

Deep Dynamic Neural Networks for Gesture Segmentation and Recognition

Di Wu and Lio

Di Wu, *Member, IEEE*, Lionel Pigou, *Member, Belgium*, and Ling Shao, *Senior Member, IEEE*

Abstract—The purpose of this paper is to describe a novel method called Deep Dynamic Neural Networks (*DDNN*) for the Track 3 of the Chalearn Looking at People 2014 challenge [1]. A generalised semi-supervised hierarchical dynamic framework is proposed for simultaneous gesture segmentation and recognition taking both skeleton and depth images as input modules. First, Deep Belief Networks (*DBN*) and 3D Convolutional Neural Networks (*3DCNN*) are adopted for skeletal and depth data accordingly to extract high level spatio-temporal features. Then the learned representations are used for estimating emission probabilities of the Hidden Markov Models to infer an action sequence. The framework can be easily extended by including an ergodic state to segment and recognise video sequences by a frame-to-frame mechanism, rendering it possible for online segmentation and recognition for diverse input modules. Some normalisation details pertaining to preprocessing raw features are also discussed. This purely data-driven approach achieves **0.8162** score in this gesture spotting challenge. The performance is on par with a variety of the state-of-the-art hand-tuned-feature approaches and other learning-based methods, opening the doors for using deep learning techniques to explore time series multimodal data.

Index Terms—Computer Society, IEEEtran, journal, LATEX, paper, template.

1 INTRODUCTION

IN recent years, human action recognition has drawn increasing attention of researchers, primarily due to its growing potential in areas such as video surveillance, robotics, human-computer interaction, user interface design, and multimedia video retrieval.

Previous works on video-based motion recognition focused on adapting handcrafted features and low-level hand-designed features [2], [3], [4] have been heavily employed with much success. These methods usually have two stages: an optional feature detection stage followed by a feature description stage. Well-known feature detection methods (“interest point detectors”) are Harris3D [5], Cuboids [6] and Hessian3D [7]. For descriptors, popular methods are Cuboids [8], HOG/HOF [5], HOG3D [9] and Extended SURF [7]. In a recent work of Wang *et al.* [10], dense trajectories with improved motion-based descriptors epitomized the pinnacle of handcrafted features and achieved state-of-the-art results on a variety of “in the wild” datasets. Given the current trends, challenges and interests in action recognition, this list would probably continue to spread out extensively. The very high-dimensional dense-trajectory features usually require advanced dimensionality reduction techniques [11], [12] to make them applicable.

In the evaluation paper of Wang *et al.* [13], one interesting finding is that there is no universally best hand-engineered feature for all datasets, suggesting that learning features directly from the dataset itself may be more advantageous.

Albeit the dominant methodology for visual recognition from images and videos relies on hand-crafted features, there has been a growing interest in methods that learn low-level and mid-level features, either in supervised, unsupervised, or semi-supervised settings [14], [15], [16].

With the recent resurgence of neural networks invoked by Hinton and others [17], deep neural architectures have been proposed as an effective solution for extracting high level features from data. Deep artificial neural networks (including the family of recurrent neural networks) have won numerous contests in pattern recognition and representation learning. Schmidhuber [18] compiled a historical survey compactly summarising relevant works with more than 850 entries of credited works. Such models have been successfully applied to a plethora of different domains: the GPU-based cuda-convnet [19] classifies 1.2 million high-resolution images into 1000 different classes; multi-column Deep Neural Networks [20] achieve near-human performance on the handwritten digits and traffic signs recognition benchmarks; 3D Convolutional Neural Networks [21], [22] recognize human actions in surveillance videos; Deep Belief Networks combining with Hidden Markov Models [23], [24] for acoustic and skeletal joints modeling outperform the decade-dominating paradigm of Gaussian Mixture Models+Hidden Markov Models. In these fields, deep architectures have shown great capacity to discover and extract higher level relevant features.

However, direct and unconstrained learning of complex problems is difficult, since (i) the amount of required training data increases steeply with the complexity of the prediction model and (ii) training highly complex models with very general learning algorithms is extremely difficult. It is therefore common practice to restrain the complexity of

- M. Shell is with the Department of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA, 30332.
E-mail: see <http://www.michaelshell.org/contact.html>
- J. Doe and J. Doe are with Anonymous University.

Manuscript received April 19, 2005; revised September 17, 2014.

the model and this is generally done by operating on small patches to reduce the input dimension and diversity [16], or by training the model in an unsupervised manner [15], or by forcing the model parameters to be identical for different input locations (as in convolutional neural networks [19], [20], [21]).

With the immense popularity of Kinect [25], [26], there has been renewed interest in developing methods for human gesture and action recognition from 3D skeletal data and depth images. A number of new datasets [27], [28], [29], [30] have provided researchers with the opportunity to design novel representations and algorithms and test them on a much larger number of sequences. It may seem that the task of action recognition given 3D joint positions is trivial, but this is not the case, largely due to the high dimensionality of the pose space. Furthermore, to achieve continuous action recognition, the sequence need to be segmented into contiguous action segments; such segmentation is as important as recognition itself and is often neglected in action recognition research.

In this paper, a data driven framework is proposed, focusing on analysis of acyclic video sequence labeling problems, *i.e.*, video sequences are non-repetitive as opposed to longer repetitive activities, *e.g.*, jogging, walking and running.

2 EXPERIMENTS AND ANALYSIS

2.1 Chalearn LAP Dataset & Evaluation Metrics

This dataset¹ is on “multiple instance, user independent learning and continuous gesture spotting” [27] of gestures. And in the 3 track, there are more than 14,000 gestures are drawn from a vocabulary of 20 Italian cultural/anthropological sign gesture categories with 700 sample sequences for training and validation and 240 sample sequences for testing.

