

# Batch-based Activity Recognition from Egocentric Photo-Streams Revisited

Alejandro Cartas    Juan Marín    Petia Radeva    Mariella Dimiccoli

May 11, 2018

## Abstract

Wearable cameras can gather large amounts of image data that provide rich visual information about the daily activities of the wearer. Motivated by the large number of health applications that could be enabled by the automatic recognition of daily activities, such as lifestyle characterization for habit improvement, context-aware personal assistance and tele-rehabilitation services, we propose a system to classify 21 daily activities from photo-streams acquired by a wearable photo-camera. Our approach combines the advantages of a Late Fusion Ensemble strategy relying on convolutional neural networks at image level with the ability of recurrent neural networks to account for the temporal evolution of high level features in photo-streams without relying on event boundaries. The proposed batch-based approach achieved an overall accuracy of 89.85%, outperforming state of the art end-to-end methodologies. These results were achieved on a dataset consists of 44,902 egocentric pictures from three persons captured during 26 days in average.

## 1 Introduction

During the last decade there has been a growing interest in analyzing human activities from sensory data [16]. More recently, the introduction of wearable cameras opened the unique opportunity to capture richer contextual information than using only traditional sensors [17]. Lifelogging cameras, as a particular case of wearable cameras, are the only tools that allow to capture image data during several days thanks to their very low

frame rate (2fpm) without recharging its battery. Consequently, they are best suited to monitor the Activities of Daily Living (ADLs) of the wearer that include, but are not limited to, the activities that an independent person performs on daily basis for living at home or in a community [15]. The monitoring of these activities has several applications in health-related research including active aging monitoring, tele-rehabilitation, frailty prevention and stroke-survivor monitoring among others [15, 21].

However, activity recognition from first-person (egocentric) photo-streams has received relatively little attention in the literature [5, 4, 18, 3]. One of its major challenges is that photo-streams are characterized by a very low frame-rate, and consequently useful important features such as optical flow cannot be reliably estimated. With the only exception of [3], all existing methods treat photo-streams as a set of unrelated images, neglecting the fact that while the user is performing a given activity such as *cooking*, the temporal coherence of environment concepts such as *food*, *kitchen*, *pan* is preserved despite of suddenly image appearance changes. This can be easily appreciated in Fig. 1, where examples of consecutive frames are shown while the user is performing different activities, including very dynamic activities such as *biking* and *walking outdoor*. Additionally, some images can contain poor information that is difficult to interpret without considering it within its temporal image context. For instance in the 4th image of the activity *walking outdoor* in Fig. 1, it would be impossible to assign the correct label without looking at neighboring images.

In this paper, we propose a new method for activity recognition from photo-streams that is based



Figure 1: Examples of consecutive frames captured while the user was performing four activities: Biking, Walking outdoor, TV, and Reading.

on our previous work [4], but extends it by taking into account the temporal coherence of high level features. More specifically, in [4] we applied a Late Fusion Ensemble (LFE) method that merges through a random decision forest the activity probabilities extracted by a Convolutional Neural Network (CNN) with contextual information embedded in a fully connected layer of the CNN. In this work, we integrate a random decision forest within a temporal framework implemented in terms of a Long Short Term Memory (LSTM) recurrent neural network that processes overlapping batches of fixed size and learns the temporal evolution of high level features.

With respect to [4], our contributions in this paper are as follows:

- we straighten our originally proposed pipeline by integrating for the first time LFE in a framework based on the temporal coherence of high-level features in low temporal resolution photo-streams instead of treating images as non-temporally related,
- we extend the annotated dataset from  $\approx 18K$  images to about  $\approx 45K$  and made the annotations publicly available<sup>1</sup>, we call it the UB Extended Annotations (UBEA) dataset,

<sup>1</sup>The annotations are publicly available at <https://www.github.com/gorayni/egocentric-photostreams>.

- we extend significantly the comparisons with state of the art techniques.
- we provide a more extensive validation of the method previously proposed by us in [4]

The reminder of the paper is organized as follows: in section 2, we discuss related work; in section 3, we detail the proposed approach. In section 4.1, we introduce the dataset used in the experiments. We detail the methodology we followed for conducting our experiments and the different combinations of networks and layers we used in section 4.3, which details the training settings. The results we obtained are discussed in section 4.4. Finally, we present our conclusions and final remarks in section 5.

## 2 Related work

### Activity recognition from egocentric video.

Recently, activity recognition from egocentric videos has become an active area of research, specially using high-temporal frame rate videos. Fathi et al. [9] proposed a probabilistic model that maps activities into a set of multiple actions, and each action is modeled per frame as a spatio-temporal relationship between the hands and the objects involved on it. They further extend their work [10] by proposing another probabilistic generative model that incorporates the gaze features and that

models an action as a sequence of frames. Pirsavash and Ramanan [20] introduced a dataset of 18 egocentric actions of daily activities performed by 20 persons in unscripted videos. On this dataset, they presented a temporal pyramid to encode spatio-temporal features along with detected active objects knowledge. These temporal pyramids are the input of support vector machines trained for action recognition. More recently, Ma et al. [14] proposed a twin stream CNN architecture for activity recognition from videos. One of the streams is used for recognizing the appearance of an object based on a hand segmentation and a region of interest. The other stream recognizes the action using an optical flow sequence. In order to recognize activities, both streams are joined and the last layers are fine-tuned.

The method proposed by Singh et al. [23] first extracts hand masks, head motion, and a saliency map from a sequence of input images. Then, it uses these features as the input of a neural network architecture that combines four streams. The first two streams capture temporal egocentric features by processing a video segment containing binary hand-arm masks. The third stream processes spatial information by using as input the full RGB frame. The last stream handles motion taken from stacked optical flow encoded in a saliency map.

**Activity recognition from egocentric photo-streams.** In contrast to all above mentioned methods that use egocentric videos, activity classification from low-temporal frame rate egocentric photo-streams captured by lifelogging devices has received comparatively little attention in the literature [1]. This may be due partially to the lack of available benchmarks, and partially to the fact that photo-streams provide less contextual action information with respect to video data since images are taken at periodic intervals of 20 or 30 seconds. Due to the lack of temporal coherence, motion features that are the most largely exploited in egocentric videos, cannot be reliably estimated. The pioneering work of Castro et al. [5] proposed a LFE method that combines, through a random decision forest, the classification probabilities of a CNN with time and global features, namely color histogram, from the input image, to classify images into 19 different activ-

ity categories. Their method has been tested on a nonpublic lifelogging dataset made of 40,103 egocentric images captured by a single person during a 6 months period. Since the user activities are performed almost daily at the same time and in the same environment, time and global image features such as color convey useful information for describing the activities in the context of the same wearer. However, the method can hardly generalize to multiple users, since the network needs to adapt or be trained again with respect to the contextual information for each new user. For instance, two distinct persons might have different daily routines, depending on their job, hobbies, age, etc. and, consequently, the system needs to deal with an increased intra-class variability. For instance the “work” activity for different persons could have a very difference appearance.

