# VPN: Learning Video-Pose Embedding for Activities of Daily Living

Srijan Das, Saurav Sharma, Rui Dai, François Brémond, and Monique Thonnat

INRIA Université Nice Côte d'Azur, France
`name.surname@inria.fr`

**Abstract.** In this paper, we focus on the spatio-temporal aspect of recognizing Activities of Daily Living (ADL). ADL have two specific properties (i) subtle spatio-temporal patterns and (ii) similar visual patterns varying with time. Therefore, ADL may look very similar and often necessitate to look at their fine-grained details to distinguish them. Because the recent spatio-temporal 3D ConvNets are too rigid to capture the subtle visual patterns across an action, we propose a novel Video-Pose Network: **VPN**. The 2 key components of this VPN are a spatial embedding and an attention network. The spatial embedding projects the 3D poses and RGB cues in a common semantic space. This enables the action recognition framework to learn better spatio-temporal features exploiting both modalities. In order to discriminate similar actions, the attention network provides two functionalities - (i) an end-to-end learnable pose backbone exploiting the topology of human body, and (ii) a coupler to provide joint spatio-temporal attention weights across a video. Experiments[1] show that VPN outperforms the state-of-the-art results for action classification on a large scale human activity dataset: **NTU-RGB+D 120**, its subset **NTU-RGB+D 60**, a real-world challenging human activity dataset: **Toyota Smarthome** and a small scale human-object interaction dataset **Northwestern UCLA**.

**Keywords:** action recognition, video, pose, embedding, attention

## 1 Introduction

Monitoring human behavior requires fine-grained understanding of actions. Activities of Daily Living (ADL) may look simple but their recognition is often more challenging than activities present in sport, movie or Youtube videos. ADL often have very low inter-class variance making the task of discriminating them from one another very challenging. The challenges characterizing ADL are illustrated in fig 1: (i) short and subtle actions like *pouring water* and *pouring grain* while *making coffee* ; (ii) actions exhibiting similar visual patterns while differing in motion patterns like *rubbing hands* and *clapping*; and finally, (iii) actions observed from different camera views. In the recent literature, the main

---

[1] Code/models: `https://github.com/srijandas07/VPN`

focus is the recognition of actions from internet videos [5,54,12,53,13] and very few studies have attempted to recognize ADL in indoor scenarios [16,4,8].
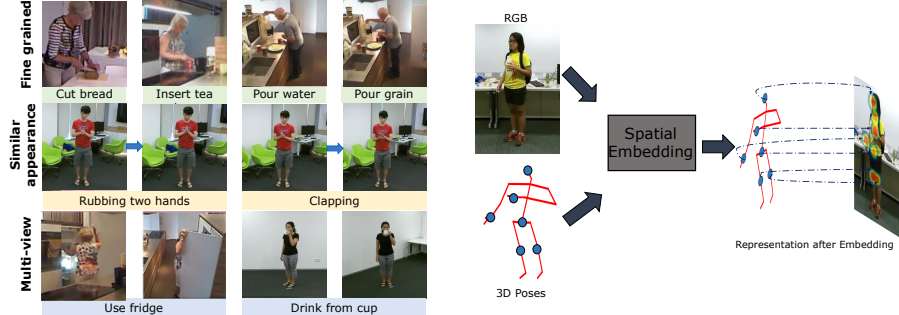


Fig. 1: Illustration of the challenges in Activities of Daily Living: fine-grained actions (top), actions with similar visual pattern (middle) and actions viewed from different cameras (below).

Fig. 2: Illustration of spatial embedding. Input is a RGB image and its corresponding 3D poses. For convenience, we only show 6 relevant human joints. The embedding enforces the human joints to represent the relevant regions in the image.

For instance, state-of-the-art 3D convolutional networks like I3D [5] pretrained on huge video datasets [17,45,21] have successfully boosted the recognition of actions from internet videos. But, these networks with similar spatio-temporal kernels applied across the whole space-time volume cannot address the complex challenges exhibited by ADL. Attention mechanisms have thus been proposed on top of these 3D convolutional networks to guide them along the regions of interest of the targeted actions [54,16,13]. Following a different direction, action recognition for ADL has been dominated by the use of human 3D poses [57,56]. They provide a strong clue for understanding the visual patterns of an action over time. 3D poses are robust to illumination changes, view adaptive and provide critical geometric information about human actions. However, they lack incorporating the appearance information which is an essential property in ADL (especially for human-object interaction).

Consequently, attempts have been made to utilize 3D poses to weight the discriminative parts of a RGB feature map [3,4,2,7,8]. These methods have improved the action recognition performance but they do not take into account the alignment of the RGB cues and the corresponding 3D poses. Therefore, we propose a spatial embedding to project the visual features and the 3D poses in the same referential. Before describing our contribution, we answer two intuitive questions below.

**First, why is spatial embedding important?** - Previous pose driven attention networks can be perceived as guiding networks to help the RGB cues focus on the salient information for action classification. For these guiding networks, it is important to have an accurate correspondences between the poses and RGB data. So, the objective of the spatial embedding is to find correspondences between the 3D human joints and the image regions representing these joints as illustrated in fig 2. This task of finding correlation between both modalities can

(i) provide informative pose aware feedback to the RGB cues, and (ii) improve the functionalities of the guiding network.

**Second, why not performing temporal embedding?** - We argue that the need of embedding is to provide proper alignment between the modalities. Across time, the 3D poses are already aligned assuming that there is a 3D pose for every images. However, even if the number of 3D poses does not correspond to the number of image frames (as in [3,4,2,7,8]), the fact that variance in poses for few consecutive frames is negligible, especially for ADL, implies temporal embedding is not needed.

We propose a recognition model based on a Video-Pose Network, **VPN** to recognize a large variety of human actions. VPN consists of a spatial embedding and an attention network. VPN exhibits the following novelties: (i) a spatial embedding learns an accurate video-pose embedding to enforce the relationships between the visual content and 3D poses, (ii) an attention network learns the attention weights with a tight spatio-temporal coupling for better modulating the RGB feature map, (iii) the attention network takes the spatial layout of the human body into account by processing the 3D poses through Graph Convolutional Networks (GCNs).

The proposed recognition model is end-to-end trainable and our proposed VPN can be used as a layer on top of any 3D ConvNets.

## 2 Related Work

Below, we discuss the relevant action recognition algorithms w.r.t. their input modalities.

**RGB** - Traditionally, image level features [49,50] have been aggregated over time using encoding techniques like Fisher Vector [34] and NetVLAD [1]. But these video descriptors do not encode long-range temporal information. Then, temporal patterns of actions have been modelled in videos using sequential networks. These sequential networks like LSTMs are fed with convolutional features from images [10] and thus, they model the temporal information based on the evolution of appearance of the human actions. However, these methods first process the image level features and then capture their temporal evolution preventing the computation of joint spatio-temporal patterns over time.

Due to this reason, Du et al. [48] have proposed 3D convolution to model the spatio-temporal patterns within an action. The 3D kernels provide tight coupling of space and time towards better action classification. Later on, holistic methods like I3D [5], slow-fast network [12], MARS [6] and two-in-one stream network [59] have been fabricated for generic datasets like Kinetics [17] and UCF-101 [45]. But these networks are trained globally over the whole 3D volume of a video and thus, are too rigid to capture salient features for subtle spatio-temporal patterns for ADL.

Recently several attention mechanisms have been proposed on top of the aforementioned 3D ConvNets to extract salient spatio-temporal patterns. For instance, Wang et al. [54] have proposed a non-local module on top of I3D which

computes the attention of each pixel as a weighted sum of the features of all pixels in the space-time volume. But this module relies too much on the appearance of the actions, i.e., pixel position within the space-time volume. As a consequence, this module though effective for the classification of actions in internet videos, fails to disambiguate ADL with similar motion and fails to address view invariant challenges.

**3D Poses** - To focus on the view-invariant challenge, temporal evolution of 3D poses have been leveraged through sequential networks like LSTM and GRU for skeleton based action recognition [57,26,56]. Taking a step ahead, LSTMs have also been used for spatial and temporal attention mechanisms to focus on the salient human joints and key temporal frames [44]. Another framework represents 3D poses as pseudo image to leverage the successful image classification CNNs for action classification [11,27]. Recently, graph-based methods model the data as a graph with joints as vertexes and bones as edges [55,47,40]. Compared to sequential networks and pseudo image based methods, graph-based methods make use of the spatial configuration of the human body joints and thus, are more effective. However, the skeleton based action recognition lacks in encoding the appearance information which is critical for ADL recognition.

