

A Deep Learning Approach for Human Action Recognition using Skeletal Information

Eirini Mathe^{1,2}, Apostolos Maniatis³, Evaggelos Spyrou^{1,3}, and Phivos Mylonas²

- ¹ Institute of Informatics and Telecommunications, National Center for Scientific Research – “Demokritos”, Athens, Greece
{emathe, espyrou}@iit.demokritos.gr,
- ² Department of Informatics, Ionian University, Corfu, Greece
fmylonas@ionio.gr,
- ³ Department of Computer Engineering T.E., Technological Education Institute of Sterea Ellada, Lamia, Greece
amaniatis@teiste.gr

Abstract. In this paper we present an approach towards human action detection for activities of daily living (ADLs) that uses a Convolutional Neural Network (CNN). The network is trained on Discrete Fourier Transform (DFT) images that result from raw sensor readings, i.e., each human action is ultimately described by an image. More specifically, we work using 3D skeletal positions of human joints, which originate from processing of raw RGB sequences enhanced by depth information. The motion of each joint may be described by a combination of 3 1D signals, representing its coefficients into a 3D Euclidean space. All such signals from a set of human joints are concatenated to form an image, which is then transformed by DFT and is used for training and evaluation of a CNN. We evaluate our approach using a publicly available challenging dataset of human actions that may involve one or more body parts simultaneously and for two sets of actions which resemble to common ADLs.

Keywords: human action recognition, activities of daily living, convolutional neural networks

1 Introduction

Human action recognition still remains one of the most challenging research areas in the field of computer vision. Several open challenges in this area include the representation, the analysis and ultimately the recognition of the human actions [2]. To this goal, machine learning approaches have been widely used. However, traditional machine learning approaches fail to show robustness when the number of possible actions increases or when the camera angle changes. During the last few years, advances in hardware have facilitated training and application of deep neural network architectures [9] which are able to learn representations from data without the need for hand-crafted rules or features, while

their accuracy may significantly increase when they are provided with more data. Moreover, recently several publicly available datasets [13] have emerged in the field of human action recognition, enabling the evaluation of novel architectures and action representations in real-like scenarios. Apart from the aforementioned challenges, the design of novel deep architectures and their application in real-life scenarios, are also among the targets of research in this field.

In this paper, we propose a novel visual representation of human actions, which based on the Discrete Fourier Transformation (DFT). More specifically, we concatenate raw signal images that result from the 3D motion of human skeletal joints. The input required for the extraction of these joints consists of aligned RGB and depth video sequences and is performed using the well-known Kinect v2 camera and its accompanying SDK. Moreover, we propose a novel Convolutional Neural Network (CNN) architecture which uses as input the DFT transformation of the aforementioned images. We evaluate the proposed approach using the challenging PKU-MMD dataset [13] consisting of 51 human actions and we demonstrate that the proposed approach may be used in real-like environments for the recognition of activities of daily living (ADLs) [10].

The rest of this paper is organized as follows: Section 2 presents state-of-the-art in the field of human action recognition using deep learning approaches and focusing on those that work on skeletal information. Section 3 presents the concepts of deep learning and Convolutional Neural Networks that have been used in the context of this work. The proposed action representation and deep network architecture are then described in section 4. Experimental results are presented in 5 and discussed in 6, which also included plans for further extensions and applications of this work.

2 Related Work

The problem of human action recognition has attracted many research efforts, which have been continuously growing during the last decade. In this section we aim to present approaches that are based on deep networks. Typically, these works do not include a feature extraction step; they are instead based upon a representation of the action. However, we should herein emphasize that there still exist approaches that propose the extraction of features [15].