To address the problem of generalizing the LFE method to multiple users, Cartas et al. [4] used the outputs of different layers of a CNN as contextual information instead of using color and time information that are too much tied to a single user context. The authors tested their approach on a dataset acquired by three different users having different lifestyles. Instead of focusing on contextual information, Oliveira et al. [18] used a gradient boosting machine approach to retrieve activities based on their estimated relations with objects in the scene. Their method was tested on an extended version of the dataset used by Cartas et al. [4], showing promising results.

While all the approaches described so far treated photo-streams as a set of temporally unrelated images, more recently Cartas et al. [3] proposed an end-to-end approach that takes into account the temporal coherence of photo-streams. Their proposed architecture consists of Long Short Term Memory units on the top of a CNN for each frame that is trained by processing the photo-streams using batches of fixed size. Experimental results over the same dataset used in [18], have shown a 6% improvement in accuracy with respect to the VGG-16 baseline.

Inspired by these results, we propose here a new approach that exploits the temporal evolution of high-order features to improve the results of the LTE method proposed in [4].

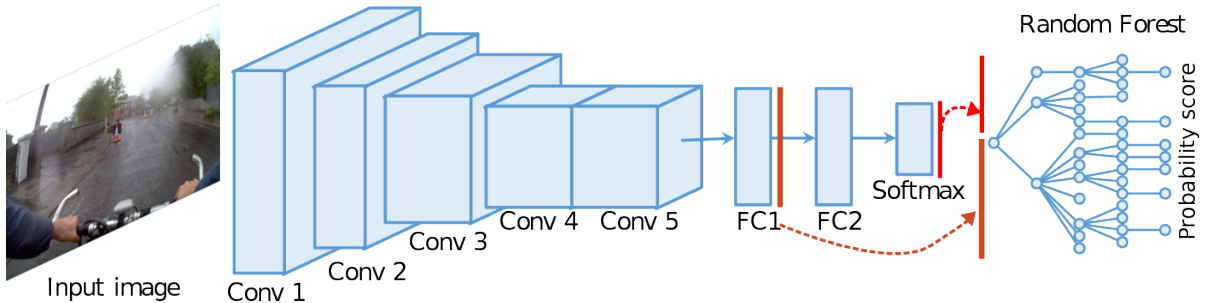


Figure 2: Activity recognition at image level. After fine-tuning a CNN, we combine the softmax probabilities and a fully connected layer into a single vector. Later on, we train a random forest using this vector as its input.

### 3 Ego-centric activity classification

Our base activity classification method is an ensemble classifier composed of a CNN and a random forest, as illustrated on Fig. 2, that acts at image level. Specifically, the random forest takes as input one or more concatenated output vectors from the final layers of a CNN. Depending on the CNN architecture, the input for the random forest can be extracted from the last convolutional layer, a fully-Connected (FC) layer, or the softmax layer. For instance, a random forest that takes as input the output of the softmax and a FC layer is shown on Fig. 2. The training of the ensemble is a two-step process. In the first step, the CNN is fine-tuned. In the second step, a random forest is trained over the output vectors of the CNN. The insight underlying this method is that a fully connected layer provides global contextual information which helps to generalize the features characterizing a given activity among different users.

With the goal of improving classification performance while keeping the advantage of the ensemble classifier, we extended the ensemble architecture to take into account the temporal information from neighbor frames. Firstly, we considered overlapping batches of fixed size and we used a classifier to make a single prediction for each batch. From an implementation point of view, this is equivalent to have a many-to-one predictor since to all frames of the batch is assigned a single label. More specifically, we employed a random forest as

classifier that takes as input the concatenation of the outputs of the first fully-connected layer of the CNN of  $n$  consecutive frames as shown in Fig. 3b. This choice was motivated by the fact that it is less prone to overfitting, does not expect linear features and has a small number of hyperparameters. However, since the Random Forest does not have an internal memory state, it does not take into account explicitly the temporal evolution of high level features, but the set of temporal predictions as a whole.

Secondly, encouraged by the promising results achieved in [3], showing that it is possible to capture the temporal evolution of high-level features extracted with a CNN without knowing event boundaries, and by the success of LSTM in video classification tasks [25, 8], we propose the architecture showed in Fig. 3b. More specifically, we introduce a many-to-many LSTM on the top of the ensemble previously described. Previous attempts to employ LSTM for modeling the temporal evolution of features over time have led to end-to-end trainable architectures [25, 8, 3]. Similarly to the architecture proposed in Fig. 3a, the one proposed in Fig. 3b is not end-to-end: the input of the LSTM is the probability classification score from the ensemble. The training is performed in two phases. We first train an ensemble as stated in the previous section. In order to reduce the computational cost in the second phase, we store the classification scores of each training image. During the second phase, the LSTM unit is trained using the classification scores of  $n$  consecutive frames. As a way of performing data augmentation, the training batches are sampled using a sliding win-

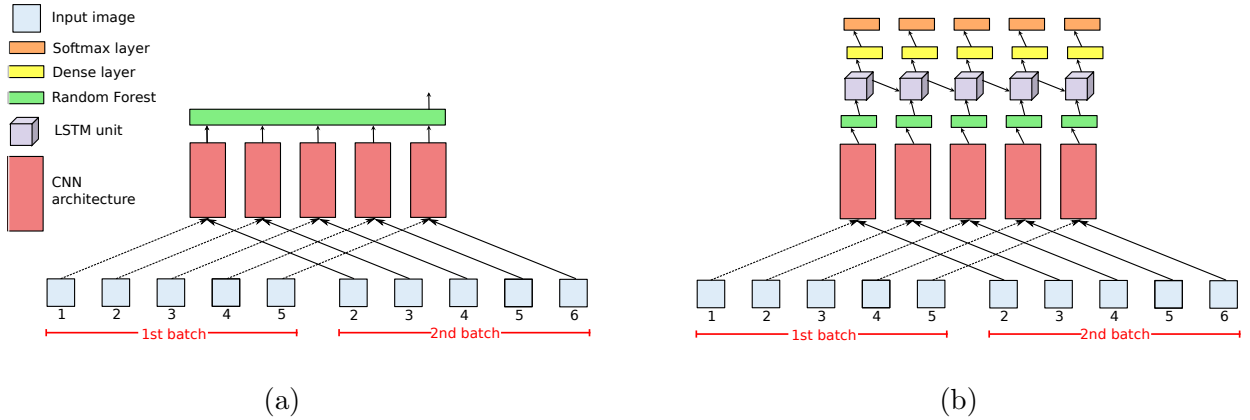


Figure 3: Activity recognition using temporal contextual information. With the activity recognition scores obtained as shown in Fig. 2, in a) we train a random forest by feeding it with overlapping batches of fixed size and a single label is assigned to each image of the batch; b) we train a many-to-many LSTM, feeding it in the same way than in (a).

dow. For example, Fig. 3b shows two consecutive batches of 5 frames that are sampled.