The evaluation criteria for this track is the *Jaccard* index (overlap) on a frame-to-frame basis.

$$J(A, B) = \frac{A \cap B}{A \cup B} \quad (1)$$

2.2 Model Architecture: Deep Dynamic Neural Networks

Inspired by the framework successfully applied to the speech recognition [23], the proposed model borrows the idea of a data driven learning system, relying on a pure learning approach in which all the knowledge in the model comes from the data without sophisticated pre-processing or dimensionality reduction. The proposed Deep Dynamic Neural Networks (DDNN) can be seen as an extension to [24] in that instead of only using the Restricted Boltzmann Machines to model human motion, various connectivity layers, *e.g.*, fully connected layers, convolutional layers, *etc.*, are stacked together to learn higher level features justified by a variational bound [17] from different input modules.

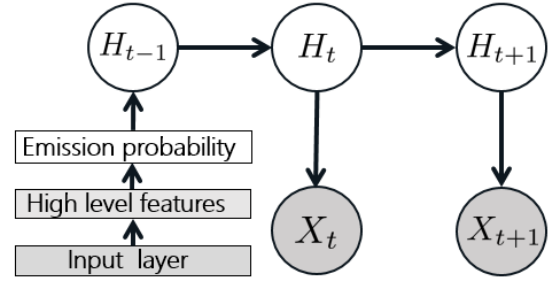


Fig. 1: Per-action model: a forward-linked chain. Inputs (skeletal features or depth image features) are first passed through Deep Neural Nets (Deep Belief Networks for skeletal modality or 3D Convolutional Neural Networks for depth modality) to extract high level features. The outputs are the emission probabilities of the hidden states.

A continuous-observation HMM with discrete hidden states is adopted for modelling higher level temporal relationships. At each time step t , we have one random observation variable X_t . Additionally we have an unobserved variable H_t taking values of a finite set $\mathcal{H} = (\bigcup_{a \in \mathcal{A}} \mathcal{H}_a)$, where \mathcal{H}_a is a set of states associated to an individual action a by force-alignment scheme defined in Sec. 2.4. The intuition motivating this construction is that an action is composed of a sequence of poses where the relative duration of each pose may vary. This variance is captured by allowing flexible forward transitions within the chain. With this definitions, the full probability model is now specified as HMM:

$$p(H_{1:T}, X_{1:T}) = p(H_1)p(X_1|H_1) \prod_{t=2}^T p(X_t|H_t)p(H_t|H_{t-1}), \quad (2)$$

where $p(H_1)$ is the prior on the first hidden state; $p(H_t|H_{t-1})$ is the transition dynamics model and $p(X_t|H_t)$ is the emission probability modelled by the deep neural nets.

The motivation for using deep neural nets to model marginal distribution is that by constructing multi-layer networks, semantically meaningful high level features will be extracted whilst learning the parametric prior of human pose from mass pool of data. In the recent work of [31], a non-parametric bayesian network is adopted for human pose prior estimation, whereas in the proposed framework, the parametric networks are incorporated. The graphical representation of a per-action model is shown as Fig. 1.

2.3 ES-HMM: Simultaneous Segmentation and Recognition

The aforementioned framework can be easily adapted for simultaneous action segmentation and recognition by adding an ergodic states- \mathcal{ES} which resembles the silence state for speech recognition. Hence, the unobserved variable H_t takes an extra finite set $\mathcal{H} = (\bigcup_{a \in \mathcal{A}} \mathcal{H}_a) \cup \mathcal{ES}$, where \mathcal{ES} is the ergodic state as the resting position between actions and we refer the model as *ES-HMM*.

Since our goal is to capture the variation in speed of performing gestures, we set the transitions in the following way: when being in a particular node n in time t , moving to time $t + 1$, we can either stay in the same node (slower

1. <http://gesture.chalearn.org/homewebsourcerefferrals>

performance), move to node $n + 1$ (the same speed of performance), or move to node $n + 2$ (faster performance). From the \mathcal{ES} we can move to the first three nodes of any gesture class, and from the last three nodes of any gesture class we can move to the \mathcal{ES} as shown in Fig. 2. The *ES-HMM* framework differs from the Firing Hidden Markov Model of [32] in that we strictly follow the temporal independent assumption, forbidding inter-states transverse, preconditioned that a non-repetitive sequence would maintain its unique states throughout its performing cycle.

The emission probability of the trained model is represented as a matrix of size $N_{\mathcal{TC}} \times N_{\mathcal{F}}$ where $N_{\mathcal{F}}$ is the number of frames in a test sequence and output target class $N_{\mathcal{TC}} = N_{\mathcal{A}} \times N_{\mathcal{H}_a} + 1$ where $N_{\mathcal{A}}$ is the number of action class and $N_{\mathcal{H}_a}$ is the number of states associated to an individual action a and one \mathcal{ES} state. Once we have the trained model, we can use the normal online or offline smoothing, inferring the hidden marginal distributions $p(H_t|X_t)$ of every node (frame) of the test video. Because the graph for the Hidden Markov Model is a directed tree, this problem can be solved exactly using the max-sum algorithm. The number of possible paths through the lattice grows exponentially with the length of the chain. The Viterbi algorithm searches this space of paths efficiently to find the most probable path with a computational cost that grows only linearly with the length of the chain [33]. We can infer the action presence in a new sequence by Viterbi decoding as:

$$V_{t,\mathcal{H}} = P(H_t|X_t) + \log(\max_{\mathcal{H} \in \mathcal{H}_a} (V_{t-1,\mathcal{H}})) \quad (3)$$

where initial state $V_{1,\mathcal{H}} = \log(P(H_1|X_1))$. From the inference results, we define the probability of an action $a \in \mathcal{A}$ as $p(y_t = a|x_{1:t}) = V_{t,\mathcal{H}}$. Result of the Viterbi algorithm is a path-sequence of nodes which correspond to hidden states of gesture classes. From this path we can infer the class of the gesture (c.f. Fig. ??). The overall algorithm for training and testing are presented in Algorithm 1 and ??.