**RGB + 3D Poses** - In order to make use of the pros of both modalities, i.e. RGB and 3D Poses, it is desirable to fuse these multi-modal information into an integrated set of discriminative features. As these modalities are heterogeneous, they must be processed by different kinds of network to show their effectiveness. This limits their performance in simple multi-modal fusion strategy [38,23,30]. As a consequence, many pose driven attention mechanisms have been proposed to guide the RGB cues for action recognition. In [2,3,4], the pose driven attention networks implemented through LSTMs, focus on the salient image features and the key frames. Then, with the success of 3D CNNs, 3D poses have been exploited to compute the attention weights of a spatio-temporal feature map. Das et al. [7] have proposed a spatial attention mechanism on top of 3D ConvNets to weight the pertinent human body parts relevant for an action. Then, authors in [8] have proposed a more general spatial and temporal attention mechanism in a dissociated manner. But these methods have the following drawbacks: (i) there is no accurate correspondence between the 3D poses and the RGB cues in the process of computing the attention weights [2,3,4,7,8]; (ii) the attention sub-networks [2,3,4,7,8] neglect the topology of the human body while computing the attention weights; (iii) the attention weights in [7,8] provide identical spatial attention along the video. As a result, action pairs with similar appearance like *jumping* and *hopping* are mis-classified.

In contrast, we propose a new spatial embedding to enforce the correspondences between RGB and 3D pose which has been missing in the state-of-the-art methods. The embedding is built upon an end-to-end learnable attention network. The attention network considers the human topology to better activate the relevant body joints for computing the attention weights. To the best of our knowledge, none of the previous action recognition methods have combined human topology with RGB cues. In addition, the proposed attention network

couples the spatial and temporal attention weights in order to provide spatial attention weights varying along time.
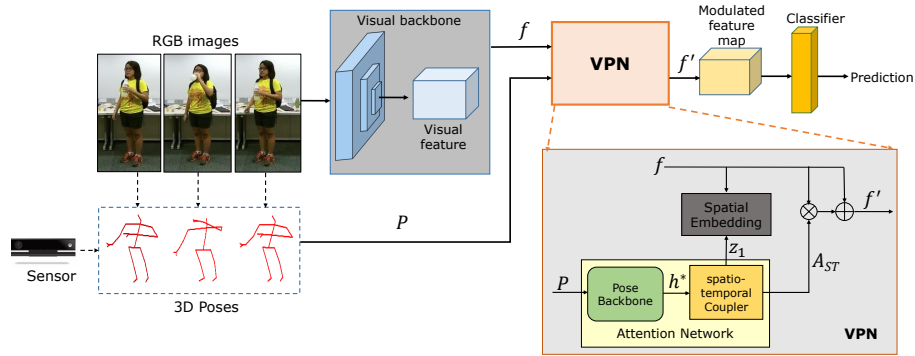


Fig. 3: **Proposed Action Recognition Model**: Our model takes as input RGB images with their corresponding 3D poses. The RGB images are processed by a visual backbone which generates a spatio-temporal feature map ($f$). The proposed **VPN** takes as input the feature map ($f$) and the 3D poses ($P$). VPN consists of two components: an attention network and a spatial embedding. The attention network further consists of a Pose Backbone and a spatio-temporal Coupler. VPN computes a modulated feature map $f'$. This modulated feature map $f'$ is then used for classification.

## 3  Proposed Action Recognition Model

Our objective is to design an accurate spatial embedding of poses and visual content to better extract the discriminative spatio-temporal patterns. As shown in fig. 3, the input of our proposed recognition model are the RGB images and their 3D poses. The 3D poses are either extracted from depth sensor or from RGB using LCRNet [37]. The proposed Video-Pose Network **VPN** takes as input the visual feature map and the 3D poses. Below, we discuss the action recognition model in details.

### 3.1  Video Representation

Taking as input a stack of human cropped images from a video clip, the spatio-temporal representation $f$ is computed by a 3D convolutional network (the visual backbone in fig. 3). $f$ is a feature map of dimension $t_c \times m \times n \times c$, where $t_c$ denotes the temporal dimension, $m \times n$ the spatial scale and $c$ the channels. Then, the feature map $f$ and the corresponding poses $P$ are processed by the proposed network.

### 3.2   VPN

VPN can be thought as a layer which can be placed on top of any 3D convolutional backbone. VPN takes as input a 3D feature map ($f$) and its corresponding 3D poses ($P$) to perform two functionalities. First, to provide an accurate alignment of the human joints with the feature map $f$. Second, to compute a modulated feature map ($f'$) which is further classified for action recognition. The modulated feature map ($f'$) is weighted along space and time as per its relevance. VPN exploits the highly informative 3D pose information to transform the visual feature map $f$ and finally, compute the attention weights. This network has two major components as shown in fig 4: (I) an attention network and (II) a spatial embedding. Though the intrinsic parameters of the attention network and the spatial embedding learns in parallel, we present these two components in the following order for better understanding.
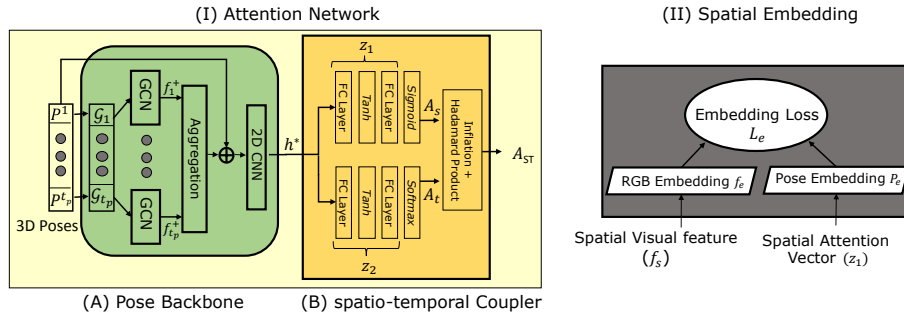


Fig. 4: The components in VPN: (I) Attention Network (left) and (II) Spatial Embedding (right). We present a zoom of the attention Network with: (A) a GCN Pose Backbone, and (B) a spatio-temporal Coupler to generate spatio-temporal attention weights $A_{ST}$

**(I) Attention Network -** The attention network consists of a Pose Backbone and a spatio-temporal Coupler. Such a framework for pose driven attention network is unique compared to the other state-of-the-art methods using poses and RGB. The proposed attention network unlike [3,4,7,8] takes into account the human spatial configuration and it also learns coupled spatio-temporal attention weights for the visual feature map $f$.

**Pose Backbone -** The input poses along the video are processed in a Pose Backbone. The pose based input of VPN are the 3D human joint coordinates $P \in \mathbb{R}^{3 \times J \times t_p}$ stacked along $t_p$ temporal dimension, where $J$ is the number of skeleton joints. The Pose Backbone processes these 3D poses to compute pose features $h^*$ which are used further in the attention network for computing

the spatio-temporal attention weights. They carry meaningful information in a compact way, so the proposed attention network can efficiently focus on salient action parts.

For the Pose Backbone, we use **GCNs** to learn the spatial relationships between the 3D human joints to provide attention weights to the visual feature map $(f)$. We aim at exploiting the graphical structure of the 3D poses. In fig. 4(I), we illustrate our GCN pose backbone (marked (A)). For each pose input $P_t \in \mathbb{R}^{3 \times J}$ with $J$ joints, we first construct a graph $\mathcal{G}_t(P_t, E)$ where $E$ is the $J \times J$ weighted adjacency matrix:

$$e_{ij} = \begin{cases} 0, & \text{if } i = j \\ \alpha, & \text{if joint i and joint j are connected} \\ \beta, & \text{if joint i and joint j are disconnected} \end{cases}$$

Each graph $\mathcal{G}_t$ at time $t$ is processed by a GCN to compute feature $f_t^+$:

$$f_t^+ = D^{-\frac{1}{2}}(E + I)D^{-\frac{1}{2}}\mathcal{G}_t W_t, \tag{1}$$

where $W_t$ is the weight matrix and $D$ is the diagonal degree matrix with $D_{ii} = \Sigma_j(E_{ij} + I_{ij})$ its diagonal elements. For all $t = 1, 2, ..., t_p$, the GCN output features $f_t^+$ are aggregated along time, resulting in a 3D tensor $[f_1^+, f_2^+, ..., f_{t_p}^+]$. Finally, the 3D pose tensor is combined with the original pose input by a residual connection followed by a set of convolutional operations. Now, the GCN pose backbone provides salient features $h^*$ because of its use of the graphical structure of the 3D joints.