## 4 Validation

The main objective of the experiments performed in this work is to prove that temporal contextual information combined with a CNN LFE improves the state of the art activity classification accuracy on egocentric photo-streams. To this goal, we performed two classes of experiments. Specifically, the goal of the first class of experiments was to determine the ensemble with the best combination of layers for performance improvement based on the proposed approach, firstly proposed in [4], illustrated in Fig. 2. In these experiments, we used three networks as the base of our ensembles, namely the VGG-16 [22], InceptionV3 [24], and ResNet 50 [12]. The goal of the second class of experiments was to determine the best approach to exploit temporal contextual information from consecutive frames. We used the same training setting and CNN as in [3] to compare and evaluate the approaches depicted in Fig. 3 directly to [3], that used the same setting, and other state of the art methods.

We describe the dataset in section 4.1 and detail the ensembles training in section 4.3. We then present the experimental results on activity recognition in section 4.4.

### 4.1 Dataset

In our experiments, we used a subset of images from the NTCIR-12 dataset [11]. This dataset consists of 89,593 egocentric pictures belonging to three persons and acquired with an OMG Autograph camera that captured two pictures per minute. Each user worn this camera in a period about three weeks, totaling 79 days.

We used a subset of 44,901 images from the NTCIR-12 dataset, that we annotated in batches with 21 activity labels, extending the dataset used in [4]. These pictures correspond to all users at different dates and times, but having a similar proportion of about 15,000 images per user. We named these additional labels the UB Extended Annotations (UBEA) dataset. The complete list of activity categories and their distribution of the number of images are shown in Fig. 4 and Fig. 5, respectively. The dataset splits used for the experiments of both methods are explained below.

**Activity recognition at image level.** For this task we used the pictures from the three users for all categories regardless of their date and time as in [4] and illustrated in Fig. 6a (a). In order to maintain the same percentage of samples for each class, we first performed a stratified 10-fold cross-validation over the images. Then a validation split for each fold containing 10% of the data was created by further making a stratified shuffle of its training split.

**Activity recognition using temporal infor-**

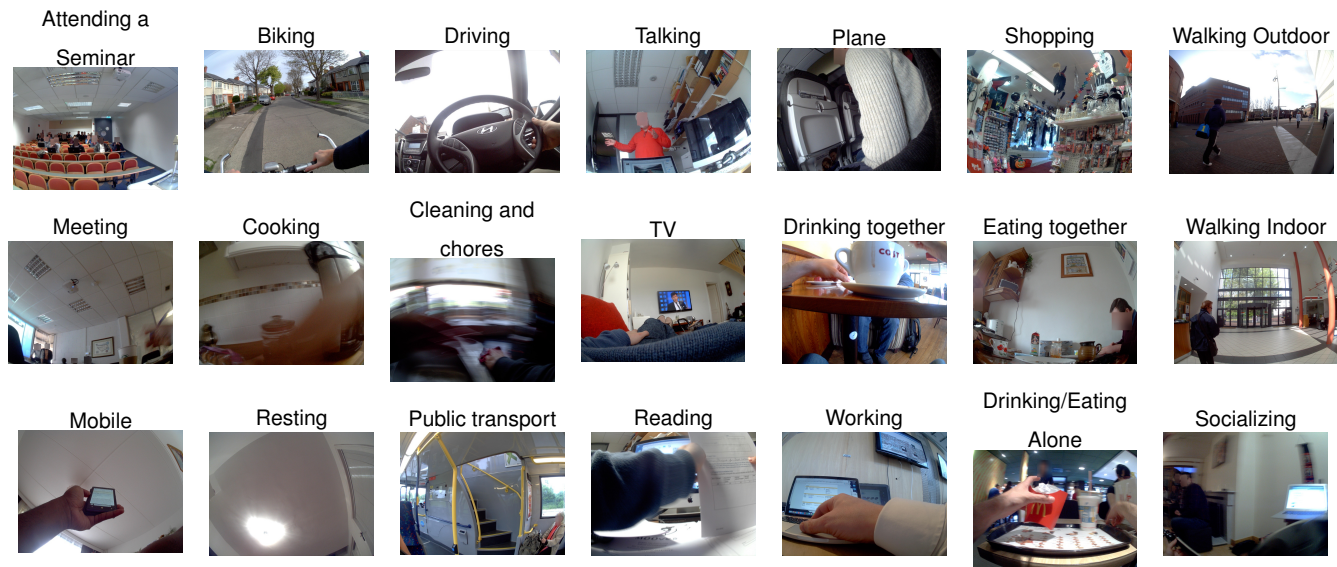


Figure 4: Examples of all activity categories from the annotated dataset.

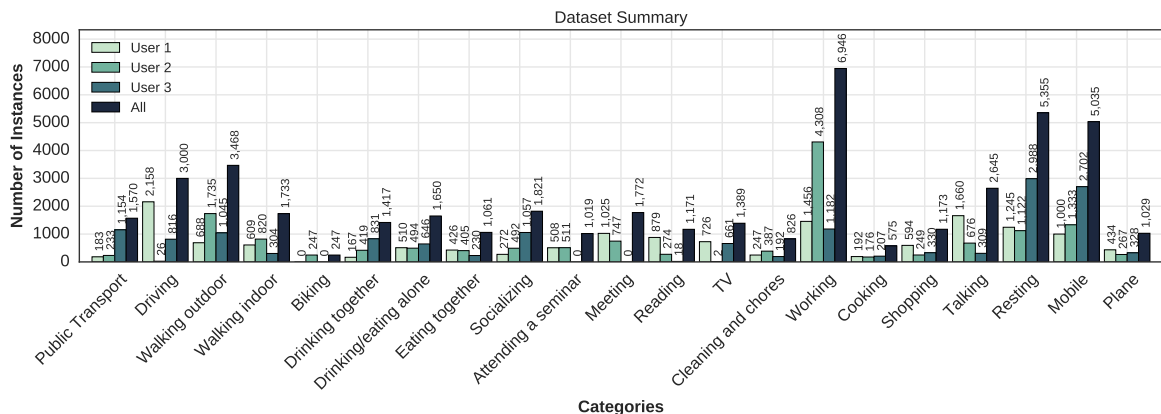


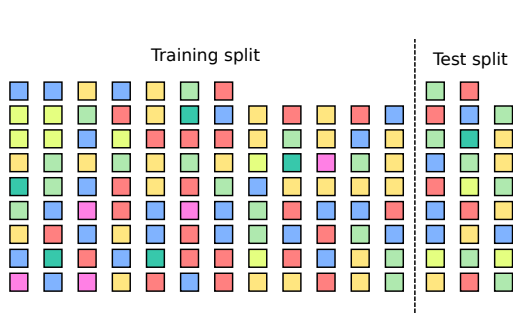
Figure 5: Per class and per user label distribution in our dataset.