Skeletal data consist of the 3D positions of human skeleton joints. These may be considered as high-level features for the recognition process. The most popular method to extract the skeleton is based on RGB sequences accompanied by corresponding depth maps, i.e., as the approach adopted by Kinect sensors. Of course, skeletons are prone to errors, due to e.g., occlusion and viewpoint changes. Moreover, certain actions may have significantly different appearance upon abrupt changes of viewpoint. There exist two major categories of tasks: a) segmented recognition; and b) continuous (online) recognition [17]. The difference between the two categories is that within the first, we assume that the input video sequence only contains the action to be recognized, i.e., frames not depicting the action (before/after the action) have been removed. Note, that for the

first category, common deep architectures used are Recurrent Neural Networks (RNNs) [6] and Convolutional Neural Networks (CNNs) [11]. For the second category RNNs are typically used. In case a CNN is used, the majority of the approaches includes a step of converting skeleton sequences to a single image, in a way that both spatial and temporal information is maintained and reflected to low-level image properties, i.e., color and/or texture may be used for the separation of classes. Note, that the proposed approach uses a CNN and a step for converting 1D skeleton sequences to a single image, as it will be described in section 4.

In the work of Du et al. [4], the authors divide the skeleton joints into five groups (arms, legs and trunk), i.e., corresponding joints are concatenated as a single vector. All five parts are then concatenated so as to capture the spatial information per frame, while x , y and z components of their 3D coordinates correspond to the R, G, B components of a color image, respectively. Then, representations of all frames of a sequence are arranged chronologically, to capture its temporal properties. A CNN architecture is used for classification. Wang et al. [18] propose the use of “joint trajectory maps,” where hue is used to capture the motion direction information of skeleton joints. Motion trajectories are projected onto three Cartesian planes (i.e., front, top and side plane) and motion magnitude is encoded by appropriately settings of saturation and brightness, so that motion changes are reflected to changes of texture. The resulting maps are classified into actions by CNNs. Similarly, Hou et al. [5] also encoded skeleton joints’ sequences into “skeleton optical spectra,” which were also color texture images. The variation of color was used to introduce the temporal information to the representation, as changes of hue.

Li et al. [12] proposed the use of “joint distance maps,” which are also texture images. Contrast to [18] and [5], projections to the three Cartesian planes are unnecessary. Instead, pair-wise distances of joints are used. 3 maps are used to encode the distances in the 3 orthogonal 2D planes and a fourth one is used to encode distances in the 3D space. Hue is used for encoding variations of distances. This way, the description is more robust to changes of viewpoint, which as we have already discussed are common in real-life applications. A CNN is then used for each map and classification is a result of a late fusion scheme which is applied. In the work of Liu et al. [14], transforms are applied to skeleton sequences in an effort to make them invariant to the position and the initial orientation of the skeleton. Skeleton data are considered as points into a 5D space; each consists of 3D space coordinates, time and joint label. They are then projected into a 2D image by selecting two of the aforementioned dimensions, while the remaining three are used as R, G, B values. This way, color images are formed and used as input to a multi-stream CNN scheme. Finally, Ke et al. [8] presented “SkeletonNet,” where contrary to the majority of the approaches, they did not extract 3D coordinates. Instead, they extracted translation, rotation and scale invariant features. More specifically, the skeleton was divided into five parts as in [4]. Then from each part they extracted vector representations which are generated from pairwise relative positions between joints. Cosine distances between

the aforementioned vectors within a specific part and normalized magnitudes of each vector are extracted. These ten representations are concatenated and then used as input to a two-stream CNN.

3 Deep Learning and Convolutional Neural Networks

Deep learning is a sub-field of machine learning which has attracted a lot of research interest during the last few years. Its main idea is the use of multiple layers to non-linearly process the network’s input, so as to “learn” to extract features. The output of each layer is fed to the next layer. Ultimately, they become able to learn multiple levels of representations which correspond to multiple levels of abstraction. Deep network architectures play a key role in several application fields such as computer vision, audio analysis, speech recognition etc., i.e., in tasks where traditional machine learning approaches fail to achieve acceptable levels of accuracy for real-life applications. It is generally accepted that the computer vision is the area that has benefited the most; a plethora of deep architectures have been proposed during the last few years and have been successfully applied to traditional computer vision problems as well as to novel applications.

