# Layout-induced Video Representation for Recognizing Agent-in-Place Actions

Ruichi Yu[1]    Hongcheng Wang[2]    Ang Li[3*]    Jingxiao Zheng[1]    Vlad I. Morariu[4*]    Larry S. Davis[1]

[1]University of Maryland, College Park        [2]Comcast Applied AI Research

[3]DeepMind        [4]Adobe Research

{yrcbsg, jxzheng, lsd}@umiacs.umd.edu, hongcheng_wang@comcast.com

anglili@google.com, morariu@adobe.com

## Abstract

*We address scene layout modeling for recognizing agent-in-place actions, which are actions associated with agents who perform them and the places where they occur, in the context of outdoor home surveillance. We introduce a novel representation to model the geometry and topology of scene layouts so that a network can generalize from the layouts observed in the training scenes to unseen scenes in the test set. This Layout-Induced Video Representation (LIVR) abstracts away low-level appearance variance and encodes geometric and topological relationships of places to explicitly model scene layout. LIVR partitions the semantic features of a scene into different places to force the network to learn generic place-based feature descriptions which are independent of specific scene layouts; then, LIVR dynamically aggregates features based on connectivities of places in each specific scene to model its layout. We introduce a new Agent-in-Place Action (APA) dataset[1] to show that our method allows neural network models to generalize significantly better to unseen scenes.*

## 1. Introduction

Recent advances in deep neural networks have brought significant improvements to many fundamental computer vision tasks, including video action recognition [8, 35, 44, 46, 22, 36, 19, 4]. Current action recognition methods are able to detect, recognize or localize general actions and identify the agents (people, vehicles, *etc*.) [8, 35, 44, 22, 36, 19, 4, 7, 28, 26, 40, 38, 18]. However, in applications such as surveillance, relevant actions often involve locations and moving directions that relate to to scene layouts –for ex-



<vehicle, move away       <pet, move along, sidewalk>   <person, move away
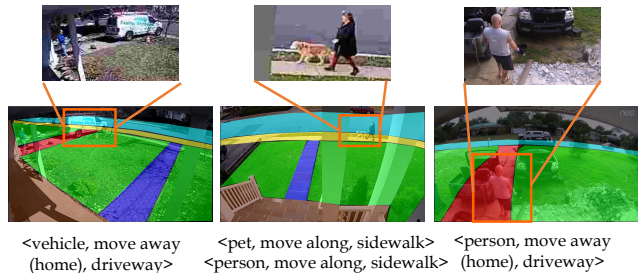(home), driveway>      <person, move along, sidewalk>     (home), driveway>

Figure 1. Example agent-in-place actions and segmentation maps. Different colors represent different places. We zoom in to the agents performing the actions for clarity. An agent-in-place action is represented as <*agent, action, place*>. Same colors indicate same place types (*e.g.,* green for lawn, blue for walkway, *etc.*).

ample, it might be of interest to detect (and issue an alert about) a person walking towards the front door of a house, but not to detect a person walking along the sidewalk. So, what makes an action "interesting" is how the it interacts with the geometry and topology of the scene.

Examples of these actions in outdoor home surveillance scenarios and the semantic segmentation maps of scenes are shown in Fig.1. We refer to these actions as "*agent-in-place*" actions to distinguish them from the widely studied generic action categories. From those examples we observe that although the types of place are limited (*e.g.,* street, walkway, lawn), the layout (*i.e.,* structure of places, including reference to location, size, appearance of places and their adjacent places) vary significantly from scene to scene. Without large-scale training data (which is hard to collect considering privacy issues), a naive method that directly learns from raw pixels in training videos without layout modeling can easily overfit to scene-specific textures and absolute pixel coordinates, and exhibit poor generalization on layouts of new scenes.

To address the generalization problem, we propose the *Layout-Induced Video Representation* (LIVR), which explicitly models scene layout for action recognition by en-

---

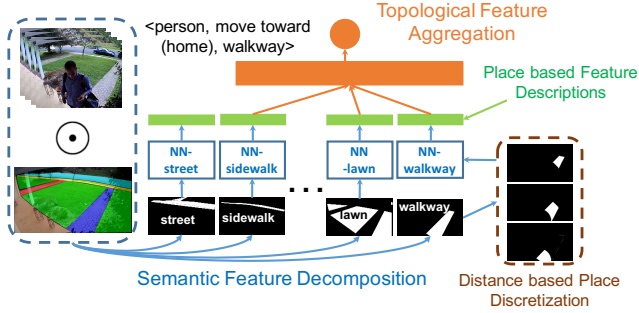[1]The dataset is pending legal review and will be released upon the acceptance of this paper.

Figure 2. Framework of LIVR. Given the segmentation map, we decompose the semantic features into different places and extract place-based feature descriptions individually, then dynamically aggregate them at inference time according to the topology of the scene. ⊙ denotes the masking operation for spatial decomposition. "NN" stands for neural network.



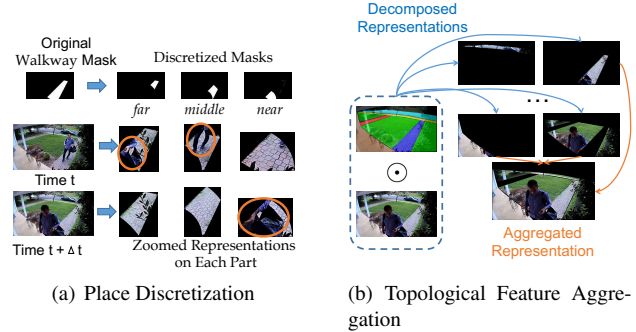(a) Place Discretization     (b) Topological Feature Aggregation

Figure 3. (a) illustrates distance-based place discretization. We segment the bit mask representing a given semantic class based on the distance transform of a second class, to explicitly represent the spatial-temporal order of moving agents which captures the moving direction w.r.t. that place. For example, this figure shows the partition of the *walkway* map into components that are "far," "middle" and "near" to the porch class. We use *move toward (home)* action as an example: we first observe the person on the part of *far* and *middle* (distance), and after some time, the person appears in the part of *near*. We use orange ellipses to highlight person parts. (b) illustrates the motivation behind topological feature aggregation. We seek a representation that covers the entire body of the person, which can be accomplished by aggregating the masked images from places that are connected to *walkway*.

