

## Network of Recurrent Neural Networks

**Authors:** Wang, Chao-Ming, Wang, Chao-Ming

**Date:** 2017-10-02T00:00:00+00:00

### Abstract

We describe a class of neural networks based on systems theory, termed “Network Of Recurrent neural networks” (NOR), which introduces a novel structural level to RNN-related models. In NOR, RNNs are regarded as high-level neurons and utilized to construct high-level layers. More specifically, we propose several methodologies for designing different NOR topologies based on the theory of system evolution. We then conduct experiments on three different tasks to evaluate our implementations. Experimental results demonstrate that our models remarkably outperform simple RNNs under the same number of parameters, and sometimes achieve even better results than GRU and LSTM.

### Full Text

#### Preamble

##### Network of Recurrent Neural Networks

Chao-Ming Wang

School of Software, Beijing Jiaotong University, Beijing, China

[oujago@gmail.com](mailto:oujago@gmail.com)

### Abstract

We describe a class of systems theory based neural networks called “Network Of Recurrent neural networks” (NOR), which introduces a new structural level to RNN-related models. In NOR, RNNs are viewed as high-level neurons and used to build high-level layers. More specifically, we propose several methodologies to design different NOR topologies according to the theory of system evolution. We then conduct experiments on three different tasks to evaluate our implementations. Experimental results show our models outperform simple RNN remarkably under the same number of parameters, and sometimes achieve even better results than GRU and LSTM.

## Introduction

In recent years, Recurrent Neural Networks (RNNs) (Elman 1990) have been widely used in Natural Language Processing (NLP). Traditionally, RNNs are directly used to build the final models. In this paper, we propose a novel idea called “Network Of Recurrent neural networks” (NOR) which utilizes existing basic RNN layers to make the structure design of the higher-level layers. From a standpoint of systems theory (Von Bertalanffy 1968; Von Bertalanffy 1972), a recurrent neural network is a group or an organization made up of a number of interacting parts, and it is actually viewed as a complex system, or a complexity. Dialectically, every system is relative: it is not only the system of its parts, but also the part of a larger system. In NOR structures, RNN is viewed as the high-level neuron, and several high-level neurons are used to build the high-level layers rather than directly used to construct the whole models.

Conventionally, there are three levels of structure in Deep Neural Networks (DNNs): neurons, layers and whole nets (or called models). From a perspective of systems theory, at each level of such increasing complexity, novel features that do not exist at lower levels emerge (Lehn 2002). For example, at the neurons level, a single neuron is simple and its generalization capability is very poor. But when a certain number of such neurons are accumulated into a certain elaborate structure by certain ingenious combinations, the layers at the higher level begin to get the unprecedented ability of classification and feature learning. More importantly, such newly gained capability or property is deducible from but not reducible to constituent neurons of lower levels. It’s not a property of the simple superposition of all constituent neurons, and the whole is greater than the sum of the parts. In systems theory, this kind of phenomenon is known as whole emergence (Wierzbicki 2015). Whole emergence often comes from the evolution of the system (Arthur and others 1993), in which a system develops from the lower level to the higher level, from simplicity to complexity. In this paper, the motivation of NOR structures is to introduce a new structure level to RNN-related networks by transferring traditional RNN from the system to the agent and from the outer dimension to the inner dimension (Fromm 2004).

In 1993, W. Brian Arthur (Arthur and others 1993) identified three mechanisms by which complexity tends to grow as systems evolve: (1) increase in co-evolutionary diversity, where the agent in the system seems to be a new instance of agent class, type or species, resulting in the system having new external agent types or capabilities; (2) increase in structural sophistication, where the individual system steadily accumulates increasing numbers of new systems or parts, thus newly formed system seems to have new internal subsystems or capabilities; and (3) increase by “capturing software”, where the system captures simpler elements and learns to “program” these as “software” to be used for its own ends.

In this paper, with the guidance of the first two mechanisms, we introduce two methodologies to NOR structures design, which are named as aggregation

and specialization. Aggregation and specialization are natural operations for increasing complexity in complex systems (Fromm 2004). The former is related to Arthur' s second mechanism, in which traditional RNNs are aggregated and accumulated into a high-level layer in accordance with a specific structure, and the latter is related to Arthur' s first mechanism, in which the RNN agent in a high-level layer is specialized as the RNN agent that performs a specific function.

We make several implementations and carry out experiments on three different tasks, including sentiment classification, question type classification and named entity recognition. Experimental results show that our models outperform constituent simple RNN remarkably with the same number of parameters and achieve even better results than GRU and LSTM sometimes.

## Background

### Systems Theory

Systems Theory was originally proposed by biologist Ludwig Von Bertalanffy (Von Bertalanffy 1968; Von Bertalanffy 1972) for biological phenomena. In biology systems, there are several different levels which begin with the smallest units of life and reach to the largest and most extensive category: molecule, cell, tissue, organ, organ system, organization, etc. Traditionally, a system could be decomposed into its individual components so that each component could be analyzed as an independent entity, and components could be added in a linear fashion to describe the totality of the system (Walonick 1993). However, Von Bertalanffy argued that we cannot fully comprehend a phenomenon by simply breaking it down into elementary parts and then reforming it. We instead need to apply a global and systematic perspective to underline its functionality (Mele, Pels, and Polese 2010), because a system is characterized by the interactions of its components and the nonlinearity of those interactions (Walonick 1993).