The most common approach when dealing with computer vision problems are the Convolutional Neural Networks (CNNs) [11]. The architecture of a CNN resembles to the one of a traditional neural network (NN), however, its goal is to learn a set of convolutional filters. Training takes place as with every other NN; a forward propagation of data and a backward propagation of error do take place to update weights. The *convolutional* layers are those that play the key role in the whole process. Their neurons are grouped in rectangular grids, so that each would perform a convolution in a part of the input image. Learning process aims to learn the parameters of this convolution. *Pooling* layers are usually placed after a single or a set of serial or parallel convolutional layers. Their input consists of small rectangular image blocks from the convolutional layer. The latter are then subsampled; a single output is produced from each block. Finally, *dense* layers (which are commonly referred to as “fully-connected” layers) are the ones that are responsible for classification, based on the features that were previously extracted by the convolutional layers and subsampled by the pooling layers. Note that each node of a dense layer is connected to all nodes of its previous layer. To avoid overfitting, one approach (which we also adopt in this work) is the use of the *dropout* regularization technique [16]. When using this technique, at each training stage several nodes are “dropped out” of the network. This way, complex co-adaptations on training data are prevented and this leads to the reduction or even total prevention of overfitting.

4 Human Action Recognition

The proposed approach uses as its input 3D skeletal data that have been captured by the Microsoft Kinect v2 sensor [19] which combines a traditional RGB and

a depth camera. Kinect is complemented by its SDK, which among others is able to provide the 3D positions of a predefined set of human skeletal joints, in real time. A graph representation has been adopted; nodes correspond to body parts (e.g., arms, legs, head etc.), edges follow the joints' structure. Note that a parent-child relationship is implied, i.e., HEAD is parent of NECK, while NECK is parent of SPINE SHOULDER, etc. A total of 25 joints are available. We should emphasize that since for each joint its x , y and z coordinates are provided. We consider each coordinate of each joint as a single 1D signal, thus 75 1D signals result, for any given video sequence and for each person. In Fig. 1 we illustrate the 25 human skeleton joints, that are extracted using the Kinect SDK.

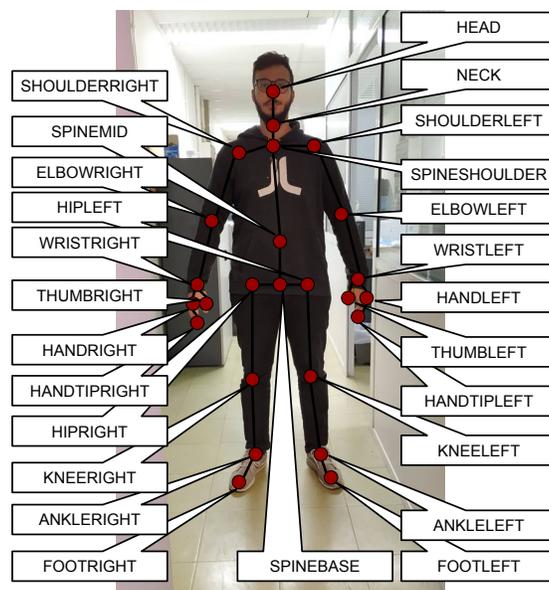


Fig. 1: Extracted human skeleton 3D joints using the Kinect SDK.

Inspired by the work of Jiang and Yin [7], we first create an activity image by concatenating the aforementioned 75 1D signals. We will refer to the result of this concatenation as “signal image.” Then, we apply the 2D Discrete Fourier Transform (DFT) to the signal image and preserve only the magnitude of the transform (i.e., the phase is discarded). The result is again an image, which we will refer to as “activity” image. In Fig. 2 we illustrate an example signal image and the corresponding activity image.

We should herein emphasize that our work focuses only on the classification of a given action into a set of predefined classes. Therefore, we should clarify that it does not perform any temporal segmentation (i.e., to detect the beginning and the ending of a possible action); instead we consider this problem as solved.

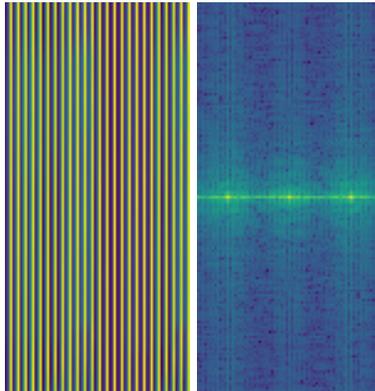


Fig. 2: Left: a signal image. Right: an activity image. For visualization purposes only, the activity image has been processed with a log transformation. Figure best viewed in color.