coding the layout in the network architecture given semantic segmentation maps. The representation has three components: *1)* A semantic component represented by a set of bitmaps used for decomposing features in different *"places"* (*e.g.,* walkway, street, *etc.*), and force the network to learn place-based features that are independent of scene layout; *2)* A geometric component represented by a set of coarsely quantized distance transforms of each semantic place incorporated into the network to model moving directions; *3)* A topological component represented through the connection structure in a dynamically gated fully connected layer of the network–essentially aggregating representations from adjacent (more generally $h$-connected for $h$ hops in the adjacency graph of the semantic map) places. By encoding layout information (class membership of places, layout geometry and topology) into the network architecture using this decompotision-aggregation framework, we encourage our model to abstract away low-level appearance variations and focus on modeling high-level scene layouts, and eliminate the need to collect massive amounts of training data.

The first two components require *semantic feature decomposition* (Fig.2). We utilize bitmaps encoded with the semantic labels of places to decompose video representations into different places and train models to learn place-based feature descriptions. This decomposition encourages the network to learn features of generic place-based motion patterns that are independent of scene layouts. As part of the semantic feature decomposition, we encode scene geometry to model moving directions by discretizing a place into parts based on a quantized distance transform w.r.t. another place. Fig.2 (brown) shows the discretized bitmaps of *walkway* w.r.t. *porch*. As illustrated in Fig.3(a), features decomposed by those discretized bitmaps capture moving agents in spatial-temporal order, which reveals the moving direction, and can be generalized to different scene layouts.

The actions occurring in one place may be projected

onto adjacent places from the camera view(see Fig.3(b)). We propose *topological feature aggregation* to dynamically aggregate the decomposed features within the place associated with that action and adjacent places. The aggregation controls the "on/off" state of neuron connections from generic place-based feature descriptions to action nodes to model scene layout based on topological connectivity among places.

We created the Agent-in-Place Action (APA) dataset, which to the best of our knowledge, is the first dataset that addresses recognizing actions associated with scene layouts. APA dataset contains over 5,000 15s videos obtained from 26 different surveillance scenes with around 7,100 actions from 15 categories. To evaluate the generalization of LIVR, we split the scenes into observed and unseen scenes. Extensive experiments show that LIVR significantly improves the generalizability of the model trained on only observed scenes and tested on unseen scenes (improving the mean average precision (mAP) from around 20% to more than 50%). Consistent improvements are observed on almost all action categories.

## 2. Related Work

**Video Action Recognition Methods and Datasets.** Recent advances in video action recognition were driven by many large scale action recognition datasets. UCF101 [36], HMDB [22] and Kinetics [19] were widely used for recog-

nizing actions in video clips [40, 29, 45, 8, 35, 44, 7, 28, 26, 38, 18, 41]; THUMOS [17], ActivityNet[4] and AVA [13] were introduced for temporal/spatial-temporal action localization [33, 48, 27, 37, 52, 53, 3, 5, 24]. Recently, significant attention has been drawn to model human-human [13] and human-object interactions in daily actions [31, 34, 42]. In contrast to these datasets that were designed to evaluate motion and appearance modeling, or human-object interactions, our Agent-in-Place Action (APA) dataset is the first one that focuses on actions that are defined with respect to scene layouts, including interaction with places and moving directions. Recognizing these actions requires the network to not only detect and recognize agent categories and motion patterns, but also how they interact with the layout of the semantic classes in the scene. With the large variations of scene layouts, it is critical to explicitly model scene layout in the network to improve generalization on unseen scenes.

**Surveillance Video Understanding.** Prior work focuses on developing robust, efficient and accurate surveillance systems that can detect and track actions or events [14, 9, 6], or detect abnormal events [54, 15, 47]. The most related work to ours is ReMotENet [50], which skips expensive object detection [30, 12, 10, 11, 23] and utilizes 3D ConvNets to detect motion of an object-of-interest in outdoor home surveillance videos. We employ a similar 3D ConvNet model as proposed in [50] as a backbone architecture for extracting place-based feature descriptions for our model.

**Knowledge Transfer.** The biggest challenge of agent-in-place action recognition is to generalize a model trained with limited scenes to unseen scenes. Previous work on knowledge transfer for both images and videos has been based on visual similarity, which requires a large amount of training data [2, 16, 25, 39, 1, 49, 51]. For trajectory prediction, Ballan *et al.*[2] transferred the priors of statistics from training scenes to new scenes based on scene similarity. Kitani *et al.*[21] extracted static scene features to learn scene-specific motion dynamics for predicting human activities. Instead of utilizing low-level visual similarity for knowledge transfer, our video representation abstracts away appearance and location variance and models geometrical and topological relationships in a scene. which are more abstract and easier to generalize from limited training scenes.

# 3. Layout-Induced Video Representation

## 3.1. Framework Overview

The network architecture of layout-induced video representation is shown in Fig.4. For each video, we stack sampled frames of a video clip into a 4-D tensor. Our backbone network is similar to the architecture of ReMotENet [50],
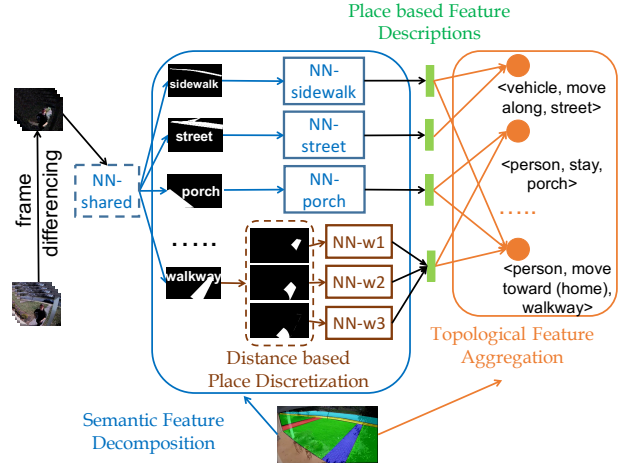


