

# Multi-Task Recurrent Modular Networks

Dongkuan Xu<sup>1\*</sup>, Wei Cheng<sup>2\*</sup>, Xin Dong<sup>3</sup>, Bo Zong<sup>2</sup>, Wenchao Yu<sup>2</sup>, Jingchao Ni<sup>2</sup>,  
Dongjin Song<sup>4</sup>, Xuchao Zhang<sup>2</sup>, Haifeng Chen<sup>2</sup>, Xiang Zhang<sup>1</sup>

<sup>1</sup>The Pennsylvania State University, {dux19, xzz89}@psu.edu

<sup>2</sup>NEC Laboratories America, Inc., {weicheng, bzong, wyu, jni, xuczhang, haifeng}@nec-labs.com

<sup>3</sup>Rutgers University, {xd48}@rutgers.edu <sup>4</sup>University of Connecticut, {dongjin.song}@uconn.edu

## Abstract

We consider the models of deep multi-task learning with recurrent architectures that exploit regularities across tasks to improve the performance of multiple sequence processing tasks jointly. Most existing architectures are painstakingly customized to learn task relationships for different problems, which is not flexible enough to model the dynamic task relationships and lacks generalization abilities to novel test-time scenarios. We propose multi-task recurrent modular networks (MT-RMN) that can be incorporated in any multi-task recurrent models to address the above drawbacks. MT-RMN consists of a shared encoder and multiple task-specific decoders, and recurrently operates over time. For better flexibility, it modularizes the encoder into multiple layers of sub-networks and dynamically controls the connection between these sub-networks and the decoders at different time steps, which provides the recurrent networks with varying degrees of parameter sharing for tasks with dynamic relatedness. For the generalization ability, MT-RMN aims to discover a set of generalizable sub-networks in the encoder that are assembled in different ways for different tasks. The policy networks augmented with the differentiable routers are utilized to make the binary connection decisions between the sub-networks. The experimental results on three multi-task sequence processing datasets consistently demonstrate the effectiveness of MT-RMN.

## Introduction

Deep learning models with recurrent architectures have been widely studied for sequence processing tasks, such as sequence labeling (Lin et al. 2018; Cui et al. 2019), time series analysis (Dennis et al. 2019; Guo, Lin, and Antulov-Fantulin 2019), spatio-temporal prediction (Wang et al. 2018; Zhou et al. 2018; Sun et al. 2019) and question answering tasks (Ke, Zolna et al. 2018). These models utilize a chain of repeating cells to encode the global information of a sequence into a word-level representation of its elements. Recent progresses have been made on jointly learning multiple tasks to improve overall performance by exploiting regularities across tasks (Chen et al. 2018; Gupta, Chakraborty, and Chakrabarti 2019; Cui et al. 2019; Lu, Bai, and Langlais 2019). Multi-task learning (MTL) not only helps reduce the risk of over-fitting to

individual tasks, but also saves computation cost by sharing model architectures and low-level representations (Ruder 2017; Zhang and Yang 2017; Ma et al. 2019; Stickland and Murray 2019; Meyerson and Miikkulainen 2019).

The multi-task architectures applicable to recurrent models, however, are underexplored in literature. Existing approaches typically adopt the general architectures of multi-task learning to address sequence processing tasks. Most of them use a set of task-specific decoders to extract task-specific knowledge (Luong et al. 2016; Subramanian et al. 2018) and an ad-hoc structure to store the knowledge shared by related tasks (Chen et al. 2018), however, the task relationships are not considered explicitly and the structure is usually pre-defined, restricting the model capacity. It is also infeasible when users have limited knowledge of the task relatedness. Moreover, we highlight the dynamics of task relationships, which indicates the strength of relatedness between tasks is not static but subject to change, depending on the input data at hand (Cui et al. 2019; Liu et al. 2019), such as the relationships between the POS tagging tasks of code-switched sentences at different token positions (Gupta, Chakraborty, and Chakrabarti 2019; Lu, Bai, and Langlais 2019). Existing approaches are not flexible enough to learn the dynamic relationships, which calls for effective multi-task architecture that is elaborately developed for recurrent models. In addition, for more complex tasks of sequence processing, we conjecture that the generalization ability is essential (Lake 2019; Goyal et al. 2019), which helps the models handle the unseen tasks that are common in the real-world applications (Chang et al. 2019; Huang, Nair et al. 2019; Purushwalkam et al. 2019). Previous studies show that recurrent neural network models generalize well when the training and test data is similar, but fail spectacularly when generalization requires systematic compositional skills (Lake and Baroni 2018; Loula, Baroni, and Lake 2018). How to build recurrent neural architectures with the compositional generalization abilities would constitute a daunting challenge.

To address these challenges, we propose multi-task recurrent modular networks (MT-RMN) that dynamically learn task relationships and accordingly learn to assemble composable modules into complex layouts to jointly solve multiple sequence processing tasks. MT-RMN consists of a shared encoder to store the knowledge shared by related tasks and multiple decoders to extract task-specific knowledge. It

\*This work was done when the first author was an intern at NEC Labs America. Wei Cheng is the corresponding author.  
Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

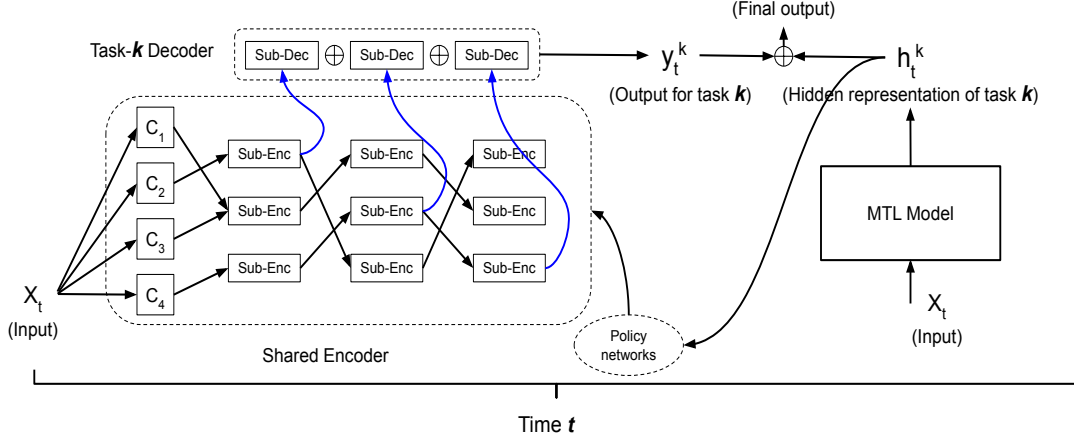


Figure 1: Architecture of MT-RMN.

recurrently operates over time. MT-RMN is general and with good compositional generalization ability. It is easy to be incorporated in any multi-task recurrent model by combining the extracted task-specific knowledge with the task-specific output of the recurrent model.

Specifically, we employ modularity as a way of improving adaptability and generalization to novel test-time scenarios (Chang et al. 2019; Lake 2019; Pathak et al. 2019). The encoder is modularized into layers of sub-networks and a task-specific decoder consists of several sub-decoders. The first layer of the encoder is a group of recurrent cells and other layers consist of sub-encoders, representing a set of generalizable sub-networks that are assembled in different ways for different tasks. More importantly, modularity provides flexibility for the recurrent model to capture the dynamics of task relationships. The recurrent cells in MT-RMN operate independently and are used to capture the dynamics originated from different feature subspaces, similar to (Goyal et al. 2019). By dynamically learning the connections between the sub-networks of the encoder and the sub-decoders at different time steps, MT-RMN selectively activates different parts of the encoder for a task over time, which provides the recurrent networks with varying degrees of parameter sharing for tasks with dynamic relatedness. The connections are learned by decision policies (Guo et al. 2019; Ahmed and Torresani 2019), which are sampled from discrete distributions indirectly parameterized by the output of lightweight neural networks called policy networks. These policies decide to connect or disconnect the route between two sub-networks, on a per-instance basis.

Intuitively, similar tasks are routed through similar paths and dissimilar tasks are routed through different paths. This enables the model to express tasks as a combination of sub-tasks and to generalize to unseen categories by dynamically routing a set of sub-encoders. As such, apart from reusability and good generalization, modularity also provides interpretability. Users can inspect how the network operates to understand which tasks are deemed more related, and at which time steps the task relatedness dramatically change, etc. As the decisions are non-differentiable, we propose the dif-

ferentiable binary routers based upon the Gumbel-Softmax reparameterization (Jang, Gu, and Poole 2017; Maddison, Mnih, and Teh 2017) to train the policy networks, which are jointly trained with other parts of the networks. MT-RMN is first validated on a task-fMRI dataset, where we analyze its ability to model task relationships dynamically and its generalization ability, then applied to a part-of-speech tagging dataset and a clinical time series dataset. Experimental results verify the effectiveness of MT-RMN.

## The MT-RMN Architecture

MT-RMN consists of a shared encoder and multiple task-specific decoders as shown in Figure 1. The connection decisions are made by policy networks augmented with straight-through routers.

### Shared Encoder

The shared encoder is motivated by modular networks (Kirsch, Kunze, and Barber 2018; Pahuja et al. 2019). We aim to discover a set of generalizable sub-networks that are assembled in different ways for different tasks. The encoder has  $m$  layers and each layer  $\ell \in \{1, 2, \dots, m\}$  contains  $n_\ell$  sub-networks. The first layer is a group of independent LSTM cells. Other layers consist of MLPs that are used as sub-encoders.

