

Skeleton-based action recognition using spatio-temporal lstm network with trust gates

Liu, Jun; Shahroudy, Amir; Xu, Dong; Kot, Alex Chichung; Wang, Gang

2018

Liu, J., Shahroudy, A., Xu, D., Kot, A. C., & Wang, G. (2018). Skeleton-based action recognition using spatio-temporal lstm network with trust gates. IEEE transactions on pattern analysis and machine intelligence, 40(12), 3007-3021.

doi:10.1109/TPAMI.2017.2771306

<https://hdl.handle.net/10356/136885>

<https://doi.org/10.1109/TPAMI.2017.2771306>

© 2018 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. The published version is available at:
<https://doi.org/10.1109/TPAMI.2017.2771306>.

Downloaded on 20 Mar 2024 20:19:53 SGT

Skeleton-Based Action Recognition Using Spatio-Temporal LSTM Network with Trust Gates

Jun Liu, Amir Shahroudy, Dong Xu, Alex C. Kot, and Gang Wang

Abstract—Skeleton-based human action recognition has attracted a lot of research attention during the past few years. Recent works attempted to utilize recurrent neural networks to model the temporal dependencies between the 3D positional configurations of human body joints for better analysis of human activities in the skeletal data. The proposed work extends this idea to spatial domain as well as temporal domain to better analyze the hidden sources of action-related information within the human skeleton sequences in both of these domains simultaneously. Based on the pictorial structure of Kinect’s skeletal data, an effective tree-structure based traversal framework is also proposed. In order to deal with the noise in the skeletal data, a new gating mechanism within LSTM module is introduced, with which the network can learn the reliability of the sequential data and accordingly adjust the effect of the input data on the updating procedure of the long-term context representation stored in the unit’s memory cell. Moreover, we introduce a novel multi-modal feature fusion strategy within the LSTM unit in this paper. The comprehensive experimental results on seven challenging benchmark datasets for human action recognition demonstrate the effectiveness of the proposed method.

Index Terms—Action Recognition, Recurrent Neural Networks, Long Short-Term Memory, Spatio-Temporal Analysis, Tree Traversal, Trust Gate, Skeleton Sequence.

1. Introduction

Human action recognition is a fast developing research area due to its wide applications in intelligent surveillance, human-computer interaction, robotics, and so on. In recent years, human activity analysis based on human skeletal data has attracted a lot of attention, and various methods for feature extraction and classifier learning have been developed for skeleton-based action recognition [1], [2], [3]. A hidden Markov model (HMM) is utilized by Xia *et al.* [4] to model the temporal dynamics over a histogram-based

representation of joint positions for action recognition. The static postures and dynamics of the motion patterns are represented via eigenjoints by Yang and Tian [5]. A Naive-Bayes-Nearest-Neighbor classifier learning approach is also used by [5]. Vemulapalli *et al.* [6] represent the skeleton configurations and action patterns as points and curves in a Lie group, and then a SVM classifier is adopted to classify the actions. Evangelidis *et al.* [7] propose to learn a GMM over the Fisher kernel representation of the skeletal quads feature. An angular body configuration representation over the tree-structured set of joints is proposed in [8]. A skeleton-based dictionary learning method using geometry constraint and group sparsity is also introduced in [9].

Recently, recurrent neural networks (RNNs) which can handle the sequential data with variable lengths [10], [11], have shown their strength in language modeling [12], [13], [14], image captioning [15], [16], video analysis [17], [18], [19], [20], [21], [22], [23], [24], [25], and RGB-based activity recognition [26], [27], [28], [29]. Applications of these networks have also shown promising achievements in skeleton-based action recognition [30], [31], [32].

In the current skeleton-based action recognition literature, RNNs are mainly used to model the long-term context information across the temporal dimension by representing motion-based dynamics. However, there is often strong dependency relations among the skeletal joints in spatial domain also, and the spatial dependency structure is usually discriminative for action classification.

To model the dynamics and dependency relations in both temporal and spatial domains, we propose a spatio-temporal long short-term memory (ST-LSTM) network in this paper. In our ST-LSTM network, each joint can receive context information from its stored data from previous frames and also from the neighboring joints at the same time frame to represent its incoming spatio-temporal context. Feeding a simple chain of joints to a sequence learner limits the performance of the network, as the human skeletal joints are not semantically arranged as a chain. Instead, the adjacency configuration of the joints in the skeletal data can be better represented by a tree structure. Consequently, we propose a traversal procedure by following the tree structure of the skeleton to exploit the kinematic relationship among the body joints for better modeling spatial dependencies.

Since the 3D positions of skeletal joints provided by depth sensors are not always very accurate, we further introduce a new gating framework, so called “trust gate”, for our ST-LSTM network to analyze the reliability of the input data

- J. Liu, A. Shahroudy, and A. C. Kot are with School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore. E-mail: {jliu029, amir3, eackot}@ntu.edu.sg.
- G. Wang is with Alibaba Group, Hangzhou, 310052, China. E-mail: wanggang@ntu.edu.sg.
- D. Xu is with School of Electrical and Information Engineering, University of Sydney, Sydney, NSW 2006, Australia. E-mail: dong.xu@sydney.edu.au.

at each spatio-temporal step. The proposed trust gate gives better insight to the ST-LSTM network about when and how to update, forget, or remember the internal memory content as the representation of the long-term context information.

In addition, we introduce a feature fusion method within the ST-LSTM unit to better exploit the multi-modal features extracted for each joint.

We summarize the main contributions of this paper as follows. (1) A novel spatio-temporal LSTM (ST-LSTM) network for skeleton-based action recognition is designed. (2) A tree traversal technique is proposed to feed the structured human skeletal data into a sequential LSTM network. (3) The functionality of the ST-LSTM framework is further extended by adding the proposed “trust gate”. (4) A multi-modal feature fusion strategy within the ST-LSTM unit is introduced. (5) The proposed method achieves state-of-the-art performance on seven benchmark datasets.

The remainder of this paper is organized as follows. In section 2, we introduce the related works on skeleton-based action recognition, which used recurrent neural networks to model the temporal dynamics. In section 3, we introduce our end-to-end trainable spatio-temporal recurrent neural network for action recognition. The experiments are presented in section 4. Finally, the paper is concluded in section 5.

2. Related Work

Skeleton-based action recognition has been explored in different aspects during recent years [33], [34], [35], [36], [37], [38], [39], [40], [41], [42], [43], [44], [45], [46], [47]. In this section, we limit our review to more recent approaches which use RNNs or LSTMs for human activity analysis.

Du *et al.* [30] proposed a Hierarchical RNN network by utilizing multiple bidirectional RNNs in a novel hierarchical fashion. The human skeletal structure was divided to five major joint groups. Then each group was fed into the corresponding bidirectional RNN. The outputs of the RNNs were concatenated to represent the upper body and lower body, then each was further fed into another set of RNNs. By concatenating the outputs of two RNNs, the global body representation was obtained, which was fed to the next RNN layer. Finally, a softmax classifier was used in [30] to perform action classification.

Veeriah *et al.* [31] proposed to add a new gating mechanism for LSTM to model the derivatives of the memory states and explore the salient action patterns. In this method, all of the input features were concatenated at each frame and were fed to the differential LSTM at each step.

Zhu *et al.* [48] introduced a regularization term to the objective function of the LSTM network to push the entire framework towards learning co-occurrence relations among the joints for action recognition. An internal dropout [49] technique within the LSTM unit was also introduced in [48].

Shahrudy *et al.* [32] proposed to split the LSTM’s memory cell to sub-cells to push the network towards learning the context representations for each body part separately.

The output of the network was learned by concatenating the multiple memory sub-cells.

Harvey and Pal [50] adopted an encoder-decoder recurrent network to reconstruct the skeleton sequence and perform action classification at the same time. Their model showed promising results on motion capture sequences.

Mahasseni and Todorovic [51] proposed to use LSTM to encode a skeleton sequence as a feature vector. At each step, the input of the LSTM consists of the concatenation of the skeletal joints’ 3D locations in a frame. They further constructed a feature manifold by using a set of encoded feature vectors. Finally, the manifold was used to assist and regularize the supervised learning of another LSTM for RGB video based action recognition.

Different from the aforementioned works, our proposed method does not simply concatenate the joint-based input features to build the body-level feature representation. Instead, the dependencies between the skeletal joints are explicitly modeled by applying recurrent analysis over temporal and spatial dimensions concurrently. Furthermore, a novel trust gate is introduced to make our ST-LSTM network more reliable against the noisy input data.

This paper is an extension of our preliminary conference version [52]. In [52], we validated the effectiveness of our model on four benchmark datasets. In this paper, we extensively evaluate our model on seven challenging datasets. Besides, we further propose an effective feature fusion strategy inside the ST-LSTM unit. In order to improve the learning ability of our ST-LSTM network, a last-to-first link scheme is also introduced. In addition, we provide more empirical analysis of the proposed framework.

3. Spatio-Temporal Recurrent Networks

In a generic skeleton-based action recognition problem, the input observations are limited to the 3D locations of the major body joints at each frame. Recurrent neural networks have been successfully applied to this problem recently [30], [32], [48]. LSTM networks [53] are among the most successful extensions of recurrent neural networks. A gating mechanism controlling the contents of an internal memory cell is adopted by the LSTM model to learn a better and more complex representation of long-term dependencies in the input sequential data. Consequently, LSTM networks are very suitable for feature learning over time series data (such as human skeletal sequences over time).

We will briefly review the original LSTM model in this section, and then introduce our ST-LSTM network and the tree-structure based traversal approach. We will also introduce a new gating mechanism for ST-LSTM to handle the noisy measurements in the input data for better action recognition. Finally, an internal feature fusion strategy for ST-LSTM will be proposed.