2.4 Experimental Setups

For input sequences, there are three modalities, i.e. skeleton, RGB and depth (with user segmentation) provided. However, only skeletal modality and the depth modality are considered (c.f. Fig. 3). In the following experiments, the first 650 sample sequences are used for training, 50 for validation and the rest 240 for testing where each sequence contains around 20 gestures with some noisy non-meaningful vocabulary tokens.

Hidden States (Alignment) is used to extract the hidden states, i.e., if a gesture token is 100 frames, the first 10 frames are assigned as hidden state 1 and the 10-20 frames are assigned as hidden state 2 and so on and so forth.

Ergodic States are extracted as 5 frames before or after a gesture tokens labelled by ground truth.

2.5 Skeleton Module & DBN training

Only upper body joints are relevant to the discriminative gesture recognition tasks. Therefore, only the 11 upper body joints are considered. The 11 upper body joints used

Algorithm 1: Multimodal Deep Dynamic Networks – training

Data:

$\mathbf{X}^1 = \{\mathbf{x}_i^1\}_{i \in [1...t]}$ - raw input(skeletal) feature sequence.

$\mathbf{X}^2 = \{\mathbf{x}_i^2\}_{i \in [1...t]}$ - raw input(depth) feature sequence in the form of

$M_1 \times M_2 \times T$, where M_1, M_2 are the height and width of the input image and T is the number of contiguous frames of the spatio-temporal cuboid.

Note that the GPU library *cuda-convnet* [19] used requires square size images and T is a multiple of 4.

$\mathbf{Y} = \{\mathbf{y}_i\}_{i \in [1...t]}$ - frame based local label (achieved by semi-supervised forced-alignment),

where $\mathbf{y}_i \in \{C * S + 1\}$ with C is the number of class, S is the

number of hidden states for each class, 1 as ergodic state.

1 **for** $m \leftarrow 1$ to 2 **do**

2 **if** m is 1 **then**

3 Preprocessing the data \mathbf{X}^1 as in Eq.4.

4 Normalizing(zero mean, unit variance per dimension) the above features and feed to to Eq.6.

5 Pre-training the networks using *Contrastive Divergence*.

6 Supervised fine-tuning the Deep Belief Networks using \mathbf{Y} by standard mini-batch SGD backpropagation.

7 **else**

8 Preprocessing the data \mathbf{X}^2 (normalizing, median filtering the depth data) Algo.2 or Algo.3.

9 Feeding the above features to Eq.9.

10 Supervised fine-tuning the Deep 3D Convolutional Neural Networks using \mathbf{Y} by standard mini-batch SGD Backpropagation.

Result:

GDBN - a gaussian bernoulli visible layer Deep Belief Network to

generate the emission probabilities for hidden markov model.

3DCNN - a 3D Deep Convolutional Neural Networks to generate the

emission probabilities for hidden markov model.

$\mathbf{p}(\mathbf{H}_1)$ - prior probability for \mathbf{Y} .

$\mathbf{p}(\mathbf{H}_t|\mathbf{H}_{t-1})$ - transition probability for \mathbf{Y} , enforcing the beginning and

ending of a sequence can only start from the first or the last state.

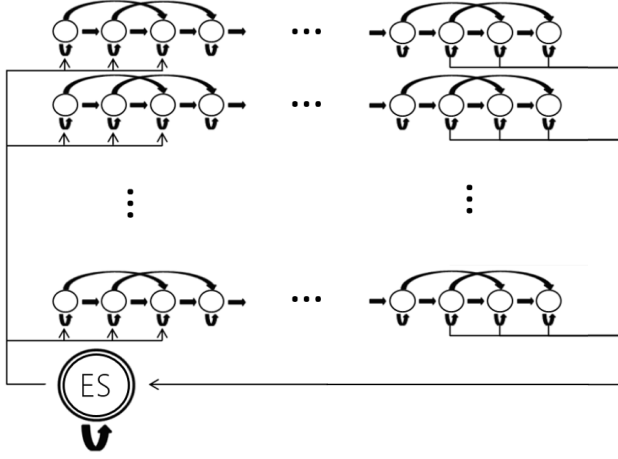


Fig. 2: State diagram of the ES-HMM model for low-latency action segmentation and recognition. An ergodic states (ES) shows the resting position between action sequence. Each node represents a single frame and each row represents a single action model. The arrows indicate possible transitions between states.

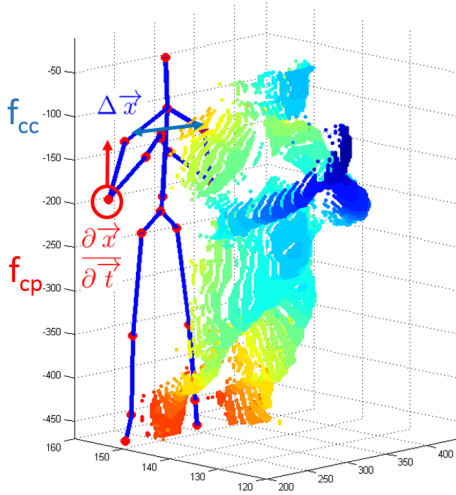


Fig. 3: Point cloud projection of depth image and the 3D positional features.

are “ElbowLeft, WristLeft, ShoulderLeft, HandLeft, ElbowRight, WristRight, ShoulderRight, HandRight, Head, Spine, HipCenter”.