**Recurrent Independent Cells.** Let  $x_t$  be the input data at time step  $t$ . We assume that the overall dynamical system of interest is divided into  $k$  small subsystems (Parascandolo et al. 2018; Goyal et al. 2019). We use  $k$  recurrent cells to model these subsystems and each cell has its own independent dynamics operating. The cell denoted by  $C_i$  is flexible to use different RNN modules and we use LSTM (Hochreiter and Schmidhuber 1997) in this paper. The hidden state  $h_{t,i}$  after cell  $C_i$  is applied is

$$h_{t,i} = \text{LSTM}_i(h_{t-1,i}, x_t; \theta_i). \quad (1)$$

**Policy Networks for Sub-Encoders.** We build a policy network (Sec. ) for each sub-encoder to make decision connections between itself and the sub-networks of the previous layer. Specifically, for sub-encoder  $j$  in layer  $\ell$ , policy network  $\mathcal{N}_j$  estimates a decision vector  $\alpha_{t,ij} \in \mathbb{R}^2$  for every

sub-network  $i$  in layer  $\ell-1$  at time step  $t$ , given the output  $\mathbf{u}_{t,i}$  of sub-network  $i$  at time step  $t$ .

$$\{\mathcal{N}_j : \mathbf{u}_{t,i} \rightarrow \alpha_{t,ij} | i \in \{1, \dots, n_{\ell-1}\}, j \in \{1, \dots, n_\ell\}\} \quad (2)$$

**Routers.** Given  $\alpha_{t,ij}$ , a straight-through router (Sec. ) is used to learn the decision policy  $\mathcal{P}_{ij}$ , which estimates a binary decision value  $\zeta_{t,ij} \in \{0, 1\}$ , indicating whether to connect ( $\zeta_{t,ij}=1$ ) or disconnect ( $\zeta_{t,ij}=0$ ) the route between sub-networks  $i$  and sub-encoder  $j$  at time step  $t$ .

$$\{\mathcal{P}_{ij} : \alpha_{t,ij} \rightarrow \zeta_{t,ij} | i \in \{1, \dots, n_{\ell-1}\}, j \in \{1, \dots, n_\ell\}\} \quad (3)$$

**Calculation of Input and Output.** Each sub-encoder  $j$  in layer  $\ell$  receives a list of  $n_{\ell-1}$  tuples of features from the sub-networks in layer  $\ell-1$ , where tuple  $(\mathbf{u}_{t,i}, \zeta_{t,ij})$  corresponds to sub-network  $i$  in layer  $\ell-1$ . The input  $\mathbf{v}_{t,j}$  and output  $\mathbf{u}_{t,j}$  of sub-network  $j$  in layer  $\ell$  are

$$\mathbf{v}_{t,j} = \sum_{i=1}^{n_{\ell-1}} \zeta_{t,ij} / \left( \sum_{i'=1}^{n_{\ell-1}} \zeta_{t,i'j} \right) \mathbf{u}_{t,i}, \quad (4)$$

$$\mathbf{u}_{t,j} = \text{MLP}_j(\mathbf{v}_{t,j}). \quad (5)$$

### Task-Specific Decoders

Task-specific decoders share a similar idea with the encoder and they selectively activate only parts of the encoder. Every task-specific decoder  $\mathcal{D}^k$  has  $m-1$  independent sub-decoder  $\mathcal{D}_\ell^k$ , which are used to extract task-related knowledge from different layers of the encoder. We use MLPs as the sub-decoders in this paper. In particular, task-specific decoder  $\mathcal{D}^k$  is used for task  $k$  and sub-decoder  $\mathcal{D}_\ell^k$  is used to extract the knowledge from the  $\ell$ -th layer of encoder. Because of the domain specialization of different tasks (Cui et al. 2019), we use different policy networks for different tasks.

The input and output calculation of the sub-decoder is similar to the one of the sub-encoder. For sub-decoder  $\mathcal{D}_\ell^k$ , we use policy network  $\hat{\mathcal{N}}^k$  to estimate a decision  $\beta_{t,j}^k \in \mathbb{R}^2$  for sub-encoder  $j$  in layer  $\ell$  of encoder, i.e.,  $\{\hat{\mathcal{N}}^k : \mathbf{u}_{t,j} \rightarrow \beta_{t,j}^k\}$ , where  $\mathbf{u}_{t,j}$  is the output of sub-encoder  $j$ . A router is further used to learn a policy  $\hat{\mathcal{P}}_j^k$  to estimates a binary decision value, i.e.,  $\{\hat{\mathcal{P}}_j^k : \beta_{t,j}^k \rightarrow \hat{\zeta}_{t,j}^k \in \{0, 1\}\}$ . Thus, sub-decoder  $\mathcal{D}_\ell^k$  receives a list of tuples  $(\mathbf{u}_{t,j}, \hat{\zeta}_{t,j}^k)$ , each of which corresponds to a sub-encoder of layer  $\ell$ . The input and output of sub-decoder  $\mathcal{D}_\ell^k$  are  $\hat{\mathbf{v}}_{t,\ell}^k = \sum_{j=1}^{n_\ell} \hat{\zeta}_{t,j}^k / (\sum_{j'=1}^{n_\ell} \hat{\zeta}_{t,j'}^k) \mathbf{u}_{t,j}$  and  $\hat{\mathbf{u}}_{t,\ell}^k = \text{MLP}_\ell(\hat{\mathbf{v}}_{t,\ell}^k)$  respectively. To consider the hierarchical structure information of the encoder, we concatenate the output of every sub-decoder  $\mathcal{D}_\ell^k$  to construct the output of  $\mathcal{D}^k$  as  $[\hat{\mathbf{u}}_{t,1}^k \oplus \dots \oplus \hat{\mathbf{u}}_{t,m-1}^k]$ .

### Policy Networks

The role of policy networks is to estimate the decision vector  $\alpha_{t,ij}$  (or  $\beta_{t,j}^k$ ), which is further fed into the routers to make the binary connection decisions. The form of policy networks is flexible. We opt to use MLP as the policy network for the encoder, which is defined as  $\alpha_{t,ij} = \mathcal{N}_j(\mathbf{u}_{t,i}) = \text{MLP}(\mathbf{u}_{t,i})$ . Note that MT-RMN is used to be incorporated

in an existing multi-task recurrent model. There are two types of multi-task recurrent models in general and they have different ways to incorporate MT-RMN. Thus, we propose two types of policy networks for the decoders. First, when the multi-task model itself generates task-specific representation  $\mathbf{r}_k$  for task  $k$ , the policy networks are defined as

$$\beta_{t,j}^k = \hat{\mathcal{N}}^k(\mathbf{u}_{t,j}) = \text{MLP}(\mathbf{u}_{t,j} \oplus \mathbf{W}_k \mathbf{r}_k), \quad (6)$$

where  $\mathbf{W}_k$  is the transformation matrix for task  $k$ . Second, when the multi-task model does not generate task-specific representations, the policy networks are defined as

$$\beta_{t,j}^k = \hat{\mathcal{N}}^k(\mathbf{u}_{t,j}) = \text{MLP}(\mathbf{u}_{t,j}). \quad (7)$$

The policy networks are jointly trained with other parts of MT-RMN. Their simple architecture makes the estimation of the decision vector fast and efficient, which is similar to the design of the policy network in (Guo et al. 2019) and the design of decision-learner in (Ahmed and Torresani 2019).

### Straight-Through Router

Given the output,  $\alpha$  (or  $\beta$ ), of a policy network, a router is applied to learn the decision policy which estimates a binary decision value  $\zeta \in \{0, 1\}$ . Conceptually, the policy  $\mathcal{P}$  can be seen as a binarization function of the decision scores  $\alpha = \{\alpha_0, \alpha_1\}$  and each value in the pair of the binary outcomes is the complement of the other. A simple way to implement the binarization function is to select the position with maximum value of  $\{\alpha_0, \alpha_1\}$ , however, this approach is non-differentiable. There are several ways that allow us to propagate gradients through the discrete nodes (Bengio, Léonard, and Courville 2013). In this work, we adopt the Gumbel-Softmax sampling approach (Jang, Gu, and Poole 2017; Maddison, Mnih, and Teh 2017).

We use Gumbel-Max trick (Maddison, Mnih, and Teh 2017). Specifically, we have two classes (disconnect and connect).  $\{\alpha_0, \alpha_1\}$  represent the log probabilities  $\{\log(p_0), \log(p_1)\}$  of the two classes. Thus, we can draw samples from a Bernoulli distribution parameterized by class probabilities  $\{p_0, p_1\}$  in the following way: we first draw i.i.d samples  $\{g_0, g_1\}$  from a Gumbel distribution, i.e.,  $g = -\log(-\log(x)) \sim \text{Gumbel}$ , where  $x \sim \text{Uniform}(0, 1)$ . Then produce the discrete sample by adding  $g$  to introduce stochasticity:

$$\zeta = \arg \max_i [\alpha_i + g_i], i \in \{0, 1\}. \quad (8)$$

The arg max operation is non-differentiable, but we can use the softmax as a continuous differentiable approximation to it (Eq. (9)).  $\tau$  is the temperature to control the discreteness. Thus, we use the arg max to make the binary connection decision on the forward pass, while approximate it with softmax on the backward pass, which is called the straight-through estimator (Jang, Gu, and Poole 2017).

$$\mu_i = \exp((\alpha_i + g_i)/\tau) / \left( \sum_{\hat{i}=0}^1 \exp((\alpha_{\hat{i}} + g_{\hat{i}})/\tau) \right), i \in \{0, 1\} \quad (9)$$

## Sparsity Constrains

We add the sparsity constraint on the connection decision values  $\zeta_{t,ij}$ . We want the proposed module to avoid connecting every pair of sub-networks in the encoder, which can be viewed as the coarse-grained neural network pruning (Frankle and Carbin 2019) and benefits the specialization of sub-networks (Goyal et al. 2019). We also add the sparsity constraint on the connection decision values  $\hat{\zeta}_{t,i}^k$ , which helps the task-specific decoders avoid focusing on every sub-network of the encoder.