3.1. Temporal Modeling with LSTM

In the standard LSTM model, each recurrent unit contains an input gate i_t , a forget gate f_t , an output gate o_t , and

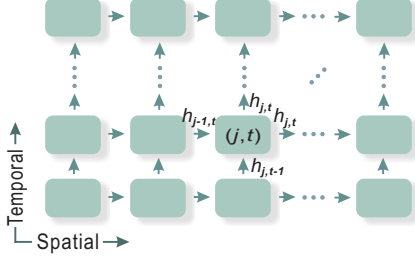


Figure 1. Illustration of the spatio-temporal LSTM network. In temporal dimension, the corresponding body joints are fed over the frames. In spatial dimension, the skeletal joints in each frame are fed as a sequence. Each unit receives the hidden representation of the previous joints and the same joint from previous frames.

an internal memory cell state c_t , together with a hidden state h_t . The input gate i_t controls the contributions of the newly arrived input data at time step t for updating the memory cell, while the forget gate f_t determines how much the contents of the previous state (c_{t-1}) contribute to deriving the current state (c_t). The output gate o_t learns how the output of the LSTM unit at current time step should be derived from the current state of the internal memory cell. These gates and states can be obtained as follows:

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ u_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left(M \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \right) \quad (1)$$

$$c_t = i_t \odot u_t + f_t \odot c_{t-1} \quad (2)$$

$$h_t = o_t \odot \tanh(c_t) \quad (3)$$

where x_t is the input at time step t , u_t is the modulated input, \odot denotes the element-wise product, and $M : \mathbb{R}^{D+d} \rightarrow \mathbb{R}^{4d}$ is an affine transformation. d is the size of the internal memory cell, and D is the dimension of x_t .

3.2. Spatio-Temporal LSTM

RNNs have already shown their strengths in modeling the complex dynamics of human activities as time series data, and achieved promising performance in skeleton-based human action recognition [30], [31], [32], [48]. In the existing literature, RNNs are mainly utilized in temporal domain to discover the discriminative dynamics and motion patterns for action recognition. However, there is also discriminative spatial information encoded in the joints' locations and posture configurations at each video frame, and the sequential nature of the body joints makes it possible to apply RNN-based modeling to spatial domain as well.

Different from the existing methods which concatenate the joints' information as the entire body's representation, we extend the recurrent analysis to spatial domain by discovering the spatial dependency patterns among different body joints. We propose a spatio-temporal LSTM (ST-LSTM) network to simultaneously model the temporal dependencies among different frames and also the spatial dependencies

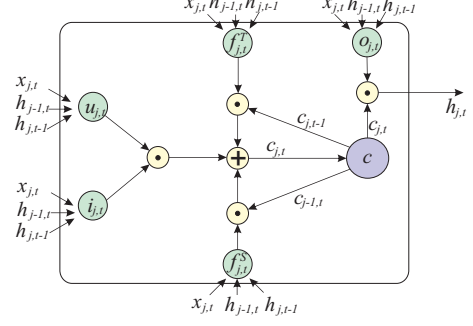


Figure 2. Illustration of the proposed ST-LSTM with one unit.

of different joints at the same frame. Each ST-LSTM unit, which corresponds to one of the body joints, receives the hidden representation of its own joint from the previous time step and also the hidden representation of its previous joint at the current frame. A schema of this model is illustrated in Figure 1.

In this section, we assume the joints are arranged in a simple chain sequence, and the order is depicted in Figure 3(a). In section 3.3, we will introduce a more advanced traversal scheme to take advantage of the adjacency structure among the skeletal joints.

We use j and t to respectively denote the indices of joints and frames, where $j \in \{1, \dots, J\}$ and $t \in \{1, \dots, T\}$. Each ST-LSTM unit is fed with the input ($x_{j,t}$, the information of the corresponding joint at current time step), the hidden representation of the previous joint at current time step ($h_{j-1,t}$), and the hidden representation of the same joint at the previous time step ($h_{j,t-1}$).

As depicted in Figure 2, each unit also has two forget gates, $f_{j,t}^T$ and $f_{j,t}^S$, to handle the two sources of context information in temporal and spatial dimensions, respectively. The transition equations of ST-LSTM are formulated as follows:

$$\begin{pmatrix} i_{j,t} \\ f_{j,t}^S \\ f_{j,t}^T \\ o_{j,t} \\ u_{j,t} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left(M \begin{pmatrix} x_{j,t} \\ h_{j-1,t} \\ h_{j,t-1} \end{pmatrix} \right) \quad (4)$$

$$c_{j,t} = i_{j,t} \odot u_{j,t} + f_{j,t}^S \odot c_{j-1,t} + f_{j,t}^T \odot c_{j,t-1} \quad (5)$$

$$h_{j,t} = o_{j,t} \odot \tanh(c_{j,t}) \quad (6)$$

3.3. Tree-Structure Based Traversal

Arranging the skeletal joints in a simple chain order ignores the kinematic interdependencies among the body joints. Moreover, several semantically false connections between the joints, which are not strongly related, are added.

The body joints are popularly represented as a tree-based pictorial structure [54], [55] in human parsing, as shown in Figure 3(b). It is beneficial to utilize the known interdependency relations between various sets of body joints as

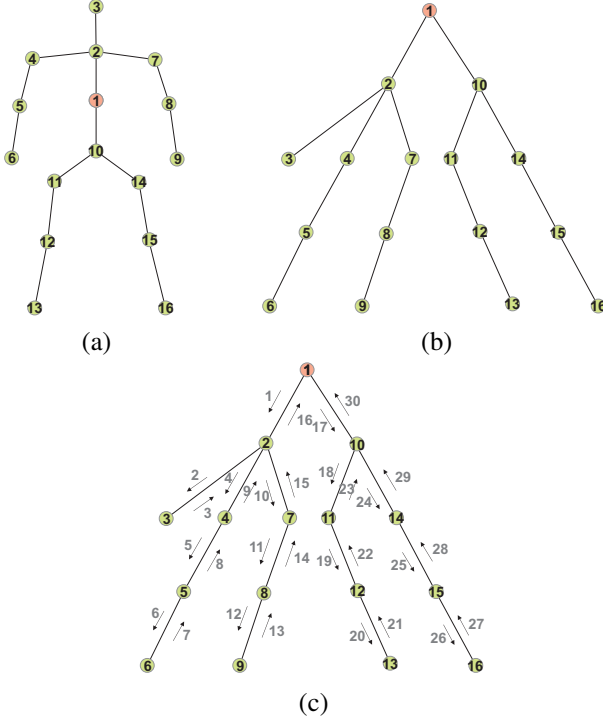


Figure 3. (a) The skeleton of the human body. In the simple joint chain model, the joint visiting order is 1-2-3-...-16. (b) The skeleton is transformed to a tree structure. (c) The tree traversal scheme. The tree structure can be unfolded to a chain with the traversal scheme, and the joint visiting order is 1-2-3-2-4-5-6-5-4-2-7-8-9-8-7-2-1-10-11-12-13-12-11-10-14-15-16-15-14-10-1.

an adjacency tree structure inside our ST-LSTM network as well. For instance, the hidden representation of the neck joint (joint 2 in Figure 3(a)) is often more informative for the right hand joints (7, 8, and 9) compared to the joint 6, which lies before them in the numerically ordered chain-like model. Although using a tree structure for the skeletal data sounds more reasonable here, tree structures cannot be directly fed into our current form of the proposed ST-LSTM network.

In order to mitigate the aforementioned limitation, a bidirectional tree traversal scheme is proposed. In this scheme, the joints are visited in a sequence, while the adjacency information in the skeletal tree structure will be maintained. At the first spatial step, the root node (central spine joint in Figure 3(c)) is fed to our network. Then the network follows the depth-first traversal order in the spatial (skeleton tree) domain. Upon reaching a leaf node, the traversal backtracks in the tree. Finally, the traversal goes back to the root node.

In our traversal scheme, each connection in the tree is met twice, thus it guarantees the transmission of the context data in both top-down and bottom-up directions within the adjacency tree structure. In other words, each node (joint) can obtain the context information from both its ancestors and descendants in the hierarchy defined by the tree structure. Compared to the simple joint chain order

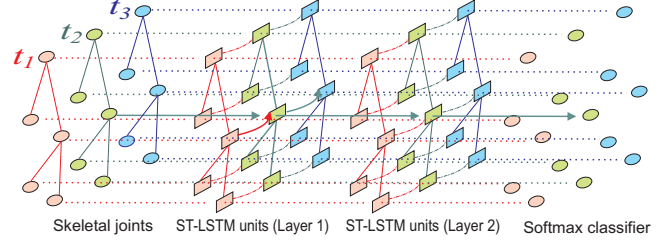


Figure 4. Illustration of the deep tree-structured ST-LSTM network. For clarity, some arrows are omitted in this figure. The hidden representation of the first ST-LSTM layer is fed to the second ST-LSTM layer as its input. The second ST-LSTM layer’s hidden representation is fed to the softmax layer for classification.

described in section 3.2, this tree traversal strategy, which takes advantage of the joints’ adjacency structure, can discover stronger long-term spatial dependency patterns in the skeleton sequence.

Our framework’s representation capacity can be further improved by stacking multiple layers of the tree-structured ST-LSTMs and making the network deeper, as shown in Figure 4.

It is worth noting that at each step of our ST-LSTM framework, the input is limited to the information of a single joint at a time step, and its dimension is much smaller compared to the concatenated input features used by other existing methods. Therefore, our network has much fewer learning parameters. This can be regarded as a weight sharing regularization for our learning model, which leads to better generalization in the scenarios with relatively small sets of training samples. This is an important advantage for skeleton-based action recognition, since the numbers of training samples in most existing datasets are limited.

3.4. Spatio-Temporal LSTM with Trust Gates

In our proposed tree-structured ST-LSTM network, the inputs are the positions of body joints provided by depth sensors (such as Kinect), which are not always accurate because of noisy measurements and occlusion. The unreliable inputs can degrade the performance of the network.