**Spatio-temporal Coupler -** The attention network in VPN learns the spatio-temporal attention weights from the output of Pose Backbone in two steps as shown in fig. 4(I)(B). In the first step, the spatial and temporal attention weights $(A_S$ and $A_T)$ are classically trained as in [44] to get the most important body part and key frames for an action. The output feature $h^*$ of Pose Backbone follows two separate non-linear mapping functions to compute the spatial and temporal attention weights. These spatial $A_S$ and temporal $A_T$ weights are defined as

$$A_S = \sigma(z_1); \qquad A_T = softmax(z_2) \tag{2}$$

where $z_r = W_{z_r} tanh(W_{h_r} h^* + b_{h_r}) + b_{z_r}$ (for $r = 1, 2$) with subscripted $W$ and $b$, the corresponding weights and biases are the latent spatial and temporal attention vectors. The dissociated attention weights $A_S$ and $A_T$ having dimension $m \times n$ and $t_c$ respectively, can undergo a linear mapping to obtain spatially and temporally modulated feature maps. The resultant model is equivalent to the separable STA model [8]. In contrast, we propose to further perform a coupling of the spatial and temporal attention weights. Thus in the second step, joint spatio-temporal attention weights are computed by performing a Hadamard product on the spatial and temporal attention weights. In order to perform this matrix

multiplication, the spatial and temporal attention weights are inflated by duplicating the same attention weights in temporal and spatial dimension respectively. Hence, the $m \times n \times t_c$ dimensional spatio-temporal attention weights $A_{ST}$ are obtained by $A_{ST} = inflate(A_S) \circ inflate(A_T)$. This two-step attention learning process enables the attention network to compute spatio-temporal attention weights in which the spatial saliency varies with time. The obtained attention weights are crucial to disambiguate actions with similar appearance as they may have dissimilar motion over time.

Finally, the spatio-temporal attention weights $A_{ST}$ are linearly multiplied with the input video feature map $f$, followed by a residual connection with the original feature map $f$ to output the modulated feature map $f'$. The residual connection enables the network to retain the properties of the original visual features.

**(II) Spatial Embedding of RGB and Pose -** The objective of the embedding model is to provide tight correspondences between both pose and RGB modalities used in VPN. The state-of-the-art methods [7,8] attempt to provide the attention weights on the RGB feature map using 3D pose information without projecting them into the same 3D referential. The mapping with the pose is only done by cropping the person within the input RGB images. The spatial attention computed through the 3D joint coordinates does not correspond to the part of the image (no pixel to pixel correspondence), although it is crucial for recognizing fine-grained actions. To correlate both modalities, an embedding technique inspired from image captioning task [32,33] is used to build an accurate RGB-Pose embedding in order to enable the poses to represent the visual content of the actions (see fig. 4(II)).

We assume that a low dimensional embedding exists for the global spatial representation of video feature map $f_s = \Sigma_{i=1}^{t_c} f(i,:,:,:)$ (a $D_v$ dimensional vector) and its corresponding pose based latent spatial attention vector $z_1$ (a $D_p$ dimensional vector). The mapping function can be derived from this embedding by

$$f_e = T_v f_s \quad and \quad P_e = T_p z_1, \tag{3}$$

where $T_v \in R^{D_e \times D_v}$ and $T_p \in R^{D_e \times D_p}$ are the transformation matrices that project the video content and the 3D poses into the common $D_e$ dimensional embedding space. This mapping function is applied on the global spatial representation of the visual feature map and the pose based features in order to attain the aforementioned objective of the spatial embedding.

To measure the correspondence between the video content and the 3D poses, we compute the distance between their mappings in the embedding space. Thus, we define an embedding loss as a hypersphere feature metric space

$$L_e = ||\widehat{T_v f_s} - \widehat{T_p z_1}||_2^2 \qquad s.t. \quad ||T_v||_2 = ||T_p||_2 = 1 \tag{4}$$

$\widehat{T_v f_s} = \frac{T_v f_s}{||T_v f_s||_2}$ and $\widehat{T_p z_1} = \frac{T_p z_1}{||T_p z_1||_2}$ are the feature representations projected to the unit hypersphere. The norm constraint $||T_v||_2 = 1$ & $||T_p||_2 = 1$ simply

prevents the trivial solution $\hat{T}_v = \hat{T}_p = 0$. This embedding loss along with the global classification loss provides a linear transformation on the RGB feature map that preserves the low-rank structure for the action representation and introduces a maximally separated features for different actions. Now, the kernels at the visual backbone are updated with a gradient proportional to $(f_e - P_e)$, which in turn transforms the visual feature map to learn pose aware characteristics. Consequently, we strengthen the correspondences between video and poses by minimizing the embedding loss. This embedding ensures that the pose information to be used for computing the spatial attention weights aligns with the content of the video.

Note that the embedding loss also provides feedback to the pose based latent spatial attention vectors $(z_1)$, which in turn transfers knowledge from the 2D image space to pose 3D referential. This allows the attention network to provide better and meaningful spatial attention weights $(A_s)$ compared to the attention network without the embedding. We will quantify this observation in the experiments.

### 3.3  Training jointly the 3D ConvNet and VPN

VPN can be trained as a layer on top of any 3D ConvNet. The 3D ConvNet can be pre-trained for the action classification task for faster convergence. Finally, VPN is plugged into the 3D ConvNet for an end-to-end training with a regularized loss $L$ formulated as

$$L = \lambda_1 L_C + (1 - \lambda_1)L_e + \lambda_2 L_a \tag{5}$$

Here, $L_C$ is the cross-entropy loss, $L_e$ is the embedding loss; the trade-off between these two losses is captured by linear fusion with a positive parameter $\lambda_1$; $L_a$ is the attention regularizer with $\lambda_2$ weighting factor. The attention regularizer consists of the spatial and temporal attention weight regularizer and is formulated as

$$L_a = \sum_{j=1}^{m \times n} \left\| A_s(j) \right\|_2 + \sum_{j=1}^{t_c} (1 - A_{t_c}(j))^2 \tag{6}$$

This additional regularization term $L_a$ ensures that the attention weights are not biased to provide extremely high values to the parts of the spatio-temporal feature map with more relevance and completely neglecting the other parts.

## 4  Experiments

We evaluate the effectiveness of our model for action classification. We consider four public datasets which are the popular datasets for ADL: NTU-60 [39], NTU-120 [25], Toyota-Smarthome [8] and Northwestern-UCLA [51].

**NTU RGB+D** (NTU-60 & NTU-120): NTU-60 is acquired with a Kinect v2 camera and consists of 56880 video samples with 60 activity classes. The activities were performed by 40 subjects and recorded from 80 viewpoints. For each frame, the dataset provides RGB, depth and a 25-joint skeleton of each subject in the frame. For evaluation, we follow the two protocols proposed in [39]: cross-subject (CS) and cross-view (CV). NTU-120 is a super-set of NTU-60 adding a lot of new similar actions. NTU-120 dataset contains 114k video clips of 106 distinct subjects performing 120 actions in a laboratory environment with 155 camera views. For evaluation, we follow a cross-subject ($CS_1$) protocol and a cross-setting ($CS_2$) protocol proposed in [25].

**Toyota-Smarthome** (Smarthome) is a recent ADL dataset recorded in an apartment where 18 older subjects carry out tasks of daily living during a day. The dataset contains 16.1k video clips, 7 different camera views and 31 complex activities performed in a natural way without strong prior instructions. This dataset provides RGB data and 3D skeletons which are extracted from LCRNet [37]. For evaluation on this dataset, we follow cross-subject (CS) and cross-view ($CV_1$ and $CV_2$) protocols proposed in [8].

**Northwestern-UCLA Multiview activity 3D Dataset** (N-UCLA) is acquired simultaneously by three Kinect v1 cameras. The dataset consists of 1194 video samples with 10 activity classes. The activities were performed by 10 subjects, and recorded from three viewpoints. We performed experiments on N-UCLA using the cross-view (CV) protocol proposed in [51]: we trained our model on samples from two camera views and tested on the samples from the remaining view. For instance, the notation $V_{1,2}^3$ indicates that we trained on samples from view 1 and 2, and tested on samples from view 3.

The presence of ADL challenges like fine-grained and similar appearance activities is in higher magnitude in NTU-120 and Smarthome datasets. So, we perform all our ablation studies on these two datasets. We abbreviate Smarthome as SH in table 1, 2, 3 and 4.

### 4.1   Implementation details

**Training.** In our experiments, the selected **visual backbone** is I3D [5] network pre-trained on ImageNet [9] and Kinetics-400 [17]. The visual backbone takes 64 video frames as input. The input of the **VPN** consists of the feature map extracted from `Mixed_5c` layer of I3D and the corresponding 3D poses.

The **pose backbone** takes as input a sequence of $t_p$ 3D poses uniformly sampled from each clip. Hyper-parameter $t_p = 20, 20, 30$ and 5 for NTU-60, NTU-120, Smarthome and N-UCLA respectively. For the pose backbone, we use $t_p$ number of GCNs, each processing a pose from the sequence. The weighting parameters $\alpha$ and $\beta$ for computing the adjacency matrix of the pose based graph are set to 5 and 2 respectively. GCN projects the input joint coordinates to a $64 - dimensional$ space. The output of the GCN is passed to a set of convolutional operations (see fig. 4(I)(A)) which consists of three 2D convolutional layers each are followed by a Batch Normalization layer and a ReLU layer. The output channels of the convolutional layers are 64, 64 and 128.