Inspired by (Ke, Zolna et al. 2018), we add two penalty terms  $\lambda_1 R_1(\zeta)$  and  $\lambda_2 R_2(\hat{\zeta})$ . They represent the penalties for connecting sub-networks in the encoder and the penalty for connections between the decoders and these sub-networks. Their definitions are

$$\lambda_1 R_1(\zeta) = \lambda_1 \text{ReLU}((\sum \zeta_{t,ij}) - \gamma_1 C_1), \quad (10)$$

$$\lambda_2 R_2(\hat{\zeta}) = \lambda_2 \text{ReLU}((\sum \hat{\zeta}_{t,i}^k) - \gamma_2 C_2), \quad (11)$$

$\lambda$  and  $\gamma$  are hyper-parameters.  $C_1$  is the number of all possible connections between sub-networks in the encoder and  $C_2$  is the number of all possible connections between the sub-networks and sub-decoders.  $\gamma$  represents the proportion of connections that are without being penalized. Intuitively, each connection above the threshold  $\gamma_1 C_1$  or  $\gamma_2 C_2$  is penalized.

**Complexity Analysis.** MT-RMN is local in time. Its complexity per time step is proportional to the number of parameters. MT-RMN includes the encoder, decoders and policy networks, and they contain  $(ns+cr)$ ,  $K(m-1)d$ ,  $(K+n)p$  parameters respectively, where  $n$ ,  $s$ ,  $c$ ,  $r$ ,  $K$ ,  $m$ ,  $d$ ,  $p$  are #sub-encoders, size of a sub-encoder, #cells, size of a cell, #tasks, #layers of an encoder, size of a sub-decoder, size of a policy network respectively. Thus, the complexity of MT-RMN is  $\mathcal{O}(ns+cr+Kms)$ .

## Related Work

**Multi-Task Learning.** A number of MTL approaches (Luo et al. 2016; Subramanian et al. 2018) learned a shared low-level representation that is followed by unshared decoders to extract task-specific representations. However, the typical MTL methods might degenerate a lot when tasks are less related (Ruder 2017; Zhang and Yang 2017; Ma et al. 2019). Several approaches have been proposed to address this challenge. (Misra et al. 2016; Ruder et al. 2019) adopted trainable parameters to control the communications between different task-specific layers to share knowledge between tasks. Some used the routing methods to achieve flexible parameter sharing (Rosenbaum, Klinger, and Riemer 2018; Ma et al. 2018, 2019). However, the routing learned by (Ma et al. 2018, 2019) was applied to all tasks, instead of task-dependent. (Rosenbaum, Klinger, and Riemer 2018) employed the multi-agent reinforcement learning approach to jointly learn the routing and function blocks, but it did not focus on recurrent models. (Chen et al. 2018; Liu et al. 2019; Gupta, Chakraborty, and Chakrabarti 2019) applied MTL to sequence learning, however, these models are limited because they used a pre-defined structure to learn task relatedness. (Cases, Rosenbaum et al. 2019; Rosenbaum et al. 2019a) are closely related to ours. They extended (Rosenbaum, Klinger,

and Riemer 2018) to natural language understanding. Our model differs from them in two crucial ways. First, our model considers layers of function modules instead of the recursive application of modules. The promise is that we can assemble the modules more flexibly to account for more complex task relationships, such as the hierarchical relationship (Huang, Nair et al. 2019; Lin et al. 2015). Second, our routing mechanism is more flexible, allowing multiple and dynamic connections between two layers of sub-encoders, which benefits the model capacity.

**Modular Networks.** We take inspiration from recent advances of modular works (Kirsch, Kunze, and Barber 2018; Pahuja et al. 2019). Modular networks are composed of modules and each module is a function with its own parameters. Learning an efficient set of such modules is akin to learning a set of functional primitives, which are evident in the natural world and can be combined to solve a given task (Meyerson and Miikkulainen 2018). The major advantage of modular networks is the compositional generalization ability, which has been studied in the context of modeling dynamical systems (Chang et al. 2019; Huang, Nair et al. 2019; Purushwalkam et al. 2019). Some recent progress has been made on the generalization of the complex sequence processing tasks (Lake and Baroni 2018; Loula, Baroni, and Lake 2018; Lake 2019). However, these methods are not designed for MTL. In addition, modular networks are clearly related to routing networks (Haimlerl, Savin, and Simoncelli 2019; Ramachandran and Le 2019; Rosenbaum, Klinger, and Riemer 2018), where each sample selectively activates only parts of the entire network. Some routing networks based on conditional computation have been proposed to improve computational efficiency (Bengio et al. 2015; Ke, Zolna et al. 2018; Guo et al. 2019). Several approaches routed the examples through a network in an MTL setting (Rosenbaum, Klinger, and Riemer 2018; Ramachandran and Le 2019). However, unlike these methods, we consider the recurrent models. (Cases, Rosenbaum et al. 2019; Rosenbaum et al. 2019a) considered the modular recurrent representation via routing. The differences between our model and them are discussed in the above paragraph. More introduction about modular and routing networks can be found in (Rosenbaum et al. 2019b).

**Neural Architecture Search.** Neural architecture search (NAS) aims to design the neural network architecture for a given task automatically (Zoph and Le 2017; Real et al. 2019). Our work can be considered as a special case of NAS, where we search for a dynamical multi-task model architecture that consists of a set of generalizable sub-networks over time. We particularly care about the parameter sharing problem dynamically. Previous NAS methods usually used an RNN-based controller to generate model architecture and then trained the target task based on this architecture (Zoph and Le 2017). Unfortunately, they are typically with high computation cost. Some efficient NSA methods have been proposed recently (Pham et al. 2018; Liu, Simonyan, and Yang 2019) and our module is of a similar spirit in the sense that the parameters of architecture and model are learned jointly, which is also shared by (Ma et al. 2019; Maziarz et al. 2019). However, all these methods focused on static

models, while we focus on the recurrent models that differ per example and per time step.

## Experiments

In our experiments, all the models are trained by Adam (Kingma and Ba 2014). The reported results were got by 5 times 5-fold cross validation. To make fair comparisons, we tuned the baselines as much as we could. We tried different parameter settings via grid-search, increased the number of parameters to be similar and utilized different tricks. The learning rate was set to  $10^{-3}$  initially and decreased during the training. We experimented with the values of  $k$ ,  $\lambda_1$ ,  $\lambda_2$ ,  $\gamma_1$ ,  $\gamma_2$  in the sets  $\{2, 3, 4\}$ ,  $\{2^{-2}, 2^{-1}, 2^0, 2^1, 2^2\}$ ,  $\{2^{-2}, 2^{-1}, 2^0, 2^1, 2^2\}$ ,  $\{0.25, 0.5, 0.75, 1\}$ ,  $\{0.25, 0.5, 0.75, 1\}$  respectively. We obtained best results using  $k=3$ ,  $\lambda_1=1$ ,  $\lambda_2=1$ ,  $\gamma_1=0.75$ ,  $\gamma_2=0.75$  in general. For the temperature  $\tau$  used in the Gumbel-Max trick, we set its value as 100 initially and divided the value by 2 at each epoch. For all other parameters in the baseline methods, we adopted the default settings according to their original papers. The experimental results on the clinic time series dataset and the analysis of different model architectures are summarized in the appendix<sup>1</sup>.

### Task-Evoked Functional Brain Activity

To test the flexibility and the generalization ability of MT-RMN in a controlled environment where we know the dynamic relationships between tasks, we create the *tfMRI* tasks based upon the task-evoked functional MRI dataset<sup>2</sup> (task-fMRI), which is widely used to analyze the relationship between brain connectivity and human behavior (Barch, Burgess et al. 2013). Figure 2 shows the visualization of the task-fMRI data of a participant who performed seven tasks successively. There are twelve time series grouped into four groups and each time series corresponds to a high-order region of human brains. The x-axis represents time and the y-axis represents the normalized activation degree (see (Barch, Burgess et al. 2013) for more details). The functional connectivity between regions varies over time, which indicates the relationships between the sequence processing tasks of brain regions are dynamic.

A *tfMRI* task is defined as a multi-class classification task of time series and we construct four tasks as shown in Figure 2. The brain regions involved in Task-1 and Task-4 are related to the cognition and sensory systems of humans respectively. The regions in Task-2 have mixed functions related to both cognition and sensory abilities. Task-3 is more related to the attention system. To extract meaningful features, each time series were further processed to be 28 time steps long with 10 feature dimensions (see supplementary material for pre-processing details). Every 4 time steps represent the time period when the participant conducted a task. The data was split into train (80%) and test (20%) sets randomly.