To circumvent this difficulty, we propose to add a novel additional gate to our ST-LSTM network to analyze the reliability of the input measurements based on the derived estimations of the input from the available context information at each spatio-temporal step. Our gating scheme is inspired by the works in natural language processing [11], which use the LSTM representation of previous words at each step to predict the next coming word. As there are often high dependency relations among the words in a sentence, this idea works decently. Similarly, in a skeletal sequence, the neighboring body joints often move together, and this articulated motion follows common yet complex patterns, thus the input data $x_{j,t}$ is expected to be predictable by using the contextual information ($h_{j-1,t}$ and $h_{j,t-1}$) at each spatio-temporal step.

Inspired by this predictability concept, we add a new mechanism to our ST-LSTM calculating a prediction of the input at each step and comparing it with the actual input. The amount of estimation error is then used to learn a new “trust gate”. The activation of this new gate can be used to assist the ST-LSTM network to learn better decisions about when and how to remember or forget the contents in the memory cell. For instance, if the trust gate learns that the current joint has wrong measurements, then this gate can block the input gate and prevent the memory cell from being altered by the current unreliable input data.

Concretely, we introduce a function to produce a prediction of the input at step (j, t) based on the available context information as:

$$p_{j,t} = \tanh \left(M_p \begin{pmatrix} h_{j-1,t} \\ h_{j,t-1} \end{pmatrix} \right) \quad (7)$$

where M_p is an affine transformation mapping the data from \mathbb{R}^{2d} to \mathbb{R}^d , thus the dimension of $p_{j,t}$ is d . Note that the context information at each step does not only contain the representation of the previous temporal step, but also the hidden state of the previous spatial step. This indicates that the long-term context information of both the same joint at previous frames and the other visited joints at the current frame are seamlessly incorporated. Thus this function is expected to be capable of generating reasonable predictions.

In our proposed network, the activation of trust gate is a vector in \mathbb{R}^d (similar to the activation of input gate and forget gate). The trust gate $\tau_{j,t}$ is calculated as follows:

$$x'_{j,t} = \tanh(M_x(x_{j,t})) \quad (8)$$

$$\tau_{j,t} = G(p_{j,t} - x'_{j,t}) \quad (9)$$

where $M_x : \mathbb{R}^D \rightarrow \mathbb{R}^d$ is an affine transformation. The activation function $G(\cdot)$ is an element-wise operation calculated as $G(z) = \exp(-\lambda z^2)$, for which λ is a parameter to control the bandwidth of Gaussian function ($\lambda > 0$). $G(z)$ produces a small response if z has a large absolute value and a large response when z is close to zero.

Adding the proposed trust gate, the cell state of ST-LSTM will be updated as:

$$\begin{aligned} c_{j,t} = & \tau_{j,t} \odot i_{j,t} \odot u_{j,t} \\ & + (1 - \tau_{j,t}) \odot f_{j,t}^S \odot c_{j-1,t} \\ & + (1 - \tau_{j,t}) \odot f_{j,t}^T \odot c_{j,t-1} \end{aligned} \quad (10)$$

This equation can be explained as follows: (1) if the input $x_{j,t}$ is not trusted (due to the noise or occlusion), then our network relies more on its history information, and tries to block the new input at this step; (2) on the contrary, if the input is reliable, then our learning algorithm updates the memory cell regarding the input data.

The proposed ST-LSTM unit equipped with trust gate is illustrated in Figure 5. The concept of the proposed trust gate technique is theoretically generic and can be used in other domains to handle noisy input information for recurrent network models.

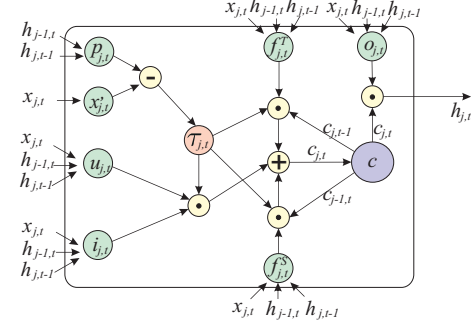


Figure 5. Illustration of the proposed ST-LSTM with trust gate.

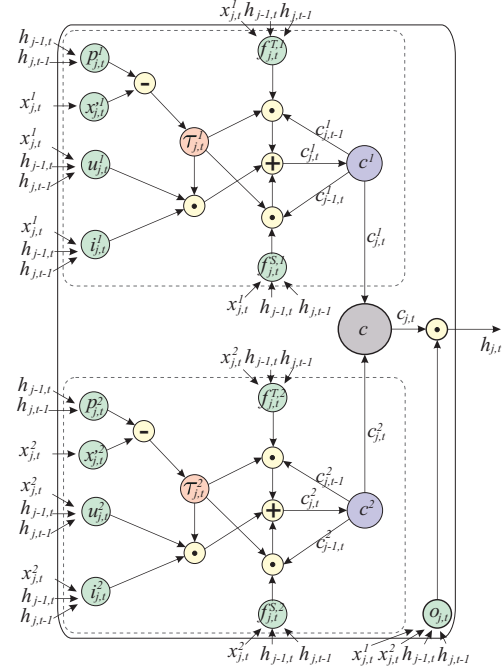


Figure 6. Illustration of the proposed structure for feature fusion inside the ST-LSTM unit.

3.5. Feature Fusion within ST-LSTM Unit

As mentioned above, at each spatio-temporal step, the positional information of the corresponding joint at the current frame is fed to our ST-LSTM network. Here we call joint position-based feature as a geometric feature. Beside utilizing the joint position (3D coordinates), we can also extract visual texture and motion features (e.g. HOG, HOF [56], [57], or ConvNet-based features [58], [59]) from the RGB frames, around each body joint as the complementary information. This is intuitively effective for better human action representation, especially in the human-object interaction scenarios.

A naive way for combining geometric and visual features for each joint is to concatenate them in the feature level and feed them to the corresponding ST-LSTM unit as network’s input data. However, the dimension of the geometric feature

is very low intrinsically, while the visual features are often in relatively higher dimensions. Due to this inconsistency, simple concatenation of these two types of features in the input stage of the network causes degradation in the final performance of the entire model.

The work in [32] feeds different body parts into the Part-aware LSTM [32] separately, and then assembles them inside the LSTM unit. Inspired by this work, we propose to fuse the two types of features inside the ST-LSTM unit, rather than simply concatenating them at the input level.

We use $x_{j,t}^{\mathcal{F}}$ ($\mathcal{F} \in \{1, 2\}$) to denote the geometric feature and visual feature for a joint at the t -th time step. As illustrated in Figure 6, at step (j, t) , the two features ($x_{j,t}^1$ and $x_{j,t}^2$) are fed to the ST-LSTM unit separately as the new input structure. Inside the recurrent unit, we deploy two sets of gates, input gates ($i_{j,t}^{\mathcal{F}}$), forget gates with respect to time ($f_{j,t}^{T,\mathcal{F}}$) and space ($f_{j,t}^{S,\mathcal{F}}$), and also trust gates ($\tau_{j,t}^{\mathcal{F}}$), to deal with the two heterogeneous sets of modality features. We put the two cell representations ($c_{j,t}^{\mathcal{F}}$) together to build up the multimodal context information of the two sets of modality features. Finally, the output of each ST-LSTM unit is calculated based on the multimodal context representations, and controlled by the output gate ($o_{j,t}$) which is shared for the two sets of features.

For the features of each modality, it is efficient and intuitive to model their context information independently. However, we argue that the representation ability of each modality-based sets of features can be strengthened by borrowing information from the other set of features. Thus, the proposed structure does not completely separate the modeling of multimodal features.

Let us take the geometric feature as an example. Its input gate, forget gates, and trust gate are all calculated from the new input ($x_{j,t}^1$) and hidden representations ($h_{j,t-1}$ and $h_{j-1,t}$), whereas each hidden representation is an associate representation of two features' context information from previous steps. Assisted by visual features' context information, the input gate, forget gates, and also trust gate for geometric feature can effectively learn how to update its current cell state ($c_{j,t}^1$). Specifically, for the new geometric feature input ($x_{j,t}^1$), we expect the network to produce a better prediction when it is not only based on the context of the geometric features, but also assisted by the context of visual features. Therefore, the trust gate ($\tau_{j,t}^1$) will have stronger ability to assess the reliability of the new input data ($x_{j,t}^1$).

The proposed ST-LSTM with integrated multimodal feature fusion is formulated as:

$$\begin{pmatrix} i_{j,t}^{\mathcal{F}} \\ f_{j,t}^{S,\mathcal{F}} \\ f_{j,t}^{T,\mathcal{F}} \\ u_{j,t}^{\mathcal{F}} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left(M^{\mathcal{F}} \begin{pmatrix} x_{j,t}^{\mathcal{F}} \\ h_{j-1,t} \\ h_{j,t-1} \end{pmatrix} \right) \quad (11)$$

$$p_{j,t}^{\mathcal{F}} = \tanh \left(M_p^{\mathcal{F}} \begin{pmatrix} h_{j-1,t} \\ h_{j,t-1} \end{pmatrix} \right) \quad (12)$$

$$x'_{j,t}^{\mathcal{F}} = \tanh \left(M_x^{\mathcal{F}} \begin{pmatrix} x_{j,t}^{\mathcal{F}} \end{pmatrix} \right) \quad (13)$$

$$\tau_{j,t}^{\mathcal{F}} = G(x'_{j,t}^{\mathcal{F}} - p_{j,t}^{\mathcal{F}}) \quad (14)$$

$$\begin{aligned} c_{j,t}^{\mathcal{F}} &= \tau_{j,t}^{\mathcal{F}} \odot i_{j,t}^{\mathcal{F}} \odot u_{j,t}^{\mathcal{F}} \\ &\quad + (1 - \tau_{j,t}^{\mathcal{F}}) \odot f_{j,t}^{S,\mathcal{F}} \odot c_{j-1,t}^{\mathcal{F}} \\ &\quad + (1 - \tau_{j,t}^{\mathcal{F}}) \odot f_{j,t}^{T,\mathcal{F}} \odot c_{j,t-1}^{\mathcal{F}} \end{aligned} \quad (15)$$

$$o_{j,t} = \sigma \left(M_o \begin{pmatrix} x_{j,t}^1 \\ x_{j,t}^2 \\ h_{j-1,t} \\ h_{j,t-1} \end{pmatrix} \right) \quad (16)$$

$$h_{j,t} = o_{j,t} \odot \tanh \begin{pmatrix} c_{j,t}^1 \\ c_{j,t}^2 \end{pmatrix} \quad (17)$$

3.6. Learning the Classifier

As the labels are given at video level, we feed them as the training outputs of our network at each spatio-temporal step. A softmax layer is used by the network to predict the action class \hat{y} among the given class set Y . The prediction of the whole video can be obtained by averaging the prediction scores of all steps. The objective function of our ST-LSTM network is as follows:

$$\mathcal{L} = \sum_{j=1}^J \sum_{t=1}^T l(\hat{y}_{j,t}, y) \quad (18)$$

where $l(\hat{y}_{j,t}, y)$ is the negative log-likelihood loss [60] that measures the difference between the prediction result $\hat{y}_{j,t}$ at step (j, t) and the true label y .

