, D. Xu, and G. Wang, “Spatio-temporal lstm with trust gates for 3d human action recognition,” in *ECCV*. Springer, 2016, pp. 816–833.
- [20] K. Hara, H. Kataoka, and Y. Satoh, “Learning spatio-temporal features with 3d residual networks for action recognition,” in *ICCV Workshops*, 2017, pp. 3154–3160.
- [21] J. Jang, D. Kim, C. Park, M. Jang, J. Lee, and J. Kim, “Etri-activity3d: A large-scale rgb-d dataset for robots to recognize daily activities of the elderly,” 2020.
- [22] A. Roitberg, N. Somani, A. Perzylo, M. Rickert, and A. Knoll, “Multimodal human activity recognition for industrial manufacturing processes in robotic workcells,” in *ICMI*, 2015.
- [23] C. R. Dreher, M. Wächter, and T. Asfour, “Learning object-action relations from bimanual human demonstration using graph networks,” *IEEE RAL*, vol. 5, no. 1, pp. 187–194, 2019.
- [24] M. Mofijul Islam and T. Iqbal, “Hamlet: A hierarchical multimodal attention-based human activity recognition algorithm,” *arXiv e-prints*, pp. arXiv-2008, 2020.
- [25] A. Roitberg, A. Perzylo, N. Somani, M. Giuliani, M. Rickert, and A. Knoll, “Human activity recognition in the context of industrial human-robot interaction,” in *APSIPA*. IEEE, 2014.
- [26] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford, *et al.*, “The limits and potentials of deep learning for robotics,” *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 405–420, 2018.
- [27] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta, “Hollywood in homes: Crowdsourcing data collection for activity understanding,” in *ECCV*. Springer, 2016, pp. 510–526.
- [28] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, *et al.*, “Scaling egocentric vision: The epic-kitchens dataset,” in *ECCV*, 2018.
- [29] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Nibbles, “Activitynet: A large-scale video benchmark for human activity understanding,” in *CVPR*, 2015, pp. 961–970.
- [30] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” in *ICCV*. IEEE, 2015.
- [31] R. Dai, S. Das, S. Sharma, L. Minciullo, L. Garattoni, F. Bremond, and G. Francesca, “Toyota smarhome untrimmed: Real-world untrimmed videos for activity detection,” *arXiv preprint arXiv:2010.14982*, 2020.
- [32] S. Das, S. Sharma, R. Dai, F. Bremond, and M. Thonnat, “Vpn: Learning video-pose embedding for activities of daily living,” in *ECCV*. Springer, 2020, pp. 72–90.
- [33] S. Das, R. Dai, D. Yang, and F. Bremond, “Vpn++: Rethinking video-pose embeddings for understanding activities of daily living,” *arXiv preprint arXiv:2105.08141*, 2021.
- [34] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, “Real-time human pose recognition in parts from single depth images,” in *CVPR*, 2011.
- [35] G. Varol, J. Romero, X. Martin, N. Mahmood, M. J. Black, I. Laptev, and C. Schmid, “Learning from synthetic humans,” in *CVPR*, 2017.
- [36] J. Liu, H. Rahmani, N. Akhtar, and A. Mian, “Learning human pose models from synthesized data for robust rgb-d action recognition,” *IJCV*, vol. 127, no. 10, pp. 1545–1564, 2019.
- [37] C. Roberto de Souza, A. Gaidon, Y. Cabon, and A. Manuel Lopez, “Procedural generation of videos to train deep action recognition networks,” in *CVPR*, 2017, pp. 4757–4767.
- [38] X. Puig, K. Ra, M. Boben, J. Li, T. Wang, S. Fidler, and A. Torralba, “Virtualhome: Simulating household activities via programs,” in *CVPR*, 2018, pp. 8494–8502.
- [39] D. Ludl, T. Gulde, and C. Curio, “Enhancing data-driven algorithms for human pose estimation and action recognition through simulation,” *TransITS*, vol. 21, no. 9, pp. 3990–3999, 2020.