For classification, a global-average pooling layer followed by a dropout [46] of 0.3 and a *softmax* layer are added at the end of the recognition model for class prediction. Our recognition model is trained with a 4-GPU machine where each GPU has 4 video clips in a mini-batch. Our model is trained for 30 epochs in total, with SGD optimizer having initial learning rate of 0.01 and decay rate of 0.1 after every 10 epochs. The trade off ($\lambda_1$) and regularizer ($\lambda_2$) parameters are set to 0.8 and 0.00001 respectively for all the experiments.

**Inference.** For the recognition model, we perform fully convolutional inference in space as in [54]. The final classification is obtained by max-pooling the softmax scores.

Table 1: Ablation study to show the effectiveness of each VPN component.

| VPN components | NTU-120 $CS_1$ | NTU-120 $CS_2$ | SH CS | SH $CV_2$ |
|---|---|---|---|---|
| $l_1$: visual backbone | 77.0 | 80.1 | 53.4 | 45.1 |
| $l_2$: $l_1$ + attention network | 85.4 | 86.9 | 56.4 | 50.5 |
| $l_3$: $l_2$ + spatial embedding | **86.3** | **87.8** | **60.8** | **53.5** |

Table 2: Performance of VPN with different choices of Attention Network.

| Model | Pose Backbone | Coupler | NTU-120 $CS_1$ | NTU-120 $CS_2$ | SH CS | SH $CV_2$ |
|---|---|---|---|---|---|---|
| $l_4$: VPN | LSTM | × | 84.7 | 83.6 | 57.1 | 50.6 |
| $l_5$: VPN | GCN | × | 85.6 | 86.8 | 60.1 | 53.1 |
| $l_6$: VPN | LSTM | ✓ | 85.3 | 84.1 | 57.6 | 51.5 |
| $l_7$: VPN | GCN | ✓ | **86.3** | **87.8** | **60.8** | **53.5** |

Table 3: Performance of VPN with different embedding losses $l_e$.

| Loss | NTU-120 $CS_1$ | NTU-120 $CS_2$ | SH CS | SH $CV_2$ |
|---|---|---|---|---|
| KL-divergence $D_{KL}(f_e||P_e)$ | 85.5 | 87.1 | 57.2 | 50.9 |
| KL-divergence $D_{KL}(P_e||f_e)$ | 85.6 | 86.9 | 57.0 | 51.1 |
| Bi-directional KL-divergence | 86.1 | 87.2 | 57.2 | 51.7 |
| Normalized Euclidean loss | **86.3** | **87.8** | **60.8** | **53.5** |

Table 4: Impact of Spatial Embedding on Spatial Attention.

| Model | Pose Backbone | Spatial Embedding | NTU-120 $CS_1$ | NTU-120 $CS_2$ | SH CS | SH $CV_2$ |
|---|---|---|---|---|---|---|
| VPN | LSTM | × | 81.7 | 81.2 | 45.5 | 50.0 |
| VPN | LSTM | ✓ | 82.7 | 82.0 | 56.5 | 52.6 |
| VPN | GCN | × | 82.6 | 84.3 | 49.1 | 51.7 |
| VPN | GCN | ✓ | **83.1** | **85.3** | **58.4** | **53.1** |

### 4.2 Ablation Study

Our model includes two novel components, the spatial embedding and the attention network. Both of them are critical for good performance on ADL recognition. We show the importance of the attention network and the spatial embedding of VPN in table 1. We also show the effectiveness of the spatial embedding with different instantiation of the attention network in table 2.

**How effective is VPN?** In order to answer this point, we show the action classification accuracy with baseline I3D ($l_1$) which is the visual backbone and then incorporate the VPN components: the attention network ($l_2$) and the spatial embedding ($l_3$) one-by-one in table 1. The attention network ($l_2$) improves significantly the classification of the actions (upto 8.4% on NTU-120 and 5.4% on Smarthome) by providing spatio-temporal saliency to the I3D feature maps. With the spatial embedding ($l_3$), the action classification further improves (upto 0.9% on NTU-120 and 4.4% on Smarthome).

**Diagnosis of the attention network -** In table 2, we further illustrate the importance of each component in the attention network, i.e. the Pose Backbone and the spatio-temporal coupler. We have designed a baseline attention network with LSTM as pose backbone following [8]. We compare the LSTM pose backbone in $l_4$ and $l_6$ with our proposed GCN instantiation in $l_5$ and $l_7$. The attention network without a spatio-temporal coupler provides dissociated spatial and temporal attention weights in $l_4$ and $l_5$ in contrast to our proposed coupler in $l_6$ and $l_7$. Firstly, we observe that the GCN pose backbone makes use of the human joint topology, thus improves the classification accuracy in all scenarios with or without the coupler. Consequently, actions like *Snapping Finger* (+24.5%) and *Apply cream on face* (+23.9%) improves significantly with GCN instantiation ($l_6$) compared to LSTM ($l_7$). Secondly, we observe that the spatio-temporal coupler provides fine spatial attention weights for the most important frames in a video, which enables the model to disambiguate actions with similar appearance but dissimilar motion. Consequently, the coupler ($l_7$) improves the classification accuracy up to 1% on NTU-120 and 0.7% on Smarthome w.r.t. dissociating the attention weights ($l_5$). For instance, with dissociation of the attention weights, *rubbing two hands* was confused with *clapping* and *flicking hair* was confused with *putting on headphone*. With VPN, the coupler improves the classification accuracy of actions *rubbing two hands* and *flicking hair* by 25% and 19.6% respectively.

**Which loss is better for learning the spatial embedding?** In this ablation study (Table 3), we compare different losses for projecting the 3D poses and RGB cues in a common semantic space. First, we compare the KL-divergence losses [19,15] ($D_{KL}(f_e||P_e)$ and $D_{KL}(P_e||f_e)$) from $P_e$ to $f_e$ and vice-versa. Then, we compare a bi-directional KL-divergence loss [58,52,29] ($D_{KL}(f_e||P_e)$ + $D_{KL}(P_e||f_e)$) to our normalized euclidean loss. We observe that (i) the loss using $D_{KL}(f_e||P_e)$ and $D_{KL}(P_e||f_e)$ deteriorates the action classification accuracy as the feedback is in one direction either towards RGB or poses, implying two-way feedback for the visual features and the attention network is necessary, (ii) our normalized euclidean loss outperforms the bi-directional KL divergence loss, exhibit its superiority.

**Impact of Embedding on Spatial attention -** In table 4, we show the impact of spatial embedding on the attention network providing spatial attention only. We perform the experiments with different choice of Pose Backbone, i.e. LSTM as discussed above and our proposed GCN. The spatial embedding provides a tight correspondence between the RGB data and poses. As a result, it boosts the classification accuracy in all the experiments. It is worth noting that the improvement is significant for Smarthome as it contains many fine-grained actions with videos captured by fixed cameras in an unconstrained Field of View. Thus, enforcing the embedding loss enhances the spatial precision during inference. As a result, the classification accuracy of fine-grained actions like *pouring water* (+77.7%), *pouring grains* (+76.1%) for making coffee, *cutting bread* (+50%), *pouring from kettle* (+42.8%) and *inserting teabag* (+35%) improves VPN with GCN pose backbone compared to its counterpart without embedding.
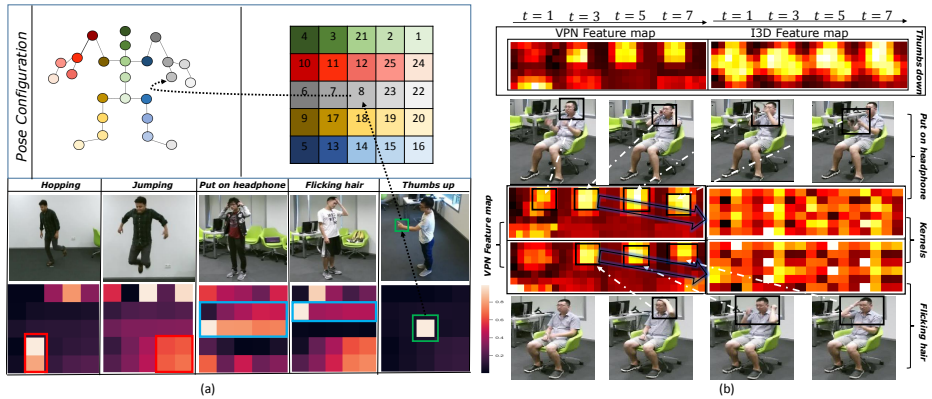
Fig. 5: (a) The heatmaps of the activations of the 3D joint coordinates (output of GCN) in the attention network of VPN. The area in the colored bounding boxes shows that different joints are activated for similar actions. (b) Heatmaps of visual feature maps & corresponding activated kernels for different time stamps. These heatmaps show that VPN has better discriminative power than I3D.

