

Learning Sequence Representations by Non-local Recurrent Neural Memory

Wenjie Pei^{1†}, Xin Feng^{1†}, Canmiao Fu², Qiong Cao³, Guangming Lu^{1*} and Yu-Wing Tai⁴

¹Department of Computer Science, Harbin Institute of Technology at Shenzhen, Shenzhen, 518057, Guangdong, China.

²Tecent, China.

³JD Explore Academy, China.

⁴Kuaishou Technology, China.

*Corresponding author(s). E-mail(s): luguangm@hit.edu.cn;

Contributing authors: wenjiecoder@outlook.com; fengx_hit@outlook.com; fcm@pku.edu.cn; mathqiong2012@gmail.com; yuwing@gmail.com;

†These authors contributed equally to this work.

Abstract

The key challenge of sequence representation learning is to capture the long-range temporal dependencies. Typical methods for supervised sequence representation learning are built upon recurrent neural networks to capture temporal dependencies. One potential limitation of these methods is that they only model one-order information interactions explicitly between adjacent time steps in a sequence, hence the high-order interactions between nonadjacent time steps are not fully exploited. It greatly limits the capability of modeling the long-range temporal dependencies since the temporal features learned by one-order interactions cannot be maintained for a long term due to temporal information dilution and gradient vanishing. To tackle this limitation, we propose the Non-local Recurrent Neural Memory (*NRNM*) for supervised sequence representation learning, which performs non-local operations by means of self-attention mechanism to learn full-order interactions within a sliding temporal memory block and models global interactions between memory blocks in a gated recurrent manner. Consequently, our model is able to capture long-range dependencies. Besides, the latent high-level features contained in high-order interactions can be distilled by our model. We validate the effectiveness and generalization of our *NRNM* on three types of sequence applications across different modalities, including sequence classification, step-wise sequential prediction and sequence similarity learning. Our model compares favorably against other state-of-the-art methods specifically designed for each of these sequence applications.

Keywords: Sequence representation learning, non-local, recurrent, neural memory, long-range temporal dependencies.

1 Introduction

Supervised sequence representation learning aims to build models to learn effective features incorporating both the single-frame information for each time step and the temporal dependencies between different

time steps from variety of sequence data such as video data, speech data or text data via supervised learning. It has extensive applications ranging from computer vision (Pei et al, 2017; Shahrudy et al, 2016) to natural language processing (Grave et al, 2017; Vaswani et al, 2017) and biological engineering like protein

structure prediction (Wang et al, 2016; Li and Yu, 2016). The key challenge in sequence representation learning is to capture the long-range temporal dependencies, which are used to further learn high-level feature representation for the whole sequence.

Most state-of-the-art methods for supervised sequence modeling are built upon the recurrent neural network (RNN) (Rumelhart et al, 1988), which has been validated its effectiveness (Sak et al, 2014; Zhang et al, 2018). One crucial limitation of the vanilla-RNN is the gradient-vanishing problem along the temporal domain, which results in the inability to model long-term dependencies. This limitation is then substantially mitigated by gated recurrent networks such as GRU (Cho et al, 2014) and LSTM (Hochreiter and Schmidhuber, 1997), which employ learnable gates to selectively retain information in the memory or hidden states. The memory-based methods for sequence representation learning (Santoro et al, 2018; Sukhbaatar et al, 2015; Weston et al, 2015) are further proposed to address the issue of limited memory of recurrent networks. However, a potential drawback of these methods is that they only model explicitly the information interactions between adjacent time steps in the sequence, whereas the high-order interactions between nonadjacent time steps are not fully exploited. This drawback gives rise to two negative consequences: 1) the high-level features contained in the high-order interactions between nonadjacent time steps cannot be distilled; 2) it greatly limits the modeling of long-range temporal dependencies since the temporal features learned by one-order interactions cannot be maintained in a long term due to the information dilution and gradient vanishing during recurrent operations in the temporal domain.

Inspired by non-local methods (Buades et al, 2005; Wang et al, 2018) which aim to explore potential interactions between all pairs of feature portions, we propose to perform non-local operations to model the high-order interactions between non-adjacent time steps in a sequence. Since the non-local operations facilitate the temporal feature propagation and thus substantially alleviate the vanishing-gradient problem, the captured high-order interactions can be used to not only distill latent high-level features which are hardly learned by typical sequence modeling methods focusing on one-order interactions, but also contribute to modeling long-range temporal dependencies. We leverage self-attention mechanism (Vaswani et al, 2017) to perform such non-local operations, which

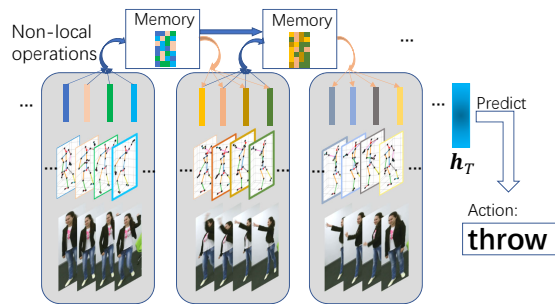


Fig. 1: The action ‘throw’ is performed slowly in this video sample. To recognize the action correctly, our proposed model (*NRNM*) performs non-local operations within each memory block to learn high-order interactions between hidden states of different time steps. Meanwhile, the global interactions between memory blocks are modeled in a gated recurrent manner. The learned memory states are in turn leveraged to refine the hidden states in future time steps. Thus, the long-range dependencies can be captured. Our model is able to predict the action correctly based on the hidden state (h_T) in the last time step. By contrast, it is challenging for typical recurrent sequence models like LSTM, which tend to perform recognition relying mainly on the frames in the end of the video due to limited capability of modeling temporal dependencies.

allows the feature learning for each time step to attend to the features of all other steps.

Exploring full-order interactions between all time steps for a long sequence is computationally expensive and also not necessary due to information redundancy, thus we model full-order interactions by non-local operations within a temporal block (a segment of sequence) and slide the block to recurrently update the extracted information. More specifically, we propose the Non-local Recurrent Neural Memory (*NRNM*) to perform the block-wise non-local operations by means of self-attention for learning full-order interactions within each memory block and capture the local but high-resolution temporal dependencies. Meanwhile, the global interactions between adjacent blocks are captured by updating the memory states in a gated recurrent manner when sliding the memory cell. Consequently, the long-range dependencies can be captured, thereby yielding effective sequence representations.

The learned memory embeddings are in turn leveraged to refine the hidden states to achieve final sequence representations in future time steps. Defining the directly involved time steps for learning the hidden states of a specified time step as the *direct interacting field* (similar to definition of the receptive field in convolutional networks), then the *direct interacting field* for learning hidden states by typical recurrent models (like LSTM) is 2 steps, including

the preceding and the current time steps. In contrast, *NRNM* enlarges such *direct interacting field* by the size of memory block when learning hidden states for a time step. Thus, another remarkable merit of the proposed *NRNM* is that it enables the learning process of hidden states for each time step to have (temporally) much longer *direct interacting field* than conventional recurrent sequence models, which potentially yields high-level features incorporating useful information from high-order interactions among all time steps of the corresponding interacting field (memory block size). Such learning process resembles the way of distilling high-level semantic features from large receptive fields by deep convolutional neural networks.

Figure 1 illustrates our method by an example of action recognition, in which the action ‘throw’ is performed slowly and thus the long-range temporal features in this video are crucial for recognizing the action. Our designed *NRNM* is able to learn effective sequence representations that captures long-range temporal dependencies and thus recognize the action ‘throw’ correctly. By contrast, typical recurrent sequence models (like LSTM) tend to perform recognition relying primarily on the frames in the end of the video due to limited capability of modeling temporal features, and therefore make false predictions.

Compared to typical supervised sequence models for sequence representation learning, especially recurrent networks with memory mechanism, our *NRNM* benefits from the following advantages:

- It is able to model 1) the local full-order interactions between all time steps within a segment of sequence (memory block) and 2) the global interactions between memory blocks. Thus, it can capture much longer temporal dependencies than existing recurrent sequence models.
- The proposed *NRNM* enlarges the *direct interacting field* for learning hidden states of each time step, which allows the model to learn high-level semantic features potentially missed by conventional sequence models.
- The *NRNM* cell can be seamlessly integrated into any existing sequence models with stepwise latent structure to enhance the power of sequence representation learning. The integrated model can be trained in an end-to-end manner.
- We evaluate the learned sequence representations by our model on three types of sequence applications across different modalities including 1)

sequence classification (action recognition and sentiment analysis), 2) step-wise sequential prediction (protein secondary structure prediction) and 3) sequence similarity learning (action similarity learning from videos). These extensive experiments demonstrate that our model compares favorably against other recurrent models for sequence representation learning.

An earlier conference version of this paper appeared in (Fu et al, 2019a). Compared to the prior version, this longer article is improved in three aspects. First, more comprehensive explanations and analysis on the proposed *NRNM* are provided in Section 1 and Section 3, including the introduction of the newly proposed multi-scale dilated *NRNM* mechanism that performs *NRNM* operations covering different scales of *direct interacting field* with different strides. Second, this article makes more detailed review of related work in Section 2, especially discussing the related methods for capturing long-range dependencies in Computer Vision. Third, substantial additional experiments are conducted to evaluate our model on more types of sequence applications across different modalities. Apart from the evaluation of our model on sequence classification in the conference version, we also validate the effectiveness of our *NRNM* on step-wise sequence prediction by performing protein secondary structure prediction (Section 4.3), and on sequence similarity learning between a pair of input sequences (Section 4.4). Besides, more experimental investigations and ablation study on *NRNM* are conducted to obtain more insights into the model.

2 Related work

In this section, we review the related work to our method. We first summarize two classical types of methods for supervised sequence modeling: graphical models and recurrent networks. Then we discuss the methods that design specialized memory structure upon recurrent networks and compare them to our method. Finally, we discuss various methods for capturing long-range dependencies in Computer Vision, including two typical ways of applying self-attention to vision models and other specialized models for modeling long-term information in vision tasks.

Graphical sequence models. The conventional graphical models for sequence modeling can be roughly divided into two categories: generative

and discriminative models. A well-known example of generative model is Hidden Markov Model (HMM) (Rabiner, 1989), which models sequence data in a chain of latent k -nomial features. Discriminative graphical models model the distribution over all class labels conditioned on the input data. Conditional Random Fields (CRF) (Lafferty et al, 2001) is a discriminative model for sequential predictions by modeling the linear mapping between observations and labels. To tackle its limitation of linear mapping, many nonlinear CRF-variants (Morency et al, 2007; Pei et al, 2018; Van Der Maaten et al, 2011) and Conditional Neural Fields (Peng et al, 2009) are proposed. The disadvantages of graphical model compared to recurrent networks lie in the hard optimization and limited capability of temporal modeling. Our model is designed based on recurrent networks.

Recurrent Networks. Recurrent Neural Network (Rumelhart et al, 1988) learns a hidden representation for each time step by taking into account both current and previous information. Benefited from its advantages such as easy training and temporal modeling, it has been successfully applied to, amongst others, handwriting recognition (Bertolami et al, 2009) and speech recognition (Sak et al, 2014). However, the key limitation of vanilla-RNN is the gradient vanishing problem during training (Hochreiter et al, 2001) and thus cannot model long-range temporal dependencies. This limitation is alleviated by gated recurrent networks such as Long Shot-Term Memory (LSTM) (Hochreiter and Schmidhuber, 1997) and Gate Recurrent Unit (GRU) (Cho et al, 2014), which selectively retain information by learnable gates. Such models is further improved either by attention mechanism (Bahdanau et al, 2015) to improve the gated structure (Tan et al, 2018; Zhang et al, 2018), or by integrating the convolutional structure to enhance the feature learning (Wang et al, 2016; Li and Yu, 2016). Nevertheless, a potential limitation of these models is that they only model explicitly one-order interactions between adjacent time steps, hence the high-order interactions between nonadjacent time steps are not fully captured. Our model is proposed to circumvent this drawback by performing non-local operations by means of self-attention to model full-order interactions in a block-wise manner. Meanwhile, the global interactions between blocks are modeled by a gated recurrent mechanism. Thus, our model is able to model long-range temporal dependencies and distill high-level features that are contained in high-order interactions.