Figure 4. Layout-induced Video Representation Network: The dashed blue box indicates a shared 3D ConvNet to extract low-level features. We utilize the segmentation maps to decompose features into different places, and the solid blue boxes indicate that we train place-based models to extract place-based feature descriptions. When relevant to the activities of interest, we conduct distance-based place discretization to model moving directions; finally, we leverage the connectivity of places to aggregate the place-based feature descriptions at inference level.

which is composed of 3D Convolution (3D-conv) blocks. A key component of our framework is semantic feature decomposition, which decomposes feature maps according to region semantics obtained from given segmentation masks. This feature decomposition can be applied after any 3D-conv layer. Spatial Global Max Pooling (SGMP) is applied to extracted features within places, allowing the network to learn abstract features independent of shapes, sizes and absolute coordinates of both places and moving agents. For predicting each action label, we aggregate features from different places based on their connectivity in the segmentation map, referred to as Topological Feature Aggregation.

## 3.2. Semantic Feature Decomposition

**Segmentation Maps** Semantic Feature Decomposition utilizes a segmentation map of each place to decompose features and cause the network to extract place-based feature descriptions individually. The segmentation maps can be manually constructed using a mobile app we developed [2] to annotate each place by drawing points to construct polygons. This is reasonable since most smart home customers have only one or two cameras in their home. We employ human annotations because automatic semantic segmentation methods segment places based on appearance (*e.g.,* color, texture, *etc.*), while our task requires differentiation between places with similar appearance based on function-

---

[2]Details for the app can be found in the supplementary materials.

ality. For example, walkway, street and driveway have different functionalities in daily life, which can be easily and efficiently differentiated by humans. However, they may confuse appearance-based methods due to their similar appearance. Furthermore, since home surveillance cameras are typically fixed, users can annotate one map per camera very efficiently. However, we will discuss the performance of our method using automatically generated segmentation maps in Sec. 5.4.

**Place-based Feature Descriptions (PD).** Given a segmentation map, we extract place-based feature descriptions as shown in the blue boxes in Fig.4. We first use the segmentation map to decompose feature maps spatially into regions, each capturing the motion occurring in a certain place. The decomposition is applied to features instead of raw inputs to retain context information[3]. Let $\mathbf{X}_L \in \mathbb{R}^{w_L \times h_L \times t_L \times c}$ be the output tensor of the $L^{th}$ conv block, where $w_L, h_L, t_L$ denote its width, height and temporal dimensions, and $c$ is the number of feature maps. The place-based feature description of a place indexed with $p$ is

$$f_{L,p}(\mathbf{X}_L) = \mathbf{X}_L \odot \mathbb{I}[\mathbf{M}_L = p] \tag{1}$$

where $\mathbf{M}_L \in \mathbb{I}^{w_L \times h_L \times 1}$ is the segmentation index map and $\odot$ is a tiled element-wise multiplication which tiles the tensors to match their dimensions. Place descriptions can be extracted from different levels of feature maps. $L = 0$ means the input level; $L > 0$ means after the $L^{th}$ 3D-conv blocks. A higher $L$ generally allows the 3D ConvNet to observe more context and to abstract features. We treat $L$ as a hyper-parameter and study its effect in Sec. 5.

**Distance-based Place Discretization (DD).** Many actions are naturally associated with moving directions *w.r.t.* some scene element (*e.g.,* the house in home surveillance). To learn general patterns of the motion direction in different scenes, we further discretize the place segmentation into several parts, and extract features from each part and aggregate them to construct the place-based feature description of this place. For illustration, we use *porch* as the anchor place (shown in Fig.5). We compute the distance between each pixel and the *porch* in a scene (distance transform), and segment a place into $k$ parts based on their distances to *porch*. The left bottom map in 5 shows the porch distance transform of a scene. Let $D_L(x)$ be the distance transform of a pixel location $x$ in the $L^{th}$ layer. The value of a pixel $x$ in the part indexing map $\mathbf{M}_L^\Delta$ is computed as

$$M_L^\Delta(x) = \left\lfloor \frac{D_L^{\max}(x) - D_L^{\min}(x)}{k(D_L(x) - D_L^{\min}(x))} \right\rfloor \tag{2}$$

---

[3]An agent can be located at one place, but with part of its body projected onto another place in the view of the camera. If we use the binary map as a hard mask at input level, then for some places such as *sidewalk*, *driveway* and *walkway*, only a small part of the moving agents will remain after the masking operation.

where $D_L^{\max}(x) = \max\{D_L(x')|M_L(x') = M_L(x)\}$ and $D_L^{\min}(x) = \min\{D_L(x')|M_L(x') = M_L(x)\}$ are the max and min of pixel distances in the same place. They can be efficiently pre-computed. The feature description corresponding to the $i^{th}$ part of $p^{th}$ place in $L^{th}$ layer is

$$f_{L,p,i}^\Delta(\mathbf{X}_L) = \mathbf{X}_L \odot \mathbb{I}[\mathbf{M}_L = p \wedge \mathbf{M}_L^\Delta = i] \tag{3}$$

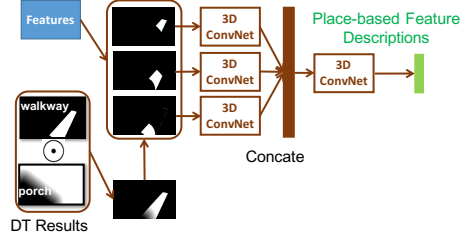where $\odot$ is the tiled element-wise multiplication.



Figure 5. The process of distance-based place discretization.

Discretizing a place into parts at different distances to the anchor place and explicitly separating their spatial-temporal features allows the representation to capture moving agents in spatial-temporal order and extract direction-related abstract features. However, not all places need to be segmented since some places (such as *sidewalk*, *street*) are not associated with any direction-related action (*e.g.,* moving toward or away from the house). For these places, we still extract the whole-place feature descriptors $f_{L,p}$. We discuss the effect of different choices of place discretization and the number of parts $k$, and show the robustness of our framework to these parameters in Sec. 5. To preserve temporal ordering, we apply 3D-conv blocks with spatial-only max pooling to extract features from each discretized place, and concatenate them channel-wise. Then, we apply 3D-conv blocks with temporal-only max pooling to abstract temporal information. Finally, we obtain a 1-D place-based feature description after applying GMP (see Fig.5). The final description obtained after distance-based place discretization has the same dimensionality as non-discretized place descriptions.