The back-propagation through time (BPTT) algorithm [60] is often effective for minimizing the objective function for the RNN/LSTM models. As our ST-LSTM model involves both spatial and temporal steps, we adopt a modified version of BPTT for training. The back-propagation runs over spatial and temporal steps simultaneously by starting at the last joint at the last frame. To clarify the error accumulation in this procedure, we use $e_{j,t}^T$ and $e_{j,t}^S$ to denote the error back-propagated from step $(j, t+1)$ to (j, t) and the error back-propagated from step $(j+1, t)$ to (j, t) , respectively. Then the errors accumulated at step (j, t) can be calculated as $e_{j,t}^T + e_{j,t}^S$. Consequently, before back-propagating the error at each step, we should guarantee both its subsequent joint step and subsequent time step have already been computed.

The left-most units in our ST-LSTM network do not have preceding spatial units, as shown in Figure 1. To update the cell states of these units in the feed-forward stage, a popular strategy is to input zero values into these nodes to substitute

the hidden representations from the preceding nodes. In our implementation, we link the last unit at the last time step to the first unit at the current time step. We call the new connection as last-to-first link. In the tree traversal, the first and last nodes refer to the same joint (root node of the tree), however the last node contains holistic information of the human skeleton in the corresponding frame. Linking the last node to the starting node at the next time step provides the starting node with the whole body structure configuration, rather than initializing it with less effective zero values. Thus, the network has better ability to learn the action patterns in the skeleton sequence.

4. Experiments

The proposed method is evaluated and empirically analyzed on seven benchmark datasets for which the coordinates of skeletal joints are provided. These datasets are NTU RGB+D, UT-Kinect, SBU Interaction, SYSU-3D, ChaLearn Gesture, MSR Action3D, and Berkeley MHAD. We conduct extensive experiments with different models to verify the effectiveness of individual technical contributions proposed, as follows:

- (1) “ST-LSTM (Joint Chain)”. In this model, the joints are visited in a simple chain order, as shown in Figure 3(a);
- (2) “ST-LSTM (Tree)”. In this model, the tree traversal scheme illustrated in Figure 3(c) is used to take advantage of the tree-based spatial structure of skeletal joints;
- (3) “ST-LSTM (Tree) + Trust Gate”. This model uses the trust gate to handle the noisy input.

The input to every unit of our network at each spatio-temporal step is the location of the corresponding skeletal joint (i.e., geometric features) at the current time step. We also use two of the datasets (NTU RGB+D dataset and UT-Kinect dataset) as examples to evaluate the performance of our fusion model within the ST-LSTM unit by fusing the geometric and visual features. These two datasets include human-object interactions (such as making a phone call and picking up something) and the visual information around the major joints can be complementary to the geometric features for action recognition.

4.1. Evaluation Datasets

NTU RGB+D dataset [32] was captured with Kinect (v2). It is currently the largest publicly available dataset for depth-based action recognition, which contains more than 56,000 video sequences and 4 million video frames. The samples in this dataset were collected from 80 distinct viewpoints. A total of 60 action classes (including daily actions, medical conditions, and pair actions) were performed by 40 different persons aged between 10 and 35. This dataset is very challenging due to the large intra-class and viewpoint variations. With a large number of samples, this dataset is highly suitable for deep learning based activity analysis. The parameters learned on this dataset can also be used to initialize the models for smaller datasets to improve

and speed up the training process of the network. The 3D coordinates of 25 body joints are provided in this dataset.

UT-Kinect dataset [4] was captured with a stationary Kinect sensor. It contains 10 action classes. Each action was performed twice by every subject. The 3D locations of 20 skeletal joints are provided. The significant intra-class and viewpoint variations make this dataset very challenging.

SBU Interaction dataset [61] was collected with Kinect. It contains 8 classes of two-person interactions, and includes 282 skeleton sequences with 6822 frames. Each body skeleton consists of 15 joints. The major challenges of this dataset are: (1) in most interactions, one subject is acting, while the other subject is reacting; and (2) the 3D measurement accuracies of the joint coordinates are low in many sequences.

SYSU-3D dataset [62] contains 480 sequences and was collected with Kinect. In this dataset, 12 different activities were performed by 40 persons. The 3D coordinates of 20 joints are provided in this dataset. The SYSU-3D dataset is a very challenging benchmark because: (1) the motion patterns are highly similar among different activities, and (2) there are various viewpoints in this dataset.

ChaLearn Gesture dataset [63] consists of 23 hours of videos captured with Kinect. A total of 20 Italian gestures were performed by 27 different subjects. This dataset contains 955 long-duration videos and has predefined splits of samples as training, validation and testing sets. Each skeleton in this dataset has 20 joints.

MSR Action3D dataset [64] is widely used for depth-based action recognition. It contains a total of 10 subjects and 20 actions. Each action was performed by the same subject two or three times. Each frame in this dataset contains 20 skeletal joints.

Berkeley MHAD dataset [65] was collected by using a motion capture network of sensors. It contains 659 sequences and about 82 minutes of recording time. Eleven action classes were performed by five female and seven male subjects. The 3D coordinates of 35 skeletal joints are provided in each frame.

4.2. Implementation Details

In our experiments, each video sequence is divided to T sub-sequences with the same length, and one frame is randomly selected from each sub-sequence. This sampling strategy has the following advantages: (1) Randomly selecting a frame from each sub-sequence can add variation to the input data, and improves the generalization strengths of our trained network. (2) Assume each sub-sequence contains n frames, so we have n choices to sample a frame from each sub-sequence. Accordingly, for the whole video, we can obtain a total number of n^T sampling combinations. This indicates that the training data can be greatly augmented. We use different frame sampling combinations for each video over different training epochs. This strategy is useful for handling the over-fitting issues, as most datasets have limited numbers of training samples. We observe this strategy achieves better performance in contrast with uniformly

sampling frames. We cross-validated the performance based on the leave-one-subject-out protocol on the large scale NTU RGB+D dataset, and found $T = 20$ as the optimum value.

We use Torch7 [66] as the deep learning platform to perform our experiments. We train the network with stochastic gradient descent, and set the learning rate, momentum, and decay rate to 2×10^{-3} , 0.9, and 0.95, respectively. We set the unit size d to 128, and the parameter λ used in $G(\cdot)$ to 0.5. Two ST-LSTM layers are used in our stacked network. Although there are variations in terms of joint number, sequence length, and data acquisition equipment for different datasets, we adopt the same parameter settings mentioned above for all datasets. Our method achieves promising results on all the benchmark datasets with these parameter settings untouched, which shows the robustness of our method.

An NVIDIA TitanX GPU is used to perform our experiments. We evaluate the computational efficiency of our method on the NTU RGB+D dataset and set the batch size to 100. On average, within one second, 210, 100, and 70 videos can be processed by using “ST-LSTM (Joint Chain)”, “ST-LSTM (Tree)”, and “ST-LSTM (Tree) + Trust Gate”, respectively.

4.3. Experiments on the NTU RGB+D Dataset

The NTU RGB+D dataset has two standard evaluation protocols [32]. The first protocol is the cross-subject (X-Subject) evaluation protocol, in which half of the subjects are used for training and the remaining subjects are kept for testing. The second is the cross-view (X-View) evaluation protocol, in which 2/3 of the viewpoints are used for training, and 1/3 unseen viewpoints are left out for testing. We evaluate the performance of our method on both of these protocols. The results are shown in TABLE 1.

TABLE 1. EXPERIMENTAL RESULTS ON THE NTU RGB+D DATASET

Method	Feature	X-Subject	X-View
Lie Group [6]	Geometric	50.1%	52.8%
Cippitelli <i>et al.</i> [67]	Geometric	48.9%	57.7%
Dynamic Skeletons [62]	Geometric	60.2%	65.2%
FTP [68]	Geometric	61.1%	72.6%
Hierarchical RNN [30]	Geometric	59.1%	64.0%
Deep RNN [32]	Geometric	56.3%	64.1%
Part-aware LSTM [32]	Geometric	62.9%	70.3%
ST-LSTM (Joint Chain)	Geometric	61.7%	75.5%
ST-LSTM (Tree)	Geometric	65.2%	76.1%
ST-LSTM (Tree) + Trust Gate	Geometric	69.2%	77.7%

In TABLE 1, the deep RNN model concatenates the joint features at each frame and then feeds them to the network to model the temporal kinetics, and ignores the spatial dynamics. As can be seen, both “ST-LSTM (Joint Chain)” and “ST-LSTM (Tree)” models outperform this method by a notable margin. It can also be observed that our approach utilizing the trust gate brings significant performance improvement, because the data provided by Kinect is often noisy and multiple joints are frequently occluded in this dataset. Note that our proposed models (such as “ST-LSTM (Tree) + Trust Gate”) reported in this table only use skeletal data as input.

We compare the class specific recognition accuracies of “ST-LSTM (Tree)” and “ST-LSTM (Tree) + Trust Gate”, as shown in Figure 7. We observe that “ST-LSTM (Tree) + Trust Gate” significantly outperforms “ST-LSTM (Tree)” for most of the action classes, which demonstrates our proposed trust gate can effectively improve the human action recognition accuracy by learning the degrees of reliability over the input data at each time step.