**Flexibility** To evaluate the flexibility, we consider two groups of *tfMRI* tasks. Group-1 includes Tasks-1,2,3 and Group-2 includes Tasks-1,3,4. The relationships between

these tasks are dynamic but the tasks of Group-1 are more related in general compared to the tasks in Group-2. We create an MTL model (*mtl*-RMN) based on our networks. The network consists of three standard LSTMs that are connected to MT-RMN at each time step. Each LSTM is used for a *tfMRI* task to extract the task-specific hidden representation. The MT-RMN used here consists of a three-layer encoder and three task-specific decoders. The first layer of the encoder includes three LSTM cells and each of other layers contains three sub-encoders. Each decoder includes two sub-decoders. We choose this architecture because of its low enough capacity that contributes to the visible competition between tasks (Ramachandran and Le 2019).

We use five baselines that use different ways to utilize task relationships. FS-MTL (Caruana 1993) uses a fully shared neural layer for all tasks along with separate task-specific output layers. SP-MTL (Liu, Qiu, and Huang 2017) uses the adversarial training to help the shared neural layer contains the task-invariant knowledge and uses the orthogonality constraint to eliminate knowledge redundancy. DC-MTL (Liu et al. 2019) has task-specific neural layers for different tasks and each task shares knowledge between any other directly. IC-MTL (Liu et al. 2019) differs from DC-MTL that uses an additional layer to share knowledge indirectly between tasks. RRNs (Cases, Rosenbaum et al. 2019) recursively chooses sub-functions from a single set of composable functions via a router for different input examples. We apply RRNs to route an LSTM input transformation.

We run these methods on Group-1 and Group-2 in the MTL setting, where the task weights are set to the same. We apply the softmax layer to the hidden representation at the last time step to get the predicted label. The results are reported in Table 1. It is observed that *mtl*-RMN achieved the best performance. When the tasks were more related (Group-1), all the methods showed high performance. However, the margin between the performance of *mtl*-RMN and the baselines was more clear when tasks were less related (Group-2). Both DC-MTL and IC-MTL outperformed FS-MTL on Group-1, but DC-MTL showed a little bit lower performance on Group-2. This is because the way of sharing knowledge in DC-MTL is more suitable for the case when tasks are more related. Compared to other baselines, RRNs showed higher performance on Group-2, because RRNs can limit task interference better for less related tasks. However, *mtl*-RMN outperformed RRNs on both groups, indicating the routing mechanism of *mtl*-RMN is more flexible. This might be because *mtl*-RMN can not only limit the shared weights for less related tasks, but also increase the shared ones for more related tasks via its dynamic/multiple connections between sub-encoder layers. We also report the loss (averaged over tasks) of FS-MTL, IC-MTL and *mtl*-RMN on the two groups over epochs, which are shown in Figure 3. We see that the average loss of *mtl*-RMN behaved in a similar way on the two groups, while FS-MTL and IC-MTL degenerated clearly when tasks were less related. This might be because there was task interference during training when tasks were less related, which resulted in negative transfer between tasks. *mtl*-RMN is more flexible and can reduce negative transfer more effectively. All the observations demonstrate that

<sup>1</sup>Data and codes can be found in the authors' website.

<sup>2</sup><https://tinyurl.com/y4hhw8ro>

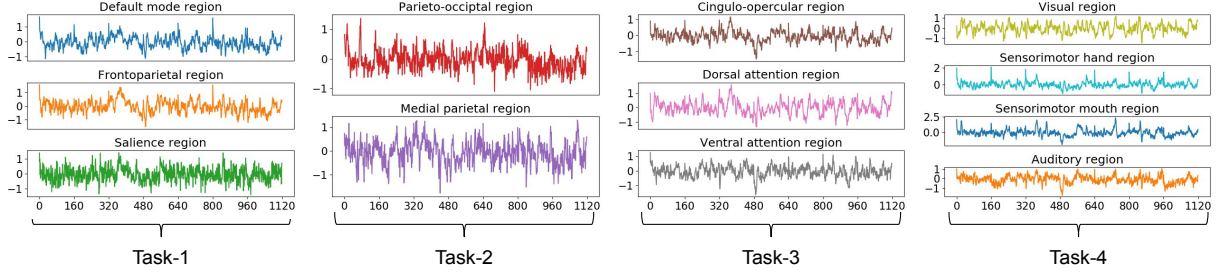
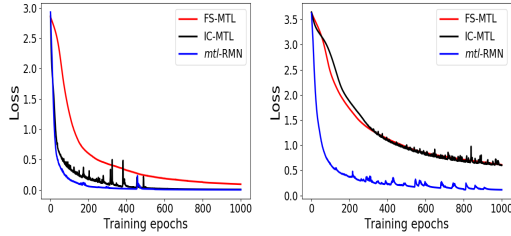


Figure 2: The task-fMRI data of twelve high-order brain regions of a participant. The regions are grouped into four groups based on their functionalities. Each group is used to construct a multi-class classification task of time series. The participant was asked to perform seven tasks successively (Barch, Burgess et al. 2013).

Table 1: Results (accuracy %) on two groups of *tfMRI* tasks.

Groups	Group-1			Group-2		
	Task-1	Task-2	Task-3	Task-1	Task-3	Task-4
FS-MTL	89.7 $\pm$ 0.6	93.5 $\pm$ 1.1	89.1 $\pm$ 0.3	73.2 $\pm$ 0.5	81.2 $\pm$ 0.6	81.3 $\pm$ 1.4
SP-MTL	88.7 $\pm$ 1.4	94.0 $\pm$ 1.9	90.7 $\pm$ 2.3	73.7 $\pm$ 1.5	78.0 $\pm$ 1.5	81.2 $\pm$ 2.1
DC-MTL	89.2 $\pm$ 0.6	95.6 $\pm$ 1.2	90.8 $\pm$ 1.0	74.5 $\pm$ 0.6	80.0 $\pm$ 0.7	80.8 $\pm$ 1.3
IC-MTL	89.6 $\pm$ 0.4	95.7 $\pm$ 1.1	91.3 $\pm$ 1.7	74.5 $\pm$ 0.4	81.4 $\pm$ 1.1	81.8 $\pm$ 2.6
RRNs	88.9 $\pm$ 3.4	95.4 $\pm$ 2.2	90.5 $\pm$ 2.8	75.4 $\pm$ 2.2	83.2 $\pm$ 3.0	82.6 $\pm$ 3.8
<i>mtl</i> -RMN	<b>90.8<math>\pm</math>2.1</b>	<b>96.7<math>\pm</math>1.6</b>	<b>92.9<math>\pm</math>2.4</b>	<b>76.1<math>\pm</math>1.8</b>	<b>84.4<math>\pm</math>2.4</b>	<b>83.5<math>\pm</math>2.2</b>



(a) Tasks are more related. (b) Tasks are less related.

Figure 3: Loss comparison when tasks are more related (*left*) and less related (*right*).

MT-RMN is able to learn task relationships flexibly.

**Generalization Ability** To evaluate the generalization ability, we create four novel test-time scenarios, *A*, *B*, *C*, *D*, based on the four *tfMRI* tasks, as shown in the left part of Table 2. Each case consists of three training tasks and a test task that is not included in the training task set. We hope to test if MT-RMN could discover a set of generalizable sub-networks to capture the sub-task patterns, which is critical for the compositional skills. We create a new model called  $RMN_{tr}$ . Specifically, we first train *mtl*-RMN on the training tasks. After training, we fix the parameters of the encoder and its policy networks. Then, we connect the fixed encoder (including its policy networks) with a new LSTM and a new decoder with a policy network. The data of test task is fed into both the new LSTM and the fixed encoder. The output of the new decoder is concatenated with the hidden representation of the new LSTM as the final representation which is further fed into a softmax layer.

We use four baselines: LSTM,  $RMN_{un}$ , RRNs and IC-MTL. We run LSTM on the test task directly.  $RMN_{un}$  is a model where an untrained encoder of MT-RMN is connected with an LSTM and a decoder. We run  $RMN_{un}$  on the test task

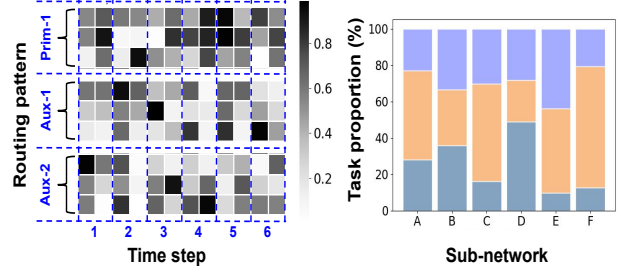


Figure 4: (*Left*) The routing patterns of examples of each task. The value represents the proportion of examples which activate the corresponding sub-network. (*Right*) The proportion of tasks assigned to each sub-network. Different colors distinguish different tasks and A-F represent the six sub-networks.

directly. For IC-MTL, we first train it on training tasks. After training, the shared LSTM in IC-MTL is fixed. Then we connect it with a new LSTM. The results are reported in the right part of Table 2.  $RMN_{tr}$  outperformed LSTM, which indicates that MT-RMN can help the base model (LSTM) handle the unseen tasks.  $RMN_{tr}$  outperformed  $RMN_{un}$ , which indicates that the regularities across tasks help improve the performance of a single task. Both  $RMN_{tr}$  and RRNs outperformed IC-MTL, indicating the modularity benefits generalization. Task-2 performed similarly in Tables 1-2, because Task-2 might be an easy task and the task relatedness does not influence its performance much. Tasks-1,3 performed worse in Table 2 than that in Group-1 of Table 1 but better than that in Group-2, because training tasks are more/less related in Group-1/2 than the ones in both scenarios B and D. All these observations demonstrate that MT-RMN can improve the systematic compositional skills of the base models and further help them generalize well for unseen tasks.



Table 2: Four test-time scenarios to evaluate generalization ability and the results (accuracy%).