### 3.3. Topological Feature Aggregation (Topo-Agg)

Semantic feature decomposition allows us to extract a feature description for each place individually. To explicitly model the layout of a scene, we need to aggregate these place features based on connectivity between places. Each action category is one-to-one mapped to a place. To predict the confidence of an action $a$ occurring in a place $p$, features from adjacent places might provide contextual information, while the ones extracted far from place $p$ are distractors. To explicitly model the topological structure of places in a scene, we propose *Topological Feature Aggregation*, which utilizes the spatial connectivity between places, to guide feature aggregation.
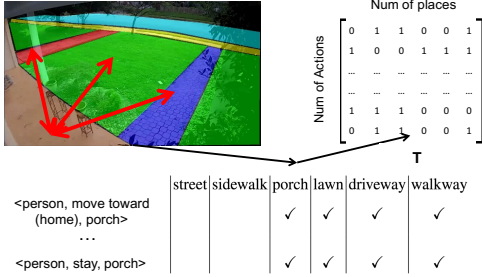
Figure 6. Topological feature aggregation which utilizes the connectivities between different places in a scene to guide the connections between the extracted place-based feature descriptions and the prediction labels. For clear visualization, we use the source places as *porch* in (b) with $h = 1$. The ✓ indicates we aggregate features from a certain place to infer the probability of an action.

Specifically, as shown in Fig.6, given a scene segmentation map, a source place $p$ and a constant $h$, we employ a Connected Component algorithm to find the $h$-connected set $C_h(p)$ which contains all places connected to place $p$ within $h$ hops. The constant $h$ specifies the minimum number of steps to walk from the source to a destination place. Given the $h$-connected place set $C_h$, we construct a binary action-place matrix ($\mathbf{T} \in \mathbb{R}^{n_a \times n_p}$) for the scene where $n_a$ is the number of actions and $n_p$ is the number of places. $\mathbf{T}_{i,j} = 1$ if and only if place $j$ is in the $C_h$ of the place corresponding to action $i$. Fig.6 shows an example segmentation map with its action-place mapping, where $C_0(porch) = \{porch\}$, $C_1(porch) = \{porch, walkway, driveway, lawn\}$, $C_2(porch)$ includes all except for *street*, and $C_3(porch)$ covers all six places. It is worth noting that since the vocabulary of our actions is closed, $\mathbf{T}$ is known at both training and testing time given the segmentation maps.

We implement topological feature aggregation using a gated fully connected layer with customized connections determined by the action-place mapping $\mathbf{T}$. Given $n_p$ $m$-D features extracted from $n_p$ places, we concatenate them to form a $(n_p \times m)$-D feature vector. We use $\mathbf{T}$ to determine the "on/off" status of each connection of a layer between the input features and the output action nodes. Let $\mathbf{T}_* = \mathbf{T} \otimes \mathbb{I}^{1 \times m}$ be the actual mask applied to the weight matrix $\mathbf{W} \in \mathbb{R}^{n_a \times n_p m}$ where $\otimes$ is the matrix Kronecker product. The final output is computed by $\mathbf{y} = (\mathbf{W} \odot \mathbf{T}_*)\mathbf{f}_*$, where $\odot$ is element-wise matrix multiplication, $\mathbf{f}_*$ is the concatenated feature vector as the input of the layer. We omit bias for simplicity. Let $J$ be the training loss function (cross-entropy loss), considering the derivative of $\mathbf{W}$, the gradient formulation is $\nabla_{\mathbf{W}} J = (\nabla_{\mathbf{y}} J \mathbf{f}_*^{\mathsf{T}}) \odot \mathbf{T}_*$, which is exactly the usual gradient ($\nabla_{\mathbf{s}} J \mathbf{f}^T$) masked by $\mathbf{T}_*$. At training time, we only back-propagate the gradients to connected neurons.

## 4. Agent-in-Place Action Dataset

We introduce a video dataset for recognizing agent-in-place actions. We collected outdoor home surveillance videos from internal donors and webcams to obtain over $7,100$ actions from around $5,000$ 15-second video clips with $1280 \times 720$ resolution. These videos are captured from 26 different outdoor cameras which cover various layouts of typical American front and back yards.

We select 15 common agent-in-place actions to label and each is represented as a tuple containing an action, the agent performing it, and the place where it occurs. The agents, actions, and places involved in our dataset are: $agent = \{person, vehicle, pet\}$; $action = \{move along, stay, move away (home), move toward (home), interact with vehicle, move across\}$; $place = \{street, sidewalk, lawn, porch, walkway, driveway\}$.

The duration of each video clip is 15s, so multiple actions can be observed involving one or more agents in one video. We formulate action recognition as a multi-label classification task. We split the 26 cameras into two sets: observed scenes (5) and unseen scenes (21) to balance the number of instances of each action in observed and unseen scenes and at the same time cover more scenes in the unseen set. We train and validate our model on observed scenes, and test its generalization capability on the unseen scenes. Details about the APA dataset including statistics can be found in the supplementary material.