Memory-based recurrent networks. Memory networks are first proposed to rectify the drawback of limited memory of recurrent networks (Sukhbaatar et al, 2015; Weston et al, 2015), which are then extended for various tasks, especially in natural language processing. Most of these models build external memory units upon a basis model to augment its memory (Graves et al, 2014; Santoro et al, 2018; Sukhbaatar et al, 2015; Weston et al, 2015). In particular, attention mechanism (Bahdanau et al, 2015) is employed to filter the information flow from memory (Grave et al, 2017; Kumar et al, 2016; Sukhbaatar et al, 2015; Xiong et al, 2016). The primary difference between these memory-based recurrent networks and our model is that these models focus on augmenting the memory size to memorize more information for reference while our model aims to model high-order interactions between different time steps in a sequence, which is not concerned by existing memory-based networks.

Capturing long-range dependencies in Computer Vision. Self-attention mechanism (Vaswani et al, 2017) has achieved great success in Natural Language Processing (Devlin et al, 2019; Dai et al, 2019) due to its merit of capturing long-range dependencies. It has been introduced into Computer Vision to enhance the modeling of long-range dependencies during feature learning either in the spatial domain for image data or in the temporal domain for video data, which is an intrinsic limitation of Convolutional Neural Networks (CNNs) (Wang et al, 2018; Ramachandran et al, 2019), the fundamental framework for modern Computer Vision systems. There are two typical ways to apply self-attention to Computer Vision. 1) Self-attention is used as an add-on to augment existing CNN models and capture long-range dependencies upon deep CNN features (Wang et al, 2018; Yin et al, 2020; Bello et al, 2019; Cao et al, 2019; Fu et al, 2019b; Srinivas et al, 2021). 2) Replacing the convolutional operation, self-attention is used as a stand-alone primitive for building vision models (Hu et al, 2019; Ramachandran et al, 2019; Zhao et al, 2020; Liu et al, 2021).

A prominent example in the first way of applying self-attention is the Non-local Neural Networks (Wang et al, 2018), which designs a specific non-local operation, using self-attention as an instantiating form, to capture the long-range dependencies in both spatial and temporal domains for video data. Such non-local operation is further extensively studied

and applied to, amongst others, video understanding (Feichtenhofer et al, 2019; Xie et al, 2018), image segmentation (Fu et al, 2019c; He et al, 2017; Carion et al, 2020) and image synthesis (Brock et al, 2019; Zhang et al, 2019). Compared with these methods, our method also performs non-local operations by means of self-attention to capture long-range temporal dependencies with two key differences. First, our method operates on sequence data instead of image data. Second, our model is built upon recurrent networks instead of CNNs.

The second classical way of applying self-attention to Computer Vision is to use (local) self-attention instead of convolution as the stand-alone fundamental building block for vision models (Hu et al, 2019; Ramachandran et al, 2019; Zhao et al, 2020; Liu et al, 2021). It has been shown that the resulting fully self-attentional models can achieve favorable performance against the convolutional counterparts in various vision tasks (Liu et al, 2021). Similar to these models, our method also performs local self-attention considering the computational complexity. Nonetheless, the primary difference is that our method focuses on capturing the temporal dependencies in sequence data rather than the spatial dependencies between pixels in image data.

Apart from self-attention, the long-range temporal information in vision tasks can be also captured by specialized models designed upon existing deep vision models. A typical example is the Long-Term Feature Bank model (LFB) (Wu et al, 2019), which proposes a long-term feature bank for video understanding to store supportive information spanning long-range temporal context. The feature bank serves as an auxiliary memory for the backbone module, which is akin to our method. Nevertheless, our model differs from LFB in that our model focuses on modeling high-order temporal interactions including both full-order intra-memory and global inter-memory interactions for arbitrary sequence data. By contrast, LFB (Wu et al, 2019) enumerates time-indexed features in the feature bank, such as 3D convolutional features, so as to provide a long-term view for its backbone video model. Thus, the high-order temporal interactions are not explicitly modeled by its feature bank. As a result, our model is potentially more advantageous than LFB in terms of capturing high-order temporal dependencies between different time steps and distilling high-level semantic features.

3 Non-local Recurrent Neural Memory

Given a sequence as input, our Non-local Recurrent Neural Memory (*NRNM*) is designed as a memory module to capture the long-term temporal dependencies and distill high-level sequence features incorporating multi-step interacting information. Specifically, *NRNM* models the high-order interactions between non-adjacent time steps of a sequence leveraging self-attention (Vaswani et al, 2017) in a non-local manner, which is in contrast to conventional recurrent sequential methods that only explicitly model one-order interactions between adjacent time steps.

In this section we will first present an overview of the whole proposed method. Then we will elaborate on the cell structure of our *NRNM*. Next we describe how *NRNM* and the LSTM backbone perform sequence modeling collaboratively. Finally, we show how to train the model in an end-to-end manner when applying the learned sequence representations by our model to different sequence applications including sequence classification, step-wise sequential prediction and sequence similarity learning.

3.1 Overview of the Method

Our proposed *NRNM* serves as an individual functional module for sequence representation learning, which can be seamlessly integrated into any existing sequential models with stepwise latent structure to enhance sequence modeling. As illustrated in Figure 2, we build our *NRNM* upon an LSTM backbone as an instantiation.

Typical recurrent models for sequence representation learning follow the similar modeling paradigm: learning hidden states as feature representation for each time step by incorporating both the input information of current time step and hidden states of the preceding time step. Consider an input sequence $\mathbf{x}_{1,\dots,T} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ of length T in which $\mathbf{x}_t \in \mathbb{R}^D$ denotes the observation at the t -th time step. The hidden state \mathbf{h}_t at the t -th time step is modeled by a recurrent model $\mathcal{F}_{\text{backbone}}$ (e.g., LSTM in the instantiation of our model) as:

$$\mathbf{h}_t = \mathcal{F}_{\text{backbone}}(\mathbf{h}_{t-1}, \mathbf{x}_t), \quad (1)$$

where \mathbf{x}_t denotes the input of the t -th time step.

All hidden states are learned sequentially in such a recurrent manner that each hidden state is expected

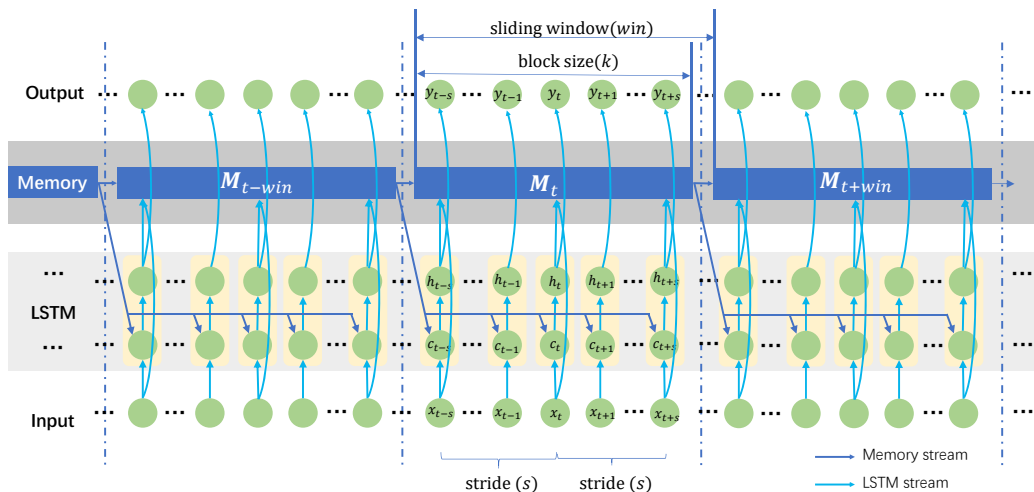


Fig. 2: Architecture of our method. Our proposed *NRNM* is built upon the LSTM backbone to learn full-order interactions between LSTM hidden states of different time steps within each memory block by non-local operations. Meanwhile, the global interactions between memory blocks are modeled in a gated recurrent manner. The learned memory states are in turn used to refine the LSTM hidden states in future time steps.

to contain both temporal dynamic features and static features in all previous steps. However, most existing recurrent methods only model the one-order interactions explicitly between two adjacent time steps as shown in Equation 1. As a result, these methods can hardly learn the long-term temporal features and thus the temporal memories are limited due to gradient vanishing problem along the temporal domain. On the other hand, as we defined before, *direct interacting field* corresponds to the directly involved time steps when learning hidden states for a time step, then the *direct interacting field* for typical recurrent models like LSTM is 2 steps including the preceding and the current time steps. Thus, the high-level semantic features covering multi-step *direct interacting field* cannot be learned by typical recurrent models.

To address these limitations, we design *NRNM* as a memory module upon an existing recurrent model (e.g., LSTM backbone in Figure 2). Specifically, it maintains a memory cell and performs non-local operations by means of self-attention on a segment of the input sequence (termed as non-local memory block) to model full-order interactions among time steps within this segment. The *NRNM* cell slides temporally along the sequence to learn memory embeddings in a block-wise manner. Such blocking design is analogous to DenseNet (Huang et al, 2017) which performs dense connection (a form of non-local operation) in blocks. Furthermore, the proposed *NRNM* captures the global interactions between adjacent memory blocks

by updating the memory state recurrently, which is consistent with the hidden state update of the LSTM backbone.

Note that the non-local operation in this work is defined consistently with Wang et al (Wang et al, 2018): the features for a time step in a sequence are learned by interacting with all time steps in a sequence. We perform non-local operations in a memory block, namely a temporal segment of the input sequence, to distill high-level features contained in full-order interactions between time steps within the memory block. Thus the feature in each time step in a memory block is learned with much larger *direct interacting field* (equal to the size of memory block) than the conventional recurrent sequence models. Formally, our proposed *NRNM* learns embeddings for the memory block at time step t by performing non-local operations:

$$\mathbf{M}_t = \mathcal{F}_{NRNM}(\mathbf{v}_{t-k+1}, \dots, \mathbf{v}_t), \quad (2)$$

where \mathcal{F}_{NRNM} is the nonlinear transformation function performed by *NRNM*, and $\mathbf{v}_{t-k+1}, \dots, \mathbf{v}_t$ are information per each time step within the memory block at t -th time step with block size k . The model structure of *NRNM* will be elaborated in Section 3.2.

The obtained memory embeddings are leveraged to perform sequence modeling by *NRNM* and LSTM backbone together to achieve the final sequence representations used for downstream tasks. Denoting the preceding memory embedding right before time step t

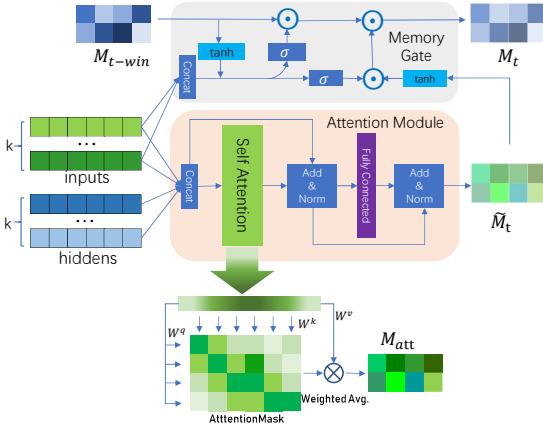


Fig. 3: Structure of *NRNM* cell. Taken the hidden states of LSTM backbone and initial input features as input, *NRNM* performs non-local interacting operations adopting self-attention mechanism to learn local full-order interactions among time steps in a memory block. The obtained memory states are updated in a gated recurrent manner to model the temporal evolution between adjacent memory blocks and thus capture global temporal dependencies.

as \mathbf{M}_{t-win} (*win* is the size of sliding window of the memory cell), the sequence representations \mathbf{r}_t at time step t is then obtained by:

$$\mathbf{r}_t = \mathcal{F}_{SM}(\mathbf{M}_{t-win}, \mathbf{h}_t). \quad (3)$$

Herein, \mathcal{F}_{SM} indicates the sequence modeling process in our model, which will be discussed in Section 3.3.

3.2 NRNM Cell

Next we elaborate on the structure of *NRNM* cell and show how to learn memory embeddings.