Scenario settings					Results of different methods				
Scenarios	Training tasks			Test task	LSTM	RMN <sub>un</sub>	RRNs	IC-MTL	RMN <sub>tr</sub>
A	Task-1	Task-2	Task-3	Task-4	82.6±1.0	82.9±1.2	83.0±2.7	81.2±1.7	<b>83.3±2.3</b>
B	Task-1	Task-2	Task-4	Task-3	82.3±0.6	84.5±1.5	84.8±2.5	82.7±1.3	<b>85.6±1.7</b>
C	Task-1	Task-3	Task-4	Task-2	93.7±0.6	95.3±0.8	95.6±1.7	95.5±0.7	<b>97.0±1.8</b>
D	Task-2	Task-3	Task-4	Task-1	71.7±1.4	73.1±1.7	78.4±4.0	77.3±2.1	<b>79.3±2.6</b>

Table 3: Results of POS tagging.

Methods	Original		DC-MTL		IC-MTL		<i>rrn</i> -RMN		MT-RMN	
	Accuracy	F <sub>1</sub>	Accuracy	F <sub>1</sub>	Accuracy	F <sub>1</sub>	Accuracy	F <sub>1</sub>	Accuracy	F <sub>1</sub>
NS-MTL	44.3±0.7	38.9±1.6	51.5±0.5	49.4±0.9	51.4±1.4	49.7±0.8	53.5±1.9	52.7±0.8	<b>55.2±1.2</b>	<b>53.4±2.4</b>
Cross-stitch.	46.0±1.9	12.7±0.8	48.5±0.9	16.1±1.3	47.7±1.7	14.3±0.7	<b>52.1±1.0</b>	19.5±0.6	51.7±1.9	<b>21.4±3.2</b>
Sluice	58.7±0.7	55.4±2.1	59.2±1.7	56.0±2.3	59.5±1.9	56.2±0.7	60.4±1.2	56.6±0.7	<b>61.5±1.7</b>	<b>57.2±1.8</b>
GIRNet	62.8±0.6	46.6±0.4	62.5±1.1	57.4±0.5	63.1±1.5	58.6±1.3	63.6±1.7	61.3±1.5	<b>64.5±2.1</b>	<b>62.6±2.4</b>

### POS Tagging of Code-Switched Sentences

To directly evaluate the ability of MT-RMN to improve the performance of multi-task models on sequence processing tasks, we test MT-RMN on the multi-task setup of part-of-speech (POS) tagging (Santos and Zadrozny 2014) of code-switched sentences, where the words are from two languages and the goal is to label the sequences from code-switched text (Gonen and Goldberg 2019). Code switching is common, especially in social media where the participants are multilingual users (Vyas et al. 2014). Given the input sentence denoted by  $(w_1, w_2, w_3, \dots)$ , where words are from languages  $A, B$ , i.e.,  $w_i \in D_A \cup D_B$ , the goal is to predict a label sequence  $(y_1, y_2, y_3, \dots)$ . We pre-process the data and create our setup exactly following (Gupta, Chakraborty, and Chakrabarti 2019). Specifically, we use the Hindi-English code-switched dataset provided in (Patra, Das, and Das 2018) for the primary task. The dataset has 19 POS tags and contains 2102 and 528 instances for the training and test sets respectively. We use the Hindi POS tagging dataset provided in (Sachdeva et al. 2014) and the English one provided in (Sang and Buchholz 2000) for the auxiliary tasks of Hindi and English respectively. The Hindi dataset includes 14084 instances with 25 POS tags and the English dataset includes 8936 instances with 45 POS tags.

We compare MT-RMN to three baselines, including DC-MTL, IC-MTL and a variant called *rrn*-RMN, where we replace the layers of sub-encoders with one set of sub-encoders in the recursive layout. We further incorporate these methods into five MTL approaches to test their ability to improve these MTL approaches. NS-MTL has an independent LSTM layer for each task and all the tasks are jointly optimized. Cross-stitch (Misra et al. 2016) has task-specific layers for different tasks. It uses parameters to control the amount of information shared between layers. Sluice (Ruder et al. 2019) splits the channels of layers and controls the information sharing by parameters. The outputs of intermediate layers are fed to the output layers. GIRNet (Gupta, Chakraborty, and Chakrabarti 2019) contains two LSTMs that are trained over auxiliary instances. The state sequences of the two LSTMs are merged into a composite state sequence for the primary task.

**Results.** We show the results in Table 3. MT-RMN improved

these base models most compared to other baselines. NS-MTL performed much better when combined with any of the four methods, indicating the necessity of exploiting regularities across tasks to improve task performance. MT-RMN and *rrn*-RMN helped the base models perform better compared to others, which verifies the better ability of modular networks with routing to exploit regularities. Both Sluice and GIRNet outperformed Cross-stitch because they share knowledge at the granularity of word. All the observations verify the advantage of MT-RMN in improving the multi-task models on sequence processing tasks.

**Sub-Encoder Specialization.** To verify that the sub-encoders can capture primitives for different tasks, we visualize the routing patterns learned by MT-RMN. We use a three-layer encoder that has 3 cells in the 1st layer and 3 sub-encoders in each of others. We incorporate this MT-RMN into GIRNet and test it on the POS tagging task. After training, examples from the test set are passed through MT-RMN. The results of the first 6 time steps are shown in Figure 4. The routing pattern of each time step is represented by a  $3 \times 2$  grid, each of which corresponds to a sub-encoder. We aggregate the connection choices (between sub-encoders and sub-decoders) for every sample and concatenate the grids of different time steps. Figure 4 (left) shows that the routing patterns varied over time and the examples from different tasks tended to be routed through different sub-encoders. Figure 4 (right) shows how different sub-encoders handle different task distributions. All the observations demonstrate that sub-encoders focus on different tasks and can be assembled for different tasks.

### Conclusion

We propose MT-RMN, the multi-task recurrent modular networks to improve the performance of multi-task recurrent models. MT-RMN includes a shared encoder and a set of task-specific decoders. It modularizes the encoder into layers of sub-encoders for flexibility and dynamically learns the connections between the sub-networks and the decoders via decision policies. MT-RMN recurrently operates over time and is easy to be incorporated in any multi-task recurrent models. Experiments demonstrate the effectiveness of MT-RMN.

## Acknowledgements

This project was partially supported by NSF projects IIS-1707548 and CBET-1638320.

## Statement of Ethical Impact

Our proposed networks are general and can be easily incorporated in any recurrent model to improve its performance via exploiting regularities across tasks. There are several impacts of using our networks. First, our networks enable models to generalize to unseen categories, which are common in real-world, via dynamically routing a set of sub-encoders. Second, our networks help models save computing resources via reusing the well-trained and generalizable sub-encoders for different tasks. Third, our networks help users inspect how the model operates to understand the dynamic relationships between tasks, such as the analysis of brain connectivity and human behavior based on the fMRI dataset. However, our networks also put deep recurrent models at a high risk of being attacked. Our networks seek a set of sub-networks that are assembled in different ways for different tasks. Disturbing these sub-networks leads to significant changes on the task performance. Besides, increasing the interpretability of deep recurrent models might bring about automation bias, such as an unwarranted trust on the models.

## References

- Ahmed, K.; and Torresani, L. 2019. STAR-Caps: Capsule networks with straight-through attentive routing. In *NeurIPS*, 9098–9107.
- Barch, D. M.; Burgess, et al. 2013. Function in the human connectome: task-fMRI and individual differences in behavior. *Neuroimage* 80: 169–189.
- Bengio, E.; Bacon, P.-L.; Pineau, J.; and Precup, D. 2015. Conditional computation in neural networks for faster models. *arXiv preprint arXiv:1511.06297*.
- Bengio, Y.; Léonard, N.; and Courville, A. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*.
- Caruana, R. 1993. Multitask learning: A knowledge-based source of inductive bias. In *ICML*, 41–48. Morgan Kaufmann.
- Cases, I.; Rosenbaum, et al. 2019. Recursive routing networks: Learning to compose modules for language understanding. In *NAACL-HLT*, 3631–3648.
- Chang, M.; Gupta, A.; Levine, S.; and Griffiths, T. L. 2019. Automatically composing representation transformations as a means for generalization. In *ICLR*.
- Chen, J.; Qiu, X.; Liu, P.; and Huang, X. 2018. Meta multi-task learning for sequence modeling. In *AAAI*.
- Cui, W.; Zheng, G.; Shen, Z.; Jiang, S.; and Wang, W. 2019. Transfer learning for sequences via learning to collocate. In *ICLR*.
- Dennis, D.; Acar, D. A. E.; Mandikal, V.; Sadasivan, V. S.; Saligrama, V.; Simhadri, H. V.; and Jain, P. 2019. Shallow RNN: accurate time-series classification on resource constrained devices. In *NeurIPS*, 12896–12906.
- Frankle, J.; and Carbin, M. 2019. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *ICLR*.
- Gonen, H.; and Goldberg, Y. 2019. Language modeling for code-switching: Evaluation, integration of monolingual data, and discriminative training. In *EMNLP-IJCNLP*.
- Gordon, E. M.; Laumann, T. O.; Adeyemo, B.; Huckins, J. F.; Kelley, W. M.; and Petersen, S. E. 2016. Generation and evaluation of a cortical area parcellation from resting-state correlations. *Cerebral cortex* 26(1): 288–303.
- Goyal, A.; Lamb, A.; Hoffmann, J.; Sodhani, S.; Levine, S.; Bengio, Y.; and Schölkopf, B. 2019. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*.
- Guo, T.; Lin, T.; and Antulov-Fantulin, N. 2019. Exploring interpretable LSTM neural networks over multi-variable data. In *ICML*.
- Guo, Y.; Shi, H.; Kumar, A.; Grauman, K.; Rosing, T.; and Feris, R. 2019. SpotTune: transfer learning through adaptive fine-tuning. In *CVPR*, 4805–4814.
- Gupta, D.; Chakraborty, T.; and Chakrabarti, S. 2019. Girnet: Interleaved multi-task recurrent state sequence models. In *AAAI*, volume 33, 6497–6504.
- Haimmerl, C.; Savin, C.; and Simoncelli, E. 2019. Flexible information routing in neural populations through stochastic comodulation. In *NeurIPS*, 14379–14388.
- Harutyunyan, H.; Khachatrian, H.; Kale, D. C.; Ver Steeg, G.; and Galstyan, A. 2019. Multitask learning and benchmarking with clinical time series data. *Scientific Data* 6(1): 96.
- Hochreiter, S.; and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8): 1735–1780.
- Huang, D.-A.; Nair, S.; et al. 2019. Neural task graphs: Generalizing to unseen tasks from a single video demonstration. In *CVPR*, 8565–8574.
- Jang, E.; Gu, S.; and Poole, B. 2017. Categorical reparameterization with gumbel-softmax. In *ICLR*.
- Ke, N. R.; Zolna; et al. 2018. Focused hierarchical rnns for conditional sequence processing. In *ICML*.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kirsch, L.; Kunze, J.; and Barber, D. 2018. Modular networks: Learning to decompose neural computation. In *NeurIPS*, 2408–2418.
- Lake, B. M. 2019. Compositional generalization through meta sequence-to-sequence learning. In *NeurIPS*.
- Lake, B. M.; and Baroni, M. 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In *ICML*.
- Lin, R.; Liu, S.; Yang, M.; Li, M.; Zhou, M.; and Li, S. 2015. Hierarchical recurrent neural network for document modeling. In *EMNLP*, 899–907.
- Lin, Y.; Yang, S.; Stoyanov, V.; and Ji, H. 2018. A multilingual multi-task architecture for low-resource sequence labeling. In *ACL*, 799–809.