## 5. Experiments

### 5.1. Implementation Details

**Network Architecture.** Unlike traditional 3D ConvNets which conduct spatial-temporal max-pooling simultaneously, we found that decoupling the pooling into spatial-only and temporal-only leads to better performance (experimental results and discussions can be found in the supplementary materials). So, for each place-specific network that extracts place-based feature descriptions, we utilize nine blocks of 3D ConvNets with the first five blocks using spatial-only max pooling and the last four blocks using temporal-only max pooling. The first two blocks have one 3D-conv layer each, and there are two convolutional (conv) layers with ReLU in between for the remaining blocks. For each place-specific network, we use 64 $3 \times 3 \times 3$ conv filters per 3D-conv layer. After conducting SGMP on features extracted by each place-specific network, the final concatenated 1-D feature dimension is $6 \times 64$ since there are 6 places in total. The inference is conducted with a gated fully connected layer, whose connections ("on/off" status) are determined by action labels and scene topology. We use the sigmoid function to obtain the predicted probability of each action. We decompose semantics after the second conv blocks ($L = 2$); we conduct distance-based place

Table 1. **The Path from Traditional 3D ConvNets to our Methods.** B/L1 and B/L2 are baseline models with raw pixels and and ConcateMap as input, respectively. For our proposed models: V1 uses segmentation maps to extract place-based feature descriptions only. V3 applies distance-based place discretization for some places. Both V1 and V3 use a FC layer to aggregate place features; V2 and V4 uses topological feature aggregation. H and FPS2 indicates using higher resolutions and FPS, and MF means using more filters per conv layer. Besides our baselines, we also compare LIVR with two state-of-the-art action recognition methods: [50, 43].

| Network Architecture | B/L1 | B/L2 | B/L2 +MF | LIVR-V1 | LIVR-V2 | LIVR-V3 | LIVR-V4 | LIVR-V4+H | LIVR-V4+MF | LIVR-V4+FPS2 | TSN [43] | ReMotENet [50] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3D ConvNet? | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | - | - |
| ConcateMap? | | ✓ | ✓ | | | | | | | | - | - |
| place-based feature description? | | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | - | - |
| distance-based place discretization? | | | | | | ✓ | ✓ | ✓ | ✓ | ✓ | - | - |
| topological feature aggregation? | | | | | ✓ | | ✓ | ✓ | ✓ | ✓ | - | - |
| higher resolutions? | | | | | | | | ✓ | | | - | - |
| more filters? | | | ✓ | | | | | | ✓ | | - | - |
| higher FPS? | | | | | | | | | | ✓ | - | - |
| Observed scenes mAP | 51.09 | 54.12 | 53.02 | 55.69 | 57.12 | 58.02 | **59.71** | 59.64 | 59.52 | 59.01 | 56.71 | 55.92 |
| Unseen scenes mAP | 19.21 | 21.16 | 20.45 | 41.57 | 43.78 | 47.76 | 50.65 | 49.03 | **50.98** | 49.56 | 23.21 | 22.05 |

discretization on $PL_{DT} = \{walkway, driveway, lawn\}$ and choose $k = 3$; for topological feature aggregation, we choose $h = 1$. The detailed place-specific network structure is shown in the supplementary materials.

**Anchor Place.** For our dataset, the directions mentioned are all relative to the house location, and *porch* is a strong indicator of the house location. So we only conduct distance transform to *porch* [4], but the distance-based place discretization method can be easily generalized to represent moving direction w.r.t any arbitrary anchor place.

**Training and Testing Details.** Our action recognition task is formulated as multi-label classification without mutual exclusion. The network is trained using the Adam optimizer [20] with 0.001 initial learning rate. For input video frames, we follow [50] to use FPS 1 and down-sample each frame to $160 \times 90$ to construct a $15 \times 160 \times 90 \times 3$ tensor for each video as input. Suggested by [50], small FPS and low resolution are sufficient to model actions for home surveillance where most agents are large and the motion patterns of actions are relatively simple. We evaluate the performance of recognizing each action independently and report Average Precision (AP) for each action and mean Average Precision (mAP) over all categories.

**Dataset Split.** We split the 26 scenes into two sets: observed scenes and unseen scenes. We further split the observed scenes into training and validation sets with a sample ratio of nearly $1 : 1$. The model is trained on observed scenes and test on unseen scenes. The validation set is used for tuning hyperparameters, which are robust with different choices (see Sec. 5.4).

### 5.2. Baseline Models

We follow [50] to employ 3D ConvNets as our baseline (B/L) model. The baseline models share the same 3D ConvNets architecture with our proposed model, except that the

last layer is fully connected instead of gated through topological feature aggregation. The difference between baselines is their input: B/L1 takes the raw frames as input; B/L2 incorporates the scene layout information by directly concatenating the 6 segmentation maps to the RGB channels in each frame (we call this method ConcateMap), resulting in an input of 9 channels per frame in total. We train the baseline models using the same setting as in the proposed model, and the performance of the baselines are shown in column 2-5 in Table 1. We observe that: 1) the testing performance gap between observed and unseen scenes is large, which reveals the poor generalization of the baseline models; 2) marginal improvements are obtained by incorporating scene layout information using ConcateMap, which suggests that it is difficulty for the network to learn the human-scene interactions directly from the raw pixels and segmentation maps. In addition, we also train a B/L2 model with $6\times$ more filters per layer to evaluate whether model size is the key factor for the performance improvement. The result of this enlarged B/L2 model is shown in column 5 of Table 1. Overall, the baseline models which directly extract features jointly from the entire video suffer from overfitting, and simply enlarging the model size or directly using the segmentation maps as features does not improve their generalization. More details about baseline models can be found in the supplementary materials.

### 5.3. Evaluation on the Proposed Method

**The path from traditional 3D ConvNets to our method.** We show the path from the baselines to our method in Table 1. In column 6-9, we report the mAP of our models on observed scene validation set and unseen scene testing set. We observe three significant performance gaps, especially on unseen scenes: 1) from B/L2 to LIVR-V1, we obtain around $20\%$ mAP improvement by applying the proposed semantic feature decomposition to extract place feature descriptions; 2) from LIVR-V1 to LIVR-V3, our model is further improved by explicitly modeling moving directions by place