The goal of *NRNM* is to distill high-level memory embeddings from a memory block (a segment of the input sequence) by modeling the full-order interactions among all time step in the block. Formally, *NRNM* cell learns memory embedding $\widetilde{\mathbf{M}}_t$ for a block with block size k covering the temporal interval $[t - k + 1, t]$ by refining the underlying information contained in this time interval. Specifically, we consider two types of source information for *NRNM* cell: 1) the learned hidden representations in this time interval $[\mathbf{h}_{t-k+1}, \dots, \mathbf{h}_t]$ by the LSTM backbone; 2) the original input features $[\mathbf{x}_{t-k+1}, \dots, \mathbf{x}_t]$. Hence Equation 2 is re-formulated as:

$$\widetilde{\mathbf{M}}_t = \mathcal{F}_{NRNM}([\mathbf{h}_{t-k+1}, \dots, \mathbf{h}_t], [\mathbf{x}_{t-k+1}, \dots, \mathbf{x}_t]), \quad (4)$$

Here we incorporate the input feature \mathbf{x} which is already assimilated in the hidden representation \mathbf{h} of the basis LSTM backbone since we aim to explore the

latent interactions between hidden representations and input features in the current block (i.e., the interval $[t - k + 1, t]$).

NRNM performs non-local operations to model full-order interactions within a memory block. There are multiple feasible ways to perform non-local operations such as the dense convolutional structure in DenseNet (Huang et al, 2017) or the global attention mechanism in Non-local Neural Networks (Wang et al, 2018). We design the *NRNM* in the similar way as Self-Attention mechanism (Vaswani et al, 2017) due to its effectiveness in language modeling. Figure 3 illustrates the structure of *NRNM*. In particular, the distilled information for a time step is constructed by attending to the source information of other time steps within the same block according to the compatibility between this step and other steps:

$$\begin{aligned} \mathbf{C} &= \text{Concat}([\mathbf{h}_{t-k+1}, \dots, \mathbf{h}_t], [\mathbf{x}_{t-k+1}, \dots, \mathbf{x}_t]), \\ \mathbf{Q}, \mathbf{K}, \mathbf{V} &= (\mathbf{W}^q, \mathbf{W}^k, \mathbf{W}^v)\mathbf{C}, \\ \mathbf{W}_{att} &= \text{softmax}(\mathbf{Q}\mathbf{K}^\top / \sqrt{m}), \\ \mathbf{M}_{att} &= \mathbf{W}_{att}\mathbf{V}. \end{aligned} \quad (5)$$

Herein, $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are queries, keys and values of Self-Attention transformed by parameters $\mathbf{W}^q, \mathbf{W}^k, \mathbf{W}^v$ from the source information \mathbf{C} respectively. \mathbf{W}_{att} is the derived attention weights calculated by dot-product attention scheme scaled by the memory hidden size m . We apply multi-head attention scheme (Vaswani et al, 2017) to learn different memory embeddings in parallel. The obtained attention embeddings \mathbf{M}_{att} is then fed into two residual-connection layers and one fully-connected layer (with non-linear activation function) to achieve the memory embedding $\widetilde{\mathbf{M}}_t$:

$$\begin{aligned} \mathbf{M}^i &= f_{\text{Norm}}(\mathbf{C}^i + \mathbf{M}_{att}^i), \\ \widetilde{\mathbf{M}}_t &= f_{\text{Norm}}\left(\sum_{i=t-k+1}^t \mathbf{M}^i + f_{\text{fc-layer}}\left(\sum_{i=t-k+1}^t \mathbf{M}^i\right)\right). \end{aligned} \quad (6)$$

Herein, f_{Norm} denotes the layer normalization (Ba et al, 2016) which is particularly designed for recurrent neural networks, and $f_{\text{fc-layer}}$ indicates the transformation by the fully-connected layer. Both the layer normalization and the residual connections are adopted to ease gradient propagation and thereby facilitate training.

Intuition. The rationale behind this design is that *NRNM* serves like a feature distiller mounted upon the hidden layer of LSTM backbone. It incorporates the useful information within a memory block, and further learns high-level features by modeling the full-order interactions with non-local attention mechanism. The source information is composed of $2k$

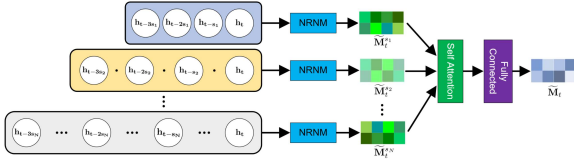


Fig. 4: Illustration of multi-scale dilated *NRNM* mechanism. Multi-scale *NRNM* operations with different strides are performed to cover different length of *direct interacting field* and thus capture temporal features from different lengths of segments in the sequence. The stride for the n -th *NRNM* $\widetilde{M}_t^{s_n}$ is s_n .

information units: k LSTM hidden states and k input features. Each information unit of the obtained memory embedding \widetilde{M}_t is constructed by attending into each of these $2k$ source information units while the size of memory embedding \widetilde{M}_t can be customized via the parametric transformation. In this way, the full-order latent interactions between the source information units are explored in a non-local way. Another benefit of such non-local operations is that it encourages latent features to be reused and thus alleviates the gradient-vanishing problem in the time domain, which is always suffered by recurrent networks.

Since the hidden states of the LSTM backbone already contain history information by recurrent structure, in practice we use a striding scheme to select hidden states as the source information for *NRNM* cell to avoid potential information redundancy and improve the modeling efficiency. For instance, we pick hidden states every s time steps in the temporal interval $[t - k + 1, t]$ for the source information, given stride $= s$.

Multi-scale dilated *NRNM*. The striding scheme for selecting the time steps for the source information within a memory block for *NRNM* cell resembles the dilated convolution (Yu and Koltun, 2015) in deep convolutional networks. Incorporating the information from the same number of time steps, larger stride leads to larger *direct interacting field* (larger block size) but coarser local sampling in the block along the temporal domain, whilst a *NRNM* cell with smaller stride focuses on learning fine-grained interacting features for shorter temporal segments. Inspired by the parallel convolutional mechanism with variable kernel size in Inception Network (Szegedy et al, 2015), we perform multi-scale dilated *NRNM* operations with different stride s in parallel to capture features from different length of segments in the input sequence. As shown in Figure 4, we learn N *NRNM* memory embeddings with increasing strides in parallel for a given time step t and fuse them by Self-Attention mechanism and a

fully-connected layer:

$$\begin{aligned} \widetilde{M}_t^{s_n} &= \mathcal{F}_{NRNM}(\mathbf{v}_{t-s_n*k+1}, \dots, \mathbf{v}_t), n = 1, 2, \dots, N, \\ \widetilde{M}_t &= f_{fc\text{-layer}}(\text{Concat}(\text{Self-Attention}(\widetilde{M}_t^{s_1}, \dots, \widetilde{M}_t^{s_N}))). \end{aligned} \quad (7)$$

Herein, Self-Attention is implemented in the similar way as Equation 5. Note that all N *NRNM* embeddings take the same number (l) of time steps of source information as input to ensure the uniform parameter size between different *NRNM* cells, thus the block size k_{s_n} for the n -th *NRNM* cell with stride s_n is $s_n * l$.

Gated recurrent update of memory state. The obtained memory embedding \widetilde{M}_t only contains information within current temporal block $([t - k + 1, t])$. To model the temporal dependencies between adjacent memory blocks, we also update memory embeddings of *NRNM* in a gated recurrent manner, which is similar to the recurrent scheme of LSTM. Specifically, the final memory state \mathbf{M}_t for current memory block is obtained by:

$$\mathbf{M}_t = \mathbf{G}_i \odot \tanh(\widetilde{M}_t) + \mathbf{G}_f \odot \mathbf{M}_{t-win}, \quad (8)$$

where win is the sliding window size of *NRNM* cell which controls the updating frequency of memory state. \mathbf{G}_i and \mathbf{G}_f are input gate and forget gate respectively to balance the memory information flow from current time step \widetilde{M}_t and the previous memory state \mathbf{M}_{t-win} . They are modeled by measuring the compatibility between current input feature and previous memory state:

$$\begin{aligned} \mathbf{G}_i &= \text{sigmoid}(\mathbf{W}_{im} \cdot [\mathbf{x}_{t-k+1}, \dots, \mathbf{x}_t, \mathbf{M}_{t-win}] + \mathbf{B}_{im}), \\ \mathbf{G}_f &= \text{sigmoid}(\mathbf{W}_{fm} \cdot [\mathbf{x}_{t-k+1}, \dots, \mathbf{x}_t, \mathbf{M}_{t-win}] + \mathbf{B}_{fm}). \end{aligned} \quad (9)$$

Herein, \mathbf{W}_{im} and \mathbf{W}_{fm} are transformation matrices while \mathbf{B}_{im} and \mathbf{B}_{fm} are bias terms.

Modeling long-range temporal dependencies. We aim to capture underlying long-range temporal dependencies in a sequence by a two-pronged strategy:

- We perform non-local operations within each temporal block by *NRNM* cell to capture full-order interactions locally between different time steps and distill high-quality memory embeddings. Hence, the local but high-resolution temporal information can be captured.
- We update the memory state in a gated recurrent manner smoothly when sliding the window of memory block temporally. It is designed to capture the global temporal dependencies between memory

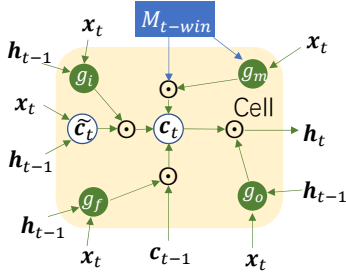


Fig. 5: The LSTM cell for current time step is updated by incorporating the memory state in the preceding time step. Thus, *NRNM* and LSTM backbone are able to perform sequence modeling collaboratively.

blocks in low resolution considering the potential information redundancy and computational efficiency.

3.3 Collaborative Sequence Modeling

Our *NRNM* can be seamlessly integrated into the LSTM backbone to enhance the power of sequence modeling. The memory state of our *NRNM* for current block is learned based on the hidden states within this block of the LSTM backbone while the obtained memory state is in turn leveraged to refine the hidden states in future time steps. Hence, our *NRNM* and the LSTM backbone are integrated seamlessly and refine each other alternately. Specifically, we incorporate the obtained memory state into the recurrent update of LSTM cell states to help refine its quality as shown in Figure 5:

$$\begin{aligned} \mathbf{v}_m &= \text{flatten}(\mathbf{M}_{t-win}), \\ \mathbf{C}_t &= \mathbf{g}_f \odot \mathbf{C}_{t-1} + \mathbf{g}_i \odot \tilde{\mathbf{C}}_t + \mathbf{g}_m \odot \mathbf{v}_m, \end{aligned} \quad (10)$$

where \mathbf{C}_{t-1} , \mathbf{C}_t and $\tilde{\mathbf{C}}_t$ are previous LSTM cell state, current cell state and candidate cell state respectively. \mathbf{v}_m is the vector flattened from the memory state \mathbf{M}_{t-win} . \mathbf{g}_f and \mathbf{g}_i are the routine forget gate and input gate of LSTM cell to balance the information flow between the current time step and previous step. All $\tilde{\mathbf{C}}_t$, \mathbf{g}_f and \mathbf{g}_i are modeled in a similar nonlinear way as a function of input feature \mathbf{x}_t and previous hidden state \mathbf{h}_{t-1} . For instance, the input gate \mathbf{g}_i is modeled as:

$$\mathbf{g}_i = \text{sigmoid}(\mathbf{W}_i \cdot \mathbf{x}_t + \mathbf{U}_i \cdot \mathbf{h}_{t-1} + \mathbf{b}_i). \quad (11)$$

In Equation 10, we enrich the modeling of the LSTM cell state \mathbf{C}_t by incorporating the cell state of *NRNM*

\mathbf{M}_{t-win} via a memory gate \mathbf{g}_m . The memory gate is constructed as a matrix to control the information flow from the memory state \mathbf{M}_{t-win} in a element-wise manner, which is derived by measuring the relevance (compatibility) between the current input features and the memory state:

$$\mathbf{g}_m = \text{sigmoid}(\mathbf{W}_m \cdot \mathbf{x}_t + \mathbf{U}_m \cdot \text{flatten}(\mathbf{M}_{t-win}) + \mathbf{b}_m), \quad (12)$$

where \mathbf{W}_m and \mathbf{U}_m are transformation matrices and \mathbf{b}_m is the bias term.

The newly constructed cell state \mathbf{C}_t are further used to derive the sequence representation at t -th time step \mathbf{r}_t prepared for downstream tasks:

$$\mathbf{r}_t = \mathbf{g}_o \odot \tanh(\mathbf{C}_t), \quad (13)$$

where \mathbf{g}_o is the output gate which is modeled in a similar way to the input gate in Equation 11.