As shown in Figure 8, a notable portion of videos in the NTU RGB+D dataset were collected in side views. Due to the design of Kinect’s body tracking mechanism, skeletal data is less accurate in side view compared to the front view. To further investigate the effectiveness of the proposed trust gate, we analyze the performance of the network by feeding the side views samples only. The accuracy of “ST-LSTM (Tree)” is 76.5%, while “ST-LSTM (Tree) + Trust Gate” yields 81.6%. This shows how trust gate can effectively deal with the noise in the input data.

To verify the performance boost by stacking layers, we limit the depth of the network by using only one ST-LSTM layer, and the accuracies drop to 65.5% and 77.0% based on the cross-subject and cross-view protocol, respectively. This indicates our two-layer stacked network has better representation power than the single-layer network.

To evaluate the performance of our feature fusion scheme, we extract visual features from several regions based on the joint positions and use them in addition to the geometric features (3D coordinates of the joints). We extract HOG and HOF [56], [57] features from a 80×80 RGB patch centered at each joint location. For each joint, this produces a 300D visual descriptor, and we apply PCA to reduce the dimension to 20. The results are shown in TABLE 2. We observe that our method using the visual features together with the joint positions improves the performance. Besides, we compare our newly proposed feature fusion strategy within the ST-LSTM unit with two other feature fusion methods: (1) early fusion which simply concatenates two types of features as the input of the ST-LSTM unit; (2) late fusion which uses two ST-LSTMs to deal with two types of features respectively, then concatenates the outputs of the two ST-LSTMs at each step, and feeds the concatenated result to a softmax classifier. We observe that our proposed feature fusion strategy is superior to other baselines.

TABLE 2. EVALUATION OF DIFFERENT FEATURE FUSION STRATEGIES ON THE NTU RGB+D DATASET. “GEOMETRIC + VISUAL (1)” INDICATES THE EARLY FUSION SCHEME. “GEOMETRIC + VISUAL (2)” INDICATES THE LATE FUSION SCHEME. “GEOMETRIC \oplus VISUAL” MEANS OUR NEWLY PROPOSED FEATURE FUSION SCHEME WITHIN THE ST-LSTM UNIT.

Feature Fusion Method	X-Subject	X-View
Geometric Only	69.2%	77.7%
Geometric + Visual (1)	70.8%	78.6%
Geometric + Visual (2)	71.0%	78.7%
Geometric \oplus Visual	73.2%	80.6%

We also evaluate the sensitivity of the proposed network with respect to the variation of neuron unit size and λ values. The results are shown in Figure 9. When trust gate is added,

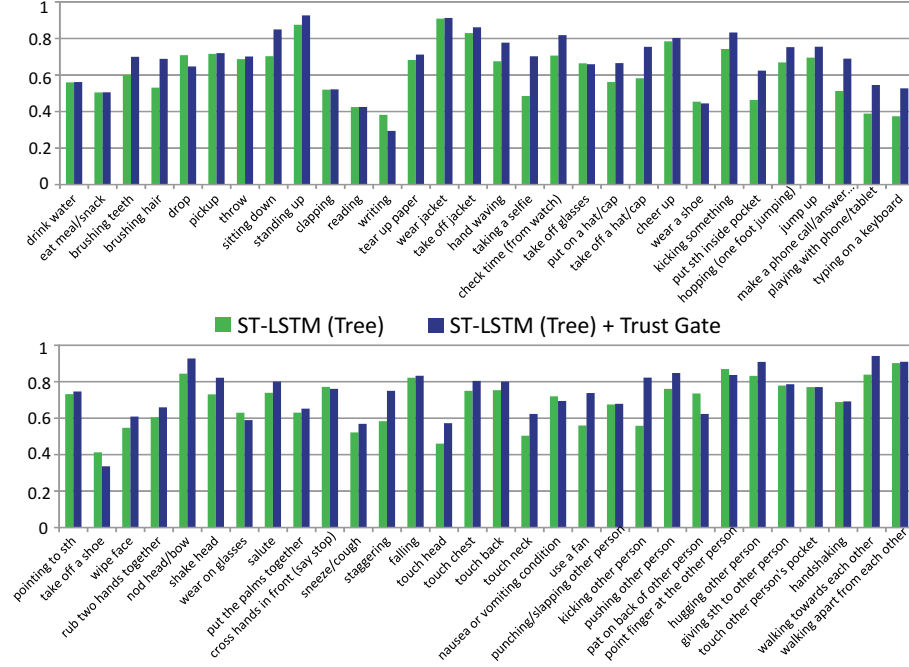


Figure 7. Recognition accuracy per class on the NTU RGB+D dataset

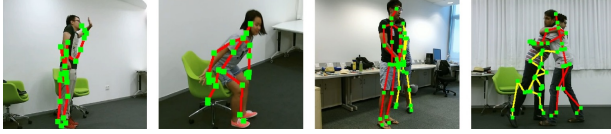


Figure 8. Examples of the noisy skeletons from the NTU RGB+D dataset.

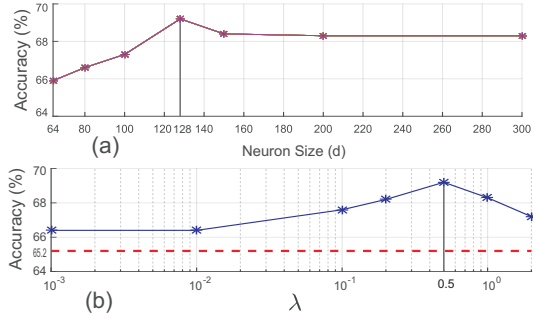


Figure 9. (a) Performance comparison of our approach using different values of neuron size (d) on the NTU RGB+D dataset (X-subject). (b) Performance comparison of our method using different λ values on the NTU RGB+D dataset (X-subject). The blue line represents our results when different λ values are used for trust gate, while the red dashed line indicates the performance of our method when trust gate is not added.

our network obtains better performance for all the λ values compared to the network without the trust gate.

Finally, we investigate the recognition performance with early stopping conditions by feeding the first p portion of the testing video to the trained network based on the cross-subject protocol ($p \in \{0.1, 0.2, \dots, 1.0\}$). The results are shown in Figure 10. We can observe that the results are

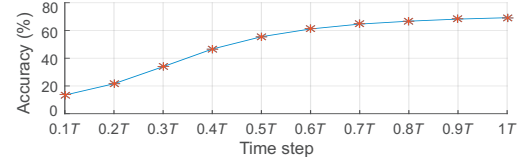


Figure 10. Experimental results of our method by early stopping the network evolution at different time steps.

improved when a larger portion of the video is fed to our network.

4.4. Experiments on the UT-Kinect Dataset

There are two evaluation protocols for the UT-Kinect dataset in the literature. The first is the leave-one-out-cross-validation (LOOCV) protocol [4]. The second protocol is suggested by [69], for which half of the subjects are used for training, and the remaining are used for testing. We evaluate our approach using both protocols on this dataset.

Using the LOOCV protocol, our method achieves better performance than other skeleton-based methods, as shown in TABLE 3. Using the second protocol (see TABLE 4), our method achieves competitive result (95.0%) to the Elastic functional coding method [70] (94.9%), which is an extension of the Lie Group model [6].

Some actions in the UT-Kinect dataset involve human-object interactions, thus appearance based features representing visual information of the objects can be complementary to the geometric features. Thus we can evaluate our proposed feature fusion approach within the ST-LSTM unit on this dataset. The results are shown in TABLE 5.

TABLE 3. EXPERIMENTAL RESULTS ON THE UT-KINECT DATASET (LOOCV PROTOCOL [4])

Method	Feature	Acc.
Grassmann Manifold [71]	Geometric	88.5%
Jetley <i>et al.</i> [72]	Geometric	90.0%
Histogram of 3D Joints [4]	Geometric	90.9%
Space Time Pose [73]	Geometric	91.5%
Riemannian Manifold [74]	Geometric	91.5%
SCs (Informative Joints) [75]	Geometric	91.9%
Chungoo <i>et al.</i> [76]	Geometric	92.0%
Key-Pose-Motifs Mining [77]	Geometric	93.5%
ST-LSTM (Joint Chain)	Geometric	91.0%
ST-LSTM (Tree)	Geometric	92.4%
ST-LSTM (Tree) + Trust Gate	Geometric	97.0%

TABLE 4. RESULTS ON THE UT-KINECT DATASET (HALF-VS-HALF PROTOCOL [69])

Method	Feature	Acc.
Skeleton Joint Features [69]	Geometric	87.9%
Chungoo <i>et al.</i> [76]	Geometric	89.5%
Lie Group [6] (reported by [70])	Geometric	93.6%
Elastic functional coding [70]	Geometric	94.9%
ST-LSTM (Tree) + Trust Gate	Geometric	95.0%

Using geometric features only, the accuracy is 97%. By simply concatenating the geometric and visual features, the accuracy improves slightly. However, the accuracy of our approach can reach 98% when the proposed feature fusion method is adopted.

TABLE 5. EVALUATION OF OUR APPROACH FOR FEATURE FUSION ON THE UT-KINECT DATASET (LOOCV PROTOCOL [4]). “GEOMETRIC + VISUAL” INDICATES WE SIMPLY CONCATENATE THE TWO TYPES OF FEATURES AS THE INPUT. “GEOMETRIC \oplus VISUAL” MEANS WE USE THE NEWLY PROPOSED FEATURE FUSION SCHEME WITHIN THE ST-LSTM UNIT.

Feature Fusion Method	Acc.
Geometric Only	97.0%
Geometric + Visual	97.5%
Geometric \oplus Visual	98.0%

4.5. Experiments on the SBU Interaction Dataset

We follow the standard evaluation protocol in [61] and perform 5-fold cross validation on the SBU Interaction dataset. As two human skeletons are provided in each frame of this dataset, our traversal scheme visits the joints throughout the two skeletons over the spatial steps.

We report the results in terms of average classification accuracy in TABLE 6. The methods in [48] and [30] are both LSTM-based approaches, which are more relevant to our method.