- [40] Y. Zhang, X. Wei, W. Qiu, Z. Xiao, G. D. Hager, and A. Yuille, “Rsa: Randomized simulation as augmentation for robust human action recognition,” *arXiv preprint arXiv:1912.01180*, 2019.
- [41] G. Varol, I. Laptev, C. Schmid, and A. Zisserman, “Synthetic humans for action recognition from unseen viewpoints,” *IJCV*, vol. 129, no. 7, pp. 2264–2287, 2021.
- [42] J. Choi, G. Sharma, S. Schulter, and J.-B. Huang, “Shuffle and attend: Video domain adaptation,” in *ECCV*. Springer, 2020, pp. 678–695.
- [43] M.-H. Chen, Z. Kira, G. AlRegib, J. Yoo, R. Chen, and J. Zheng, “Temporal attentive alignment for large-scale video domain adaptation,” in *ICCV*, 2019, pp. 6321–6330.
- [44] M.-H. Chen, B. Li, Y. Bao, and G. AlRegib, “Action segmentation with mixed temporal domain adaptation,” in *WACV*, 2020, pp. 605–614.
- [45] P. P. Busto, A. Iqbal, and J. Gall, “Open set domain adaptation for image and action recognition,” *T-PAMI*, vol. 42, pp. 413–429, 2018.
- [46] A. Jamal, V. P. Namboodiri, D. Deodhare, and K. Venkatesh, “Deep domain adaptation in action space,” in *BMVC*, 2018.
- [47] J. Choi, G. Sharma, M. Chandraker, and J.-B. Huang, “Unsupervised and semi-supervised domain adaptation for action recognition from drones,” in *WACV*, 2020, pp. 1717–1726.
- [48] B. Pan, Z. Cao, E. Adeli, and J. C. Nibbles, “Adversarial cross-domain action recognition with co-attention,” in *AAAI*, 2020.
- [49] S. Reiß, A. Roitberg, M. Haurilet, and R. Stiefelhagen, “Deep classification-driven domain adaptation for cross-modal driver behavior recognition,” in *IV*. IEEE, 2020, pp. 1042–1047.
- [50] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, “RMPE: Regional multi-person pose estimation,” in *ICCV*, 2017.
- [51] J. Li, C. Wang, H. Zhu, Y. Mao, H.-S. Fang, and C. Lu, “Crowdpose: Efficient crowded scenes pose estimation and a new benchmark,” *arXiv preprint arXiv:1812.00324*, 2018.
- [52] Y. Xiu, J. Li, H. Wang, Y. Fang, and C. Lu, “Pose Flow: Efficient online pose tracking,” in *BMVC*, 2018.
- [53] G. Farnebäck, “Two-frame motion estimation based on polynomial expansion,” in *Scandinavian conference on Image analysis*. Springer, 2003, pp. 363–370.
- [54] M. Cuturi, “Sinkhorn distances: Lightspeed computation of optimal transport,” *NeurIPS*, vol. 26, 2013.
- [55] K. Hara, H. Kataoka, and Y. Satoh, “Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?” in *CVPR*, June 2018.
- [56] G. Varol, I. Laptev, and C. Schmid, “Long-term temporal convolutions for action recognition,” *T-PAMI*, 2017.
- [57] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [58] D. Ulyanov, A. Vedaldi, and V. Lempitsky, “Instance normalization: The missing ingredient for fast stylization,” *arXiv preprint arXiv:1607.08022*, 2016.
- [59] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [60] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [61] D. Kim, I. Lee, D. Kim, and S. Lee, “Action recognition using close-up of maximum activation and etri-activity3d livinglab dataset,” *Sensors*, vol. 21, no. 20, p. 6774, 2021.
- [62] M. Martin, A. Roitberg, M. Haurilet, M. Horne, S. Reiß, M. Voit, and R. Stiefelhagen, “Drive&Act: A Multi-modal Dataset for Fine-grained Driver Behavior Recognition in Autonomous Vehicles,” in *ICCV*. IEEE, October 2019.
- [63] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [64] Z. Yao, Y. Wang, J. Wang, P. Yu, and M. Long, “Videodg: Generalizing temporal relations in videos to novel domains,” *T-PAMI*, 2021.