### 4.3   Qualitative Analysis

Fig. 5(a) visualizes the activation of the human joints at the output of pose backbone (with GCNs) in VPN. The figure depicts the activations of the 3D joints. They are presented in a sequence of the human body topological order (follow first row of fig. 5(a)) for convenient visualization. VPN is able to disambiguate actions with similar appearance like *hopping* and *jumping* due to high order activation at relevant joints of the human legs. The discriminative leg joints with high activation have been marked with a red bounding box in fig. 5(a) (third row). Similarly, for actions like *put on headphone* with two hands and *flicking hair* with one hand, the blue bounding boxes demonstrate high activation of both the hand joints for the former action as compared to high activation of a single hand joints for the latter. For a very fine-grained action like *thumbs up*, the thumb joint is highly activated as compared to the other joints. This shows that the GCN pose backbone in VPN is a crucial ingredient for better action recognition.

In fig. 5(b), we compare the heatmap of the VPN and I3D feature maps for different time stamps. We observe the sharpness in the VPN feature maps compared to that of I3D for *thumbs down* action which is localized over a small space. For similar actions like *put on headphone* and *flicking hair*, along with salience precision of the VPN feature map, the activations of their corresponding receptive fields show the discriminative power of VPN.

In fig. 6(a), we illustrate the performance of VPN w.r.t. I3D baseline for the dynamicity of an action along the videos. This dynamicity is computed by averaging the Procrustes distance [20] between subsequent 3D poses along the videos. If the average distance is large, it means the poses change a lot in an action. VPN significantly improves for actions with subtle motion like *hush* (+52.7%),

*staple book* $(+40.7\%)$ and *reading* $(+36.2\%)$ which indicates the efficacy of VPN for fine-grained actions. The degradation of the VPN performance for high action dynamicity is negligible(-0.8%). In fig. 6(b), we show the t-SNE plots of the feature spaces produced by I3D and VPN for some selected actions with similar appearance. It clearly shows the discriminative power of VPN for actions with similar appearance which is a frequent challenge in ADL.
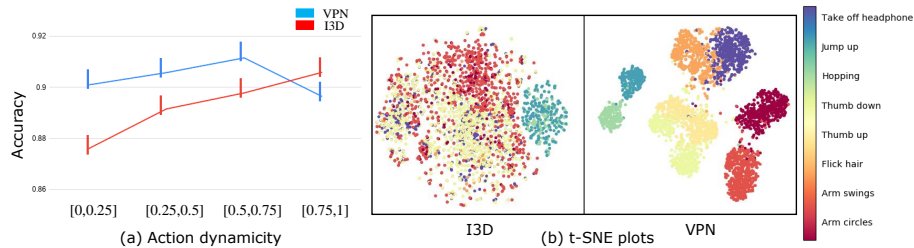


Fig. 6: (a) We compare our model against baseline I3D across action dynamicity. Our model significantly improves for most actions. (b) t-SNE plots of feature spaces produced by I3D and VPN for similar appearance actions.

Table 5: Results (accuracies in %) on NTU-60 with cross-subject (CS) and cross-view (CV) settings (at left) and NTU-120 with cross-subject $(CS_1)$ and cross-setup $(CS_2)$ settings (at right); Att indicates attention mechanism, $\circ$ indicates that the modality has only been used for training, the methods indicated with $^*$ are reproduced on this dataset. 3D ResNeXt-101 is abbreviated as RNX3D101.

| Methods | Pose | RGB | Att | CS | CV |
|---|---|---|---|---|---|
| AGC-LSTM [42] | ✓ | × | ✓ | 89.2 | 95.0 |
| DGNN [40] | ✓ | × | × | 89.9 | 96.1 |
| STA-Hands [2] | ✓ | ✓ | ✓ | 82.5 | 88.6 |
| altered STA-Hands [3] | ✓ | ✓ | ✓ | 84.8 | 90.6 |
| Glimpse Cloud [4] | ∘ | ✓ | ✓ | 86.6 | 93.2 |
| PEM [28] | ✓ | ✓ | × | 91.7 | 95.2 |
| Separable STA [8] | ✓ | ✓ | ✓ | 92.2 | 94.6 |
| P-I3D [7] | ✓ | ✓ | ✓ | 93 | 95.4 |
| **VPN** | ✓ | ✓ | ✓ | **93.5** | **96.2** |
| **VPN (RNX3D101)** | ✓ | ✓ | ✓ | **95.5** | **98.0** |

| Methods | Pose | RGB | Att | $CS_1$ | $CS_2$ |
|---|---|---|---|---|---|
| ST-LSTM [26] | ✓ | × | ✓ | 55.7 | 57.9 |
| Two stream Att LSTM [24] | ✓ | × | ✓ | 61.2 | 63.3 |
| Multi-Task CNN [18] | ✓ | × | × | 62.2 | 61.8 |
| PEM [28] | ✓ | × | ✓ | 64.6 | 66.9 |
| 2s-AGCN [41] | ✓ | × | ✓ | 82.9 | 84.9 |
| Two-streams [43] | × | ✓ | × | 58.5 | 54.8 |
| I3D$^*$ [5] | × | ✓ | × | 77.0 | 80.1 |
| Two-streams + ST-LSTM [25] | ✓ | ✓ | × | 61.2 | 63.1 |
| Separable STA$^*$ [8] | ✓ | ✓ | ✓ | 83.8 | 82.5 |
| **VPN** | ✓ | ✓ | ✓ | **86.3** | **87.8** |

### 4.4    Comparison with the state-of-the-art

We compare VPN to the state-of-the-art (SoA) on NTU-60, NTU-120, Smarthome and N-UCLA in table 5, 6 and 7. VPN outperforms on each of them. In table 5 (at left), for input modality RGB+Poses, VPN improves the SoA [7] by up to 0.8% on NTU-60 even by using one-third parameters compared to [7]. The SoA using Poses only [40] yields classification accuracy near to VPN for cross-view protocol (with 0.1% difference) due to their robustness to view changes. However,

Table 6: Results on Smarthome dataset with cross-subject (CS) and cross-view ($CV_1$ and $CV_2$) settings (accuracies in %). Att indicates attention mechanism.

| Methods | Pose | RGB | Att | CS | $CV_1$ | $CV_2$ |
|---|---|---|---|---|---|---|
| DT [49] | × | ✓ | × | 41.9 | 20.9 | 23.7 |
| LSTM [31] | ✓ | × | × | 42.5 | 13.4 | 17.2 |
| I3D [5] | × | ✓ | × | 53.4 | 34.9 | 45.1 |
| I3D+NL [54] | × | ✓ | ✓ | 53.6 | 34.3 | 43.9 |
| P-I3D [7] | ✓ | ✓ | ✓ | 54.2 | 35.1 | 50.3 |
| Separable STA [8] | ✓ | ✓ | ✓ | 54.2 | 35.2 | 50.3 |
| **VPN** | ✓ | ✓ | ✓ | **60.8** | **43.8** | **53.5** |

Table 7: Results on N-UCLA dataset with cross-view $V^3_{1,2}$ settings (accuracies in %); $\overline{Pose}$ indicate its usage only in the training phase.

| Methods | Data | Att | $V^3_{1,2}$ |
|---|---|---|---|
| HPM+TM [36] | Depth | × | 91.9 |
| Ensemble TS-LSTM [22] | Pose | × | 89.2 |
| NKTM [35] | RGB | × | 85.6 |
| Glimpse Cloud [4] | RGB+ $\overline{Pose}$ | ✓ | 90.1 |
| Separable STA [8] | RGB+Pose | ✓ | 92.4 |
| P-I3D [7] | RGB+Pose | ✓ | 93.1 |
| **VPN** | RGB+Pose | ✓ | **93.5** |

the lack of appearance information restricts these methods [40,42] to disambiguate actions with similar visual appearance, thus resulting in lower accuracy for cross-subject protocol. We have also tested VPN with 3D ResNeXt-101 [14] on NTU-60 dataset. The results in table 5 show that VPN can be adapted with other existing video backbones.

Compared to the SoA results, the improvement by 3.9% and 4.9% (averaging over the protocols) on NTU-120 and Smarthome respectively are significant. It is worth noting that VPN improves further the classification of actions with similar appearance as compared to Separable STA [8]. For example, actions like *clapping* (+44.3%) and *flicking hair* (+19.1%) are now discriminated with better accuracy. In addition, the superior performance of VPN in cross-view protocol for both NTU-120 and Smarthome implies that it provides better view-adaptive characterization compared to all the prior methods.