**mation.** For this experiment we used the same data split as in [3], illustrated in Fig. 6b. The purpose of this split was to separate full days of activities rather than single frames. This separation made harder the classification task since similar consecutive frames only appear in one split. Additionally, the splits proportionally maintained the class imbalance of the dataset distribution shown in Fig. 5. In order to create the splits, all the combinations of training and testing day splits were enumerated. Then, the Bhattacharyya distance between the normalized distributions of the whole dataset and each split for all combinations was calculated. Finally, the selected combination was the one with the minimal sum of Bhattacharyya dis-

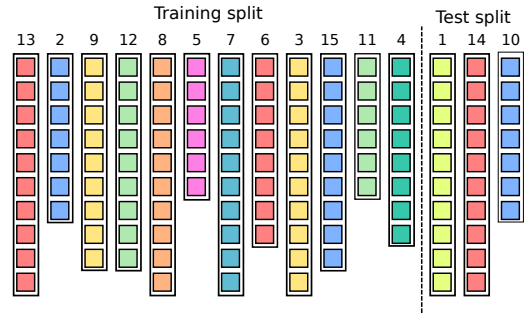
tances.

## 4.2 Evaluation metrics

Since the dataset is highly imbalanced, the classification performance of all methods was assessed by not only using the accuracy, but also macro metrics for precision, recall, and F1-score. Since macro metrics are an unweighted average of the metrics taken separately for each class, they do not take into account of the number of instances available for each class. Moreover, the results on the activity recognition at the image level experiments are cross-validated. Considering that the photo-streams lack of event boundaries, the predictions



(a) Splits used in [4] for activity recognition at image level



(b) Splits used in [3] and this work for activity recognition at batch level.

Figure 6: Visual description of the dataset splits for the experiments carried out. Images of the same color belong to the same day. In (a) the split into training and test does not take into account time metadata, so that images of the same day can be found in both training and test, whereas in (b) images of the same day, can be found or in training or in test even of the split does not take into account the temporal adjacency of days so that temporally adjacent days, say 1 and 2 can be found in different sets (training or test).

done using temporal contextual information were calculated per frame.

### 4.3 Implementation

#### 4.3.1 Activity recognition at image level

We used three different CNN architectures as base models for our ensembles, specifically the VGG-16, InceptionV3, and ResNet 50. All these networks were pre-trained on ImageNet using the Keras framework [6]. Moreover, we employed a class weighting scheme for all CNN models based on [13] to handle class imbalance during training. The random forests in our ensembles were trained using the Scikit Learning framework [19]. As input of the random forest we considered the output of each and all fully connected layers following the last convolutional layer. The training was done on the dataset split described in the previous subsection. First, we trained each CNN on all the cross-validation folds. In order to find the number of trees of each random forest, we trained all the random forests combinations of each CNN using a different number of trees. The criteria used for training them was the Gini impurity[2] and their nodes were expanded until all the leaves were pure. All these random forest combinations were tested on each validation split and the mean accuracy for the corresponding number of trees was plotted, as

shown in Fig. 7. Based on this plot, the number of trees to be used was determined as the lowest one after which the performance did not improve significantly. The training details of each ensemble configuration are detailed below.

**VGG-16.** We fine-tuned a VGG-16 network [22] in two phases. During the first phase, only the top layers were back-propagated with the objective of initialize their weights. The optimization method used was the Stochastic Gradient Descent (SGD) for 10 epochs for all folds, a learning rate  $\alpha = 1 \times 10^{-5}$ , a batch size of 1, a momentum  $\mu = 0.9$ , and a weight decay equal to  $5 \times 10^{-6}$ . In the second phase, the last three convolutional layers were also fine-tuned and the initial weights were obtained from the best epoch of the first phase. Moreover, the SGD ran for another 10 epochs for each fold and set with the same parameters except the learning rate  $\alpha = 4 \times 10^{-5}$ . Additionally, dropout layers were added after each fully-connected layer as a mechanism for regularization.

**VGG-16+RF.** A set of random forests was trained using the distinct layers of the best epoch from the fine-tuned VGG-16 for each fold. The different combinations of layers for the ensembles were: FC1, FC2, FC1+FC2, FC1+softmax, and FC2+softmax. Their corresponding number of trees of each combination were 400, 500, 600, 300, and 200. The maximum depth of each combina-

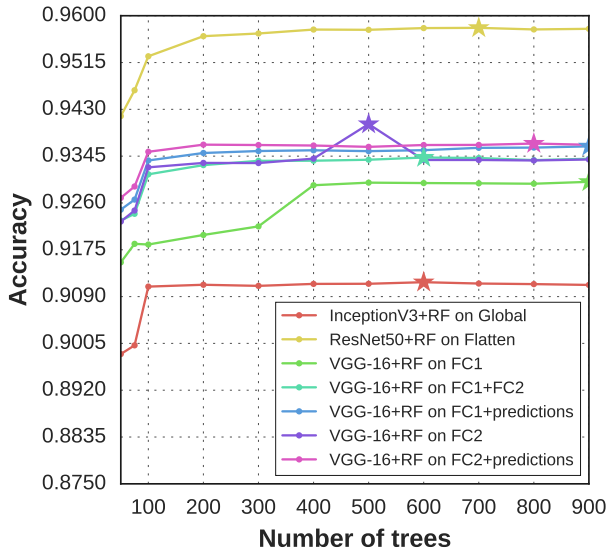


Figure 7: Mean accuracy on the validation set of each fold using a different number of trees. The star mark points at the maximum value. This figure is best seen in color.

tion was 49, 47, 53, 58, and 48, respectively.

**InceptionV3.** InceptionV3 was also fine-tuned in two phases. During the first phase, the last fully-connected layer was optimized using SGD for 10 epochs for all folds, a learning rate  $\alpha = 1 \times 10^{-5}$ , a batch size of 32, a momentum  $\mu = 0.9$ , and a weight decay equal to  $5 \times 10^{-6}$ . During the second phase, the last inception block was added to the optimization process and the network was optimized for another 10 epochs setting the learning rate  $\alpha = 4 \times 10^{-5}$  and the batch size to 10. Moreover, dropout layers were added before and after the last global average pooling layer.