The 3D coordinates of N joints of frame c are given as: $X_c = \{x_1^c, x_2^c, \dots, x_N^c\}$. 3D positional pairwise differences of joints [24] are deployed for observation domain \mathcal{X} . They capture posture features, motion features by direction concatenation: $\mathcal{X} = [f_{cc}, f_{cp}]$ as demonstrated in Eq 4. Note that offset features f_{ci} used in [24] depend on the first frame, if the initialization fails which is a very common scenario, the feature descriptor will be generally very noisy. Hence, the offset features f_{ci} are discarded and only the two more robust features $[f_{cc}, f_{cp}]$ (as shown in Fig. 3) are kept:

$$f_{cc} = \{x_i^c - x_j^c | i, j = 1, 2, \dots, N; i \neq j\} \quad (4)$$

$$f_{cp} = \{x_i^c - x_j^p | x_i^c \in X_c; x_j^p \in X_p\} \quad (5)$$

This results in a raw dimension of $N_{\mathcal{X}} = N_{joints} * (N_{joints} - 1)/2 + N_{joints}^2 * 3$ where N_{joints} is the number of joints. Therefore, in the experiment with $N_{joints} = 11$, $N_{\mathcal{X}} = 528$. Admittedly, we do not completely neglect human prior knowledge about information extraction for relevant static postures, velocity and offset overall dynamics of motion data. Nevertheless, the aforementioned three attributes are all very crude pairwise features without any tweak into the dataset or handpicking the most relevant pairwise, triple wise, etc., designed features.

A similar data driven approach has been adopted in [28] where random forest classifiers were adapted to the problem of recognizing gestures using a bundle of 35 frames. These sets of feature extraction processes resemble the *Mel Frequency Cepstral Coefficients (MFCCs)* for the speech recognition community [23].

2.5.1 Gaussian Bernoulli Restricted Boltzmann machines

Because input skeletal features (a.k.a. observation domain \mathcal{X}) are continuous instead of binomial features, we use the Gaussian RBM (GRBM) to model the energy term of first visible layer:

$$E(v, h; \theta) = - \sum_{i=1}^D \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_{i=1}^D \sum_{j=1}^F W_{ij} h_j \frac{v_i}{\sigma_i} - \sum_{j=1}^F a_j h_j \quad (6)$$

The conditional distributions needed for inference and generation are given by:

$$P(h_{j=1} | \mathbf{v}) = g(\sum_i W_{ij} v_i + a_j); \quad (7)$$

$$P(v_{i=1} | \mathbf{h}) = \mathcal{N}(v_i | \mu_i, \sigma_i^2). \quad (8)$$

where $\mu_i = b_i + \sigma_i^2 \sum_j W_{ij} h_j$ and \mathcal{N} is normal distribution. In general, we normalize the data (mean subtraction and standard deviation division) in the preprocessing phase. Hence, in practice, instead of learning σ_i^2 , one would typically use a fixed, predetermined unit value $\mathbf{1}$ for σ_i^2 .

For high level skeleton feature extraction, two network architectures, i.e., a smaller one and a larger one were experimented: $[N_{\mathcal{X}}, 1000, 1000, 500, N_{\mathcal{T}C}]$ and $[N_{\mathcal{X}}, 2000, 2000, 1000, N_{\mathcal{T}C}]$, where $N_{\mathcal{X}} = 528$ is the observation domain dimension; $N_{\mathcal{T}C}$ is the output target class. Because in all our experiments the number of states associated to an individual action $N_{\mathcal{H}_a}$ is chosen as 10 for modeling the states of an action class, therefore $N_{\mathcal{T}C} = 20 + 1 = 201$.

In the training set, there are in total 400,117 frames. During the training of DBN, 90% is used for training, 8% for validation (for the purpose of early stopping) 2% is used for test evaluation. The feed forward networks are pre-trained with a fixed recipe using stochastic gradient decent with a mini-batch size of 200 training cases. Unsupervised initializations tend to avoid local minima and increase the networks performance stability and we have run 100 epochs for unsupervised pre-training. For Gaussian-binary RBMs, learning rate is fixed at 0.001 while for binary-binary RBMs as 0.01 (note in generally training GRBM requires smaller learning rate). For fine-tuning, the learning rate starts at 1 with 0.99999 mini-batch scaling. Maximum number of fine-tuning epoch is 500 with early stopping strategy and during the experiments, early stopping occurs around 440

epoch. Optimization complete with best validation score (the frame based prediction error rate) of 38.19%, with test performance 38.11%.

Though we believe further carefully choosing network architecture would lead to more competitive results, in order not to “creeping overfitting”, as algorithms over time become too adapted to the dataset, essentially memorizing all its idiosyncrasies, and losing ability to generalize [34], we would like to treat the model as the aforementioned more generic approach. Since a completely new approach will initially have a hard time competing against established, carefully fine-tuned methods. More fundamentally, it may be that the right way to treat dataset performance numbers is not as a competition for the top place. This way, fundamentally new approaches will not be forced to compete for top performance right away, but will have a chance to develop and mature.

The performance of skeleton module is shown in Tab 1. And it can be seen that larger net (Net2) will generally perform better than smaller net (Net1), averaging multi-column nets almost will certainly further improve the performance [20]. Hence, in the following experiments, only the multi-column averaging results are reported.

2.6 Depth 3D Module

2.6.1 Preprocessing & Normalizing: shifting, scaling and resizing.

Working directly with raw input Kinect recorded data frames, which are 480×640 pixel images, can be computationally demanding. Deepmind technology [35] presented the first deep learning model to successfully learn control policies directly from high-dimensional sensory input using deep reinforcement learning.

Our first step in the preprocessing stage is cropping the highest hand and the upper body using the given joint information. We discovered that the highest hand is the most interesting. If both hands are used, they perform the same (mirrored) movement. If one hand is used, it is always the highest one. If the left hand is used, the videos are mirrored. This way, the model only needs to learn one side.

The preprocessing results in four video samples (body and hand with gray-scale and depth) of resolution $4 \times 64 \times 64$ (4 frames of size 64×64). Furthermore, the noise in the depth maps is reduced with thresholding, background removal using the user index, and median filtering. The outcome is shown in Figure 4.