For N-UCLA which is a small-scale dataset, we pre-train the visual backbone with NTU-60 for a fair comparison with [4,8,7]. We also outperform the SoA [7] by 0.4% on this dataset.

## 5    Conclusion

This paper addresses the challenges of ADL classification. We have proposed a novel Video-Pose Network **VPN** which provides an accurate video-pose embedding. We show that the embedding along with attention network yields a more discriminative feature map for action classification. The attention network leverages the topology of the human joints and with the coupler provides precise spatio-temporal attention weights along the video.

Our recognition model outperforms the state of-the-art results for action classification on 4 public datasets. This is a first step towards combining RGB and Pose through an explicit embedding. A future perspective of this work is to exploit this embedding even in case of noisy 3D poses in order to also boost action recognition for internet videos. This embedding could even help to refine these noisy 3D poses in a weakly supervised manner.

## Acknowledgement

We are grateful to INRIA Sophia Antipolis - Mediterranean "NEF" computation cluster for providing resources and support.

## References

1. Arandjelović, R., Gronat, P., Torii, A., Pajdla, T., Sivic, J.: NetVLAD: CNN architecture for weakly supervised place recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
2. Baradel, F., Wolf, C., Mille, J.: Human action recognition: Pose-based attention draws focus to hands. In: 2017 IEEE International Conference on Computer Vision Workshops (ICCVW). pp. 604–613 (Oct 2017). https://doi.org/10.1109/ICCVW.2017.77
3. Baradel, F., Wolf, C., Mille, J.: Human activity recognition with pose-driven attention to rgb. In: The British Machine Vision Conference (BMVC) (September 2018)
4. Baradel, F., Wolf, C., Mille, J., Taylor, G.W.: Glimpse clouds: Human activity recognition from unstructured feature points. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
5. Carreira, J., Zisserman, A.: Quo vadis, action recognition? a new model and the kinetics dataset. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4724–4733. IEEE (2017)
6. Crasto, N., Weinzaepfel, P., Alahari, K., Schmid, C.: MARS: Motion-Augmented RGB Stream for Action Recognition. In: CVPR (2019)
7. Das, S., Chaudhary, A., Bremond, F., Thonnat, M.: Where to focus on for human action recognition? In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 71–80 (Jan 2019). https://doi.org/10.1109/WACV.2019.00015
8. Das, S., Dai, R., Koperski, M., Minciullo, L., Garattoni, L., Bremond, F., Francesca, G.: Toyota smarthome: Real-world activities of daily living. In: ICCV (2019)
9. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR09 (2009)
10. Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2015)
11. Du, Y., Fu, Y., Wang, L.: Skeleton based action recognition with convolutional neural network. In: 2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR). pp. 579–583 (Nov 2015). https://doi.org/10.1109/ACPR.2015.7486569
12. Feichtenhofer, C., Fan, H., Malik, J., He, K.: Slowfast networks for video recognition. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
13. Girdhar, R., Carreira, J., Doersch, C., Zisserman, A.: Video action transformer network. CoRR **abs/1812.02707** (2018), http://arxiv.org/abs/1812.02707
14. Hara, K., Kataoka, H., Satoh, Y.: Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)

15. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network (2015)
16. Hussein, N., Gavves, E., Smeulders, A.W.M.: Timeception for complex action recognition. CoRR **abs/1812.01289** (2018), `http://arxiv.org/abs/1812.01289`
17. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
18. Ke, Q., Bennamoun, M., An, S., Sohel, F., Boussaid, F.: Learning clip representations for skeleton-based 3d action recognition. IEEE Transactions on Image Processing **27**(6), 2842–2855 (June 2018). https://doi.org/10.1109/TIP.2018.2812099
19. Kim, S., Seltzer, M., Li, J., Zhao, R.: Improved training for online end-to-end speech recognition systems. In: Proc. Interspeech 2018. pp. 2913–2917 (2018). https://doi.org/10.21437/Interspeech.2018-2517, `http://dx.doi.org/10.21437/Interspeech.2018-2517`
20. Krzanowski, W.J.: Principles of Multivariate Analysis: A Users Perspective. Oxford University Press, Inc., USA (1988)
21. Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: Hmdb: a large video database for human motion recognition. In: 2011 International Conference on Computer Vision. pp. 2556–2563. IEEE (2011)
22. Lee, I., Kim, D., Kang, S., Lee, S.: Ensemble deep learning for skeleton-based action recognition using temporal sliding lstm networks. In: Proceedings of the IEEE International Conference on Computer Vision (2017)
23. Liu, G., Qian, J., Wen, F., Zhu, X., Ying, R., Liu, P.: Action recognition based on 3d skeleton and rgb frame fusion. In: 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 258–264 (Nov 2019). https://doi.org/10.1109/IROS40897.2019.8967570
24. Liu, J., Wang, G., Hu, P., Duan, L., Kot, A.C.: Global context-aware attention lstm networks for 3d action recognition. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3671–3680 (July 2017). https://doi.org/10.1109/CVPR.2017.391
25. Liu, J., Shahroudy, A., Perez, M., Wang, G., Duan, L.Y., Kot, A.C.: Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding. IEEE Transactions on Pattern Analysis and Machine Intelligence (2019). https://doi.org/10.1109/TPAMI.2019.2916873
26. Liu, J., Shahroudy, A., Xu, D., Wang, G.: Spatio-temporal lstm with trust gates for 3d human action recognition. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) Computer Vision – ECCV 2016. pp. 816–833. Springer International Publishing, Cham (2016)
27. Liu, M., Liu, H., Chen, C.: Enhanced skeleton visualization for view invariant human action recognition. Pattern Recognition **68**, 346 – 362 (2017). https://doi.org/https://doi.org/10.1016/j.patcog.2017.02.030, `http://www.sciencedirect.com/science/article/pii/S0031320317300936`
28. Liu, M., Yuan, J.: Recognizing human actions as the evolution of pose estimation maps. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
29. Liu, Y., Guo, Y., Bakker, E.M., Lew, M.S.: Learning a recurrent residual fusion network for multimodal matching. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 4127–4136 (Oct 2017). https://doi.org/10.1109/ICCV.2017.442

30. Luo, Z., Hsieh, J.T., Jiang, L., Carlos Niebles, J., Fei-Fei, L.: Graph distillation for action detection with privileged modalities. In: The European Conference on Computer Vision (ECCV) (September 2018)

31. Mahasseni, B., Todorovic, S.: Regularizing long short term memory with 3d human-skeleton sequences for action recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3054–3062 (2016)

32. Miech, A., Laptev, I., Sivic, J.: Learning a text-video embedding from incomplete and heterogeneous data. CoRR **abs/1804.02516** (2018), `http://arxiv.org/abs/1804.02516`

33. Pan, Y., Mei, T., Yao, T., Li, H., Rui, Y.: Jointly modeling embedding and translation to bridge video and language. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)

34. Perronnin, F., Sánchez, J., Mensink, T.: Improving the fisher kernel for large-scale image classification. In: European conference on computer vision. pp. 143–156. Springer (2010)

35. Rahmani, H., Mian, A.: Learning a non-linear knowledge transfer model for cross-view action recognition. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2458–2466 (June 2015). https://doi.org/10.1109/CVPR.2015.7298860

36. Rahmani, H., Mian., A.: 3d action recognition from novel viewpoints. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1506–1515 (June 2016). https://doi.org/10.1109/CVPR.2016.167

37. Rogez, G., Weinzaepfel, P., Schmid, C.: LCR-Net++: Multi-person 2D and 3D Pose Detection in Natural Images. IEEE Transactions on Pattern Analysis and Machine Intelligence (2019)

38. Shahroudy, A., Wang, G., Ng, T.: Multi-modal feature fusion for action recognition in rgb-d sequences. In: 2014 6th International Symposium on Communications, Control and Signal Processing (ISCCSP). pp. 1–4 (May 2014). https://doi.org/10.1109/ISCCSP.2014.6877819

39. Shahroudy, A., Liu, J., Ng, T.T., Wang, G.: Ntu rgb+d: A large scale dataset for 3d human activity analysis. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)

40. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Skeleton-based action recognition with directed graph neural networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)

41. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In: CVPR (2019)

42. Si, C., Chen, W., Wang, W., Wang, L., Tan, T.: An attention enhanced graph convolutional lstm network for skeleton-based action recognition. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1227–1236 (June 2019). https://doi.org/10.1109/CVPR.2019.00132

43. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: Advances in neural information processing systems. pp. 568–576 (2014)