---
[4]If there is no *porch* in a scene, we let the user to draw a line (click to generate two endpoints) to indicate its location.
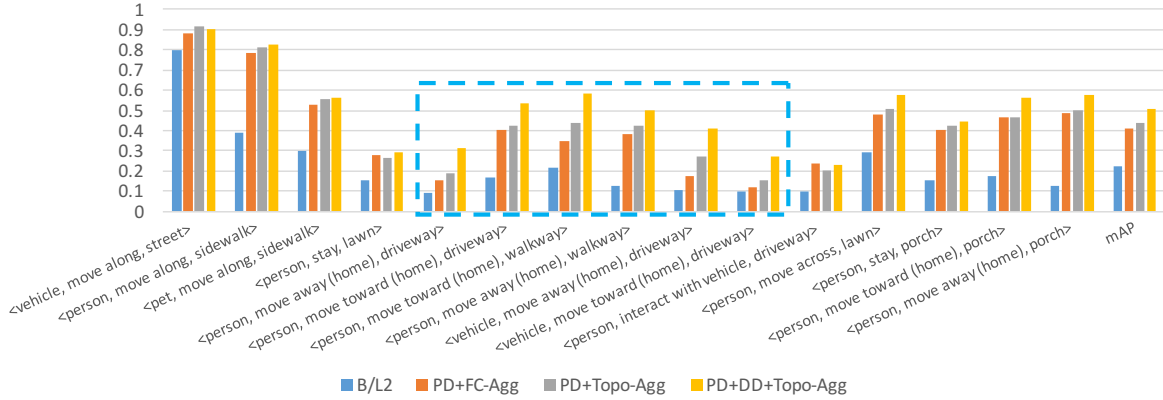
Figure 7. Per-category average precision of the baseline 3 and our methods on unseen scenes. The blue dashed box highlights actions which require modeling moving directions. We observe that the proposed place-based feature descriptions (PD), distance-based place discretization (DD) and topological feature aggregation (Topo-Agg) significantly improve the average precision on almost all action categories. FC-Agg stands for using a FC layer to aggregate place descriptions.



Figure 8. Qualitative examples: The predicted confidences of groundtruth actions using different methods. We use 3 frames to visualize a motion and orange ellipses to highlight moving agents.

discretization; 3) when compared to using a fully connected layer for feature aggregation (V1 and V3), our topological method (V2 and V4) leads to another significant improvement, which shows the efficacy of feature aggregation based on scene layout connectivity. We also evaluate the effect of resolutions, FPS and number of filters using our best model (LIVR-V4). Doubling the resolution ($320 \times 180$), FPS (2) and number of filters (128) only results in a slight change of the model's accuracy (columns 10-12 in Table 1). Besides our baselines, we also apply other state-of-the-art video action recognition methods (TSN [43] and ReMotENet [50]) on our dataset. LIVR outperforms them by a large margin, especially on the unseen scenes.

**Per-category Performance.** Fig.7 shows the average precision for each action on unseen scenes. LIVR outperforms the baseline methods by a large margin on almost all action categories. When comparing the orange and green bars in Fig.7, we observe that the proposed topological feature

aggregation (Topo-Agg) leads to consistently better generalization for almost all actions. The blue dashed box highlights the actions that include moving directions, and consistent improvements are brought by distance-based place discretization (DD). For some actions, especially the ones occurring on street and sidewalk, since they are relatively easy to recognize, adding DD or Topo-Agg upon the place-based feature descriptions (PD) does not help much. Overall, LIVR improves the generalization capability of the network, especially for actions that are more challenging, and are associated with moving directions.

**Qualitative Results.** Some example actions are visualized using three frames in temporal order and the predicted probabilities of the groundtruth actions using different methods are reported in Fig.8. It is observed that for relatively easy actions such as *<vehicle, move along, street>*, performance is similar across approaches. However, for challenging actions, especially ones requiring modeling moving directions such as *<person, move toward (home), walkway>*, our method outperforms baselines significantly.

## 5.4. Ablation Analysis on Unseen Scenes

**Place-based Feature Description.** The hyper-parameter for PD is the level $L$, controlling when to decompose semantics in different places. Fig.9(a) and 9(c) show that the generalization capability of our model is improved when we allow the network to observe the entire video at input level, and decompose semantics at feature level (after the 2nd conv blocks). Generally, the improvements of PD are robust across different feature levels.

**Distance-based Place Discretization.** We study different strategies for determining $PL_{DT}$ and the number of parts to discretize ($k$) per place. From our observations, including the anchor place–*porch*, the six places in our dataset can be

(a) Decomposition Level $L$     (b) Place Discretization     (c) Topological Aggregation     (d) Segmentation: GT v.s. Auto
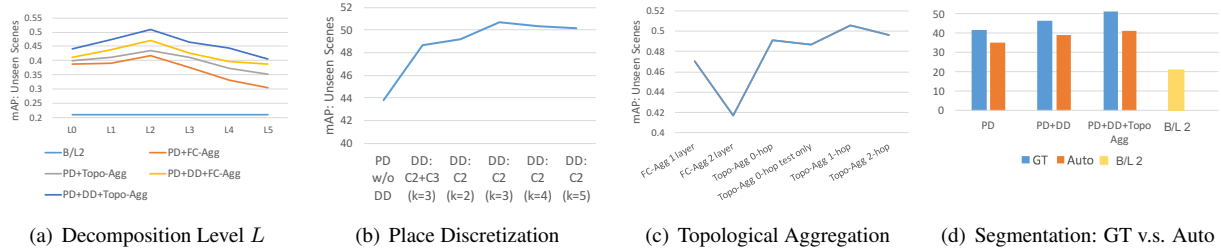
Figure 9. Evaluation: (a) The effect of extracting place-based feature descriptions (PD) at different levels using different variants of our proposed model. (b) Different strategies for distance-based place discretization. (c) Different feature aggregation approaches on unseen scenes. (d) Performance of LIVR using groundtruth (GT) and automatically generated (Auto) segmentation map.
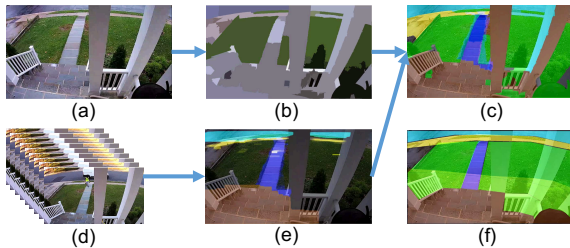


Figure 10. Process of Automatically Generating Segmentation Maps. (a) is the input camera image. (b) is the output of normalized cut method. (d) is the set of all videos captured by this camera. (e) shows the heatmaps we obtained by analyzing the patterns of moving objects from the videos. (c) is the generated segmentation map. (f) is the ground truth map.

clustered into three categories with regard to the distance to camera: C1 includes only *porch*, which is usually the closest place to camera; C2 includes *lawn*, *walkway*, *driveway*, and actions occurring in those places usually require modeling the moving direction directly; C3 includes *sidewalk* and *street*, which are usually far away from the house, and actions on them are not sensitive to directions (*e.g.,* "*move along*"). We evaluate our method with two strategies to apply DD on: 1) all places belong to C2 and C3; 2) only places in C2. The results are shown in Fig.9(b). We observe that applying DD on C3 dose not help much, but if we only apply DD on places in C2, our method achieves the best performance. In terms of the number of discretized parts $k$, we evaluate $k$ from 2 to 5 and observe from Fig.9(b) that the performance is robust when $k \geqslant 3$.