- Liu, H.; Simonyan, K.; and Yang, Y. 2019. Darts: Differentiable architecture search. In *ICLR*.
- Liu, P.; Fu, J.; Dong, Y.; Qiu, X.; and Cheung, J. C. K. 2019. Learning multi-task communication with message passing for sequence learning. In *AAAI*, volume 33, 4360–4367.
- Liu, P.; Qiu, X.; and Huang, X. 2017. Adversarial multi-task learning for text classification. In *ACL*.
- Loula, J.; Baroni, M.; and Lake, B. M. 2018. Rearranging the familiar: Testing compositional generalization in recurrent networks. In *EMNLP Workshop*.
- Lu, P.; Bai, T.; and Langlais, P. 2019. SC-LSTM: Learning task-specific representations in multi-task learning for sequence labeling. In *NAACL-HLT*, 2396–2406.
- Luong, M.-T.; Le, Q. V.; Sutskever, I.; Vinyals, O.; and Kaiser, L. 2016. Multi-task sequence to sequence learning. In *ICLR*.
- Ma, J.; Zhao, Z.; Chen, J.; Li, A.; Hong, L.; and Chi, E. H. 2019. SNR: Sub-Network Routing for Flexible Parameter Sharing in Multi-task Learning. In *AAAI*, 216–223.
- Ma, J.; Zhao, Z.; Yi, X.; Chen, J.; Hong, L.; and Chi, E. H. 2018. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *KDD*, 1930–1939.
- Maddison, C. J.; Mnih, A.; and Teh, Y. W. 2017. The concrete distribution: A continuous relaxation of discrete random variables. *ICLR*.
- Maziarz, K.; Kokiopoulou, E.; Gesmundo, A.; Sbaiz, L.; Bartok, G.; and Berent, J. 2019. Gumbel-Matrix Routing for Flexible Multi-task Learning. *arXiv:1910.04915*.
- Meyerson, E.; and Miikkulainen, R. 2018. Beyond shared hierarchies: deep multitask learning through soft layer ordering. In *ICLR*.
- Meyerson, E.; and Miikkulainen, R. 2019. Modular universal reparameterization: deep multi-task learning across diverse domains. In *NeurIPS*, 7901–7912.
- Misra, I.; Shrivastava, A.; Gupta, A.; and Hebert, M. 2016. Cross-stitch networks for multi-task learning. In *CVPR*, 3994–4003.
- Pahuja, V.; Fu, J.; Chandar, S.; and Pal, C. J. 2019. Structure learning for neural module networks. *arXiv preprint arXiv:1905.11532*.
- Parascandolo, G.; Kilbertus, N.; Rojas-Carulla, M.; and Schölkopf, B. 2018. Learning independent causal mechanisms. In *ICML*, 4036–4044.
- Pathak, D.; Lu, C.; Darrell, T.; Isola, P.; and Efros, A. A. 2019. Learning to control self-assembling morphologies: a study of generalization via modularity. In *NeurIPS*.
- Patra, B. G.; Das, D.; and Das, A. 2018. Sentiment Analysis of Code-Mixed Indian Languages: An Overview of SAIL-Code-Mixed Shared Task@ ICON-2017. *arXiv preprint arXiv:1803.06745*.
- Pham, H.; Guan, M. Y.; Zoph, B.; Le, Q. V.; and Dean, J. 2018. Efficient neural architecture search via parameter sharing. In *ICML*.
- Purushwalkam, S.; Nickel, M.; Gupta, A.; and Ranzato, M. 2019. Task-Driven Modular Networks for Zero-Shot Compositional Learning. In *ICCV*, 3593–3602.
- Ramachandran, P.; and Le, Q. V. 2019. Diversity and depth in per-example routing models. In *ICLR*.
- Real, E.; Aggarwal, A.; Huang, Y.; and Le, Q. V. 2019. Regularized evolution for image classifier architecture search. In *AAA*, volume 33, 4780–4789.
- Rosenbaum, C.; Cases, I.; Riemer, M.; Geiger, A.; Karttunen, L.; Greene, J. D.; Jurafsky, D.; and Potts, C. 2019a. Dispatched routing networks. *Technical Report Stanford AI Lab, NLP Group Tech Report 2019-1*.
- Rosenbaum, C.; Cases, I.; Riemer, M.; and Klinger, T. 2019b. Routing networks and the challenges of modular and compositional computation. *arXiv:1904.12774*.
- Rosenbaum, C.; Klinger, T.; and Riemer, M. 2018. Routing networks: Adaptive selection of non-linear functions for multi-task learning. In *ICLR*.
- Ruder, S. 2017. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*.
- Ruder, S.; Bingel, J.; Augenstein, I.; and Søgaard, A. 2019. Latent multi-task architecture learning. In *AAAI*, 4822–4829.
- Sachdeva, K.; Srivastava, R.; Jain, S.; and Sharma, D. M. 2014. Hindi to English Machine Translation: Using Effective Selection in Multi-Model SMT. In *LREC*, 1807–1811.
- Sang, E. F.; and Buchholz, S. 2000. Introduction to the CoNLL-2000 shared task: Chunking. In *CONLL Workshop*.
- Santos, C. D.; and Zadrozny, B. 2014. Learning character-level representations for part-of-speech tagging. In *ICML*, 1818–1826.
- Shickel, B.; Tighe, P. J.; Bihorac, A.; and Rashidi, P. 2017. Deep EHR: a survey of recent advances in deep learning techniques for electronic health record (EHR) analysis. *IEEE journal of biomedical and health informatics* 22(5): 1589–1604.
- Stickland, A. C.; and Murray, I. 2019. BERT and PALs: Projected attention layers for efficient adaptation in multi-task learning. In *ICML*.
- Subramanian, S.; Trischler, A.; Bengio, Y.; and Pal, C. J. 2018. Learning general purpose distributed sentence representations via large scale multi-task learning. In *ICLR*.
- Sun, C.; Shrivastava, A.; Vondrick, C.; Sukthankar, R.; Murphy, K.; and Schmid, C. 2019. Relational action forecasting. In *CVPR*, 273–283.
- Vyas, Y.; Gella, S.; Sharma, J.; Bali, K.; and Choudhury, M. 2014. Pos tagging of english-hindi code-mixed social media content. In *EMNLP*, 974–979.
- Wang, Y.; Gao, Z.; Long, M.; Wang, J.; and Yu, P. S. 2018. Predrnn++: Towards a resolution of the deep-in-time dilemma in spatiotemporal predictive learning. In *ICML*.
- Zhang, Y.; and Yang, Q. 2017. A survey on multi-task learning. *arXiv preprint arXiv:1707.08114*.

Zhou, Y.; Li, Z.; Xiao, S.; He, C.; Huang, Z.; and Li, H. 2018. Auto-conditioned recurrent networks for extended complex human motion synthesis. In *ICLR*.

Zoph, B.; and Le, Q. V. 2017. Neural architecture search with reinforcement learning. In *ICLR*.

# Appendices: Multi-Task Recurrent Modular Networks

## A. The Task-Evoked Functional MRI Data

**Data Description.** The task-evoked functional MRI data (task-fMRI) is used to analyze the relationship between brain connectivity and human behavior (Barch, Burgess et al. 2013). We use the public task-fMRI dataset<sup>3</sup> to conduct experiments. For this dataset, the participants performed seven tasks successively and the duration of each task is 160 time steps (each time step represents 0.72 seconds), which is shown in Figure 5. There are 460 participants in total and we extract 12 time series for each participant from the original dataset. Each time series corresponds to a high-order region of human brains. Figure 6 shows the 12 high-order regions (Gordon et al. 2016) and Figure 7 shows the extracted time series data of four participants.