For evaluation purposes we work on pre-segmented sequences of videos, aiming to only recognize the performed actions within each segment. We also assume that each segment contains exactly one action. We should highlight that human-performed actions may typically vary in terms of duration, even when performed by the same user, thus an interpolation step should be necessary. To tackle this issue and upon experimentation we decided to set a threshold T_s for the duration of each action. Signals resulting from all actions with a duration $T_a < T_s$ are padded with zeros, while the length of those with $T_a > T_s$ is reduced upon a linear interpolation step. This way, all signal images have a fixed length of $T_s \times 75$.

The architecture of our proposed CNN is presented in detail in Fig. 3. The first convolutional layer filters the 159×75 input activity image with 32 kernels of size 3×3 . Then the first pooling layer uses “max-pooling” to perform 2×2 subsampling. A second convolutional layer filters the 36×78 resulting image with 64 kernels of size 3×3 . Then a second pooling layer uses “max-pooling” to perform 2×2 subsampling. The third convolutional layer filters the 17×38 resulting image with 128 kernels of size 3×3 . A third pooling layer uses “max-pooling” to perform 2×2 subsampling. Then, a flatten layer transforms the output image of size 7×18 of the last pooling to a vector, which is then used as input to a dense layer using dropout. Finally, a second dense layer produces the output of the network.

5 Experimental Results

For the experimental evaluation of the proposed approach we used the PKU-MMD dataset [13]. This dataset aims to provide a large scale benchmark, focusing on 3D human action understanding. It contains approx. 20K action instances spanning into 5.4M video frames and belonging to 51 action categories. A total of 66 human subjects have been involved, while video recordings have been

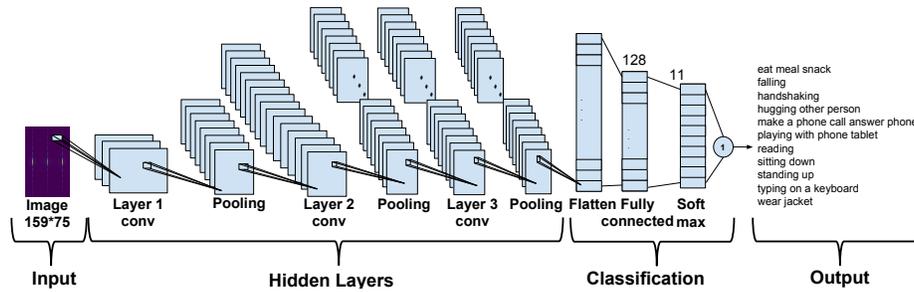


Fig. 3: The proposed CNN architecture.

captured from 3 camera angles, using the Microsoft Kinect v2 camera. Provided modalities are raw RGB video, depth sequences, infrared radiation captured by the Kinect and extracted 3D positions of skeletons.

From this dataset we decided to use 11 classes which in our opinion are the most close to ADLs or events that should be recorded at a use case of e.g., home monitoring. More specifically, these classes were: *eat meal snack*, *falling*, *handshaking*, *hugging other person*, *make a phone call*, *answer phone*, *playing with phone tablet*, *reading*, *sitting down*, *standing up*, *typing on a keyboard* and *wear jacket*. As described in section 4, we worked with the provided skeleton positions. For further evaluation, we also include experiments in all 51 classes of PKU-MMD. Sample signal and activity images from the 11 classes are illustrated in Fig. 4. Note the visual difference of these images which may not be significant, yet allows the CNN to learn the differences between two classes.