The results show that the proposed “ST-LSTM (Tree) + Trust Gate” model outperforms all other skeleton-based methods. “ST-LSTM (Tree)” achieves higher accuracy than “ST-LSTM (Joint Chain)”, as the latter adds some false links between less related joints.

Both Co-occurrence LSTM [48] and Hierarchical RNN [30] adopt the Svaitzky-Golay filter [80] in the temporal

TABLE 6. EXPERIMENTAL RESULTS ON THE SBU INTERACTION DATASET

Method	Feature	Acc.
Yun <i>et al.</i> [61]	Geometric	80.3%
Ji <i>et al.</i> [78]	Geometric	86.9%
CHARM [79]	Geometric	83.9%
Hierarchical RNN [30]	Geometric	80.4%
Co-occurrence LSTM [48]	Geometric	90.4%
Deep LSTM [48]	Geometric	86.0%
ST-LSTM (Joint Chain)	Geometric	84.7%
ST-LSTM (Tree)	Geometric	88.6%
ST-LSTM (Tree) + Trust Gate	Geometric	93.3%

domain to smooth the skeletal joint positions and reduce the influence of noise in the data collected by Kinect.

The proposed “ST-LSTM (Tree)” model without the trust gate mechanism outperforms Hierarchical RNN, and achieves comparable result (88.6%) to Co-occurrence LSTM. When the trust gate is used, the accuracy of our method jumps to 93.3%.

4.6. Experiments on the SYSU-3D Dataset

We follow the standard evaluation protocol in [62] on the SYSU-3D dataset. The samples from 20 subjects are used to train the model parameters, and the samples of the remaining 20 subjects are used for testing. We perform 30-fold cross validation and report the mean accuracy in TABLE 7.

TABLE 7. EXPERIMENTAL RESULTS ON THE SYSU-3D DATASET

Method	Feature	Acc.
LAFF (SKL) [81]	Geometric	54.2%
Dynamic Skeletons [62]	Geometric	75.5%
ST-LSTM (Joint Chain)	Geometric	72.1%
ST-LSTM (Tree)	Geometric	73.4%
ST-LSTM (Tree) + Trust Gate	Geometric	76.5%

The results in TABLE 7 show that our proposed “ST-LSTM (Tree) + Trust Gate” method outperforms all the baseline methods on this dataset. We can also find that the tree traversal strategy can help to improve the classification accuracy of our model. As the skeletal joints provided by Kinect are noisy in this dataset, the trust gate, which aims at handling noisy data, brings significant performance improvement (about 3% improvement).

There are large viewpoint variations in this dataset. To make our model reliable against viewpoint variations, we adopt a similar skeleton normalization procedure as suggested by [32] on this dataset. In this preprocessing step, we perform a rotation transformation on each skeleton, such that all the normalized skeletons face to the same direction. Specifically, after rotation, the 3D vector from “right shoulder” to “left shoulder” will be parallel to the X axis, and the vector from “hip center” to “spine” will be aligned to the Y axis (please see [32] for more details about the normalization procedure).

We evaluate our “ST-LSTM (Tree) + Trust Gate” method by respectively using the original skeletons without rotation and the transformed skeletons, and report the results in

TABLE 8. The results show that it is beneficial to use the transformed skeletons as the input for action recognition.

TABLE 8. EVALUATION FOR SKELETON ROTATION ON THE SYSU-3D DATASET

Method	Acc.
With Skeleton Rotation	76.5%
Without Skeleton Rotation	73.0%

4.7. Experiments on the ChaLearn Gesture Dataset

We follow the evaluation protocol adopted in [82], [83] and report the F1-score measures on the validation set of the ChaLearn Gesture dataset.

TABLE 9. EXPERIMENTAL RESULTS ON THE CHALEARN GESTURE DATASET

Method	Feature	F1-Score
Portfolios [84]	Geometric	56.0%
Wu <i>et al.</i> [85]	Geometric	59.6%
Pfister <i>et al.</i> [86]	Geometric	61.7%
HiVideoDarwin [82]	Geometric	74.6%
VideoDarwin [83]	Geometric	75.2%
Deep LSTM [32]	Geometric	87.1%
ST-LSTM (Joint Chain)	Geometric	89.1%
ST-LSTM (Tree)	Geometric	89.9%
ST-LSTM (Tree) + Trust Gate	Geometric	92.0%

As shown in TABLE 9, our method surpasses the state-of-the-art methods [32], [82], [83], [84], [85], [86], which demonstrates the effectiveness of our method in dealing with skeleton-based action recognition problem.

Compared to other methods, our method focuses on modeling both temporal and spatial dependency patterns in skeleton sequences. Moreover, the proposed trust gate is also incorporated to our method to handle the noisy skeleton data captured by Kinect, which can further improve the results.

4.8. Experiments on the MSR Action3D Dataset

We follow the experimental protocol in [30] on the MSR Action3D dataset, and show the results in TABLE 10.

On the MSR Action3D dataset, our proposed method, “ST-LSTM (Tree) + Trust Gate”, achieves 94.8% of classification accuracy, which is superior to the Hierarchical RNN model [30] and other baseline methods.

TABLE 10. EXPERIMENTAL RESULTS ON THE MSR ACTION3D DATASET

Method	Feature	Acc.
Histogram of 3D Joints [4]	Geometric	79.0%
Joint Angles Similarities [8]	Geometric	83.5%
SCs (Informative Joints) [75]	Geometric	88.3%
Oriented Displacements [87]	Geometric	91.3%
Lie Group [6]	Geometric	92.5%
Space Time Pose [73]	Geometric	92.8%
Lillo <i>et al.</i> [88]	Geometric	93.0%
Hierarchical RNN [30]	Geometric	94.5%
ST-LSTM (Tree) + Trust Gate	Geometric	94.8%

4.9. Experiments on the Berkeley MHAD Dataset

TABLE 11. EXPERIMENTAL RESULTS ON THE BERKELEY MHAD DATASET

Method	Feature	Acc.
Ofli <i>et al.</i> [89]	Geometric	95.4%
Vantigodi <i>et al.</i> [90]	Geometric	96.1%
Vantigodi <i>et al.</i> [91]	Geometric	97.6%
Kapsouras <i>et al.</i> [92]	Geometric	98.2%
Hierarchical RNN [30]	Geometric	100%
Co-occurrence LSTM [48]	Geometric	100%
ST-LSTM (Tree) + Trust Gate	Geometric	100%

We adopt the experimental protocol in [30] on the Berkeley MHAD dataset. 384 video sequences corresponding to the first seven persons are used for training, and the 275 sequences of the remaining five persons are held out for testing. The experimental results in TABLE 11 show that our method achieves very high accuracy (100%) on this dataset. Unlike [30] and [48], our method does not use any preliminary manual smoothing procedures.

4.10. Visualization of Trust Gates

In this section, to better investigate the effectiveness of the proposed trust gate scheme, we study the behavior of the proposed framework against the presence of noise in skeletal data from the MSR Action3D dataset. We manually rectify some noisy joints of the samples by referring to the corresponding depth images. We then compare the activations of trust gates on the noisy and rectified inputs. As illustrated in Figure 11(a), the magnitude of trust gate’s output (l_2 norm of the activations of the trust gate) is smaller when a noisy joint is fed, compared to the corresponding rectified joint. This demonstrates how the network controls the impact of noisy input on its stored representation of the observed data.

In our next experiment, we manually add noise to one joint for all testing samples on the Berkeley MHAD dataset, in order to further analyze the behavior of our proposed trust gate. Note that the Berkeley MHAD dataset was collected with motion capture system, thus the skeletal joint coordinates in this dataset are much more accurate than those captured with Kinect sensors.

We add noise to the right foot joint by moving the joint away from its original location. The direction of the translation vector is randomly chosen and the norm is a random value around 30cm, which is a significant noise in the scale of human body. We measure the difference in the magnitudes of trust gates’ activations between the noisy data and the original ones. For all testing samples, we carry out the same operations and then calculate the average difference. The results in Figure 11(b) show that the magnitude of trust gate is reduced when the noisy data is fed to the network. This shows that our network tries to block the flow of noisy input and stop it from affecting the memory. We also observe that the overall accuracy of our network does not drop after adding the above-mentioned noise to the input data.

TABLE 12. PERFORMANCE COMPARISON OF DIFFERENT SPATIAL SEQUENCE MODELS

Dataset	NTU (X-Subject)	NTU (X-View)	UT-Kinect	SBU Interaction	ChaLearn Gesture
ST-LSTM (Joint Chain)	61.7%	75.5%	91.0%	84.7%	89.1%
ST-LSTM (Double Joint Chain)	63.5%	75.6%	91.5%	85.9%	89.2%
ST-LSTM (Tree)	65.2%	76.1%	92.4%	88.6%	89.9%

TABLE 13. PERFORMANCE COMPARISON OF TEMPORAL AVERAGE, LSTM, AND OUR PROPOSED ST-LSTM

Dataset	NTU (X-Subject)	NTU (X-View)	UT-Kinect	SBU Interaction	ChaLearn Gesture
Temporal Average	47.6%	52.6%	81.9%	71.5%	77.9%
LSTM	62.0%	70.7%	90.5%	86.0%	87.1%
LSTM + Trust Gate	62.9%	71.7%	92.0%	86.6%	87.6%
ST-LSTM	65.2%	76.1%	92.4%	88.6%	89.9%
ST-LSTM + Trust Gate	69.2%	77.7%	97.0%	93.3%	92.0%

TABLE 14. EVALUATION OF THE LAST-TO-FIRST LINK IN OUR PROPOSED NETWORK

Dataset	NTU (X-Subject)	NTU (X-View)	UT-Kinect	SBU Interaction	ChaLearn Gesture
Without last-to-first link	68.5%	76.9%	96.5%	92.1%	90.9 %
With last-to-first link	69.2%	77.7%	97.0%	93.3%	92.0 %

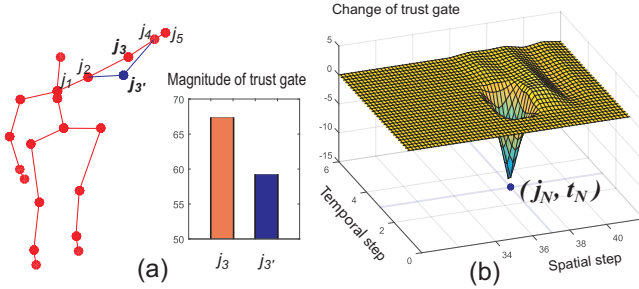


Figure 11. Visualization of the trust gate’s behavior when inputting noisy data. (a) $j_{3'}$ is a noisy joint position, and j_3 is the corresponding rectified joint location. In the histogram, the blue bar indicates the magnitude of trust gate when inputting the noisy joint $j_{3'}$. The red bar indicates the magnitude of the corresponding trust gate when $j_{3'}$ is rectified to j_3 . (b) Visualization of the difference between the trust gate calculated when the noise is imposed at the step (j_N, t_N) and that calculated when inputting the original data.