44. Song, S., Lan, C., Xing, J., Zeng, W., Liu, J.: An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In: AAAI Conference on Artificial Intelligence. pp. 4263–4270 (2017)

45. Soomro, K., Roshan Zamir, A., Shah, M.: Ucf101: A dataset of 101 human actions classes from videos in the wild (12 2012)

46. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. J. Mach. Learn. Res. **15**(1), 1929–1958 (Jan 2014), `http://dl.acm.org/citation.cfm?id=2627435.2670313`

47. Tang, Y., Tian, Y., Lu, J., Li, P., Zhou, J.: Deep progressive reinforcement learning for skeleton-based action recognition. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5323–5332 (June 2018). https://doi.org/10.1109/CVPR.2018.00558

48. Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M.: Learning spatiotemporal features with 3d convolutional networks. In: Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV). pp. 4489–4497. ICCV '15, IEEE Computer Society, Washington, DC, USA (2015). https://doi.org/10.1109/ICCV.2015.510, `http://dx.doi.org/10.1109/ICCV.2015.510`

49. Wang, H., Kläser, A., Schmid, C., Liu, C.L.: Action Recognition by Dense Trajectories. In: IEEE Conference on Computer Vision & Pattern Recognition. pp. 3169–3176. Colorado Springs, United States (Jun 2011), `http://hal.inria.fr/inria-00583818/en`

50. Wang, H., Schmid, C.: Action recognition with improved trajectories. In: IEEE International Conference on Computer Vision. Sydney, Australia (2013), `http://hal.inria.fr/hal-00873267`

51. Wang, J., Nie, X., Xia, Y., Wu, Y., Zhu, S.C.: Cross-view action modeling, learning, and recognition. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition. pp. 2649–2656 (June 2014). https://doi.org/10.1109/CVPR.2014.339

52. Wang, L., Li, Y., Lazebnik, S.: Learning deep structure-preserving image-text embeddings. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5005–5013 (June 2016). https://doi.org/10.1109/CVPR.2016.541

53. Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Val Gool, L.: Temporal segment networks: Towards good practices for deep action recognition. In: ECCV (2016)

54. Wang, X., Girshick, R.B., Gupta, A., He, K.: Non-local neural networks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition pp. 7794–7803 (2018)

55. Yan, S., Xiong, Y., Lin, D.: Spatial temporal graph convolutional networks for skeleton-based action recognition. In: AAAI (2018)

56. Zhang, P., Lan, C., Xing, J., Zeng, W., Xue, J., Zheng, N.: View adaptive recurrent neural networks for high performance human action recognition from skeleton data. In: The IEEE International Conference on Computer Vision (ICCV) (Oct 2017)

57. Zhang, S., Liu, X., Xiao, J.: On geometric features for skeleton-based action recognition using multilayer lstm networks. In: 2017 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 148–157 (March 2017). https://doi.org/10.1109/WACV.2017.24

58. Zhang, Y., Lu, H.: Deep cross-modal projection learning for image-text matching. In: ECCV (2018)

59. Zhao, J., Snoek, C.G.M.: Dance with flow: Two-in-one stream action detection. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)

**Appendix overview**

We provide in section 1 computational details regarding the normalization of Euclidean loss provided in Spatial Embedding of RGB and Pose (section 3.2 (II)). Section 2 provides the details of the baseline with LSTM pose backbone with or without coupler in Table 2 & 4 from the ablation studies. Section 3 provides the details of the divergence losses used for comparing with Normalized Euclidean loss in Table 3 from ablation studies. Finally, we provide some more insights about VPN in section 4 to illustrate its effectiveness.

For convenience, we use the same notation as in the main paper for this supplementary material.

## 1    Details on normalization of Euclidean loss

In equation (4), $\widehat{T_v f_s} = \frac{T_v f_s}{||T_v f_s||_2} = \frac{f_e}{||f_e||_2}$ and $\widehat{T_p z_1} = \frac{T_p z_1}{||T_p z_1||_2} = \frac{P_e}{||P_e||_2}$ are the feature representations projected to the unit hypersphere. Here, we compute the norm $||f_e||_2$ and $||P_e||_2$ using

$$||f_e||_2 = \sqrt{\Sigma_i f_{e_i}^2 + \epsilon} \qquad \& \qquad ||P_e||_2 = \sqrt{\Sigma_i P_{e_i}^2 + \epsilon} \qquad (7)$$

where $\epsilon$ is a small positive value to prevent dividing zero.

## 2    LSTM Pose backbone with or without coupler baselines

For the LSTM Pose Backbone in Table 2 & 4, we use a 3-layer stacked LSTM, pre-trained for action classification, as a Pose Backbone by freezing the weights of their cell gates following [8]. The output feature vector $h^*$ is computed by concatenating all the LSTM output features over time. To have a fair comparison with our GCN Pose Backbone, we also introduced residual connections between the original pose input and the LSTM output tensor. However, these residual connections do not improve the action classification accuracy.

For the experiments in Table 2 to implement the attention network without the coupler, we do not compute $A_{ST}$. Instead, we multiply the attention weights $inflate(A_S)$ and $inflate(A_T)$ separately with the RGB feature map $f$ in two streams following [8]. Finally, the modulated feature maps from both the streams are concatenated to classify the actions.

## 3    Baselines with KL divergence loss

In Table 3, we compare different forms of KL divergence loss with normalized euclidean loss for spatial embedding of RGB and 3D poses. The KL-divergence losses $D_{KL}(f_e||P_e)$ and $D_{KL}(P_e||f_e)$ for n samples are computed by

$$D_{KL}(f_e||P_e) = \sum_{i=1}^{n} f_e^i log(\frac{f_e^i}{P_e^i}) \tag{8}$$

$$D_{KL}(P_e||f_e) = \sum_{i=1}^{n} P_e^i log(\frac{P_e^i}{f_e^i}) \tag{9}$$

where $f_e^i$ and $P_e^i$ are visual and pose embedding of the $i^{th}$ input sample.

Finally, the bi-directional KL-divergence loss is given by $D_{KL}(f_e||P_e) + D_{KL}(P_e||f_e)$.

## 4    Detailed qualitative analysis of VPN

In this section, we provide illustrations to show the impact of each VPN components in section 4.1, superiority of VPN compared to other representative baselines in section 4.2, and some result visualization to highlight the solved and remaining challenges in ADL.

### 4.1    Illustration to show the impact of VPN components

In fig. 7, we illustrate a set of graphs showing the top-5 improvement of action classification accuracy using different components of VPN compared to I3D baseline. As discussed in the ablation studies of the primary paper, each component in VPN is critical for good performance on ADL recognition.

- The spatial embedding provides an accurate alignment of the RGB images and the 3D poses. As a result, the recognition performance of the fine-grained actions improves compared to its counterpart without embedding (see fig. 7 (a)).
- The GCN pose backbone of the attention network, not only provides a strategy to globally optimize the recognition model but also takes the human joint configuration into account for computing the attention weights. This further boosts the action classification performance (see fig. 7 (b)).

– The spatio-temporal coupler of the attention network provides discriminative spatio-temporal attention weights which enables the recognition model to better disambiguate the actions with similar appearance (see fig. 7 (c)).
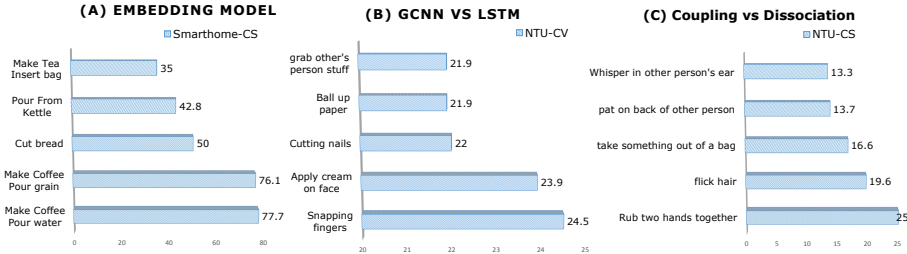


Fig. 7: Graphs illustrating the superiority of each component of VPN compared to their counterparts (without the respective components). We present the Top-5 per class improvement for (a) VPN with embedding vs without embedding (only Spaital Attention), (b) VPN with GCN vs LSTM Pose Backbone, and (c) attention in VPN with vs without spatio-temporal coupler.