**InceptionV3+RF.** A random forest was trained using the global average pooling layer from InceptionV3 pre-trained on ImageNet. This random forest had 100 trees as estimators and a maximum depth of 58.

**ResNet50.** ResNet50 was fine-tuned in two phases. The last fully-connected layer was optimized in the first phase using SGD for 10 epochs for all folds, a learning rate  $\alpha = 1 \times 10^{-3}$ , a batch size of 32, a momentum  $\mu = 0.9$ , and a weight decay equal to  $5 \times 10^{-6}$ . During the last phase, the last residual block was also optimized using SGD with same learning rate and a batch size of 10 for

three additional epochs.

**ResNet50+RF.** A random forest was trained using the last activation layer from an InceptionV3 network pre-trained on ImageNet. This random forests had 400 trees as estimators and a maximum depth of 49.

### 4.3.2 Activity recognition using temporal information

**VGG-16.** The VGG-16 was trained for 14 epochs using the SGD algorithm. During the first 10 epochs only the fully connected layers were optimized using a learning rate  $\alpha = 1 \times 10^{-5}$ , a batch size of 1, a momentum  $\mu = 0.9$ , and a weight decay equal to  $5 \times 10^{-6}$ . During the last 4 epochs, the last 3 convolutional layers were also fine-tuned and the learning rate changed to  $\alpha = 1 \times 10^{-5}$ . In comparison with the previous experiment, no class weighting scheme was used for training.

**VGG-16+RF.** Five different combinations of random forests using 500 estimators were trained for this experiment. The first four models were previous combinations of layers from the first experiment, specifically, FC1, FC2, FC1+Prediction, and FC2+Prediction. The max. depth of each combination was 55, 52, 60, and 51, correspondingly. Another model was trained for comparison purposes following the description in [5]. In other words, its input was a feature vector obtained by concatenating the softmax probability scores, the day of the week, the time of the day, and 10-bin size histogram for each color channel. Its resulting max. depth was 47. The last trained model used the FC1 layer, the softmax probability scores, the day of the week, the hour, and the time. Its maximum depth was 57.

**VGG-16+RF+LSTM.** For this configuration we trained a LSTM using the predictions from the ensemble of VGG-16 plus the RF on FC1. The LSTM had 32 units and its output was connected to a fully-connected layer. Dropout layers were added between the input and output of the LSTM layer as a way to perform regularization. We trained three configurations using batches of consecutive frames, obtained by using a temporal sliding window of size (or *timestep*) 5, 10, and 15.

**VGG-16+RF many-to-one.** Four different random forests were trained by concatenating the



Table 1: Comparison of the ensembles of CNN+Random forest on different combinations of layers. Upper table shows the recall per class and the lower table shows the performance metrics.

Activity	InceptionV3	InceptionV3+RF on GAP	ResNet50	ResNet50+RF on AP	VGG-16	VGG-16+RF on FC1	VGG-16+RF on FC1+FC2	VGG-16+RF on FC1+Pred	VGG-16+RF on FC2	VGG-16+RF on FC2+Pred
Public Transport	89.23	91.46	<b>92.41</b>	91.59	89.17	90.87	91.14	91.53	91.27	91.59
Driving	99.60	99.83	<b>99.90</b>	99.83	99.50	99.74	99.77	99.77	99.73	99.77
Walking outdoor	90.95	97.35	94.61	<b>98.27</b>	90.89	96.96	97.03	96.94	97.00	96.83
Walking indoor	79.17	95.15	92.73	<b>95.85</b>	79.63	94.81	94.92	94.75	95.04	94.87
Biking	91.93	98.80	99.17	<b>99.20</b>	92.32	98.20	97.57	97.98	97.17	97.98
Drinking together	80.74	95.63	94.93	95.49	89.42	<b>96.48</b>	95.98	96.05	95.91	95.84
Drinking/eating alone	75.94	<b>90.85</b>	88.25	90.37	74.92	88.70	89.77	89.71	89.65	89.71
Eating together	80.58	96.50	96.42	<b>98.11</b>	88.78	96.75	97.26	97.26	96.89	97.36
Socializing	89.13	97.37	97.47	98.63	91.60	98.47	98.46	98.57	<b>98.73</b>	
Attending a seminar	78.70	78.70	<b>80.57</b>	78.31	76.93	77.75	78.21	78.11	78.21	78.11
Meeting	80.92	93.29	96.61	<b>96.78</b>	90.00	95.36	95.49	95.49	95.54	95.43
Reading	81.40	97.35	96.83	<b>97.69</b>	88.23	96.49	96.67	96.67	96.59	96.59
TV	93.74	96.55	96.04	96.54	92.59	<b>97.21</b>	96.91	96.69	97.05	96.91
Cleaning and chores	57.14	92.61	92.22	92.39	61.84	92.06	92.13	92.62	92.26	<b>92.74</b>
Working	95.82	98.81	97.34	<b>99.38</b>	96.27	98.18	97.98	97.91	97.96	97.91
Cooking	81.98	95.10	91.69	96.15	84.78	96.48	95.97	95.79	95.47	<b>96.49</b>
Shopping	81.07	95.40	93.18	<b>96.68</b>	82.44	94.98	95.57	95.14	95.57	95.65
Talking	72.36	<b>95.54</b>	89.64	94.10	76.25	93.53	93.50	93.50	93.65	93.61
Resting	87.97	91.07	91.54	91.02	92.87	94.56	94.90	<b>95.14</b>	94.90	95.11
Mobile	89.59	97.22	95.27	<b>97.58</b>	93.32	97.55	97.46	97.36	97.30	97.34
Plane	<b>87.67</b>	87.16	82.61	75.70	84.94	77.21	80.66	81.15	81.54	83.58
<b>Accuracy</b>	87.07	95.27	94.08	95.39	89.46	95.26	95.41	95.41	95.41	<b>95.50</b>
<b>Macro precision</b>	83.85	94.96	92.60	<b>96.14</b>	87.06	94.63	94.59	94.62	94.57	94.61
<b>Macro recall</b>	84.08	94.37	93.31	94.27	86.51	93.92	94.16	94.19	94.15	<b>94.39</b>
<b>Macro F1-score</b>	83.77	94.53	92.78	<b>95.00</b>	86.45	94.11	94.23	94.28	94.23	94.38