**Topological Feature Aggregation.** We evaluate different $h$ values to determine the $h$-connected set and different strategies to construct and utilize the action-place mapping **T**. The results are shown in Fig.9(c). We set $L = 2$, and use both PD and DD. We observe that Topo-Agg achieves its best performance when $h = 1$, *i.e.,* for an action occurring in place $P$, we aggregate features extracted from place $P$ and its directly connected places. In addition, we compare Topo-Agg to the naive fully connected inference layer (FC-Agg: 1 layer) and two fully-connected layers with 384 neu-

rons each and a ReLU layer in between (FC-Agg: 2 layers). Unsurprisingly, we observe that the generalizability drops significantly with an extra fully-connected layer, which reflects overfitting. Our Topo-Agg outperforms both methods. We also conduct an experiment where we train a fully connected inference layer and only aggregate features based on topology at testing time ("Topo-Agg: 1-hop test only") and it shows worse performance.

**LIVR with Automatically Generated Segmentation Maps.** To evaluate the performance of LIVR using imperfect segmentation maps, we developed an algorithm to automatically generate the segmentation maps. As shown in Fig.10, we first apply normalized cut [32] on the camera images to obtain segments (Fig.10 (b))[5]. Then, to further differentiate different places with similar appearance (*e.g.,* walkway and street), we developed an algorithm to utilize the historical statistics obtained from previous videos (Fig.10 (d)) of a scene to generate heatmaps of some specific places[6] (Fig.10 (e)). Then, the two results are combined to obtain final segmentation maps ((Fig.10 (c))). Our method can generate reasonably good segmentation maps when compared to the groundtruth maps obtained manually (Fig. 10 (f)). We evaluate LIVR using the imperfect maps and observe some performance degradations (around 10%), but LIVR still outperforms the baselines by a large margin (around 20%), which demonstrate the effectiveness of our method even if the segmentation maps are imperfect. Details of how we generate the segmentation maps can be found in the supplementary materials.

## 6. Conclusion

To improve the generalization of a deep network that learns from limited training scenes, we explicitly model scene layout in a network by using layout-induced video

---

[5]We also tried deep learning based methods trained on semantic segmentation datasets, but they perform poorly on our camera images. Details can be found in the supplementary materials.

[6]We utilize the patterns of moving objects to differentiate places. For example, street is a place where vehicles move with limited scale changes (from the camera perspective), and walkway is a place where people with notably large scale changes walk.

representations, which abstracts away low-level appearance variance but encodes the semantics, geometry and topology of scene layouts. An interesting future directions would be to include integrate the estimation of the semantic maps into the network architecture, which may require collecting more scenes for training.

# References

[1] L. Ballan, M. Bertini, G. Serra, and A. Del Bimbo. A data-driven approach for tag refinement and localization in web videos. *Computer Vision and Image Understanding*, 140:58–67, Nov. 2015. 3

[2] L. Ballan, F. Castaldo, A. Alahi, F. Palmieri, and S. Savarese. Knowledge transfer for scene-specific motion prediction. In *European Conference on Computer Vision (ECCV)*, 2016. 3

[3] S. Buch, V. Escorcia, B. Ghanem, L. Fei-Fei, and J. C. Niebles. End-to-end, single-stream temporal action detection in untrimmed videos. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2017. 3

[4] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 1, 3

[5] Y.-W. Chao, S. Vijayanarasimhan, B. Seybold, D. A. Ross, J. Deng, and R. Sukthankar. Rethinking the faster R-CNN architecture for temporal action localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3

[6] R. Collins, A. Lipton, T. Kanade, H. Fujiyoshi, D. Duggins, Y. Tsin, D. Tolliver, N. Enomoto, and O. Hasegawa. A system for video surveillance and monitoring. Technical Report CMU-RI-TR-00-12, Carnegie Mellon University, Pittsburgh, PA, May 2000. 3

[7] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, T. Darrell, and K. Saenko. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, pages 2625–2634. IEEE Computer Society, 2015. 1, 3

[8] C. Feichtenhofer, A. Pinz, and A. Zisserman. Convolutional two-stream network fusion for video action recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 1933–1941, 2016. 1, 3

[9] F. Fusier, V. Valentin, F. Brémond, M. Thonnat, M. Borg, D. Thirde, and J. Ferryman. Video understanding for complex activity recognition. *Machine Vision and Applications*, 18(3):167–188, Aug 2007. 3

[10] M. Gao, A. Li, R. Yu, V. I. Morariu, and L. S. Davis. C-wsl: Count-guided weakly supervised localization. *arXiv preprint arXiv:1711.05282*, 2017. 3

[11] M. Gao, R. Yu, A. Li, V. I. Morariu, and L. S. Davis. Dynamic zoom-in network for fast object detection in large images. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3

[12] R. Girshick. Fast R-CNN. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015. 3

[13] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar, C. Schmid, and J. Malik. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3

[14] I. Haritaoglu, D. Harwood, and L. S. Davis. W4: Real-time surveillance of people and their activities. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22:809–830, 2000. 3

[15] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis. Learning temporal regularity in video sequences. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 3

[16] J. Hays and A. A. Efros. Scene completion using millions of photographs. *ACM Trans. Graph.*, 26(3), July 2007. 3