3.4 End-to-end Parameter Learning

The learned sequence representations $\{\mathbf{r}_t\}_{t=1, \dots, T}$ in Equation 13 for a sequence with length T can be used for any sequence prediction task. In subsequent experiments, we validate our model in three types of sequence applications with different modalities: sequence classification, step-wise sequential prediction and sequence similarity learning. Equipped with corresponding loss functions, our model can be readily trained for each of these different sequence applications in an end-to-end manner.

Sequence classification. Given a training set $\mathcal{D} = \{\mathbf{x}_{1, \dots, T^n}^n, y^n\}_{n=1, \dots, N}$ containing N sequences of length T^n and their associated labels y^n . We learn our *NRNM* and the LSTM backbone jointly in an end-to-end manner by minimizing the conditional negative log-likelihood of the training data with respect to the parameters:

$$\mathcal{L}_{\text{cls}} = - \sum_{n=1}^N \log P(y^n | \mathbf{x}_{1, \dots, T^n}^n), \quad (14)$$

where the probability of the predicted label y^n among K classes is calculated by the learned sequence representations in the last time step:

$$P(y^n | \mathbf{x}_{1, \dots, T^n}^n) = \frac{\exp(\mathbf{W}_y^\top \mathbf{r}_{T^n}^n + \mathbf{b}_{y^n})}{\sum_{i=1}^K \exp(\mathbf{W}_i^\top \mathbf{r}_{T^n}^n + \mathbf{b}_i)}. \quad (15)$$

Herein, \mathbf{W}^\top and \mathbf{b} is the parameters for linear transformation and the bias term.

Step-wise sequential prediction. Unlike the task of sequence classification which predicts a label for the whole sequence, the task of step-wise sequential prediction makes prediction for each time step (or every time stride) of sequence. Thus the loss function for this task is equal to the sum of the classification losses of predictions at all time steps:

$$\begin{aligned}\mathcal{L}_{\text{step-pre}} &= - \sum_{n=1}^N \sum_{t=1}^T \log P(y_t^n | \mathbf{x}_{1,\dots,t}^n), \\ &= - \sum_{n=1}^N \sum_{t=1}^T \frac{\exp(\mathbf{W}_{y_t^n}^\top \mathbf{r}_t^n + \mathbf{b}_{y_t^n})}{\sum_{i=1}^K \exp(\mathbf{W}_i^\top \mathbf{r}_t^n + \mathbf{b}_i)},\end{aligned}\quad (16)$$

where y_t^n is the groundtruth label for the prediction at the t -th step for the n -th samples and \mathbf{r}_t^n is learned sequence representation of the n -th samples at the t -th step .

Sequence similarity learning. The task of sequence similarity learning takes two sequences as input and aims to measure the similarity between them. To this end, we design a siamese-*NRNM* structure, which employs two *NRNMs* with shared parameters to learn sequence representations for two input sequences respectively. Then the similarity calculated between two learned sequence representations.

Suppose we are given a training set \mathcal{D} containing N pairs of sequences with associated binary similarity labels. We learn the parameters of our model by minimizing the binary cross-entropy loss (conditional negative log-likelihood):

$$\begin{aligned}\mathcal{L}_{\text{similarity}} &= - \sum_{(n_1, n_2) \in \mathcal{D}} \log P(y^{n_1, n_2} | (\mathbf{x}_{1,\dots, T^{n_1}}^{n_1}, \mathbf{x}_{1,\dots, T^{n_2}}^{n_2})), \\ &= - \sum_{(n_1, n_2) \in \mathcal{D}} \log \frac{1}{1 + \exp^{-(\mathbf{v}^\top \text{Concat}(\mathbf{r}_{T^{n_1}}^{n_1}, \mathbf{r}_{T^{n_2}}^{n_2}) + b)}},\end{aligned}\quad (17)$$

where $\mathbf{r}_{T^{n_1}}^{n_1}$ and $\mathbf{r}_{T^{n_2}}^{n_2}$ are learned sequence representations for two input sequences by our model respectively. The similarity is defined as a linear transformation performed on the concatenation of two vectorial representations, parameterized by \mathbf{v} and b .

4 Experiments

We conduct four sets of experiments on three types of sequence applications across different modalities to evaluate the performance of our *NRNM* model extensively. We first perform experiments to assess the performance of our model on sequence classification which predicts a label for the whole sequence. In particular, experiments are performed on two sequence

applications including Action Recognition and Sentiment Analysis to evaluate the robustness of our model across different modalities. Then we apply *NRNM* to the task of step-wise sequential prediction, i.e., making predictions for each time step sequentially. In the last set of experiments, we investigate the performance of our model on the task of sequence similarity learning, which involves two input sequences and aims to measure the similarity between them. Besides, investigations on the effect of hyper-parameters and different configurations of *NRNM* are also included in the experiments on action recognition (Section 4.1). Code reproducing the results of our experiments is available.¹

4.1 Experiment 1 on Sequence Classification: Action Recognition

To evaluate the performance of our proposed *NRNM* model on sequence classification, we first consider the task of action recognition in which the temporal dependencies between frames in a video are quite discriminative features.

4.1.1 Dataset and Evaluation Protocol

We evaluate our method on action recognition using two datasets: the NTU dataset (Shahroudy et al, 2016) and the Charades dataset (Sigurdsson et al, 2016). We first perform ablation study on the NTU dataset, and then compare our model with the state-of-the-art methods for action recognition on both two datasets. We use the NTU dataset with only skeleton joint information in each frame to evaluate the capability of our *NRNM* to capture the long-range temporal dependencies, while the experiments on the Charades dataset with RGB information in each frame are conducted to investigate the effectiveness of the *NRNM* on generic visual data.

NTU Dataset. the NTU dataset (Shahroudy et al, 2016) is a popular and large action recognition dataset. It is a RGB+D-based dataset containing 56,880 video sequences and 4 million frames collected from 40 distinct subjects. The dataset includes 60 action categories. 3D skeleton data (i.e. 3D coordinates of 25 body joints) is provided using Microsoft Kinect.

Since our primary goal is to evaluate the capability of *NRNM* to capture the long-range temporal dependencies in sequences in experiments, we opt for

¹<https://github.com/F-Frida/NRNM>

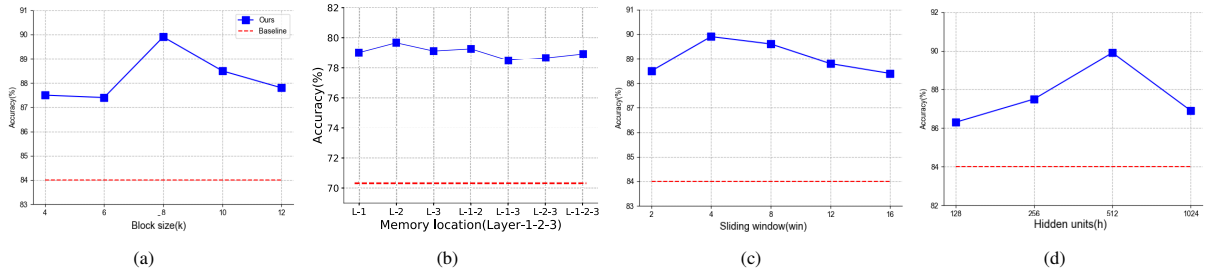


Fig. 6: Investigation on *NRNM* by performing experiments on NTU dataset to explore the effect of (a) the block size k , (b) the integrated location of *NRNM* on the LSTM backbone, (c) the sliding window size win and (d) the number of hidden units in *NRNM*. The performance of the LSTM baseline is presented in dotted red lines for reference.

NTU dataset using only 3D skeleton joint information rather than Kinetics (Kay et al, 2017) based on RGB information for action recognition since single-frame RGB information already provides much implication for action recognition and weakens the importance of temporal dependencies (Qiu and et al., 2017). Discarding RGB-D information enforces our model to recognize actions relying on temporal information of joints instead of rich single-frame information like RGB-D.

Two standard evaluation settings (Shahroudy et al, 2016) are used on the NTU dataset: Cross-Subject (CS) and Cross-View (CV). CS setting splits 40 subjects equally into training and test sets consisting of 40,320 and 16,560 samples respectively. In CV setting, samples of camera 1 are used for testing and samples from cameras 2 and 3 for training. The dropout value is set to 0.5 to prevent potential overfitting. Adam (Kingma and Ba, 2015) is used with the initial learning rate of 0.001 for gradient descent optimization.

Charades Dataset. The Charades dataset (Sigurdsson et al, 2016) is composed of 9,848 videos of human interacting actions with objects. The RGB information is provided for each frame at 24 FPS. The actions in the Charades dataset have relatively long duration, spanning around 30 seconds on average, which is challenging and particularly suitable for evaluating the long-range temporal modeling for action recognition. This dataset is a multi-class and multi-label dataset with a total of 66,500 annotations from 157 classes, i.e., one video sample could contain more than one category of actions.

Considering that our model cannot learn features from RGB data directly, we extract features for the Charades dataset from pre-trained I3D (Carreira and Zisserman, 2017) and I3D-NL (Wang et al, 2018) (I3D

integrated with non-local operations) respectively as two types of input features for our *NRNM* and the baseline LSTM. Specifically, we first pre-train I3D and I3D-NL on the Kinetics-400 dataset (Carreira and Zisserman, 2017), and fine-tune them on the Charades dataset. Then we freeze their parameters and extract features of the Charades dataset to feed into our *NRNM* and the baseline LSTM for training. Stochastic gradient descent (SGD) is used for optimization on the Charades datasets.

4.1.2 Implementation

Our *NRNM* is built on a 3-layer LSTM backbone unless otherwise specified. The number of hidden units of all recurrent networks mentioned in this work (vanila-RNN, GRU, LSTM) is tuned on a validation set by selecting the best configuration from the option set $\{128, 256, 512\}$. We employ 4-head attention scheme in practice. The size of memory state is set to be same as the combined size of input hidden states, i.e., the dimensions are $[\text{block size } (k)/\text{stride } (s), \text{dim}(\mathbf{h}_t)]$. Following Tu et al. (Tu et al, 2018), Zoneout (Krueger et al, 2017) is employed for network regularization. We tune the set of stride size to be $\{1, 3, 5\}$, based on the validation results, to apply the Multi-scale dilated *NRNM* mechanism.

4.1.3 Investigation on the Configuration of NRNM

We first perform experiments on NTU to investigate our proposed *NRNM* systematically, including the study of effect of hyper-parameters and different model configurations on the performance.

Effect of the block size k . We first investigate the performance of *NRNM* as a function of the block size k .

Concretely, we evaluate our method using an increasing number of block sizes: 4, 6, 8, 10, and 12 while fixing other hyper-parameters.

Figure 6(a) shows that the accuracy initially increases as the increase of the block size, which is reasonable since larger block size allows *NRNM* to incorporate information of more time steps in memory and thus enables *NRNM* to capture longer temporal dependencies. As the block size increases further after the saturated state at the block size of 8, the performance starts to decrease. We surmise that the non-local operations on a long block of sequence result in overfitting on the training data and information redundancy.

Effect of the integrated location of *NRNM* on the LSTM backbone. Next we study the effect of integrating the *NRNM* into different layers of the 3-layer LSTM backbone. Figure 6(b) presents the results, from which we can conclude: 1) integrating *NRNM* at any layer of LSTM outperforms the standard LSTM; 2) only integrating *NRNM* once at one layer performs better than applying *NRNM* at multiple layers which would lead to information redundancy and overfitting; 3) integrating *NRNM* at the middle layer achieves the best performance, which is probably because the layer-2 hidden states of LSTM are more suitable for *NRNM* to distill information than the low-level and high-level features learned by layer-1 and layer-3 hidden states.

Effect of sliding window size *win*. We further investigate the effect of sliding window size *win* of *NRNM*, which is used to control the updating frequency of memory state. Theoretically, smaller sliding window size implies more overlap between two adjacent memory blocks and thus tends to lead to more information redundancy. On the other hand, larger sliding window size results in larger non-accessed temporal interval between two adjacent memory blocks and would potentially miss information in the interval. In this set of experiments, we test the performance of *NRNM* with different sliding window size while fixing the block size k to be 8 (time steps). The results in Figure 6(c) indicate that the model performs well when the sliding window is around 4 to 8 while the performance decreases at other values, which validates our analysis.