Similarly, basic preprocessing steps are adopted aimed at reducing the input dimensionality from the original 480×640 pixels to 90×90 pixels. The square-sized of the final image is required because the used GPU implementation from [19] expects square inputs and the input channel should be in the set of $[1, 3, 4x]$. There are two normalization schemes implemented as 2 using depth template matching method and 3 using skeletal joints as assistant normalization. Note that the scheme 2 depends heavily on the provided maximum depth from the recording scene and scheme 3 depends on the accurate detection of skeleton joints, and both scheme require the performer remains a roughly static position (though the max pooling scheme in 3DCNN to some extent overcome the problem of position



Fig. 4: Preprocessing.

Algorithm 2: Normalization scheme 1: template matching

Data:

\mathbf{T} - exemplary template with original scale of size 320×320 ,

(Sample0003 is chosen as the exemplary template, shown in 5a).

$\mathbf{R}_{\text{depth}}$ - reference depth, fixed to **1941** (acquired from the above

exemplary template \mathbf{T}).

$\hat{\mathbf{T}}$ - test image, as shown in 5b.

\mathbf{M} - user foreground segmented mask.

- 1 Apply a 5×5 aperture median filter to test depth frame $\hat{\mathbf{T}}$ as in [36] to reduce the salt and pepper noise.
- 2 Multiply test depth frame $\hat{\mathbf{T}}$ with the user segmented mask \mathbf{M} : $\hat{\mathbf{T}} = \hat{\mathbf{T}} \times \mathbf{M}$.
- 3 Template matching test image $\hat{\mathbf{T}}$ with \mathbf{T} using normalized cross-correlation [37], the response score \mathbf{R} is shown in 5c.
- 4 Shift the image according to the maximum response \mathbf{R} to its centre applying affine transformation [38].
- 5 Scale the image according to reference depth $\mathbf{R}_{\text{depth}}$ and the median depth of a bounding box in the centre of the image with 25×25 size as shown as the green bounding box in 5d.
- 6 Resize the image from 320×320 to 90×90 .

Result:

$\hat{\mathbf{T}}$ - Resize-normalized image shown in the yellow bounding box of 5d.

shifting). Generally, scheme 3 is more robust than scheme 2 because the provided maximum depth can sometimes be very noisy, *e.g.*, Sample0671, Sample0692, Sample0699, *etc.* After the normalization and resizing, a cuboid of 4 frames, hence, size $90 \times 90 \times 4$, is extracted as a spatio-temporal unit.

2.6.2 3DCNN Architecture & Details of Learning.

The 3D convolution is achieved by convolving a 3D kernel to the cuboid formed by stacking multiple contiguous

Algorithm 3: Normalization scheme 2: skeleton normalization**Data:**

S_{spine} - Skeleton Spine joints pixel coordinates.
 S_{shoulder} - Skeleton Shoulder joints pixel coordinates.
 \hat{T} - test image.
 M - user foreground segmented mask.
 R_{length} - reference length of shoulder to spine, fixed to 100 (1 meter).

- 1 Apply a 5×5 aperture median filter to test depth frame \hat{T} .
 - 2 Multiply test depth frame \hat{T} with the user segmented mask M .
 - 3 Shift the image according to the centroid of Spine joint S_{spine} .
 - 4 Scale the image according to the $R_{\text{length}} / (S_{\text{spine}} - S_{\text{shoulder}})$.
- Result:**
 \tilde{T} - Resize the shifted-scaledp image to 90×90 .

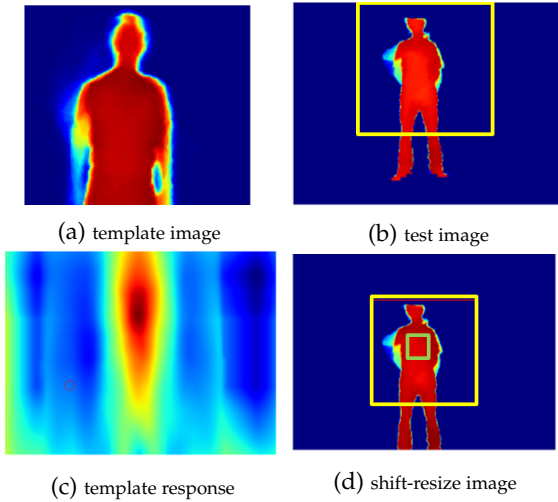


Fig. 5: Illustration of normalization scheme 1: template matching.

frames together. We follow the nomenclature as in [22]. However, instead of using *tanh* unit as in [22], the Rectified Linear Units (*ReLU*s) [19] were adopted where trainings are several times faster than their equivalents with *tanh* units. Formally, the value of a unit at position (x, y, z) (z here corresponds the time-axis) in the j th feature map in the i th layer, denoted as v_{ij}^{xyz} , is given by:

$$v_{ij}^{xyz} = \max(0, (b_{ij} + \sum_m \sum_{p=0}^{P_i-1} \sum_{q=0}^{Q_i-1} \sum_{r=0}^{R_i-1} w_{ijm}^{pqr} v_{(i-1)m}^{(x+p)(y+q)(t+r)})) \quad (9)$$

The 3DCNN architecture is depicted as Fig. 6: the input contextual frames are stacked as size $90 \times 90 \times 4$ subtracting the mean activity over the training set from each pixel, the first layer contains 16 maps of $7 \times 7 \times 4$ 3D kernel followed by local response normalization layer [19] and stride 2 max pooling; the second convolutional layers has 32 maps of 5×5 kernel followed by local response normalization layer and stride 2 max pooling; the third convolution layer is