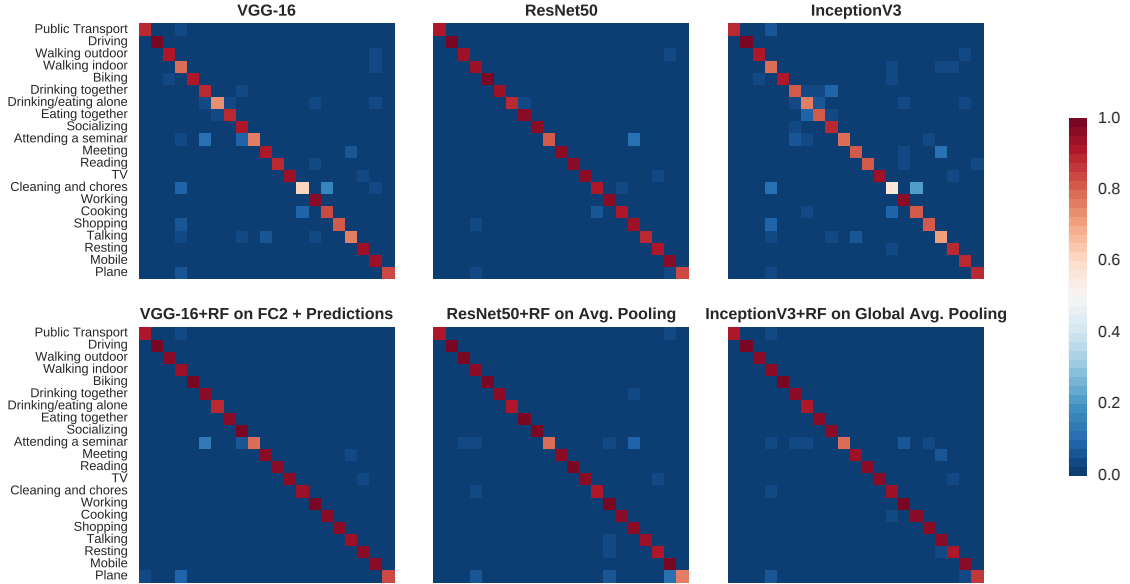


Figure 8: Normalized confusion matrices of the best combination of layers for each baseline convolutional neural network. This figure is best seen in color.

output vectors of the FC1, FC2, and predictions of the VGG-16 network of a consecutive number of frames. This configuration employed the same sampling window strategy for training. For example, Fig. 3a shows a 5 timestep configuration training on two consecutive batches and predicting the activity from the last frame. On this architecture, the same label is assigned to all frames considered. We considered configurations with timestep 5 and

10. The number of estimators for FC1, FC1 + Prediction, FC2, and FC2 + Prediction were 700. The resulting maximum depth of each combination for timestep 5 was 57, 56, 54, and 53; and for timestep 10 was 53, 60, 53, and 53.

**InceptionV3.** This network was trained for 12 epochs using the SGD algorithm. During the first 10 epochs only the last fully-connected layer was optimized using a learning rate  $\alpha = 1 \times 10^{-5}$ , a

batch size of 10, a momentum  $\mu = 0.9$ , and a weight decay equal to  $5 \times 10^{-6}$ . During the last 2 epochs, the last inception block was also fine-tuned and the learning rate changed to  $\alpha = 4 \times 10^{-5}$ . No class weighting scheme was used for training.

**InceptionV3+RF.** The global average pooling (GAP) layer was used as the input of a random forests using 100 estimators and its maximum depth was 72.

**InceptionV3+RF+LSTM.** We trained a LSTM using the predictions from the ensemble of InceptionV3 plus the RF on GAP layer. The LSTM had 32 units and its output was connected to a fully-connected layer. Dropout layers were added between the input and output of the LSTM layer as a way to perform regularization. We trained three configurations using batches of consecutive frames, obtained by using a temporal sliding window of *timestep* 5, 10, and 15.

**InceptionV3+RF many-to-one.** One random forest was trained by joining the output vector of the GAP layer of the InceptionV3 network of a consecutive number of frames. This configuration employed the same sampling window strategy for training. We considered configurations with timestep 5 and 10. The number of estimators for both cases was 700. The maximum depth obtained for both timesteps was 62 and 66, correspondingly.

**ResNet50.** This network was trained for 4 epochs using the SGD algorithm. The last fully-connected layer was optimized in the first phase using SGD for 2 epochs for all folds, a learning rate  $\alpha = 1 \times 10^{-3}$ , a batch size of 10, a momentum  $\mu = 0.9$ , and a weight decay equal to  $5 \times 10^{-6}$ . During the last 2 epochs, the last inception block was also fine-tuned using the same values. No class weighting scheme was used for training.

**ResNet50+RF.** A random forest was trained on the AP layer using 500 estimators and its resulting max. depth was 53.

**ResNet50+RF+LSTM.** We trained a LSTM using the predictions from the trained ensemble. The LSTM had 32 units and its output was connected to a fully-connected layer. Dropout layers were added between the input and output of the LSTM layer as a way to perform regularization. We trained three configurations using batches of consecutive frames, obtained by using *timesteps*

of 5, 10, and 15.

**ResNet50+RF many-to-one.** One random forest was trained by concatenating the output vector of the AP layer of the ResNet50 network of a consecutive number of frames. This configuration employed the same sampling window strategy for training. We considered configurations with timestep 5 and 10. The number of estimators for both cases was 700. Moreover, the maximum depth for the configurations with timestep 5 and 10 were 65 and 63, respectively.

## 4.4 Results and discussion

**Activity recognition at image level** The results of adding contextual information from fully connected layers are presented on Table 1. They show that CNN LFE improves the performance for all baseline CNN on different ensembles. Moreover, Table 1 shows that the best ensembles were the *InceptionV3+RF on Global Avg. Pooling*, *VGG-16+RF on FC2 + Predictions*, and *ResNet+RF on average pooling*. Furthermore, these ensembles improved the baseline accuracy by 8.2%, 6.04%, and 1.31% for InceptionV3, VGG-16, and ResNet, respectively. Specifically, the recall of some categories with fewer learning instances improved significantly, such as *Cooking*, and *Cleaning and Choring*. Additionally, high overlapping classes such as *Drinking/Eating alone* and *Eating together* also improved their accuracy. The improvement over the overlapping classes can also be seen on the confusion matrices shown in Fig. 8. This means that the random forest improved the classification of images belonging to categories that score similar probabilities. Moreover, the only decrease on accuracy is presented on the class *Plane*. Since its accuracy on the baseline CNN is very high (87.67%) considering the small number of learning instances (1,026), we believe this decrease is a consequence of the random forest trying to balance the prediction error among classes. Some classification examples are shown in Fig. 9.