Effect of varying the number of hidden units. Figure 6(d) shows the effect of varying the number of hidden units in the *NRNM* cell. The results reveal that our model achieves the best performance when

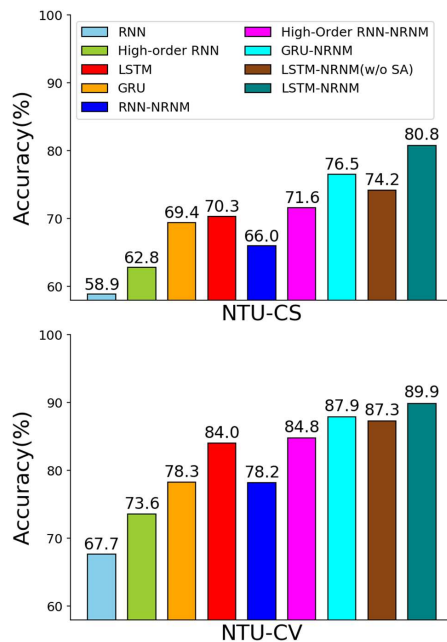


Fig. 7: Comparison of our model with other basic recurrent models in terms of classification accuracy (%) on NTU in both Cross-Subject (CS) and Cross-View (CV) settings. *NRNM*(w/o SA) denotes the *NRNM* using simple linear transformation instead of self-attention for non-local operation. We instantiate our *NRNM* with various recurrent backbones including RNN, high-order RNN, GRU and LSTM, and compare the performance between these backbones and their *NRNM* counterparts.

equipped with 512 hidden units, whilst other options lead to inferior performances due to potential underfitting (with less hidden units) and overfitting (with more hidden units).

4.1.4 Investigation on the effectiveness of *NRNM* Cell

In this set of experiments on the NTU dataset, we perform investigation on the effectiveness of the proposed *NRNM* cell. Our *NRNM* can be seamlessly integrated into any sequence models with stepwise latent states. To investigate the effectiveness and generalizability of the proposed *NRNM*, we first instantiate our model with various recurrent models as backbones and compare the resulting models to the corresponding backbones. Then we investigate the functionality of self-attention in the *NRNM*. Next, we compare our *NRNM* with Transformer, which also employs self-attention for sequence modeling, especially for NLP tasks. Finally, we perform analysis on model complexity.

Comparison with recurrent baseline models. To compare with classical recurrent models, we instantiate our *NRNM* with various recurrent models as backbones, including vanilla-RNN, high-order RNN, GRU and LSTM. Then we compare the performance between these backbones and the *NRNM* counterparts on the NTU dataset in two evaluation settings: Cross-Subject (CS) and Cross-View (CV). We make following observations from the results presented in Figure 7. 1) All recurrent models with memory or gated structure outperforms vanilla-RNN and high-order RNN by a large margin, which indicates the advantages of memory and gated structure for controlling information flow. 2) High-order RNN performs better than vanilla-RNN which reveals the necessity of the non-local interacting operations since high-order connections can be considered as a simple form of non-local operation in a local area. It is also consistent with the existing conclusions (Soltani and Jiang, 2016; Zhang and Woodland, 2018). 3) Our *NRNM* outperforms the corresponding backbone in all cases of instantiations significantly, which demonstrates the superiority of our model over these typical recurrent baseline models.

Functionality of Self-Attention in *NRNM*. We opt for Self-Attention mechanism to perform non-local operations within a memory block of *NRNM*. To investigate the effectiveness of Self-Attention mechanism in our *NRNM*, in this experiment we replace the Self-Attention with simple concatenation together with linear transformation to carry out the non-local operations among different time steps within a memory block of *NRNM*. Figure 7 presents the performance of such simple fusion scheme (*NRNM*(w/o SA)), which shows that it outperforms other recurrent baseline models including LSTM, GRU and vanilla-RNN distinctly but performs worse than our *NRNM* using self-attention operations. These results demonstrate 1) the superiority of the non-local operations in our *NRNM* compared to typical recurrent models, and 2) the effectiveness of self-attention scheme in performing non-local operations.

Comparison with Transformer. Transformer (Vaswani et al, 2017) also employs Self-Attention Mechanism to extract features and can be applied to sequence representation learning. Thus we compare its performance with our *NRNM* on the NTU dataset in Table 1. The results show that Transformer performs much worse than our model, although the model configuration of Transformer

Table 1: Classification accuracy (%) by Transformer on NTU in the Cross-View (CV) setting. d_{ff} represents the dimensions of the Feed-Forward layer whilst d_k, d_v, l, h denote the dimensions of the key, value, number of layers and heads, respectively.

Configurations of Transformer	CV
$d_{ff}=256, d_k=d_v=64, l=3, h=8$	74.3
$d_{ff}=256, d_k=d_v=64, l=6, h=4$	78.5
$d_{ff}=256, d_k=d_v=64, l=6, h=8$	70.4
$d_{ff}=512, d_k=d_v=32, l=3, h=4$	79.7
$d_{ff}=1024, d_k=d_v=32, l=3, h=4$	81.3
$d_{ff}=2048, d_k=d_v=32, l=3, h=4$	80.7
<i>NRNM</i> (ours)	89.9

Table 2: Classification accuracy (%) on NTU by different methods with different model complexities in the Cross-View (CV) setting. Here 3-LSTM (256) denotes the LSTM equipped with 3 hidden layers comprising 256 hidden units. Note that all results are reported from our implementations optimized on a held-out validation set.

Models	CV(%)	#Parameters
3-LSTM (256)	83.9	1.5M
3-LSTM (512)	84.0	5.6M
5-LSTM (512)	83.1	9.8M
3-ElAtt-LSTM (256)	85.5	1.8M
6-ElAtt-LSTM (256)	82.7	3.8M
4-ElAtt-LSTM (512)	83.4	8.9M
3-ElAtt-GRU (100)	87.1	0.3M
3-ElAtt-GRU (256)	85.4	1.4M
5-ElAtt-GRU (256)	85.0	2.5M
<i>NRNM</i> (ours)	89.9	3.8M

is carefully tuned. We surmise that the worse performance of Transformer results from two potential factors: 1) Transformer relies on the position encoding to incorporate the temporal information while our *NRNM* utilizes recurrent structure to model global temporal dependencies explicitly. It seems that the recurrent structure is more favourable than the position encoding to sequence applications like action recognition; 2) the training data in NTU is not sufficient for Transformer since Transformer typically requires a large amount of samples in language processing tasks.

Analysis on model complexity. To compare the model complexity between our model and other recurrent baselines and investigate whether the performance of our model is resulted from by the augmented

Table 3: Classification accuracy (%) on NTU by different methods in both Cross-Subject (CS) and Cross-View (CV) settings.

Models	CS	CV
HBRNN-L (Du et al, 2015)	59.1	64.0
Part-aware LSTM (Shahroudy et al, 2016)	62.9	70.3
Trust Gate ST-LSTM (Liu et al, 2016)	69.2	77.7
Two-stream RNN (Wang and Wang, 2017)	71.3	79.5
Ensemble TS-LSTM (Lee et al, 2017)	74.6	81.3
VA-LSTM (Zhang et al, 2017)	79.4	87.6
STA-LSTM (Song et al, 2018)	73.4	81.2
EleAtt-LSTM (Zhang et al, 2018)	78.4	85.0
EleAtt-GRU (Zhang et al, 2018)	79.8	87.1
LSTM (baseline)	70.3	84.0
<i>NRNM</i> (ours)	80.8	89.9

Table 4: Comparison between LFB (Wu et al, 2019) using RGB information and our *NRNM* using skeleton joint information for action recognition on NTU dataset. Performance is evaluated in terms of classification accuracy (%). The memory usage and inference time are measured for a batch of (64) test samples using 1 NVIDIA GeForce RTX 3090 GPU.

Models		LFB	<i>NRNM</i> (ours)
Frame Information		RGB	Skeleton
Performance	CS	83.5	80.8
	CV	93.4	89.9
Model Complexity	#Parameters (M)	82.3	3.8
	Memory Usage (GB)	13.59	1.86
	Inference Time (s)	4.60	0.11

model complexity, we evaluate the performance of the state-of-the-art recurrent baselines with different model complexities (configurations) on the NTU dataset, shown in Table 2. Our model substantially outperforms other baselines under optimized configurations, including many baselines with much more complexities than our model. These results demonstrate that the performance superiority of our model is not resulted from the increased capacity by the extra parameters.

4.1.5 Comparison with State-of-the-art Methods

In this set of experiments, we compare our model with the state-of-the-art methods for action recognition on both the NTU dataset and the Charades dataset.

Results on NTU. We perform experiments in both Cross-Subject (CS) and Cross-View (CV) settings on the NTU dataset. It should be noted that we do not

compare with the methods which employ extra information or prior knowledge such as joint connections for each part of body or human body structure modeling (Si et al, 2018; Yan et al, 2018) for a fair comparison.

Table 3 reports the experimental results. Our model achieves the best performance in both CS and CV settings, which demonstrates the superiority of our model over other recurrent networks, especially those with memory or gated structures. While our model outperforms the standard LSTM model substantially, the methods based on LSTM (Song et al, 2018; Zhang et al, 2018) boost the performance over LSTM by introducing extra attention mechanisms. However, all these methods focus on enhancing the effectiveness of controlling information flow by improving the gate structure or the cell structure of recurrent models. In contrast, our *NRNM* aims to capture the long-range temporal dependencies and distill high-level semantic features by modeling the high-order interactions between non-adjacent time steps with designed non-local interacting operations.

***NRNM* using skeleton joint information vs LFB (Wu et al, 2019) using RGB information.** Similar to our *NRNM*, LFB designs a long-term feature bank, which stores long-range temporal context by enumerating 3D convolutional features for detected objects extracted from RGB data, to augment existing video models. LFB and our model adopt two different kinds of pipelines, and both methods have their own merits for action recognition on video data. On the one hand, LFB is able to store much richer information from RGB data for each time step in the constructed feature bank than our proposed *NRNM* memory distilled from skeleton joint information, potentially leading to favorable performance of LFB compared with our model. On the other hand, LFB contains much more parameters and also consumes much more memory than our model due to substantially more parameters introduced by the 3D convolutional operations of LFB for learning features from RGB data and more memory space demanded by the feature bank of LFB than the *NRNM* memory of our model. Besides, high computational complexity of 3D convolutional operations of LFB also leads to more inference time than our model.

To validate above analysis, we conduct experiments to evaluate LFB on NTU dataset based on RGB information and compare its performance with that of our model on NTU dataset with only skeleton

Table 5: Action recognition accuracy (mAP in %) using RGB data on Charades by different methods.

Methods	Feature backbone	mAP
2-Strm (Simonyan and Zisserman, 2014)	VGG16	18.6
Multiscale TRN (Zhou et al, 2018a)	Inception	25.2
I3D (Carreira and Zisserman, 2017)	R101-I3D	35.5
I3D+NL (Wang et al, 2018)	R101-I3D-NL	37.5
STRG (Wang and Gupta, 2018)	R101-I3D-NL	39.7
PoTion (Choutas et al, 2018)	GCN+I3D-NL+I3D	40.8
PA3D (Yan et al, 2019)	GCN+I3D-NL+I3D	41.0
SlowFast+NL (Feichtenhofer et al, 2019)	3D ResNet	42.5
LFB+NL (Wu et al, 2019)	R101-I3D-NL	42.5
LSTM (Baseline)	R101-I3D	35.7
<i>NRNM</i> (ours)	R101-I3D	40.1
LSTM (Baseline)	R101-I3D-NL	38.9
<i>NRNM</i> (ours)	R101-I3D-NL	41.7

joint information. Table 4 shows that LFB achieves better performance than our model in both Cross-Subject (CS) and Cross-View (CV) settings. However, it contains significantly (an order of magnitude) more parameters, and also demands a lot more memory and inference time than our model. These results are consistent with our theoretical analysis. Overall, our *NRNM* focuses on modeling high-order temporal interactions between non-adjacent time steps while LFB seeks to distill rich context information per frame and obtain a long-term view by constructing a feature bank from RGB information.