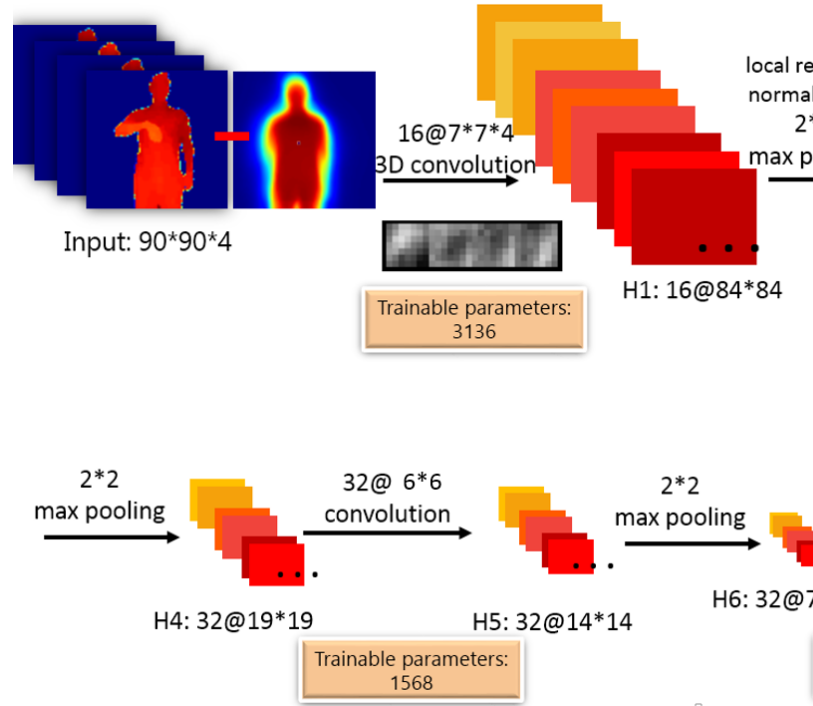


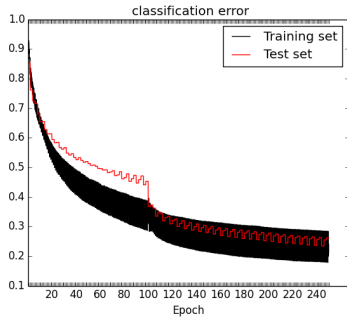
Fig. 6: An illustration of the architecture of the 3DCNN architecture.

composed of 32 maps of 6×6 kernel followed by max pooling; then we have one fully connected layer of size 1000; the output layer N_{TC} is of size $201 = 10 \times 20 + 1$ (number of hidden states for each class \times number of classes plus one ergodic state).

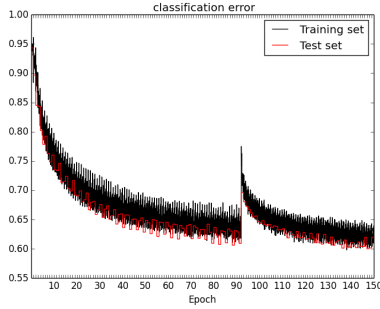
The training set is roughly of 400,000 frames and is divided into 33 mini-batches with first 30 batches for training and the rest 3 batches for validation. Standard *SGD* is run for the first 100 epochs with learning rate of 0.1 and the weight learning rate as 0.001 and weight bias learning rate 0.002 both momentum are fixed as 0.9, weight decay is fixed to 0.0005, the next 100 epochs with $0.1 \times$ learning rate. Another network trained by randomly cropping 82×82 pixels on the flight as [19] is also implemented to enhance the model's robustness. During the test time, the centre part and other 4 corner parts are averaged to obtain the final score, *c.f.* Fig 7c. Due to the time constraint, only 150 epochs are trained with the learning rate reduced to one tenth at the 92th epoch. The training frame based classification error for the aforementioned two networks are shown in 7a and 7b. One interesting observation is that for the network with uncropped input, reducing the learning rate at 100 epoch, the frame-based classification rate reduces drastically whereas for the network with cropped input, reducing the learning rate results in a spike increase of frame-based classification error rate. The reason for this discrepancy worths further investigation. The performance of depth module is shown in Tab 1.

2.6.3 Looking into the networks-visualization of filter banks.

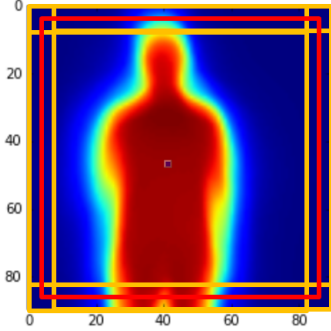
The weight filters of the first *conv1* layer are illustrated in Fig 8 and it can be seen that both shape pattern filters and motion filters are learnt effectively and the filters/weights of the cropped input trained networks are smoother than the



(a) Frame based classification error with uncropped input.



(b) Frame based classification error with cropped input.



(c) Cropped images to enhance model's robustness.

Fig. 7: Visualization of the first filters and training statistics for 3DCNN.

uncropped one. Interestingly, the 3DCNN is able to learn the most informative motion part of the body effectively (highest response parts are the arms/hands areas), albeit no signal was explicitly given during training instructing which body parts the gesture recognition tasks should focus on.

2.7 Post-Processing

The predicted token less than 20 frames are discarded as noisy tokens. Note that there are many noisy gesture tokens predicted by viterbi decoding. One way to sift through the noisy tokens is to discard the token path log probability small than certain threshold. However, because the metric of this challenge: *Jaccard index* strongly penalizes false negatives, experiments show that it's better to have more false positives than to miss true positives. Effective ways