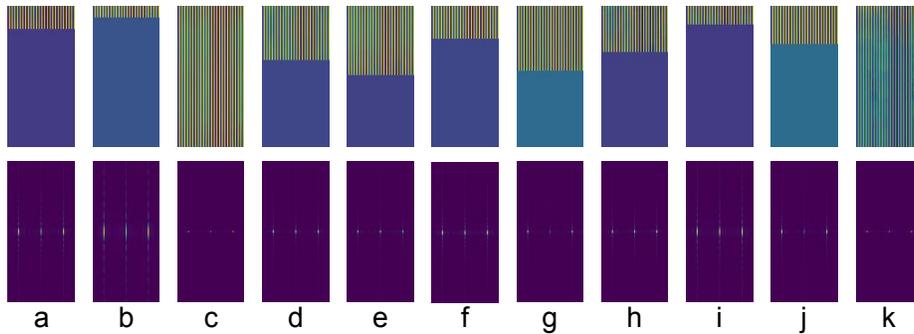


Fig. 4: Sample images from 11 actions of the PKU-MMD dataset that have been used throughout our experiments: upper row: signal images; lower row: activity images. (a) eat meal snack; (b) falling; (c) hand shaking; (d) hugging other person; (e) make answer phone call; (f) playing phone tablet; (g) reading; (h) sitting down; (i) standing up; (j) typing on keyboard; (k) wearing jacket. Figure best viewed in color.

We have set $T_s=158$ frames so as to prevent significant loss of information upon interpolation. The evaluation protocol we followed is as follows: We first performed experiments per camera position, namely *middle* (M), *left* (L) and *right* (R). In this case, both training and testing sets derived from the same position. Then, we performed cross-view experiments, where 1 position was used for training, while the other two were used for testing. The goal therein was to test the robustness of the proposed approach in abrupt changes of camera angle, which are expected to happen in real-life scenarios. In all cases we measured the accuracy of classification. Detailed results are depicted in Table 1.

As it may be observed, the proposed approach in the aforementioned case of 11 classes is able to achieve accuracy ranging from 0.75 to 0.85 when the camera angle remains unchanged (i.e., when samples from the same angle have been used for both training and testing). Also, it achieves adequate performance in case two “neighboring” angles are used, e.g., M for training, L for testing etc. In this case, accuracy ranges from 0.56 to 0.64. We noticed that when samples from the right camera position are used, a significant drop of performance is observed. Moreover, and as it has been expected, dramatic changes of camera angle, e.g., when L is used for training, R for testing or vice versa, performance ranges between 0.35 and 0.40. Finally, when the whole set of 51 classes has been used, performance is acceptable only in cases where the same angle has been used for both training and testing; corresponding accuracies range from 0.55 to 0.73. In all other cases, a strong drop of performance is observed.

For the implementation of the CNN we have used Keras [3] running on top of Tensorflow [1]. All data pre-processing and processing steps have been implemented in Python 3.6 using NumPy (<http://www.numpy.org/>) and SciPy (<https://www.scipy.org/>).

Table 1: Experimental results of the proposed approach. M, L and R denote the middle, left and right camera angles, respectively. 11 and 51 are the numbers of classes considered for evaluation. Results indicate the achieved accuracy.

Experiment	Train	M	M	M	L	L	L	R	R	R
	Test	M	L	R	M	L	R	M	R	L
Dataset	11	0.82	0.56	0.64	0.61	0.85	0.40	0.56	0.75	0.35
	51	0.73	0.29	0.28	0.25	0.55	0.11	0.29	0.73	0.12

6 Discussion

In this paper we presented a methodology for the recognition of human actions which was based on a novel image representation of 3D human skeletal information and a novel convolutional neural network architecture. We used an image representation of a human action, which resulted upon the concatenation of raw

1D signals corresponding to 3D motion of skeletal joints' coefficients and the application of the Discrete Fourier Transform to the created image.

We evaluated the proposed approach using a state-of-the-art and challenging dataset, which consisted of sequences corresponding to 51 human actions. These sequences had been captured with 3 Kinect v2 cameras, under different camera angles and the skeletal joints of the human actors involved had been extracted. We performed experiments involving either only one or two cameras (cross-view). We mainly focused on a subset of 11 actions which in our opinion are the most close to real-life ADLs. However, we also experimented with the whole dataset. Our initial results indicate that the proposed approach may be successfully applied to human action recognition in real-like conditions, yet a drop of performance is expected when camera angle would change.

Among our plans for future are the following: a) investigation on methods for creating the signal image, possibly with the use of other types of sensor measurements such as wearable accelerometers, gyroscopes etc.; b) investigation on image processing methods for transforming the signal image to the activity image. To this goal transforms such as wavelets, discrete cosine transformation (DCT) etc. may be used; c) exploitation of other types of visual modalities in the process, such as RGB and depth data; d) evaluation of the proposed approach on several other public datasets; and e) application into a real-like or even real-live assistive living environment.