Results on Charades. We conduct experiments on the Charades dataset to investigate whether our model is effective compared to the baseline (LSTM) on the generic video data based on RGB information, although our model is not specifically designed towards visual RGB data. Note that we do not compare with the methods based on Neural Architecture Search (NAS) since these methods introduce the additional NAS stage, which is not a fair comparison w.r.t. the algorithmic complexity. We evaluate the performance of our *NRNM* and LSTM using two backbones for input features: R101-I3D and R101-I3D-NL.

We can make following observations from the experimental results in Table 5. First, our model substantially improves the performance over LSTM using each type of input features (R101-I3D or R101-I3D-NL), which demonstrates the effectiveness of the proposed *NRNM* on the RGB data. Besides, the performance gain using R101-I3D features is larger than that using R101-I3D-NL features. It is reasonable since R101-I3D-NL also performs non-local operations temporally in the convolutional feature space,

which shares similar function with our *NRNM*. Nevertheless, our model still boosts the performance over LSTM by 2.8% using R101-I3D-NL features. Second, our model outperforms most of the specialized methods for action recognition except for SlowFast (Feichtenhofer et al, 2019) and LFB (Wu et al, 2019). SlowFast designs two pathways to model the spatial semantics and temporal dynamics respectively, which is algorithmically optimized for video data. LFB designs a long-term feature bank to store temporal context in the form of 3D convolutional features, which contains much higher dimensions of features and richer information in each step than the *NRNM* memory distilled from the hidden states of LSTM. Both models are designed specifically towards RGB information of video data while our model focuses on sequence modeling, thus it is acceptable that our model performs slightly worse than these two models.

4.1.6 Qualitative Analysis

To qualitatively illustrate the advantages of the proposed *NRNM*, figure 8 presents a concrete video example from the NTU dataset with the groundtruth action label “walking towards each other”. In this example, it is quite challenging to recognize the action since it can only be inferred by the temporal variations of the relative distance between two persons in the scene. Hence, capturing the long-range dependencies is crucial for recognizing the action correctly. The standard LSTM misclassifies it as “punching/slapping other person” while our model is able to correctly recognize it due to the capability of modeling long-range temporal information by our designed *NRNM*.

Figure 8 visualizes two blocks of memory states, each of which is learned by *NRNM* cell via incorporating information of multiple frames within the corresponding block, including input features x_i and the hidden states h_i of LSTM backbone. To obtain more insights into the non-local operations of *NRNM*, we visualize the attention weights \mathbf{W}_{att} in Equation 5 to show that each unit of memory state is calculated by attending to all units of source information.

4.2 Experiment 2 on Sequence Classification: Sentiment Analysis

Next we perform experiments on sentiment analysis, another task of sequence classification, to evaluate our model on the text modality. Specifically, we aim to identify online movie reviews as positive or negative.

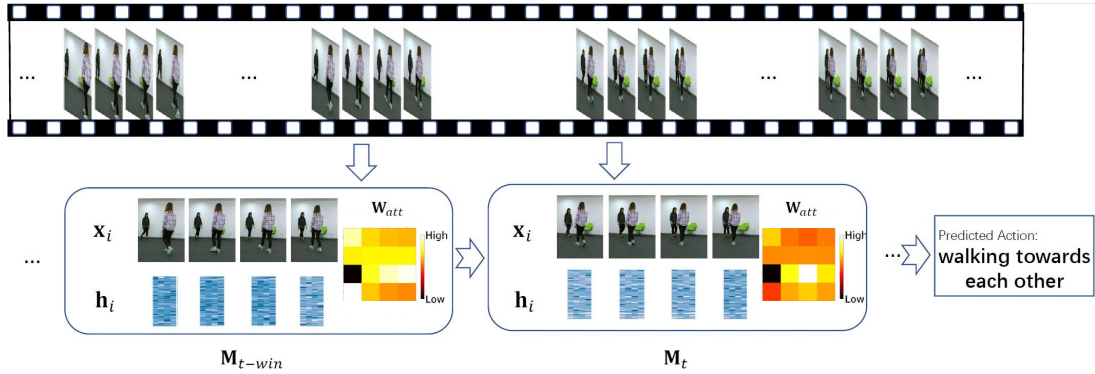


Fig. 8: Visualization of an example with labeled action “walking towards each other” from NTU. Our model is able to correctly recognize it while LSTM misclassifies it as “punching/slapping other person”. The temporal variations of relative distance between two persons are key to recognize the action. Our model can successfully capture it while LSTM fails. Two blocks of memory states and the attention weights \mathbf{W}_{att} in Equation 5 are visualized.

4.2.1 Dataset and Evaluation Protocol

In this set of experiments, we perform evaluation on IMDB Review dataset (Maas et al, 2011) which is a standard benchmark for sentiment analysis. It contains 50,000 labeled reviews among which 25,000 samples are used for training and the rest for testing. The average length of reviews is 241 words and the maximum length is 2526 words (Dai and Le, 2015). Note that the IMDB dataset also provides additional 50,000 unlabeled reviews, which are used by several customized semi-supervised learning methods (Dai and Le, 2015; Dieng et al, 2017; Johnson and Zhang, 2016; Miyato et al, 2017; Radford et al, 2017). Since we only use labeled data for supervised training, we compare our methods with those methods based on supervised learning using the same set of training data for a fair comparison.

The torchtext² is used for data preprocessing. Following the training strategy in Dai et al. (Dai and Le, 2015), we pretrain a language model for extracting word embeddings.

4.2.2 Comparison with Recurrent Baseline Models

We first conduct a set of experiments to compare our model to the basic recurrent networks including vanilla-RNN, GRU, LSTM and high-order RNN. Figure 9 shows that our model outperforms all other baselines significantly which reveals the remarkable advantages of our *NRNM*. Besides, while LSTM and

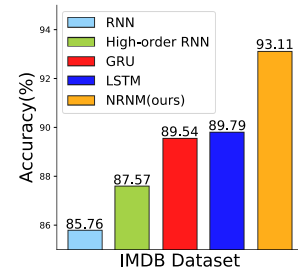


Fig. 9: Comparison of our model with other basic recurrent models for sentiment analysis in terms of classification accuracy (%) on IMDB dataset.

GRU perform much better than vanilla-RNN, high-order RNN also boosts the performance by a large margin compared to vanilla-RNN. It again demonstrates the benefits of high-order connections which are a simple form of non-local operations in local area.

4.2.3 Comparison with State-of-the-art Methods

Next we compare our *NRNM* with the state-of-the-art methods including LSTM (Xia et al, 2018), oh-CNN (Johnson and Zhang, 2014) and oh-2LSTMp (Johnson and Zhang, 2016) which learn the word embeddings by customized CNN or LSTM instead of using existing pretrained word embedding vocabulary, DSL (Xia et al, 2018) and MLDL (Xia et al, 2018) which perform a dual learning between language modeling and sentiment analysis, GLoMo (Yang et al, 2018) which is a transfer learning framework, and BCN+Char+CoVe (McCann et al, 2017) which trains a machine translation model to encode the word embeddings to improve the performance of sentiment analysis.

²<https://github.com/pytorch/text>

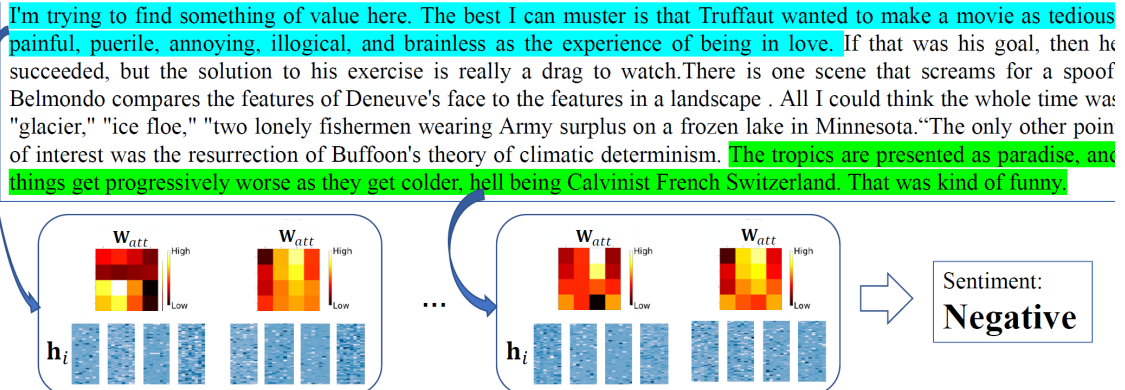


Fig. 10: Visualization of an example of movie review with the groundtruth label “negative”. Our model is able to correctly classify it while LSTM fails. The last sentence (in green color) which seems positive tends to misguide models. The first sentence is an important cue for negative sentiment, which is hardly captured by LSTM since it is easily forgotten by the hidden state h_T in the last time step.

Table 6: Classification accuracy (%) for sentiment analysis on IMDB dataset by different methods. Note that GLoMo (Yang et al, 2018) and BCN+Char+CoVe (McCann et al, 2017) use additional data for either transfer learning or training an individual machine translation model.

Methods	Accuracy
LSTM (Xia et al, 2018)	89.9
MLDL (Xia et al, 2018)	92.6
GLoMo (Yang et al, 2018)	89.2
oh-2LSTMp (Johnson and Zhang, 2016)	91.9
DSL (Xia et al, 2018)	90.8
oh-CNN (Johnson and Zhang, 2014)	91.6
BCN+Char+CoVe (McCann et al, 2017)	92.1
LSTM (baseline)	89.8
<i>NRNM</i> (ours)	93.1

Table 6 shows that our model achieves the best performance among all methods. It is worth mentioning that our model even performs better than GLoMo (Yang et al, 2018) and BCN+Char+CoVe (McCann et al, 2017), which employ additional data for either transfer learning or training an individual machine translation model.

4.2.4 Qualitative Analysis

Figure 10 illustrates an example of sentiment analysis from IMDB dataset. This example of movie review is fairly challenging since the last sentence of the review seems to be positive which is prone to misguide models, especially when we use the hidden state of last time step h_T for prediction. However, the groundtruth sentiment label for this example is negative, which is evidently implied by the first sentence. Our model

is able to correctly classify it as “negative” while LSTM fails since LSTM can hardly capture the long-term features incorporating the information of the first sentence. We also visualize the attention weights of non-local operations (W_{att} Equation 5) in two blocks of *NRNM* states to show the attendance of each information units of source information for calculating the *NRNM* states. The first memory block corresponds to the first sentence which is an important cue of negative sentiment while the second memory block corresponds to the last sentence.

4.3 Experiment 3 on Step-wise Sequential Prediction

In this set of experiments, we evaluate the performance of our model on step-wise sequence prediction, which makes predictions sequentially for each time step (or every some steps). To be specific, we conduct experiments on the task of protein secondary structure (PSS) prediction, which has extensive applications such as analyzing protein function and drug design (Zhou et al, 2018b). Given an input protein primary sequence, this task aims to predict a label out of 8 categories for each time step and the predicted sequence labels constitute the PSS result for the input protein sequence. The crux of this task stems from its various patterns of dependencies of output labels on the input sequences. To have a more distinct comparison between different methods, we opt for the more challenging problem setting, namely the 8-category classification for predictions at each time step, rather than the easier setting of 3-category classification.

4.3.1 Dataset and Evaluation Protocol

We perform experiments following the routine settings for protein secondary structure prediction (Wang et al, 2016; Li and Yu, 2016). Particularly, we adopt CB6133 dataset as the training data and perform test on three popular test datasets including CB513, CASP10 and CASP11 datasets. The details of these datasets are presented as follows.

CB6133 dataset (Zhou and Troyanskaya, 2014) is produced by PISCES CullPDB (Wang and Dunbrack Jr, 2003) and contains 6128 protein sequences. It is a non-homologous dataset, and contains 6128 protein sequences with given protein secondary structure for each sequence as groundtruth labels. Following (Wang et al, 2016; Li and Yu, 2016), we filter out the overlapped data with CB513 dataset (used for test) from CB6133 to obtain unbiased evaluation of all methods. The remaining CB6133 dataset comprising 5534 protein sequences is used as the training set in this set of experiments.

CB513 dataset (Zhou and Troyanskaya, 2014) contains 514 protein sequences and is a widely used test dataset for the task of protein secondary structure prediction. To test the generalization of methods for step-wise sequential prediction across different datasets, we also use CASP10 and CASP11 datasets (Zhou et al, 2018b; Li and Yu, 2016) as test datasets, which contain 123 and 105 protein sequences respectively.