[17] H. Idrees, A. R. Zamir, Y. Jiang, A. Gorban, I. Laptev, R. Sukthankar, and M. Shah. THE THUMOS challenge on action recognition for videos "in the wild". *CoRR*, abs/1604.06182, 2016. 3

[18] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(1):221–231, Jan. 2013. 1, 3

[19] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The kinetics human action video dataset. *CoRR*, abs/1705.06950, 2017. 1, 2

[20] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 6

[21] K. M. Kitani, B. D. Ziebart, J. A. Bagnell, and M. Hebert. Activity forecasting. In *Proceedings of the 12th European Conference on Computer Vision - Volume Part IV*, ECCV'12, pages 201–214, Berlin, Heidelberg, 2012. Springer-Verlag. 3

[22] H. Kuhne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. Hmdb: A large video database for human motion recognition. In *IEEE International Conference on Computer Vision (ICCV)*, 2011. 1, 2

[23] A. Li, J. Sun, J. Y.-H. Ng, R. Yu, V. I. Morariu, and L. S. Davis. Generating holistic 3d scene abstractions for text-based image retrieval. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3

[24] T. Lin, X. Zhao, and Z. Shou. Single shot temporal action detection. In *Proceedings of the 25th ACM International Conference on Multimedia*, MM '17, pages 988–996, New York, NY, USA, 2017. ACM. 3

[25] C. Liu, J. Yuen, and A. Torralba. *Nonparametric Scene Parsing via Label Transfer*, pages 207–236. Springer International Publishing, Cham, 2016. 3

[26] C. Ma, M. Chen, Z. Kira, and G. AlRegib. TS-LSTM and temporal-inception: Exploiting spatiotemporal dynamics for activity recognition. *CoRR*, abs/1703.10667, 2017. 1, 3

[27] C.-Y. Ma, A. Kadav, I. Melvin, Z. Kira, G. AlRegib, and H. Peter Graf. Attend and interact: Higher-order object interactions for video understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3

[28] J. Y.-H. Ng, M. J. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. Beyond short snippets: Deep networks for video classification. In *CVPR*, pages 4694–4702. IEEE Computer Society, 2015. 1, 3

[29] Z. Qiu, T. Yao, and T. Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *ICCV*, 2017. 3

[30] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 91–99. Curran Associates, Inc., 2015. 3

[31] M. Rohrbach, A. Rohrbach, M. Regneri, S. Amin, M. Andriluka, M. Pinkal, and B. Schiele. Recognizing fine-grained and composite activities using hand-centric features and script data. *International Journal of Computer Vision*, pages 1–28, 2015. 3

[32] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 22(8):888–905, Aug. 2000. 8

[33] Z. Shou, J. Chan, A. Zareian, K. Miyazawa, and S.-F. Chang. Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In *CVPR*, 2017. 3

[34] 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*, 2016. 3

[35] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1*, NIPS'14, pages 568–576, Cambridge, MA, USA, 2014. MIT Press. 1, 3

[36] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *CoRR*, abs/1212.0402, 2012. 1, 2

[37] C. Sun, A. Shrivastava, C. Vondrick, K. Murphy, R. Sukthankar, and C. Schmid. Actor-centric relation network. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XI*, pages 335–351, 2018. 3

[38] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler. Convolutional learning of spatio-temporal features. In *Proceedings of the 11th European Conference on Computer Vision: Part VI*, ECCV'10, pages 140–153, Berlin, Heidelberg, 2010. Springer-Verlag. 1, 3

[39] J. Tighe and S. Lazebnik. Superparsing: Scalable nonparametric image parsing with superpixels. In K. Daniilidis, P. Maragos, and N. Paragios, editors, *Computer Vision – ECCV 2010*, pages 352–365, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg. 3

[40] D. Tran, L. D. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 1, 3

[41] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri. A closer look at spatiotemporal convolutions for action recognition. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 6450–6459, 2018. 3

[42] G. Vaquette, A. Orcesi, L. Lucat, and C. Achard. The daily home life activity dataset: A high semantic activity dataset for online recognition. *2017 12th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2017)*, pages 497–504, 2017. 3

[43] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. V. Gool. Temporal segment networks: Towards good practices for deep action recognition. *CoRR*, abs/1608.00859, 2016. 6, 7

[44] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Val Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. 1, 3

[45] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. *CVPR*, 2018. 3

[46] C. Xu and J. J. Corso. Actor-action semantic segmentation with grouping-process models. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 1

[47] D. Xu, E. Ricci, Y. Yan, J. Song, and N. Sebe. Learning deep representations of appearance and motion for anomalous event detection. In M. W. J. Xianghua Xie and G. K. L. Tam, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 8.1–8.12. BMVA Press, September 2015. 3

[48] H. Xu, A. Das, and K. Saenko. R-c3d: Region convolutional 3d network for temporal activity detection. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017. 3

[49] X. Xu, T. M. Hospedales, and S. Gong. Discovery of shared semantic spaces for multiscene video query and summarization. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(6):1353–1367, 2017. 3

[50] R. Yu, H. Wang, and L. Davis. Remotenet: Efficient relevant motion event detection for large-scale home surveillance videos. In *IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2018. 3, 6, 7

[51] J. Yuen and A. Torralba. A data-driven approach for event prediction. In K. Daniilidis, P. Maragos, and N. Paragios, editors, *Computer Vision – ECCV 2010*, pages 707–720, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg. 3

[52] D. Zhang, X. Dai, and Y. Wang. Dynamic temporal pyramid network: A closer look at multi-scale modeling for activity detection. *CoRR*, abs/1808.02536, 2018. 3

[53] Y. S. Zhao, Y. Xiong, L. Wang, Z. Wu, D. Lin, and X. Tang. Temporal action detection with structured segment networks. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2933–2942, 2017. 3

[54] S. Zhou, W. Shen, D. Zeng, M. Fang, Y. Wei, and Z. Zhang. Spatial-temporal convolutional neural networks for anomaly detection and localization in crowded scenes. *Image Commun.*, 47(C):358–368, Sept. 2016. 3