**Pre-Processing Details.** To extract features and capture temporal patterns, we divided the time series data that corresponds to a task into 4 sub-time series of equal length. Thus, each sub-time series is with 40 time steps (each time step represents 0.72 seconds). We further applied the same PCA ( $n\_components = 10$ ) to all the sub-time series of all participants. As a result, each time series were processed to be 28 time steps (each time step represents  $0.72 \times 4$  seconds) long with 10 feature dimensions and every 4 time steps correspond to the time period when the participant conducted a task.

## B. Experiments on Clinical Time Series Data

**Experimental Setup.** We test our proposed MT-RMN networks on the multi-task setup of clinic prediction tasks. Health care is one of the frontiers in machine learning research and the widely used electronic health record (EHR) system generates the digital clinical data available for analysis (Shickel et al. 2017). We create the setup following (Harutyunyan et al. 2019), where we use the publicly available Medical Information Mart for Intensive Care (MIMIC-III) dataset<sup>4</sup> to create four clinical prediction tasks. Specifically, we use a subset of MIMIC-III dataset. It contains the multi-variate time series data (17 variables) from 42276 intensive care unit (ICU) stays of 33798 patients. The in-hospital mortality prediction task is a binary classification task, which aims to predict the in-hospital mortality based on the ICU stay of the first 48 hours. The decompensation prediction task is to predict if the health of a patient will deteriorate in the next 24 hours. The length-of-stay prediction task is a classification problem with 10 classes. Its goal is to predict the remaining time spent in ICU at each hour of a stay. The phenotype classification task is to classify which of 25 acute care conditions are present for an ICU stay. To evaluate MT-RMN, we use the two multi-task methods in (Harutyunyan et al. 2019), *i.e.*, multitask standard LSTM (MS) and

multitask channel-wise LSTM (MC), as the base models to incorporate MT-RMN. For both of the two methods, we use the instances from an ICU stay and predict all the targets for the ICU stay jointly. Note that different tasks share the same input data in this experiment, however, each task has its specific input data for the POS tagging task.

**Results.** We incorporate a MT-RMN without sparsity constraints ( $\gamma_i = 1$ ) and a MT-RMN with 75% sparsity constraints ( $\gamma_i = 0.75$ ) into these baselines. We show the results in Table 4. MAD and Kappa represent the mean absolute difference and the Cohen’s kappa score respectively. For all the metrics except MAD, larger values are better. It is observed that MT-RMN improved the performance of the two base models on all tasks, which demonstrates the effectiveness of MT-RMN. MT-RMN ( $\gamma_i=0.75$ ) outperformed MT-RMN ( $\gamma_i=1$ ) generally, which verifies the effectiveness of the sparsity constraints for the classification tasks on clinical time series data.

## C. Architecture Analysis

**Architectures Used for Different Datasets.** Table 5 describes the architectures of the encoder and decoders of MT-RMN on different datasets. Table 6 describes the architectures of policy networks for encoder and decoders on different datasets. The number represents the neuron number of each layer. For example, according to Table 5, the encoder consists of three layers for the task-fMRI data. The first layer consists of four LSTMs and the hidden dimension is 5. The second layer includes three MLPs, each MLP has two neuron layers and each neuron layer contains five neurons.

**Influence of Different Architectures.** We explored the influence of different architectures of MT-RMN. Table 7 shows the POS tagging results of GIRNet when different MT-RMNs were incorporated. For ease of comparison, we keep the architectures of decoders and policy networks as the same as the ones for the POS tagging data in Tables 5-6, while change the number of layers and the number of MLPs in each layer. We use the optimal hyperparameter setting (please see the section of Experiments) and set the MT-RMN with 75% sparsity constraints ( $\gamma_i = 0.75$ ) for all architectures. From Table 7, it is observed that the architecture of (L1:  $3 \times$  LSTM; L2:  $3 \times$  MLP; L3:  $3 \times$  MLP) achieved the best performance in general. One possible reason is that this architecture achieved the good trade-off between the routing complexity and the ease of optimization. As the number of layers or the number of MLPs in each layer increases, more hard decisions are need to make, which makes it more challenging to optimize the routing model.

**Analysis of The Number of Cells.** We also studied the influence of different number of recurrent cells, *i.e.*,  $k$ , in the first layer of encoder. We assume that the overall dynamical system of interest is divided into  $k$  small subsystems (Paras-

<sup>3</sup><https://tinyurl.com/y4hhw8ro>

<sup>4</sup><https://mimic.physionet.org/>

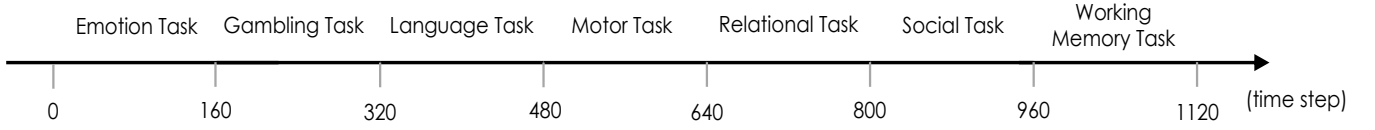
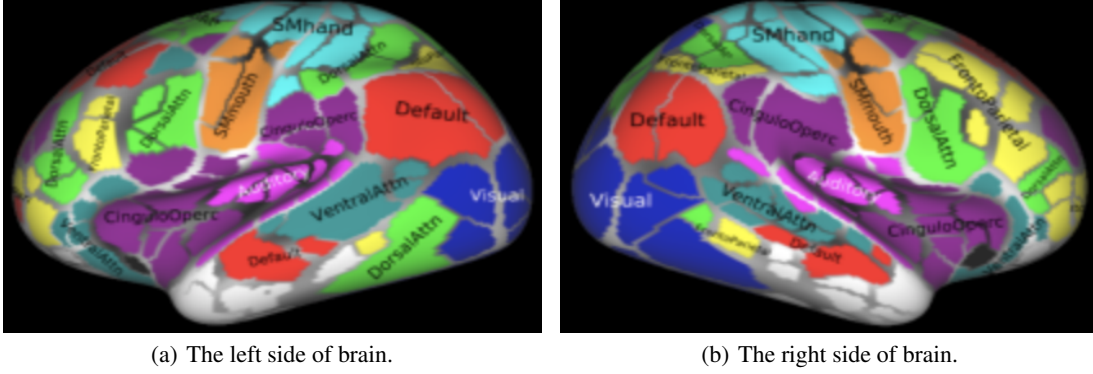


Figure 5: Seven tasks and each task lasts 160 time steps.



(a) The left side of brain.

(b) The right side of brain.

Figure 6: The 12 high-order regions of human brains, including *the default mode region, the parieto-occipital region, the frontoparietal region, the salience region, the cingulo-opercular region, the medial parietal region, the dorsal attention region, the ventral attention region, the visual region, the sensorimotor hand region, the sensorimotor mouth region and the auditory region.*

candolo et al. 2018; Goyal et al. 2019). The recurrent cells are used to model these subsystems and each cell has its own independent dynamics operating.  $k$  is a hyperparameter determined by grid-search. The best way to find it is to combine grid-search with domain knowledge. We observed that a too large  $k$  tended to attenuate the improvements to generalization and a too small  $k$  tended to hurt optimization.