4.11. Evaluation of Different Spatial Joint Sequence Models

The previous experiments showed how “ST-LSTM (Tree)” outperforms “ST-LSTM (Joint Chain)”, because “ST-LSTM (Tree)” models the kinematic dependency structures of human skeletal sequences. In this section, we further analyze the effectiveness of our “ST-LSTM (Tree)” model and compare it with a “ST-LSTM (Double Joint Chain)” model.

The “ST-LSTM (Joint Chain)” has fewer steps in the spatial dimension than the “ST-LSTM (Tree)”. One question that may rise here is if the advantage of “ST-LSTM (Tree)” model could be only due to the higher length and redundant sequence of the joints fed to the network, and not because of the proposed semantic relations between the joints. To answer this question, we evaluate the effect of using a double chain scheme to increase the spatial steps

of the “ST-LSTM (Joint Chain)” model. Specifically, we use the joint visiting order of 1-2-3-...-16-1-2-3-...-16, and we call this model as “ST-LSTM (Double Joint Chain)”. The results in TABLE 12 show that the performance of “ST-LSTM (Double Joint Chain)” is better than “ST-LSTM (Joint Chain)”, yet inferior to “ST-LSTM (Tree)”.

This experiment indicates that it is beneficial to introduce more passes in the spatial dimension to the ST-LSTM for performance improvement. A possible explanation is that the units visited in the second round can obtain the global level context representation from the previous pass, thus they can generate better representations of the action patterns by using the context information. However, the performance of “ST-LSTM (Double Joint Chain)” is still weaker than “ST-LSTM (Tree)”, though the numbers of their spatial steps are almost equal.

The proposed tree traversal scheme is superior because it connects the most semantically related joints and avoids false connections between the less-related joints (unlike the other two compared models).

4.12. Evaluation of Temporal Average, LSTM and ST-LSTM

To further investigate the effect of simultaneous modeling of dependencies in spatial and temporal domains, in this experiment, we replace our ST-LSTM with the original LSTM which only models the temporal dynamics among the frames without explicitly considering spatial dependencies. We report the results of this experiment in TABLE 13. As can be seen, our “ST-LSTM + Trust Gate” significantly outperforms “LSTM + Trust Gate”. This demonstrates that the proposed modeling of the dependencies in both temporal and spatial dimensions provides much richer representations than the original LSTM.

The second observation of this experiment is that if we add our trust gate to the original LSTM, the performance of LSTM can also be improved, but its performance gain

is less than the performance gain on ST-LSTM. A possible explanation is that we have both spatial and temporal context information at each step of ST-LSTM to generate a good prediction of the input at the current step (see Eq. (7)), thus our trust gate can achieve a good estimation of the reliability of the input at each step by using the prediction (see Eq. (9)). However, in the original LSTM, the available context at each step is from the previous temporal step, i.e., the prediction can only be based on the context in the temporal dimension, thus the effectiveness of the trust gate is limited when it is added to the original LSTM. This further demonstrates the effectiveness of our ST-LSTM framework for spatio-temporal modeling of the skeleton sequences.

In addition, we investigate the effectiveness of the LSTM structure for handling the sequential data. We evaluate a baseline method (called “Temporal Average”) by averaging the features from all frames instead of using LSTM. Specifically, the geometric features are averaged over all the frames of the input sequence (i.e., the temporal ordering information in the sequence is ignored), and then the resultant averaged feature is fed to a two-layer network, followed by a softmax classifier. The performance of this scheme is much weaker than our proposed ST-LSTM with trust gate, and also weaker than the original LSTM, as shown in TABLE 13. The results demonstrate the representation strengths of the LSTM networks for modeling the dependencies and dynamics in sequential data, when compared to traditional temporal aggregation methods of input sequences.

4.13. Evaluation of the Last-to-first Link Scheme

In this section, we evaluate the effectiveness of the last-to-first link in our model (see section 3.6). The results in TABLE 14 show the advantages of using the last-to-first link in improving the final action recognition performance.

5. Conclusion

In this paper, we have extended the RNN-based action recognition method to both spatial and temporal domains. Specifically, we have proposed a novel ST-LSTM network which analyzes the 3D locations of skeletal joints at each frame and at each processing step. A skeleton tree traversal method based on the adjacency graph of body joints is also proposed to better represent the structure of the input sequences and to improve the performance of our network by connecting the most related joints together in the input sequence. In addition, a new gating mechanism is introduced to improve the robustness of our network against the noise in input sequences. A multi-modal feature fusion method is also proposed for our ST-LSTM framework. The experimental results have validated the contributions and demonstrated the effectiveness of our approach which achieves better performance over the existing state-of-the-art methods on seven challenging benchmark datasets.

Acknowledgement

This work was carried out at Rapid-Rich Object Search (ROSE) Lab, Nanyang Technological University. ROSE Lab is supported by the National Research Foundation, Singapore, under its IDM Strategic Research Programme.

References

- [1] F. Zhu, L. Shao, J. Xie, and Y. Fang, “From handcrafted to learned representations for human action recognition: a survey,” *Image and Vision Computing*, 2016.
- [2] L. L. Presti and M. La Cascia, “3d skeleton-based human action classification: A survey,” *Pattern Recognition*, 2016.
- [3] F. Han, B. Reily, W. Hoff, and H. Zhang, “Space-time representation of people based on 3d skeletal data: a review,” *arXiv*, 2016.
- [4] L. Xia, C. Chen, and J. Aggarwal, “View invariant human action recognition using histograms of 3d joints,” in *CVPRW*, 2012.
- [5] X. Yang and Y. Tian, “Effective 3d action recognition using eigen-joints,” *Journal of Visual Communication and Image Representation*, 2014.
- [6] R. Vemulapalli, F. Arrate, and R. Chellappa, “Human action recognition by representing 3d skeletons as points in a lie group,” in *CVPR*, 2014.
- [7] G. Evangelidis, G. Singh, and R. Horaud, “Skeletal quads: Human action recognition using joint quadruples,” in *ICPR*, 2014.
- [8] E. Ohn-Bar and M. Trivedi, “Joint angles similarities and hog² for action recognition,” in *CVPRW*, 2013.
- [9] J. Luo, W. Wang, and H. Qi, “Group sparsity and geometry constrained dictionary learning for action recognition from depth maps,” in *ICCV*, 2013.
- [10] A. Graves, A.-r. Mohamed, and G. Hinton, “Speech recognition with deep recurrent neural networks,” in *ICASSP*, 2013.
- [11] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *NIPS*, 2014.
- [12] T. Mikolov, S. Kombrink, L. Burget, J. H. Černocký, and S. Khudanpur, “Extensions of recurrent neural network language model,” in *ICASSP*, 2011.
- [13] M. Sundermeyer, R. Schlüter, and H. Ney, “Lstm neural networks for language modeling,” in *INTERSPEECH*, 2012.
- [14] G. Mesnil, X. He, L. Deng, and Y. Bengio, “Investigation of recurrent-neural-network architectures and learning methods for spoken language understanding,” in *INTERSPEECH*, 2013.
- [15] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” in *CVPR*, 2015.
- [16] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *ICML*, 2015.
- [17] N. Srivastava, E. Mansimov, and R. Salakhudinov, “Unsupervised learning of video representations using lstms,” in *ICML*, 2015.
- [18] B. Singh, T. K. Marks, M. Jones, O. Tuzel, and M. Shao, “A multi-stream bi-directional recurrent neural network for fine-grained action detection,” in *CVPR*, 2016.
- [19] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena, “Structural-rnn: Deep learning on spatio-temporal graphs,” in *CVPR*, 2016.
- [20] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, “Social lstm: Human trajectory prediction in crowded spaces,” in *CVPR*, 2016.
- [21] Z. Deng, A. Vahdat, H. Hu, and G. Mori, “Structure inference machines: Recurrent neural networks for analyzing relations in group activity recognition,” in *CVPR*, 2016.