Acknowledgment

We acknowledge support of this work by the project SYNTELEESIS “Innovative Technologies and Applications based on the Internet of Things (IoT) and the Cloud Computing” (MIS 5002521) which is implemented under the “Action for the Strategic Development on the Research and Technological Sector”, funded by the Operational Programme “Competitiveness, Entrepreneurship and Innovation” (NSRF 2014-2020) and co-financed by Greece and the European Union (European Regional Development Fund).

References

1. M. Abadi et al., *TensorFlow: A system for Large-Scale Machine Learning*. In Proc. of the USENIX Symposium on Operating Systems Design and Implementation (OSDI), 2016.
2. S. Berretti, M. Daoudi, P. Turaga, and A. Basu, *Representation, Analysis, and Recognition of 3D Humans: A Survey*. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), 14(1s), 16, 2018.
3. F. Chollet, *Keras*, <https://github.com/fchollet/keras>, 2015.
4. Y. Du, Y. Fu, Y. and L. Wang, *Skeleton based action recognition with convolutional neural network*. In Proc. of 3rd IAPR Asian Conference on Pattern Recognition (ACPR), IEEE, 2015.

5. Y. Hou, Z. Li, P. Wang and W. Li, *Skeleton optical spectra-based action recognition using convolutional neural networks*. IEEE Transactions on Circuits and Systems for Video Technology, 28(3), 807-811, 2018.
6. A. Graves, A.R. Mohamed and G. Hinton, *Speech recognition with deep recurrent neural networks*. In Proc. of IEEE Int'l Conf. of Acoustics, speech and signal processing (ICASSP), 2013.
7. W. Jiang and Z. Yin, *Human activity recognition using wearable sensors by deep convolutional neural networks*. In Proc. of ACM Int'l Conf. on Multimedia (MM), 2015.
8. Q. Ke, S. An, M. Bennamoun, F. Sohel and F. Boussaid, *Skeletonnet: Mining deep part features for 3-d action recognition*. IEEE signal processing letters, 24(6), 731-735, 2017.
9. A. Krizhevsky, I. Sutskever and G.E. Hinton, *Imagenet classification with deep convolutional neural networks*. Advances in neural information processing systems (pp. 1097-1105), 2012.
10. M.P. Lawton and E.M. Brody, *Assessment of older people: self-maintaining and instrumental activities of daily living*. The gerontologist, 9(3_Part.1), 179-186, 1969.
11. Y. LeCun, L. Bottou, Y. Bengio and P. Haffner, *Gradient-based learning applied to document recognition*. Proceedings of the IEEE, 86(11), pp.2278-2324, 1998.
12. C. Li, Y. Hou, P. Wang and W. Li, *Joint distance maps based action recognition with convolutional neural networks*. IEEE Signal Processing Letters, 24(5), 624-628, 2017.
13. C. Liu, Y. Hu, Y. Li, S. Song and J. Liu, *PKU-MMD: A Large Scale Benchmark for Continuous Multi-Modal Human Action Understanding*. In Proc. of ACM Multimedia Workshop (MM), 2017.
14. M. Liu, H. Liu and C. Chen, *Enhanced skeleton visualization for view invariant human action recognition*. Pattern Recognition, 68, 346-362, 2017.
15. E. Mathe, A. Mitsou, E. Spyrou and Ph. Mylonas, *Arm Gesture Recognition using a Convolutional Neural Network*. In Proc. of Int'l Workshop on Semantic and Social Media Adaptation and Personalization (SMAP) 2018.
16. N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov, *Dropout: A simple way to prevent neural networks from overfitting*. The Journal of Machine Learning Research, 15(1), pp.1929-1958, 2014.
17. P. Wang, W. Li, P. Ogunbona, J. Wan and S. Escalera, *RGB-D-based human motion recognition with deep learning: A survey*. Computer Vision and Image Understanding, 2018.
18. P. Wang, W. Li, C. Li and Y. Hou, *Action recognition based on joint trajectory maps with convolutional neural networks*. Knowledge-Based Systems, Vol. 158, pp. 43-53, 2018.
19. Z. Zhang, *Microsoft Kinect sensor and its effect*. IEEE multimedia, 19(2), pp.4-10, 2012.