Each protein sequence in above datasets is described by 55 channels of information per protein residue, among which 21 channels are sequence features for specifying the category of the amino acid and 21 channels are the sequence profile. These 42 channels of information are used as the input features for each time step of sequences. Besides, 8 channels (out of 55 channels) of information are used to indicate the category labels of the protein secondary structure. Note that the left 5 channels of information are not used in this set of experiments. Considering the convenience of implementation, we pre-process all protein sequences to make them have equal length (700 time steps) by truncating or padding operations.

As a commonly used metric for protein secondary structure prediction, Q8 accuracy (Pollastri et al, 2002; Wang et al, 2010; Peng et al, 2009; Wang et al, 2016; Zhou and Troyanskaya, 2014; Li and Yu, 2016) is used for evaluation in our experiment. It measures the prediction accuracy of the amino-acid residues.

Table 7: Comparison between our model and other basic recurrent models for protein secondary structure prediction on three test datasets in terms of Q8 accuracy (%).

	CB513	CASP10	CASP11
RNN	57.1	61.1	59.9
LSTM	55.2	60.6	60.1
GRU	57.8	61.4	61.2
Bi-RNN	66.1	70.6	69.9
Bi-LSTM	66.4	70.8	70.4
Bi-GRU	67.1	71.5	71.0
<i>NRNM</i> (ours)	69.6	74.1	72.0

Since the bi-directional temporal features is crucial for protein secondary structure prediction, we perform bi-directional modeling for our method, which is straightforward by incorporating the bi-directional hidden states of LSTM backbone into the input source information \mathbf{C} (in Equation 5) for our *NRNM* cell:

$$\mathbf{C} = \text{Concat}([\vec{\mathbf{h}}_{t-k+1}, \dots, \vec{\mathbf{h}}_t], [\overleftarrow{\mathbf{h}}_{t-k+1}, \dots, \overleftarrow{\mathbf{h}}_t], [\mathbf{x}_{t-k+1}, \dots, \mathbf{x}_t]). \quad (18)$$

4.3.2 Comparison with Recurrent Baseline Models

As we did in previous experiments on other sequence applications, in this set of experiments for step-wise sequential prediction we first compare our model with other recurrent baseline models including the vanilla-RNN, LSTM and GRU. Besides, we also evaluate the performance of bi-directional version of these models.

Table 7 presents the experimental results, from which we make two observations. Firstly, the large performance gap between the normal recurrent models and their respective bi-directional counterparts reveal that the bi-directional temporal features is indeed distinctly beneficial for the task of protein secondary structure prediction. Secondly, our model outperforms all other recurrent models by a large margin, which demonstrates the superiority of our model over these recurrent baselines on the task of protein secondary structure prediction.

4.3.3 Comparison with State-of-the-art Methods

Next we conduct experiments to compare our model with state-of-the-art approaches for protein secondary

structure prediction, which include: **SSpro8** (Pollastri et al, 2002), which uses ensembles of bi-directional recurrent neural network architectures and PSI-BLAST-derived profiles to improve the performance of protein secondary structure prediction; **CNF** (Wang et al, 2010), which presents a probabilistic method based on Conditional Neural Field (Peng et al, 2009) for secondary structure prediction. It not only models the relationship between sequence features and secondary structures, but also exploits the inter-dependencies among secondary structures; **DeepCNF** (Wang et al, 2016), which is extended from **CNF** using deep convolutional networks; **GSN** (Zhou and Troyanskaya, 2014), which proposes a method based on generative stochastic network (GSN) to predict local secondary structure with deep hierarchical representation; **DCRNN** (Li and Yu, 2016), which proposes an end-to-end deep network that predicted protein secondary structure from integrated local and global features between amino-acid residues.

From the experimental results presented in Table 8, we observe that our model achieves the best results among all methods. In particular, our model performs much better than **SSpro8** which employs ensembles of bi-directional recurrent networks, which implies that using ensemble of recurrent networks cannot achieve the similar effect as our model. Another surprising observation is that leveraging the (1-D) convolutional networks to capture temporal features by both **DeepCNF** and **DCRNN** achieves excellent performance. We surmise that this is largely because the protein secondary structure prediction is determined not only by the long-term temporal dependencies, but also the high-level semantic features distilled from the time steps near the predicted time step, which are the strengths of convolutional networks. Nevertheless, our model also has these two key advantages (comparing to other recurrent models) and thus compares favorably to both **DeepCNF** and **DCRNN**. Besides, it is worth noting that all the state-of-the-art methods in Table 8 are specifically designed for protein secondary structure prediction whilst our model is generally applicable to any sequence task and is not modified to adapt the task of protein secondary structure prediction.

Table 8: Q8 accuracy (%) of our method and state-of-the-art methods on three test benchmarks for protein secondary structure prediction.

	CB513	CASP10	CASP11
SSpro8 (Pollastri et al, 2002)	63.5	64.9	65.6
CNF (Wang et al, 2010)	64.9	64.8	65.1
GSN (Zhou and Troyanskaya, 2014)	66.4	-	-
DeepCNF (Wang et al, 2016)	68.3	71.8	72.3
DCRNN (Li and Yu, 2016)	69.4	-	-
<i>NRNM</i> (ours)	69.6	74.1	72.0

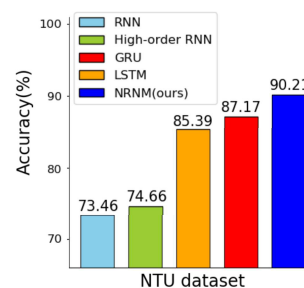


Fig. 11: Comparison of our model with other basic recurrent models for sequence similarity learning in terms of classification accuracy (%).

4.4 Experiment 4 on Sequence Similarity Learning

In the last set of experiments, we evaluate our model on sequence similarity learning, which aims to predict two input sequences are similar or not. To predict the similarity precisely, this task demands the methods for it to learn effective sequence representations that incorporate the temporal dependencies covering the whole sequence for both input sequences. To this end, we design a siamese-*NRNM* structure, which employs two *NRNMs* with shared parameters to learn sequence representations for two input sequences in the same feature space. Then the similarity between two sequences is measured between two learned representations. The whole model is trained in end-to-end manner according to the loss function in Equation 17.

4.4.1 Dataset and Evaluation Protocol

We reuse the 3D skeleton data of NTU dataset (Shahroudy et al, 2016) introduced in Section 4.1.1, which is a dataset for action recognition, to design the experiments of sequence similarity learning. In particular, we define similar sequences

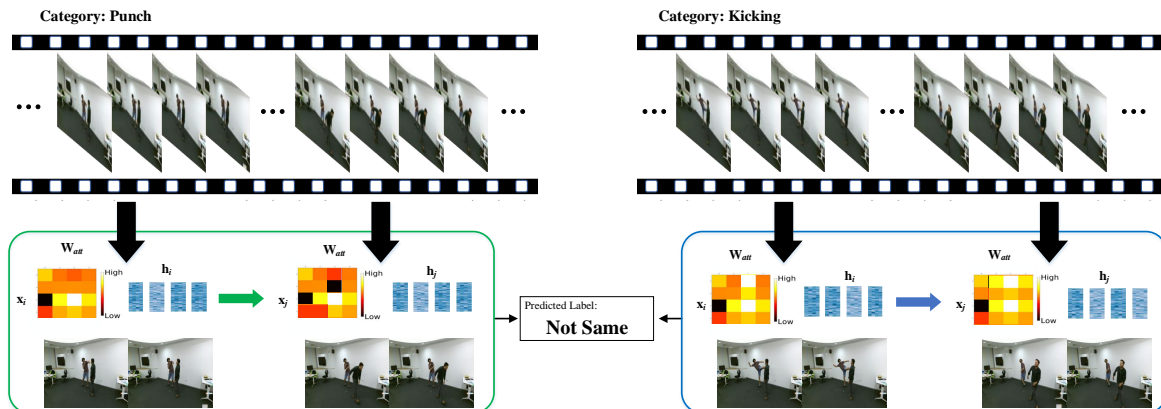


Fig. 12: Visualization of an example of sequence similarity learning with a pair of action videos. One video is labeled with action ‘Punch’ while the label of the other one is ‘Kicking’. Two videos begin with different action scenes but end with similar scenes. LSTM mis-predicts them as the same category due to its limited capability of modeling temporal features. By contrast, our model is able to correctly distinguish them owing to our designed *NRNM*.

Table 9: Classification accuracy (%) of different methods for sequence similarity learning on the constructed dataset from NTU dataset.

	Accuracy
EleAtt-LSTM (Zhang et al, 2018)	83.4
EleAtt-GRU (Zhang et al, 2018)	86.0
MaLSTM (Mueller and Thyagarajan, 2016)	82.4
MwAN (Tan et al, 2018)	84.9
RE2 (Yang et al, 2019)	89.0
<i>NRNM</i> (ours)	90.2

to be those with the same category label and the dissimilar sequences to be those with different category labels. According to such definition, we construct the training set and test set from NTU dataset to perform action similarity learning. Specifically, we randomly select 18,823 pairs of sequences from the training split of NTU dataset to be the training set for sequence similarity learning, in which half of data is similar pairs of action sequences and the other half is dissimilar pairs. We construct the test set consisting of 5,000 pairs of action sequences in the same way. Note that we keep the amount of similar and dissimilar pairs of action sequences to be balanced for each of total 60 action categories.

4.4.2 Comparison with Recurrent Baseline Models

We first conduct experiments to compare our model with recurrent baseline models including vanilla-RNN, high-order RNN, GRU and LSTM. Same as our

siamese-NRNM structure, We build siamese structure for each of them to perform sequence similarity learning and employ the same loss function shown in Equation 17 for training models.

Figure 11 presents the experimental results. We can draw similar conclusions as the previous experiments for sequence classification: 1) the models with memory or gated structures like GRU, LSTM and our model consistently perform better than the vanilla-RNN, which results from the strong capability of sequence representation learning of these models; 2) our model outperforms LSTM and GRU by a large margin, which indicates the merits of our *NRNM* module.

4.4.3 Comparison with State-of-the-art Methods

We then compare our model with state-of-the-art methods for sequence similarity learning, which includes: 1) **MaLSTM** (Mueller and Thyagarajan, 2016), which employs siamese-LSTM to learn the sequence representations for two input sequences in the same feature space, and measure the sequence similarity between the learned representations. Similar idea was also investigated in Siamese Recurrent Networks (SRN) (Pei et al, 2016). These two methods are the first to design siamese recurrent structure to model sequence similarity. 2) **MwAN** (Tan et al, 2018), which applies multiple attention functions to model the similarity between a pair of sequences. 3) **RE2** (Yang et al, 2019), which focuses on modeling multiple alignments for learning sequences similarity. Above three methods are all designed specifically

for sequence similarity learning. Additionally, we also construct the siamese structure for **EleAtt-LSTM** and **EleAtt-GRU** (Zhang et al, 2018) as a baseline model, which employs attention mechanism to explore element-wise potential of the input vector.

The experimental results are reported in Table 9. It shows that our *NRNM* achieves the best result among all methods, including the methods specifically designed for sequence similarity learning. **RE2** also performs well due to its various kinds of alignments carefully modeled for precise sequence matching.

4.4.4 Qualitative Analysis

To qualitatively validate the effectiveness of our proposed *NRNM*, we visualize a real example which is to predict the similarity between a pair of action sequences with similar content but different labels in Figure 12. One action sequence is labeled with ‘Punch’ while the other one is about ‘Kicking’. While two videos begin with different action scenes, the ending part of two videos are very similar to each other which tends to misguide the prediction. LSTM recognizes these two actions as the same category due to its limited capability of capturing temporal features. By contrast, our model is able to distinguish them correctly.

5 Conclusion

In this work, we have presented the Non-local Recurrent Neural Memory (*NRNM*) for sequence representation learning. We perform non-local operations within each memory block to model full-order interactions between non-adjacent time steps and model the global interactions between memory blocks in a gated recurrent manner. Thus, our model is able to capture long-range temporal dependencies. Furthermore, the proposed *NRNM* allows learning high-level semantic features within a memory block by non-local operations due to much longer *direct interacting field*. The proposed *NRNM* cell can be seamlessly integrated into any existing recurrent-based sequence model to enhance the power of sequence representation learning. Our method achieves the state-of-the-art performance in three types of sequence applications across different modalities including sequence classification, step-wise sequence prediction and sequence similarity learning.