- [22] M. S. Ibrahim, S. Muralidharan, Z. Deng, A. Vahdat, and G. Mori, "A hierarchical deep temporal model for group activity recognition," in *CVPR*, 2016.
- [23] S. Ma, L. Sigal, and S. Sclaroff, "Learning activity progression in lstms for activity detection and early detection," in *CVPR*, 2016.
- [24] B. Ni, X. Yang, and S. Gao, "Progressively parsing interactional objects for fine grained action detection," in *CVPR*, 2016.
- [25] Y. Li, C. Lan, J. Xing, W. Zeng, C. Yuan, and J. Liu, "Online human action detection using joint classification-regression recurrent neural networks," *arXiv*, 2016.
- [26] J. Yue-Hei Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici, "Beyond short snippets: Deep networks for video classification," in *CVPR*, 2015.
- [27] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *CVPR*, 2015.
- [28] Q. Li, Z. Qiu, T. Yao, T. Mei, Y. Rui, and J. Luo, "Action recognition by learning deep multi-granular spatio-temporal video representation," in *ICMR*, 2016.
- [29] Z. Wu, X. Wang, Y.-G. Jiang, H. Ye, and X. Xue, "Modeling spatial-temporal clues in a hybrid deep learning framework for video classification," in *ACM MM*, 2015.
- [30] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *CVPR*, 2015.
- [31] V. Veeriah, N. Zhuang, and G.-J. Qi, "Differential recurrent neural networks for action recognition," in *ICCV*, 2015.
- [32] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "Ntu rgb+d: A large scale dataset for 3d human activity analysis," in *CVPR*, 2016.
- [33] M. Meng, H. Drira, M. Daoudi, and J. Boonaert, "Human-object interaction recognition by learning the distances between the object and the skeleton joints," in *FG*, 2015.
- [34] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Learning actionlet ensemble for 3d human action recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2014.
- [35] A. Shahroudy, T. T. Ng, Q. Yang, and G. Wang, "Multimodal multipart learning for action recognition in depth videos," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016.
- [36] J. Wang and Y. Wu, "Learning maximum margin temporal warping for action recognition," in *ICCV*, 2013.
- [37] R. Vemulapalli and R. Chellapa, "Rolling rotations for recognizing human actions from 3d skeletal data," in *CVPR*, 2016.
- [38] H. Rahmani, A. Mahmood, D. Q. Huynh, and A. Mian, "Real time action recognition using histograms of depth gradients and random decision forests," in *WACV*, 2014.
- [39] A. Shahroudy, G. Wang, and T.-T. Ng, "Multi-modal feature fusion for action recognition in rgb-d sequences," in *ISCCSP*, 2014.
- [40] H. Rahmani and A. Mian, "Learning a non-linear knowledge transfer model for cross-view action recognition," in *CVPR*, 2015.
- [41] I. Lillo, A. Soto, and J. Carlos Nibbles, "Discriminative hierarchical modeling of spatio-temporally composable human activities," in *CVPR*, 2014.
- [42] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black, "Towards understanding action recognition," in *ICCV*, 2013.
- [43] C. Chen, R. Jafari, and N. Kehtarnavaz, "Fusion of depth, skeleton, and inertial data for human action recognition," in *ICASSP*, 2016.
- [44] Z. Liu, C. Zhang, and Y. Tian, "3d-based deep convolutional neural network for action recognition with depth sequences," *Image and Vision Computing*, 2016.
- [45] X. Cai, W. Zhou, L. Wu, J. Luo, and H. Li, "Effective active skeleton representation for low latency human action recognition," *IEEE Transactions on Multimedia*, 2016.
- [46] A. S. Al Alwani and Y. Chahir, "Spatiotemporal representation of 3d skeleton joints-based action recognition using modified spherical harmonics," *Pattern Recognition Letters*, 2016.
- [47] L. Tao and R. Vidal, "Moving poselets: A discriminative and interpretable skeletal motion representation for action recognition," in *ICCVW*, 2015.
- [48] W. Zhu, C. Lan, J. Xing, W. Zeng, Y. Li, L. Shen, and X. Xie, "Co-occurrence feature learning for skeleton based action recognition using regularized deep lstm networks," in *AAAI*, 2016.
- [49] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, 2014.
- [50] F. G. Harvey and C. Pal, "Semi-supervised learning with encoder-decoder recurrent neural networks: Experiments with motion capture sequences," *arXiv*, 2016.
- [51] B. Mahasseni and S. Todorovic, "Regularizing long short term memory with 3d human-skeleton sequences for action recognition," in *CVPR*, 2016.
- [52] J. Liu, A. Shahroudy, D. Xu, and G. Wang, "Spatio-temporal lstm with trust gates for 3d human action recognition," in *ECCV*, 2016.
- [53] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, 1997.
- [54] B. Zou, S. Chen, C. Shi, and U. M. Providence, "Automatic reconstruction of 3d human motion pose from uncalibrated monocular video sequences based on markerless human motion tracking," *Pattern Recognition*, 2009.
- [55] Y. Yang and D. Ramanan, "Articulated pose estimation with flexible mixtures-of-parts," in *CVPR*, 2011.
- [56] N. Dalal, B. Triggs, and C. Schmid, "Human detection using oriented histograms of flow and appearance," in *ECCV*, 2006.
- [57] H. Wang, A. Kläser, C. Schmid, and C.-L. Liu, "Action recognition by dense trajectories," in *CVPR*, 2011.
- [58] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [59] G. Chéron, I. Laptev, and C. Schmid, "P-cnn: Pose-based cnn features for action recognition," in *ICCV*, 2015.
- [60] A. Graves, *Supervised Sequence Labelling with Recurrent Neural Networks*. Springer, 2012.
- [61] K. Yun, J. Honorio, D. Chattopadhyay, T. L. Berg, and D. Samaras, "Two-person interaction detection using body-pose features and multiple instance learning," in *CVPRW*, 2012.
- [62] J.-F. Hu, W. Zheng, J. Lai, and J. Zhang, "Jointly learning heterogeneous features for rgb-d activity recognition," in *CVPR*, 2015.
- [63] S. Escalera, J. González, X. Baró, M. Reyes, O. Lopes, I. Guyon, V. Athitsos, and H. Escalante, "Multi-modal gesture recognition challenge 2013: Dataset and results," in *ICMI*, 2013.
- [64] W. Li, Z. Zhang, and Z. Liu, "Action recognition based on a bag of 3d points," in *CVPRW*, 2010.
- [65] F. Offi, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, "Berkeley mhad: A comprehensive multimodal human action database," in *WACV*, 2013.
- [66] R. Collobert, K. Kavukcuoglu, and C. Farabet, "Torch7: A matlab-like environment for machine learning," in *NIPS Workshop*, 2011.
- [67] E. Cippitelli, E. Gambi, S. Spinsante, and F. Flórez-Revuelta, "Evaluation of a skeleton-based method for human activity recognition on a large-scale rgb-d dataset," in *TechAAL*, 2016.
- [68] H. Rahmani and A. Mian, "3d action recognition from novel viewpoints," in *CVPR*, 2016.
- [69] Y. Zhu, W. Chen, and G. Guo, "Fusing spatiotemporal features and joints for 3d action recognition," in *CVPRW*, 2013.

- [70] R. Anirudh, P. Turaga, J. Su, and A. Srivastava, "Elastic functional coding of human actions: from vector-fields to latent variables," in *CVPR*, 2015.
- [71] R. Slama, H. Wannous, M. Daoudi, and A. Srivastava, "Accurate 3d action recognition using learning on the grassmann manifold," *Pattern Recognition*, 2015.
- [72] S. Jetley and F. Cuzzolin, "3d activity recognition using motion history and binary shape templates," in *ACCVW*, 2014.
- [73] M. Devanne, H. Wannous, S. Berretti, P. Pala, M. Daoudi, and A. Del Bimbo, "Space-time pose representation for 3d human action recognition," in *ICIAP*, 2013.
- [74] M. Devanne, H. Wannous, S. Berretti, P. Pala, M. Daoudi, and A. Bimbo, "3-d human action recognition by shape analysis of motion trajectories on riemannian manifold," *IEEE Transactions on Cybernetics*, 2015.
- [75] M. Jiang, J. Kong, G. Bebis, and H. Huo, "Informative joints based human action recognition using skeleton contexts," *Signal Processing: Image Communication*, 2015.
- [76] A. Chrungoo, S. Manimaran, and B. Ravindran, "Activity recognition for natural human robot interaction," in *ICSR*, 2014.
- [77] C. Wang, Y. Wang, and A. L. Yuille, "Mining 3d key-pose-motifs for action recognition," in *CVPR*, 2016.
- [78] Y. Ji, G. Ye, and H. Cheng, "Interactive body part contrast mining for human interaction recognition," in *ICMEW*, 2014.
- [79] W. Li, L. Wen, M. Choo Chuah, and S. Lyu, "Category-blind human action recognition: a practical recognition system," in *ICCV*, 2015.
- [80] A. Savitzky and M. Golay, "Smoothing and differentiation of data by simplified least squares procedures," *Analytical chemistry*, 1964.
- [81] J.-F. Hu, W.-S. Zheng, L. Ma, G. Wang, and J. Lai, "Real-time rgb-d activity prediction by soft regression," in *ECCV*, 2016.
- [82] H. Wang, W. Wang, and L. Wang, "Hierarchical motion evolution for action recognition," in *ACPR*, 2015.
- [83] B. Fernando, E. Gavves, J. Oramas, A. Ghodrati, and T. Tuytelaars, "Modeling video evolution for action recognition," in *CVPR*, 2015.
- [84] A. Yao, L. Van Gool, and P. Kohli, "Gesture recognition portfolios for personalization," in *CVPR*, 2014.
- [85] J. Wu, J. Cheng, C. Zhao, and H. Lu, "Fusing multi-modal features for gesture recognition," in *ICMI*, 2013.
- [86] T. Pfister, J. Charles, and A. Zisserman, "Domain-adaptive discriminative one-shot learning of gestures," in *ECCV*, 2014.
- [87] M. A. Gowayyed, M. Torki, M. E. Hussein, and M. El-Saban, "Histogram of oriented displacements (hod): Describing trajectories of human joints for action recognition," in *IJCAI*, 2013.
- [88] I. Lillo, J. Carlos Niebles, and A. Soto, "A hierarchical pose-based approach to complex action understanding using dictionaries of actionlets and motion poselets," in *CVPR*, 2016.
- [89] F. Ofli, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, "Sequence of the most informative joints (smij): A new representation for human skeletal action recognition," *Journal of Visual Communication and Image Representation*, 2014.
- [90] S. Vantigodi and R. V. Babu, "Real-time human action recognition from motion capture data," in *NCVPRIPG*, 2013.
- [91] S. Vantigodi and V. B. Radhakrishnan, "Action recognition from motion capture data using meta-cognitive rbf network classifier," in *ISSNIP*, 2014.
- [92] I. Kapsouras and N. Nikolaidis, "Action recognition on motion capture data using a dynemes and forward differences representation," *Journal of Visual Communication and Image Representation*, 2014.