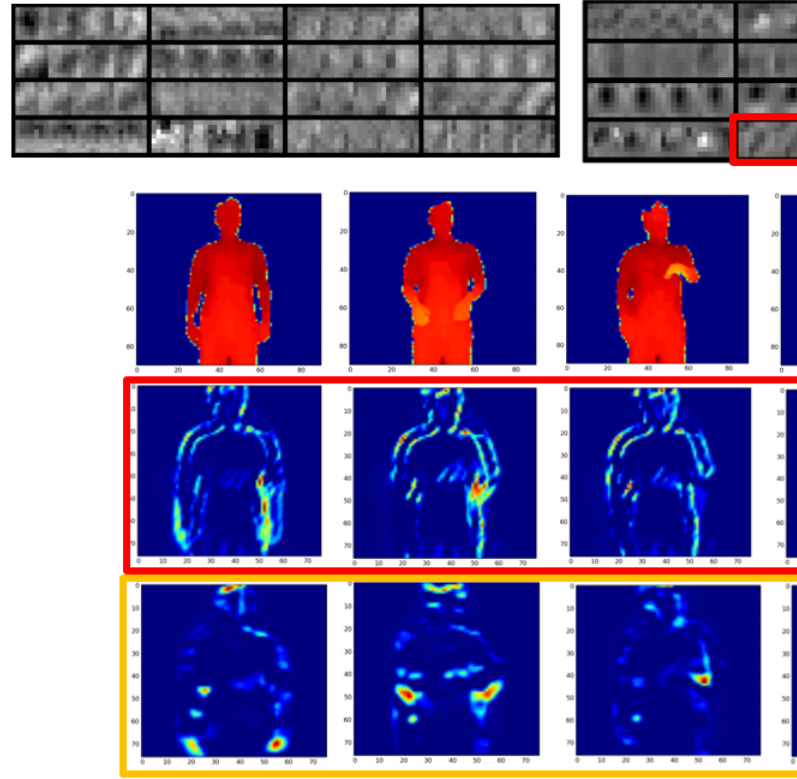


Fig. 8: Top left: the *conv1* weights of the 3DCNN learnt with uncropped input; top right: the *conv1* weights of the 3DCNN learnt with cropped input. It can be seen that filters/weights of the cropped input trained networks are smoother. Bottom: visualization of sample frames after *conv1* layer (Sample0654, 264-296 frames, sampled every 8 frames). It can be seen that the filters of the first convolutional layer are able to learn both shape pattern (red bounding box) and motion (yellow bounding box). Also note that the high response maps correspond to the most informative part of the body, even though during the training process, all local patches are learned indiscriminately regardless of its location.

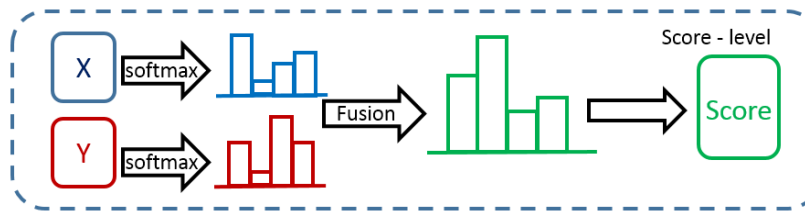


Fig. 9: Illustration of descriptor fusion.

to detect false positives should be an interesting aspect of future works.

2.8 Score Fusion

To fuse the dual model prediction, the strategy shown as Fig 9 is adopted. The complementary properties of both modules can be seen from the Viterbi path decoding plot in Fig ?? . Note that the skeleton module generally performs better than the depth module, one reason could be that the skeleton joints learnt from [25] lie in success of utilizing huge and highly varied training data: from both realistic

Module \ Evaluation Set	Validation	Test
Skeleton-DBDN Net1	0.7468	-
Skeleton-DBDN Net2	0.8017	-
Skeleton-DBDN MultiNet	0.8236	0.7873
Depth-3DCNN Norm1 2	0.6378	-
Depth-3DCNN Norm2 3	0.6924	0.6371
Score Fusion	0.8045	0.8162

TABLE 1: Comparison of results in terms of Jaccard index between different network structures and various modules. DBDN Net1 corresponds to network structure of [528, 1000, 1000, 500, 201] and DBDN Net2 [528, 2000, 2000, 1000, 201], DBDN MultiNet is the average of 3 Nets (2 Net1 and 1 Net2 with different initializations). It can be seen that larger net has better performance and multi-column net will further improve the classification rate. Norm1 corresponds to the normalization scheme 2 and Norm2 corresponds to the scheme 3.

and synthetic depth images, a total number of 1 million images were used to train the deep randomized decision forest classifier in order to avoid overfitting. Hence skeleton data are more robust.

2.9 Computational complexity

Though learning the Deep Neural Networks using stochastic gradient descent is tediously lengthy, once the model finishes training, with a low inference cost, our framework can perform real-time video sequence labeling. Specifically, a single feed forward neural network incurs trivial computational time ($\mathcal{O}(T)$) and is fast because it requires only matrix products and convolution operations. The complexity of Viterbi algorithm is $\mathcal{O}(T * |S|^2)$ with number of frames T and state number S .

3 CONCLUSION AND DISCUSSION

Hand-engineered, task-specific features are often less adaptive and time-consuming to design. This difficulty is more pronounced with multimodal data as the features have to relate multiple data sources. In this paper, we presented a novel Deep Dynamic Neural Networks (DDNN) framework that utilizes Deep Belief Networks and 3D Convolutional Neural Networks for learning contextual frame-level representations and modeling emission probabilities for Markov Field. The heterogeneous inputs from skeletal joints and depth images require different feature learning methods and the late fusion scheme is adopted at the score level. The experimental results on bi-modal time series data show that the multimodal DDNN framework can learn a good model of the joint space of multiple sensory inputs, and is consistently as good as/better than the unimodal input, opening the door for exploring the complementary representation among multimodal inputs. It also suggests that learning features directly from data is a very important research direction and with more and more data and flops-free computational power, the learning-based methods are not only more generalizable to many domains, but also are powerful in combining with other well-studied probabilistic graphical models for modeling and reasoning dynamic sequences. Future works include learning the share representation amongst the heterogeneous inputs at the penultimate

layer and backpropagating the gradient in the share space in a unified representation.

APPENDIX A

PROOF OF THE FIRST ZONKLAR EQUATION

Appendix one text goes here.

APPENDIX B

Appendix two text goes here.