**Activity recognition at batch level** The results of using temporal contextual information from batch sequences is shown on Table 2. The fact that accuracy performance is much higher when using the splits defined as in Fig. 6b (a)



Figure 9: Classification activity examples. On top of each image is shown its true activity label and on bottom its top 5 predictions by VGG-16, ResNet50 and InceptionV3. Additionally, the result of the ensembles *VGG-16+RF on FC1+Softmax*, *ResNet50+RF on Avg. Pooling*, and *InceptionV3+RF on Global Avg. Pooling* is highlighted on color in its corresponding table. The green and red colors means true positive and false positive classification, respectively.

for the VGG-16 can be easily understood when looking at consecutive frames in Fig. 1. Specially in the case of static activities such as *reading* or *watching TV*, consecutive frames are very similar and putting some of them in training and other in test is almost equivalent at doing the test including training images, which is unfair.

Comparing the performances of VGG-16+LSTM with VGG-16+RF+LSTM, our experiments demonstrated that the activity scores obtained with CNN LFE are better features for further temporal analysis with respect to features computed by a end-to-end architecture. In both cases, the best performances were achieved by a larger timestep.

**State of the art comparisons.** In Table 2, we compared the proposed method with the CNN baselines, with our previous method [4], with the work of Castro et al. [5] and with [3], that exploits temporal information. Our proposed approach *CNN+RF+LSTM* gives the best results for each CNN baseline when using 10 as timestep

value. Specifically, the *InceptionV3+RF+LSTM* with a 10 timestep achieved 89.85% accuracy. It is followed by the end-to-end *VGG-16+LSTM* proposed in [3] that also takes into account temporal information and achieves approximately 6% less in terms of accuracy. On the contrary, with *CNN+RF*, being a many-to-one architecture, the best results are achieved when using a small timestep. We can also observe that, in general, methods that take into account temporal information achieve better performance with respect to methods that act at image level. Among these latter methods, not all the fine-tuned CNNs were improved by all LFE methods used for comparison. On this dataset, the methods [4] and [5] are almost equivalent. To explain this result we performed further experiments to investigate the contribution of time metadata and we could double check that indeed they help in improving the performances (*VGG-16 + RF on FC1 + softmax + date&time*), a result that would unlikely be generalizable to a dataset with more individuals hav-

ing a very different lifestyle. The fact that color histogram as contextual feature leads to slightly better performances might be due to the different light conditions imposed the activity categories.

In this study, we limited our state of the art comparisons to activity recognition methods conceived specifically for egocentric images or photo-streams. Indeed, state of the art activity recognition methods developed for videos often rely on features that cannot be computed in a lifelogging setting. For example, the method proposed by Fathi et al. [9] relies on a hierarchy between actions and activities, and the sequential order of the former; but, given the frame-rate of photo-streams, our labels do not consider actions. Another method presented by Fathi et al. [10] relies on the egocentric gaze provided by the camera, which is not available for a chest-mounted camera as the one used for capturing our dataset. In addition, the classification methods proposed in [14, 23] depend on motion features that can not be computed on photo-streams because of the very low frame rate. Finally, the temporal pyramid approach proposed in [20] strongly relies on shot detection which is in general difficult to compute in photo-streams [7], and highly depends on additional cues based on object detection and human-object interaction.

Although our results are encouraging, we recognize that assessing the generalization capability of our system would require a dataset captured by a large number of persons having a wide variety of lifestyles. Given that lifelogging is a relatively new area of research, and that building such a benchmark would be very expensive in terms of human resources and annotation effort, we leaved this for future work.

## 5 Conclusion

In this paper, we proposed a new pipeline for activity recognition from egocentric photo-streams that relies on a two-phases approach. First, a CNN late fusion ensemble classifier, which combines different layers of a CNN through a random forest is used to predict the activity probabilities on each image. Specifically, the random forest takes as input a vector containing the output of the softmax

probability layer and a fully connected layer encoding global image features. Second, these vector probabilities are used in batches of temporally adjacent images, to train a many-to-many LSTM. In addition, we extended a subset of the NTCIR-12 egocentric dataset consisting of 44,902 images by annotating it with 21 different activity labels. Experimental results on this subset demonstrated that the proposed approach achieve better performances than a end-to-end architecture and state of the art methods. The proposed method achieved an overall accuracy of 89.85%.

## Acknowledgements

A.C. was supported by a doctoral fellowship from the Mexican Council of Science and Technology (CONACYT) (grant-no. 366596). This work was partially founded by TIN2015-66951-C2, SGR 1219, CERCA, *ICREA Academia'2014* and 20141510 (Marató TV3). The funders had no role in the study design, data collection, analysis, and preparation of the manuscript. M.D. is grateful to the NVIDIA donation program for its support with GPU card.

## Notes

A final version of this manuscript was published in the *Pattern Analysis and Applications* journal under the name *Batch-Based Activity Recognition from Egocentric Photo-Streams Revisited*.

## References

- [1] Bolaños, M., Dimiccoli, M., Radeva, P.: Toward storytelling from visual lifelogging: An overview. *IEEE Transactions on Human-Machine Systems* **47**(1), 77–90 (2017)
- [2] Breiman, L., Friedman, J., Stone, C.J., Olshen, R.A.: *Classification and regression trees*. CRC press (1984)
- [3] Cartas, A., Dimiccoli, M., Radeva, P.: Batch-based activity recognition from egocentric