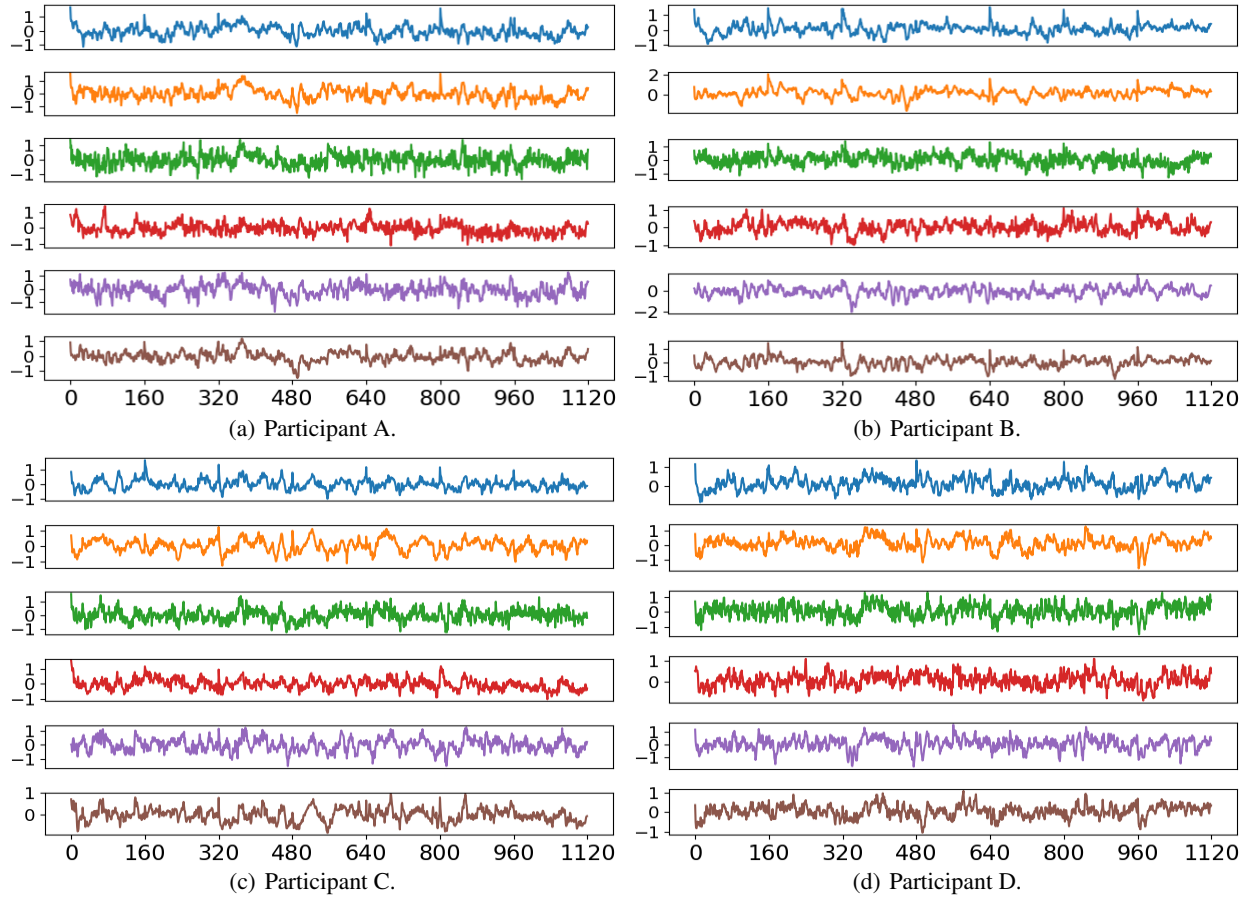


Figure 7: The extracted time series data of four participants. From top to down, the time series correspond to six high-order regions of human brains: *the default mode region, the parieto-occipital region, the frontoparietal region, the salience region, the cingulo-opercular region and the medial parietal region.*

Table 4: Results on four clinic prediction tasks.

Category	Methods	In-hospital mortality		Decompensation		Length-of-stay		Phenotype
		AUC-ROC	AUC-PR	AUC-ROC	AUC-PR	Kappa	MAD	AUC-ROC
Original	MS	85.82	55.69	88.33	31.51	39.36	93.78	75.39
	MC	86.93	58.01	89.57	32.64	36.40	98.76	77.87
MT-RMN ( $\gamma_i=1$ )	MS	87.01	<b>59.64</b>	90.60	33.04	43.32	<b>89.01</b>	77.51
	MC	87.90	56.49	<b>91.84</b>	32.92	39.35	89.47	78.70
MT-RMN ( $\gamma_i=0.75$ )	MS	<b>87.31</b>	56.07	<b>91.02</b>	<b>33.55</b>	<b>44.40</b>	92.85	<b>78.88</b>
	MC	<b>88.13</b>	<b>58.27</b>	91.57	<b>33.34</b>	<b>46.67</b>	<b>88.83</b>	<b>78.90</b>

Table 5: The architectures of encoder and decoders of MT-RMN on different datasets. ‘ $L_1$ ’ represents the first layer of encoder. ‘ $4 \times \text{LSTM} (h=5)$ ’ represent four LSTMs with the hidden dimension 5. ‘ $3 \times \text{MLP} [5 \times 5]$ ’ represent three MLPs with two layers that contain 5 neurons respectively.  $\mathcal{D}_1^2$  represents the first sub-decoder for the second task. MT-RMN adopts the activation relu for all architectures.

Dataset	Shared Encoder	Task Decoder
Task-fMRI	$L_1: 4 \times \text{LSTM} (h=5); L_2: 3 \times \text{MLP} [5 \times 5]; L_3: 3 \times \text{MLP} [5 \times 5]$	$\mathcal{D}_1^1: \text{MLP} [5 \times 5]; \mathcal{D}_2^1: \text{MLP} [5 \times 5]$
		$\mathcal{D}_1^2: \text{MLP} [5 \times 5]; \mathcal{D}_2^2: \text{MLP} [5 \times 5]$
		$\mathcal{D}_1^3: \text{MLP} [5 \times 5]; \mathcal{D}_2^3: \text{MLP} [5 \times 5]$
		$\mathcal{D}_1^4: \text{MLP} [5 \times 5]; \mathcal{D}_2^4: \text{MLP} [5 \times 5]$
POS Tagging	$L_1: 3 \times \text{LSTM} (h=50); L_2: 3 \times \text{MLP} [32 \times 32]; L_3: 3 \times \text{MLP} [32 \times 32]$	$\mathcal{D}_1^1: \text{MLP} [32 \times 32]; \mathcal{D}_2^1: \text{MLP} [32 \times 32]$
		$\mathcal{D}_1^2: \text{MLP} [32 \times 32]; \mathcal{D}_2^2: \text{MLP} [32 \times 32]$
		$\mathcal{D}_1^3: \text{MLP} [32 \times 32]; \mathcal{D}_2^3: \text{MLP} [32 \times 32]$
Clinic	$L_1: 4 \times \text{LSTM} (h=256); L_2: 3 \times \text{MLP} [128 \times 128]; L_3: 3 \times \text{MLP} [128 \times 128]$	$\mathcal{D}_1^1: \text{MLP} [128 \times 128]; \mathcal{D}_2^1: \text{MLP} [128 \times 128]$
		$\mathcal{D}_1^2: \text{MLP} [128 \times 128]; \mathcal{D}_2^2: \text{MLP} [128 \times 128]$
		$\mathcal{D}_1^3: \text{MLP} [128 \times 128]; \mathcal{D}_2^3: \text{MLP} [128 \times 128]$
		$\mathcal{D}_1^4: \text{MLP} [128 \times 128]; \mathcal{D}_2^4: \text{MLP} [128 \times 128]$

Table 6: The architectures of policy networks for encoder and decoders on different datasets. ‘ $\mathcal{N}$ ’ and ‘ $\hat{\mathcal{N}}^1$ ’ represent the policy network for encoder and the policy network for the decoder of the 1st task respectively. ‘ $\text{MLP} [5 \times 5 \times 2]$ ’ represent a MLP of three layers. MT-RMN adopts the activation relu for all architectures.

Dataset	Shared Encoder	Task Decoder
Task-fMRI	$\mathcal{N}: \text{MLP} [5 \times 5 \times 2]$	$\hat{\mathcal{N}}^1: \text{MLP} [5 \times 5 \times 2]$
		$\hat{\mathcal{N}}^2: \text{MLP} [5 \times 5 \times 2]$
		$\hat{\mathcal{N}}^3: \text{MLP} [5 \times 5 \times 2]$
		$\hat{\mathcal{N}}^4: \text{MLP} [5 \times 5 \times 2]$
POS Tagging	$\mathcal{N}: \text{MLP} [32 \times 32 \times 2]$	$\hat{\mathcal{N}}^1: \text{MLP} [32 \times 32 \times 2]$
		$\hat{\mathcal{N}}^2: \text{MLP} [32 \times 32 \times 2]$
		$\hat{\mathcal{N}}^3: \text{MLP} [32 \times 32 \times 2]$
Clinic	$\mathcal{N}: \text{MLP} [128 \times 128 \times 2]$	$\hat{\mathcal{N}}^1: \text{MLP} [128 \times 128 \times 2]$
		$\hat{\mathcal{N}}^2: \text{MLP} [128 \times 128 \times 2]$
		$\hat{\mathcal{N}}^3: \text{MLP} [128 \times 128 \times 2]$
		$\hat{\mathcal{N}}^4: \text{MLP} [128 \times 128 \times 2]$

Table 7: The POS tagging results of GIRNet with different MT-RMN. The hidden dimension of all LSTMs is 50. All MLPs are  $[32 \times 32]$  with the activation relu.

Different Architecture	Accuracy	F <sub>1</sub>	Precision	Recall
Original GIRNet	63.13	45.81	52.80	43.42
MT-RMN (L1: $3 \times \text{LSTM}$ ; L2: $2 \times \text{MLP}$ ; L3: $2 \times \text{MLP}$ )	64.65	62.04	61.84	61.16
MT-RMN (L1: $3 \times \text{LSTM}$ ; L2: $3 \times \text{MLP}$ ; L3: $3 \times \text{MLP}$ )	64.16	<b>62.62</b>	<b>62.58</b>	<b>62.22</b>
MT-RMN (L1: $3 \times \text{LSTM}$ ; L2: $4 \times \text{MLP}$ ; L3: $4 \times \text{MLP}$ )	63.67	60.29	61.43	61.94
MT-RMN (L1: $3 \times \text{LSTM}$ ; L2: $3 \times \text{MLP}$ ; L3: $3 \times \text{MLP}$ ; L4: $3 \times \text{MLP}$ )	<b>64.84</b>	62.48	62.04	62.12
MT-RMN (L1: $3 \times \text{LSTM}$ ; L2: $3 \times \text{MLP}$ ; L3: $3 \times \text{MLP}$ ; L4: $3 \times \text{MLP}$ ; L5: $3 \times \text{MLP}$ )	63.57	61.24	61.58	61.05
MT-RMN (L1: $2 \times \text{LSTM}$ ; L2: $3 \times \text{MLP}$ ; L3: $3 \times \text{MLP}$ )	63.43	61.19	61.02	61.16
MT-RMN (L1: $4 \times \text{LSTM}$ ; L2: $3 \times \text{MLP}$ ; L3: $3 \times \text{MLP}$ )	64.68	62.10	62.38	62.10