ACKNOWLEDGMENTS

The authors would like to thank...

sed to reserve space for the reference number labels box)

REFERENCES

- [1] S. Escalera, X. Bar, J. Gonzalez, M. Bautista, M. Madadi, M. Reyes, V. Ponce, H. Escalante, J. Shotton, and I. Guyon, "Chalearn looking at people challenge 2014: Dataset and results," in *European Conference on Computer Vision workshop*, 2014.
- [2] L. Liu, L. Shao, F. Zheng, and X. Li, "Realistic action recognition via sparsely-constructed gaussian processes," *Pattern Recognition*, doi: 10.1016/j.patcog.2014.07.006., 2014.
- [3] L. Shao, X. Zhen, and X. Li, "Spatio-temporal laplacian pyramid coding for action recognition," *IEEE Transactions on Cybernetics*, vol. 44, no. 6, pp. 817-827, 2014.
- [4] D. Wu and L. Shao, "Silhouette analysis-based action recognition via exploiting human poses," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 2, pp. 236-243, 2013.
- [5] I. Laptev, "On space-time interest points," *International Journal of Computer Vision*, 2005.
- [6] P. Dollár, V. Rabaud, G. Cottrell, and S. Belongie, "Behavior recognition via sparse spatio-temporal features," in *Visual Surveillance and Performance Evaluation of Tracking and Surveillance*. IEEE, 2005.
- [7] G. Willems, T. Tuytelaars, and L. V. Gool, "An efficient dense and scale-invariant spatio-temporal interest point detector," in *European Conference on Computer Vision*. Springer, 2008.
- [8] P. Scovanner, S. Ali, and M. Shah, "A 3-dimensional sift descriptor and its application to action recognition," in *International Conference on Multimedia*. ACM, 2007.
- [9] A. Klaser, M. Marszalek, and C. Schmid, "A Spatio-Temporal Descriptor Based on 3D-Gradients," in *British Machine Vision Conference*, 2008.
- [10] H. Wang, A. Klaser, C. Schmid, and C.-L. Liu, "Dense trajectories and motion boundary descriptors for action recognition," *International Journal of Computer Vision*, 2013.
- [11] T. Zhou and D. Tao, "Double shrinking sparse dimension reduction," *IEEE Transactions on Image Processing* 22(1): 244-257, 2013.
- [12] C. Xu and D. Tao, "Large-margin multi-view information bottleneck," *IEEE Trans. Pattern Anal. Mach. Intell.* 36(8): 1559-1572, 2014.
- [13] H. Wang, M. M. Ullah, A. Klaser, I. Laptev, C. Schmid et al., "Evaluation of local spatio-temporal features for action recognition," in *British Machine Vision Conference*, 2009.
- [14] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler, "Convolutional learning of spatio-temporal features," in *European Conference on Computer Vision*. Springer, 2010.
- [15] Q. V. Le, W. Y. Zou, S. Y. Yeung, and A. Y. Ng, "Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2011.
- [16] M. Baccouche, F. Mamalet, C. Wolf, C. Garcia, and A. Baskurt, "Spatio-temporal convolutional sparse auto-encoder for sequence classification," in *British Machine Vision Conference*, 2012.
- [17] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, 2006.
- [18] J. Schmidhuber, "Deep learning in neural networks: An overview," *arXiv preprint arXiv:1404.7828*, 2014.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Neural Information Processing Systems*, 2012.

- [20] D. Cireşan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
- [21] M. Y. Shuiwang Ji, Wei Xu and K. Yu, "3d convolutional neural networks for human action recognition," in *International Conference on Machine Learning*. IEEE, 2010.
- [22] S. Ji, W. Xu, M. Yang, and K. Yu, "3d convolutional neural networks for human action recognition," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 2013.
- [23] A. Mohamed, G. E. Dahl, and G. Hinton, "Acoustic modeling using deep belief networks," *Audio, Speech, and Language Processing, IEEE Transactions on*, 2012.
- [24] D. Wu and L. Shao, "Leveraging hierarchical parametric networks for skeletal joints based action segmentation and recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [25] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2011.
- [26] J. Han, L. Shao, and J. Shotton, "Enhanced computer vision with microsoft kinect sensor: A review," *IEEE Transactions on Cybernetics*, vol. 43, no. 5, pp. 1317-1333, 2013.
- [27] S. Escalera, J. Gonzalez, X. Bar, M. Reyes, O. Lops, I. Guyon, V. Athitsos, and H. J. Escalante, "Multi-modal gesture recognition challenge 2013: Dataset and results." in *ACM ChaLearn Multi-Modal Gesture Recognition Grand Challenge and Workshop*, 2013. [Online]. Available: <http://gesture.chalearn.org/>
- [28] S. Fothergill, H. M. Mentis, P. Kohli, and S. Nowozin, "Instructing people for training gestural interactive systems," in *ACM Computer Human Interaction*, 2012.
- [29] I. Guyon, V. Athitsos, P. Jangyodsuk, B. Hamner, and H. J. Escalante, "Chalearn gesture challenge: Design and first results," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2012.
- [30] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Mining actionlet ensemble for action recognition with depth cameras," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
- [31] A. Lehrmann, P. Gehler, and S. Nowozin, "A non-parametric bayesian network prior of human pose," in *International Conference on Computer Vision*, 2013.
- [32] S. Nowozin and J. Shotton, "Action points: A representation for low-latency online human action recognition," Tech. Rep., 2012.
- [33] C. Bishop, *Pattern recognition and machine learning*. Springer, 2006.
- [34] A. Torralba and A. A. Efros, "Unbiased look at dataset bias," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2011.
- [35] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [36] D. Wu, F. Zhu, and L. Shao, "One shot learning gesture recognition from rgb-d images," in *International Conference on Computer Vision and Pattern Recognition Workshops*, 2012.
- [37] J. Lewis, "Fast normalized cross-correlation," in *Vision interface*, vol. 10, no. 1, 1995, pp. 120-123.
- [38] G. Bradski, *Dr. Dobb's Journal of Software Tools*.

Jane Doe Biography text here.



Michael Shell Biography text here.

John Doe Biography text here.