Acknowledgements. This work was supported in part by the NSFC fund (U2013210, 62006060,

62176077), in part by the Guangdong Basic and Applied Basic Research Foundation under Grant (2019B1515120055, 2022A1515010306), in part by the Shenzhen Key Technical Project under Grant 2020N046, in part by the Shenzhen Fundamental Research Fund under Grant (JCYJ20210324132210025), in part by the Shenzhen Stable Support Plan Fund for Universities (GXWD20201230155427003-20200824125730001), in part by the Medical Biometrics Perception and Analysis Engineering Laboratory, Shenzhen, China, and in part by Guangdong Provincial Key Laboratory of Novel Security Intelligence Technologies (2022B1212010005).

References

- Ba JL, Kiros JR, Hinton GE (2016) Layer normalization. arXiv preprint arXiv:160706450
- Bahdanau D, Cho K, Bengio Y (2015) Neural machine translation by jointly learning to align and translate. In: ICLR
- Bello I, Zoph B, Vaswani A, et al (2019) Attention augmented convolutional networks. In: ICCV
- Bertolami R, Bunke H, Fernandez S, et al (2009) A novel connectionist system for improved unconstrained handwriting recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence 31(5)
- Brock A, Donahue J, Simonyan K (2019) Large scale gan training for high fidelity natural image synthesis. In: ICLR
- Buades A, Coll B, Morel JM (2005) A non-local algorithm for image denoising. In: CVPR
- Cao Y, Xu J, Lin S, et al (2019) Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In: ICCV Workshops
- Carion N, Massa F, Synnaeve G, et al (2020) End-to-end object detection with transformers. In: ECCV
- Carreira J, Zisserman A (2017) Quo vadis, action recognition? a new model and the kinetics dataset. In: "CVPR"
- Cho K, Van Merriënboer B, Gulcehre C, et al (2014) Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:14061078

- Choutas V, Weinzaepfel P, Revaud J, et al (2018) Potion: Pose motion representation for action recognition. In: "CVPR"
- Dai AM, Le QV (2015) Semi-supervised sequence learning. In: NeurIPS
- Dai Z, Yang Z, Yang Y, et al (2019) Transformer-xl: Attentive language models beyond a fixed-length context. In: ACL
- Devlin J, Chang MW, Lee K, et al (2019) Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL
- Dieng AB, Wang C, Gao J, et al (2017) Topicrnn: A recurrent neural network with long-range semantic dependency. In: ICLR
- Du Y, Wang W, Wang L (2015) Hierarchical recurrent neural network for skeleton based action recognition. In: CVPR
- Feichtenhofer C, Fan H, Malik J, et al (2019) Slowfast networks for video recognition. In: ICCV
- Fu C, Pei W, Cao Q, et al (2019a) Non-local recurrent neural memory for supervised sequence modeling. In: ICCV
- Fu J, Liu J, Tian H, et al (2019b) Dual attention network for scene segmentation. In: CVPR
- Fu J, Liu J, Tian H, et al (2019c) Dual attention network for scene segmentation. In: CVPR
- Grave E, Joulin A, Usunier N (2017) Improving neural language models with a continuous cache. ICLR
- Graves A, Wayne G, Danihelka I (2014) Neural Turing machines. arXiv preprint arXiv:14105401
- He K, Gkioxari G, Dollár P, et al (2017) Mask r-cnn. In: ICCV
- Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural computation* 9(8):1735–1780
- Hochreiter S, Bengio Y, Frasconi P, et al (2001) Gradient flow in recurrent nets: the difficulty of learning long-term dependencies
- Hu H, Zhang Z, Xie Z, et al (2019) Local relation networks for image recognition. In: ICCV
- Huang G, Liu Z, Van Der Maaten L, et al (2017) Densely connected convolutional networks. In: CVPR
- Johnson R, Zhang T (2014) Effective use of word order for text categorization with convolutional neural networks. arXiv preprint arXiv:14121058
- Johnson R, Zhang T (2016) Supervised and semi-supervised text categorization using lstm for region embeddings. In: ICML
- Kay W, Carreira J, Simonyan K, et al (2017) The kinetics human action video dataset. arXiv preprint arXiv:170506950
- Kingma DP, Ba J (2015) Adam: A method for stochastic optimization. ICLR
- Krueger D, Maharaj T, Kramár J, et al (2017) Zoneout: Regularizing rnns by randomly preserving hidden activations. ICLR
- Kumar A, Irsoy O, Ondruska P, et al (2016) Ask me anything: Dynamic memory networks for natural language processing. In: ICML
- Lafferty JD, McCallum A, Pereira FC (2001) Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In: ICML
- Lee I, Kim D, Kang S, et al (2017) Ensemble deep learning for skeleton-based action recognition using temporal sliding lstm networks. In: ICCV
- Li Z, Yu Y (2016) Protein secondary structure prediction using cascaded convolutional and recurrent neural networks. arXiv preprint arXiv:160407176
- Liu J, Shahroudy A, Xu D, et al (2016) Spatio-temporal lstm with trust gates for 3d human action recognition. In: ECCV
- Liu Z, Lin Y, Cao Y, et al (2021) Swin transformer: Hierarchical vision transformer using shifted windows. In: ICCV
- Maas AL, Daly RE, Pham PT, et al (2011) Learning word vectors for sentiment analysis. In: ACL
- McCann B, Bradbury J, Xiong C, et al (2017) Learned in translation: Contextualized word vectors. In: NeurIPS
- Miyato T, Dai AM, Goodfellow I (2017) Adversarial training methods for semi-supervised text classification. In: ICLR
- Morency LP, Quattoni A, Darrell T (2007) Latent-dynamic discriminative models for continuous gesture recognition. In: CVPR
- Mueller J, Thyagarajan A (2016) Siamese recurrent architectures for learning sentence similarity. In: AAAI
- Pei W, Tax DM, van der Maaten L (2016) Modeling time series similarity with siamese recurrent

- networks. arXiv preprint arXiv:160304713
- Pei W, Baltrusaitis T, Tax DM, et al (2017) Temporal attention-gated model for robust sequence classification. In: CVPR
- Pei W, Dibeklioğlu H, Tax DM, et al (2018) Multivariate time-series classification using the hidden-unit logistic model. *IEEE transactions on neural networks and learning systems* 29(4):920–931
- Peng J, Bo L, Xu J (2009) Conditional neural fields. In: *NeurIPS*
- Pollastri G, Przybylski D, Rost B, et al (2002) Improving the prediction of protein secondary structure in three and eight classes using recurrent neural networks and profiles. *Proteins: Structure, Function, and Bioinformatics* 47(2):228–235
- Qiu Z, et al. (2017) Learning spatio-temporal representation with pseudo-3d residual networks. In: *ICCV*
- Rabiner LR (1989) A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE* 77(2):257–286
- Radford A, Jozefowicz R, Sutskever I (2017) Learning to generate reviews and discovering sentiment. arXiv preprint arXiv:170401444
- Ramachandran P, Parmar N, Vaswani A, et al (2019) Stand-alone self-attention in vision models. *NeurIPS*
- Rumelhart DE, Hinton GE, Williams RJ, et al (1988) Learning representations by back-propagating errors. *Cognitive modeling* 5(3):1
- Sak H, Senior A, Beaufays F (2014) Long short-term memory recurrent neural network architectures for large scale acoustic modeling. In: Fifteenth annual conference of the international speech communication association
- Santoro A, Faulkner R, Raposo D, et al (2018) Relational recurrent neural networks. In: *NeurIPS*
- Shahroudy A, Liu J, Ng TT, et al (2016) Ntu rgb+d: A large scale dataset for 3d human activity analysis. In: *CVPR*
- Si C, Jing Y, Wang W, et al (2018) Skeleton-based action recognition with spatial reasoning and temporal stack learning. In: *ECCV*
- Sigurdsson GA, Varol G, Wang X, et al (2016) Hollywood in homes: Crowdsourcing data collection for activity understanding. In: "ECCV"
- Simonyan K, Zisserman A (2014) Two-stream convolutional networks for action recognition in videos. "NIPS"
- Soltani R, Jiang H (2016) Higher order recurrent neural networks. arXiv preprint arXiv:160500064
- Song S, Lan C, Xing J, et al (2018) Spatio-temporal attention-based lstm networks for 3d action recognition and detection. *IEEE Transactions on Image Processing* 27(7):3459–3471
- Srinivas A, Lin TY, Parmar N, et al (2021) Bottleneck transformers for visual recognition. In: *CVPR*
- Sukhbaatar S, Weston J, Fergus R, et al (2015) End-to-end memory networks. In: *NeurIPS*
- Szegedy C, Liu W, Jia Y, et al (2015) Going deeper with convolutions. In: *CVPR*
- Tan C, Wei F, Wang W, et al (2018) Multiway attention networks for modeling sentence pairs. In: *IJCAI*
- Tu J, Liu H, Meng F, et al (2018) Spatial-temporal data augmentation based on lstm autoencoder network for skeleton-based human action recognition. In: *ICIP*
- Van Der Maaten L, Welling M, Saul L (2011) Hidden-unit conditional random fields. In: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*
- Vaswani A, Shazeer N, Parmar N, et al (2017) Attention is all you need. In: *NeurIPS*
- Wang G, Dunbrack Jr RL (2003) Pisces: a protein sequence culling server. *Bioinformatics* 19(12):1589–1591
- Wang H, Wang L (2017) Modeling temporal dynamics and spatial configurations of actions using two-stream recurrent neural networks. In: *CVPR*
- Wang S, Peng J, Ma J, et al (2016) Protein secondary structure prediction using deep convolutional neural fields. *Scientific reports* 6(1):1–11
- Wang X, Gupta A (2018) Videos as space-time region graphs. In: "ECCV"
- Wang X, Girshick R, Gupta A, et al (2018) Non-local neural networks. In: *CVPR*
- Wang Z, Zhao F, Peng J, et al (2010) Protein 8-class secondary structure prediction using conditional neural fields. In: *2010 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp 109–114

- Weston J, Chopra S, Bordes A (2015) Memory networks. In: ICLR
- Wu CY, Feichtenhofer C, Fan H, et al (2019) Long-term feature banks for detailed video understanding. In: CVPR
- Xia Y, Tan X, Tian F, et al (2018) Model-level dual learning. In: ICML
- Xie S, Sun C, Huang J, et al (2018) Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In: ECCV
- Xiong C, Merity S, Socher R (2016) Dynamic memory networks for visual and textual question answering. In: ICML
- Yan A, Wang Y, Li Z, et al (2019) Pa3d: Pose-action 3d machine for video recognition. In: "CVPR"
- Yan S, Xiong Y, Lin D (2018) Spatial temporal graph convolutional networks for skeleton-based action recognition. In: AAAI
- Yang R, Zhang J, Gao X, et al (2019) Simple and effective text matching with richer alignment features. In: ACL, pp 4699–4709
- Yang Z, Dhingra B, He K, et al (2018) Glomo: Unsupervisedly learned relational graphs as transferable representations. In: NeurIPS
- Yin M, Yao Z, Cao Y, et al (2020) Disentangled non-local neural networks. In: ECCV
- Yu F, Koltun V (2015) Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:151107122
- Zhang C, Woodland PC (2018) High order recurrent neural networks for acoustic modelling. In: ICASSP
- Zhang H, Goodfellow I, Metaxas D, et al (2019) Self-attention generative adversarial networks. In: ICML
- Zhang P, Lan C, Xing J, et al (2017) View adaptive recurrent neural networks for high performance human action recognition from skeleton data. In: ICCV
- Zhang P, Xue J, Lan C, et al (2018) Adding attentiveness to the neurons in recurrent neural networks. In: ECCV
- Zhao H, Jia J, Koltun V (2020) Exploring self-attention for image recognition. In: CVPR
- Zhou B, Andonian A, Oliva A, et al (2018a) Temporal relational reasoning in videos. In: "ECCV"
- Zhou J, Troyanskaya O (2014) Deep supervised and convolutional generative stochastic network for protein secondary structure prediction. In: ICML, pp 745–753
- Zhou J, Wang H, Zhao Z, et al (2018b) Cnnh.pss: protein 8-class secondary structure prediction by convolutional neural network with highway. BMC bioinformatics 19(4):99–109