- photo-streams. In: proceedings of the International Conference on Computer Vision (ICCV), Workshop on Egocentric Perception, Interaction and Computing. IEEE (2017)
- [4] Cartas, A., Marín, J., Radeva, P., Dimiccoli, M.: Recognizing activities of daily living from egocentric images. In: Proceedings of the Iberian Conference on Pattern Recognition and Image Analysis (IbPRIA), pp. 87–95. Springer, Cham (2017)
- [5] Castro, D., Hickson, S., Bettadapura, V., Thomaz, E., Abowd, G., Christensen, H., Essa, I.: Predicting daily activities from egocentric images using deep learning. In: proceedings of the 2015 ACM International symposium on Wearable Computers, pp. 75–82. ACM (2015)
- [6] Chollet, F., et al.: Keras. <https://github.com/fchollet/keras> (2015)
- [7] Dimiccoli, M., Bolaños, M., Talavera, E., Aghaei, M., Nikolov, S.G., Radeva, P.: Sr-clustering: Semantic regularized clustering for egocentric photo streams segmentation. *Computer Vision and Image Understanding* **155**, 55–69 (2016)
- [8] Donahue, J., Hendricks, L.A., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: CVPR (2015)
- [9] Fathi, A., Farhadi, A., Rehg, J.M.: Understanding egocentric activities. In: 2011 International Conference on Computer Vision, pp. 407–414. IEEE (2011)
- [10] Fathi, A., Li, Y., Rehg, J.M.: Learning to recognize daily actions using gaze. In: European Conference on Computer Vision, pp. 314–327. Springer (2012)
- [11] Gurrin, C., Joho, H., Hopfgartner, F., Zhou, L., Albatal, R.: Overview of ntcir-12 lifelog task. In: NTCIR-12. National Institute of Informatics (NII). URL <http://eprints.gla.ac.uk/119206/>
- [12] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
- [13] King, G., Zeng, L.: Logistic regression in rare events data. *Political Analysis* **9**(2), 137–163 (2001). DOI 10.1093/oxfordjournals.pan.a004868
- [14] Ma, M., Fan, H., Kitani, K.M.: Going deeper into first-person activity recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
- [15] Martin-Lesende, I., Vrotsou, K., Vergara, I., Bueno, A., Diez, A., et al.: Design and validation of the vida questionnaire, for assessing instrumental activities of daily living in elderly people. *J Gerontol Geriat Res* **4**(214), 2 (2015)
- [16] Mukhopadhyay, S.C.: Wearable sensors for human activity monitoring: A review. *IEEE sensors journal* **15**(3), 1321–1330 (2015)
- [17] Nguyen, T.H.C., Nebel, J.C., Florez-Revuelta, F.: Recognition of activities of daily living with egocentric vision: A review. *Sensors (Basel)* **16**(1), 72 (2016). DOI 10.3390/s16010072. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4732105/>. Sensors-16-00072[PII]
- [18] Oliveira-Barra, G., Dimiccoli, M., Radeva, P.: Leveraging activity indexing for egocentric image retrieval. In: Iberian Conference on Pattern Recognition and Image Analysis, pp. 295–303. Springer (2017)
- [19] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
- [20] Pirsiavash, H., Ramanan, D.: Detecting activities of daily living in first-person camera

views. In: Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on, pp. 2847–2854. IEEE (2012)

- [21] Schüssler-Fiorenza Rose, S.M., Stineman, M.G., Pan, Q., Bogner, H., Kurichi, J.E., Streim, J.E., Xie, D.: Potentially avoidable hospitalizations among people at different activity of daily living limitation stages. Health services research (2016)
- [22] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR [abs/1409.1556](#) (2014)
- [23] Singh, S., Arora, C., Jawahar, C.V.: First person action recognition using deep learned descriptors. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
- [24] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
- [25] Yue-Hei Ng, J., Hausknecht, M., Vijayanarasimhan, S., Vinyals, O., Monga, R., Toderici, G.: Beyond short snippets: Deep networks for video classification. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4694–4702 (2015)

Method	Accuracy	Macro Precision	Macro Recall	Macro F1-score
ResNet50	78.44	72.44	70.62	69.85
ResNet50+RF on GAP	74.62	76.64	62.61	64.94
ResNet50+RF on GAP timestep 5	73.09	71.97	59.52	62.51
ResNet50+RF on GAP timestep 10	71.29	72.12	57.23	60.16
ResNet50+RF+LSTM timestep 5	81.10	75.69	72.90	72.41
ResNet50+RF+LSTM timestep 10	85.12	81.48	78.98	78.76
ResNet50+RF+LSTM timestep 15	83.29	80.75	78.04	77.24
InceptionV3	78.31	73.06	69.94	70.60
InceptionV3+RF on AP	77.08	73.68	67.24	68.54
InceptionV3+RF timestep 5	78.82	75.62	69.11	69.93
InceptionV3+RF timestep 10	78.55	75.93	69.12	69.89
InceptionV3+RF+LSTM on AP timestep 5	86.04	85.58	83.74	83.56
<b>InceptionV3+RF+LSTM on AP timestep 10</b>	<b>89.85</b>	<b>89.58</b>	<b>87.20</b>	<b>87.22</b>
InceptionV3+RF+LSTM on AP timestep 15	87.72	86.95	84.44	84.37
VGG-16	75.97	68.50	67.49	66.80
VGG-16+RF on FC1	74.22	70.77	65.06	65.80
VGG-16+RF on FC1+Pred[4]	76.80	71.93	66.86	67.51
VGG-16+RF on FC2	75.02	69.12	65.60	65.60
VGG-16+RF on FC2+Pred	75.87	69.78	66.08	66.28
VGG-16+RF on Softmax + Date&Time + Color [5]	77.39	69.66	67.79	66.99
VGG-16+RF on FC1 + Softmax + Date&Time	76.91	72.07	66.77	67.27
VGG-16+RF on FC1 timestep 5	76.04	74.44	65.98	66.93
VGG-16+RF on FC1 timestep 10	75.67	75.20	64.93	66.26
VGG-16+RF on FC1+Pred timestep 5	77.43	74.54	67.13	67.57
VGG-16+RF on FC1+Pred timestep 10	76.66	74.10	65.91	66.67
VGG-16+RF on FC2 timestep 5	76.83	73.48	67.92	67.92
VGG-16+RF on FC2 timestep 10	76.49	73.40	67.72	67.61
VGG-16+RF on FC2+Pred timestep 5	77.79	74.60	68.31	68.53
VGG-16+RF on FC2+Pred timestep 10	77.75	74.82	67.89	68.37
VGG-16+LSTM [3] on FC1 timestep 5	79.68	72.96	71.36	70.87
VGG-16+LSTM [3] timestep 10	80.39	75.25	71.86	71.97
VGG-16+LSTM [3] timestep 15	81.73	76.68	74.04	74.16
VGG-16+RF+LSTM FC1 timestep 5	85.96	79.81	81.36	80.00
VGG-16+RF+LSTM FC1 timestep 10	86.87	80.45	81.36	80.39
VGG-16+RF+LSTM FC1 timestep 15	85.55	80.00	79.34	78.45
VGG-16+RF+LSTM FC1+Pred timestep 5	85.45	79.29	79.11	78.05
VGG-16+RF+LSTM FC1+Pred timestep 10	87.71	81.55	81.88	81.10
VGG-16+RF+LSTM FC1+Pred timestep 15	86.24	80.05	80.79	79.56
VGG-16+RF+LSTM FC2 timestep 5	85.73	81.22	80.56	79.96
VGG-16+RF+LSTM FC2 timestep 10	86.78	81.37	80.97	80.37
VGG-16+RF+LSTM FC2 timestep 15	85.38	79.71	79.80	78.63
VGG-16+RF+LSTM FC2+Pred timestep 5	85.89	79.62	80.35	78.82
VGG-16+RF+LSTM FC2+Pred timestep 10	86.85	80.22	81.19	80.23
VGG-16+RF+LSTM FC2+Pred timestep 15	86.35	80.11	80.74	79.51

Table 2: Performance summary of the experiments on activity recognition using temporal contextual information.