## 4.2    Illustration to show the superiority of VPN

We illustrate in fig. 8, the top-5 per-class classification improvement compared to baseline I3D [5] and to an attention mechanism (Separable STA [8]) from the state-of-the-art, utilizing 3D poses. The significant accuracy improvements for actions with subtle motion like *hush* (+52.7%), *staple book* (+40.7%) and *reading* (+36.2%) as depicted in fig. 8 (a) illustrate the efficacy of VPN for fine-grained actions. It is worth noting that VPN improves further the classification of actions possessing similar appearance as compared to separable STA in fig. 8 (b). For example, actions like *clapping* (+44.3%) and *flicking hair* (+19.1%) are now discriminated with better accuracy. Further, in fig. 8 (c) we present a radar for the average mis-classification score of few action-pairs. The smaller area under the curve for VPN compared to I3D baseline and Separable STA shows that it is able to better disambiguate the action-pairs even with low inter-class variation.
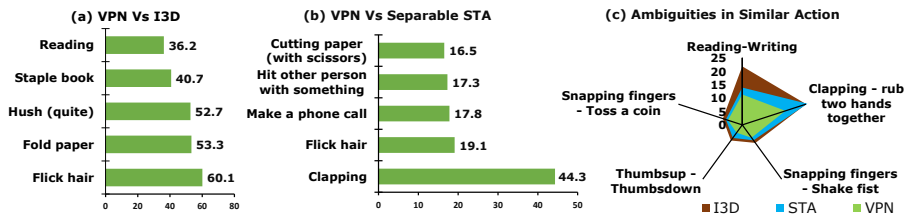
**(a) VPN Vs I3D**

| | |
|---|---|
| Reading | 36.2 |
| Staple book | 40.7 |
| Hush (quite) | 52.7 |
| Fold paper | 53.3 |
| Flick hair | 60.1 |

**(b) VPN Vs Separable STA**

| | |
|---|---|
| Cutting paper (with scissors) | 16.5 |
| Hit other person with something | 17.3 |
| Make a phone call | 17.8 |
| Flick hair | 19.1 |
| Clapping | 44.3 |

**(c) Ambiguities in Similar Action**

Reading-Writing · Clapping – rub two hands together · Snapping fingers – Shake fist · Thumbsup – Thumbsdown · Snapping fingers – Toss a coin

■ I3D    ■ STA    ■ VPN

Fig. 8: Graphs illustrating the superiority of VPN compared to the state-of-the-art methods. We present the Top-5 per class improvement for VPN over (a) I3D baseline and (b) Separable STA. In (c), we present a radar for the average mis-classification score of few action-pairs: lower scores indicate lesser ambiguities between the action-pairs.

### 4.3   Result visualization

In this section, we provide the confusion matrix for action classification on NTU RGB+D 120 and Toyota Smarthome using VPN. In fig 9, we present the confusion matrix of VPN on NTU RGB+D (on right) and a zoom of it around the red bounding box (on left). We also present the corresponding zoom of the confusion matrix of I3D. We are particularly interested in the mis-classifications performed by VPN and thus, we zoom into the region with relatively low classification accuracy. We observe that actions like *staple book* and *taking something out of bag* were confused with *cutting papers* and *put something into a bag* respectively when classified with I3D. However, with VPN these actions with similar motion are now better discriminated, improving their classification accuracy by approximately 42% and 27% respectively.

Similarly, in fig. 10 (a), we present the confusion matrix of VPN on Toyota Smarthome dataset. In fig. 10 (b), we show the poses for some images belonging to action videos mis-classified by I3D. Thanks to the high quality 3D poses for these videos, now VPN can correctly classify these actions taking the human topology of the 3D poses into account. We provide some visual results in fig. 11 where VPN outperforms I3D baseline. We notice that actions like *Drink from glass* are not recognized due to extremely low number of training samples. We further notice that actions like *using tablet* are recognized with low accuracy of 13% and largely confused with *using laptop*. However, I3D completely mis-classifies the action *using tablet*. We also observe that still few action classes are recognized with extremely low classification accuracy. We infer that these poor classification

results on certain videos are due to occlusion, low resolution of the actions and low quality poses as illustrated in fig  12.
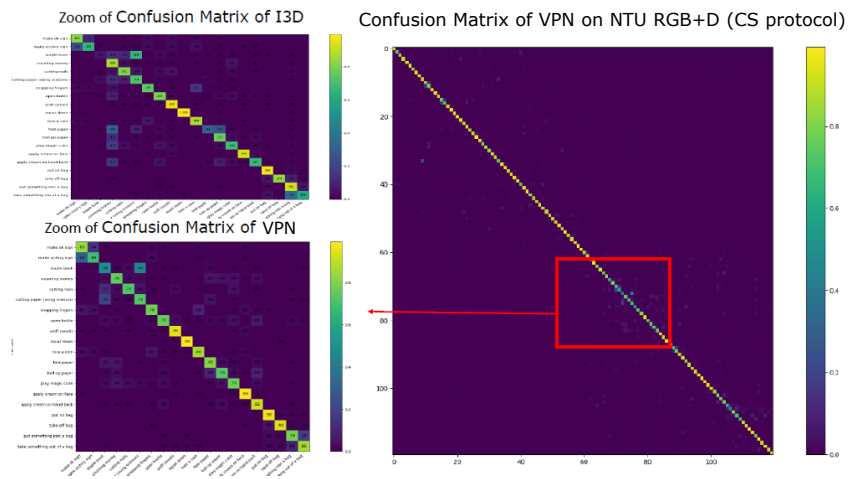


Fig. 9: Confusion matrix of VPN on NTU RGB+D (CS Protocol) on the right. Zoom of the red bounding box on the left along with the corresponding confusion matrix of I3D.
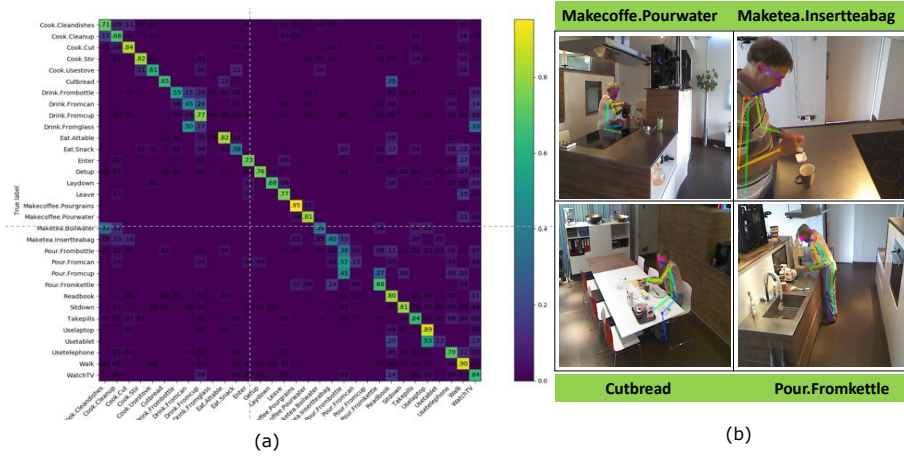
Fig. 10: (a) Confusion matrix of VPN on Toyota Smarthome (CS protocol) (b) Illustration of poses for activities mis-classified with I3D but correctly classified with VPN.
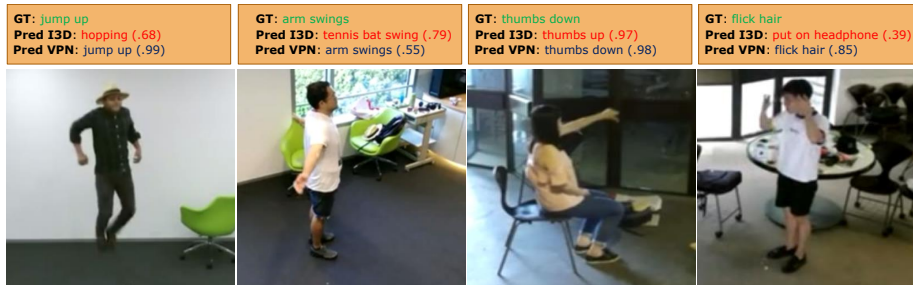


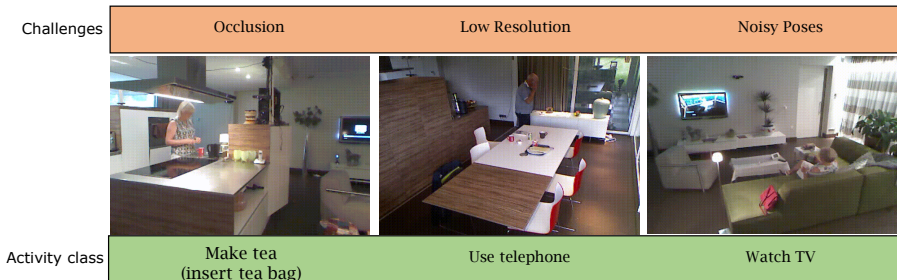Fig. 11: Visual results from NTU RGB+D 120 where VPN outperforms I3D.



Fig. 12: Illustration of the remaining challenges in Toyota Smarthome with images from activities (indicated below) and their corresponding challenges (indicated on the top)