### Whole Emergence

In systems theory, the phenomenon (i.e., the whole is irreducible to its parts) is known as emergence or whole emergence (Wierzbicki 2015). Emergence can be qualitatively described as: the whole is greater than the sum of the parts (Upton, Janeka, and Ferraro 2014). Or it can also be quantitatively expressed as:  $\bigoplus$  where  $W$  is the whole of the system and consists of  $n$  parts, and  $\{p_i\}_{i=1..n}$  is the  $i$ -th part. In 1972, Philip W. Anderson highlighted the idea of emergence in his article "More is Different" (Anderson 1972) in which he stated that a change of scale very often causes a qualitative change in the behavior of the system. For example, in human brains, when one examines a single neuron, there is nothing that suggests consciousness. But a collection of millions of neurons is clearly able to produce wonderful consciousness.

The mechanisms behind the emergence of complexity can be used to design neural network structures. One of the widely accepted reasons is the repeated

application and combination of two complementary forces or operations: stretching and folding (in Physics term (Thompson and Stewart 2002)), splitting and merging (in Computer Science term (Hannebauer 2002)), or specialization and cooperation (in Sociology term). Merging or aggregating of agents means generally a number of (sub-)agents is aggregated or conglomerated into a single agent. Splitting or specializing means the agents are clearly separated from each other and each agent is constrained to a certain class or role (Fromm 2004).

## Recurrent Neural Networks at the Edge of Chaos

Recurrent Neural Networks (RNNs) (Werbos 1988; Elman 1990) are a class of deep neural networks that possess internal short-term memory due to recurrent feed-back connections between units, which makes them able to process arbitrary sequences of inputs. Formally, given a sequence of vectors  $\{x_t\}_{t=1\dots T}$ , the equation of Simple RNN (Elman 1990) is:  $h_t = f(Wx_t + Uh_{t-1})$  where  $W$  and  $U$  are parameter matrices, and  $f$  denotes a nonlinearity function such as tanh or ReLU. For simplicity, the neuron biases are omitted from the equation.

Actually, RNNs can behave chaotically. There have been some works analyzing RNNs theoretically or experimentally from the perspective of systems theory. (Sontag 1997) provided an exposition of research regarding system-theoretic aspects of RNNs with sigmoid activation functions. (Bertschinger and Natschläger 2004) analyzed the computation at the edge of chaos in RNNs and calculated the critical boundary in parameter space where the transition from ordered to chaotic dynamics takes place. (Pascanu, Mikolov, and Bengio 2013) employed a dynamical systems perspective to understand the exploding gradients and vanishing gradients problems in RNNs. In this paper, we obtain methodologies from systems theory to conduct structure designs of RNN related models.

## Network of Recurrent Neural Networks

### Overall Architecture

[Figure 1: see original paper] Overview of NOR structure.

As the high-level illustration shown in Figure 1, NOR architecture is a three-dimensional spatial-temporal structure. We summarize NOR architecture as four components: I, M, S and O, in which component I (input) and O (output) control the head and tail of NOR layer, component S (subnetworks) is in charge of the spatial extension and component M (memories) is responsible for the temporal extension of the whole structure. We describe each component as follows:

**Component I:** Component I controls the head of NOR architecture. It does data preprocessing tasks and distributes processed input data to subnetworks. At each time-step  $t$ , the form of upcoming input data  $x_t$  may be various, such as one single vector, or several vectors with the multi-granularity information, even the feature vectors with noise. One single vector may be the simplest situation,

and the common solution is copying this vector into  $n$  duplicates and feeding each of them into one single subnetwork in the component S. In this paper, the copying method meets our needs and we formalize it as:  $\{x_t^i\}_{i=1\dots n} = C(x_t)$  where  $C$  means copy function, and  $x_t^i$  will be fed into  $i$ -th subnetwork.

**Component M:** Component M manages all memories over the whole layer, not only internal but also external memories (Weston, Chopra, and Bordes 2014). But in this paper, component M only considers internal memory and does not apply any extra processing to the individual memory of each RNN neuron. That is:  $m_t^j = I(o_t^j)$  where  $I$  means identity function, the superscript  $j$  is the identifier of  $j$ -th RNN neuron,  $m_t^j$  is the memory of  $j$ -th RNN neuron at time-step  $t$  and  $o_{t-1}^j$  is the transformation output of  $j$ -th RNN neuron at time-step  $t - 1$ .

**Component S:** Component S is made up of several different or same subnetworks. Interaction may exist in these subnetworks. The responsibility of component S is to manage the logic of each subnetwork and handle the interaction between them. Suppose component S has  $n$  in-degree and  $m$  out-degree, i.e., component S receives  $n$  inputs and produces  $m$  outputs, the  $k$ -th output is generated by necessary inputs and memories:  $s_t^k = f([X, M])$  where  $s_t^k$  is the  $k$ -th output at time-step  $t$ ,  $X$  and  $M$  are needed inputs and memories, and  $f$  is the nonlinear function which can be one-layer RNN or two-layer RNN etc.

**Component O:** To form a layer we need a certain amount of neurons. So one of the NOR properties is multiple RNNs. A natural approach to integrate multiple RNN neurons' signals is collecting all outputs first and then using a MLP layer to measure the weights of each output. Traditional neuron outputs a single real value, so the collection method is directly arranging them into a vector. But RNN neurons are different, as each of them outputs a vector not a value. A simple method is concatenating all vectors and then connecting the concatenated vector to the next MLP. Another is pooling each RNN output vector into a real value, then arranging all these real values into a vector, which seems same as traditional neurons. In this paper, the former solution is used and formalized as:  $s_t = [s_t^1; s_t^2; \dots; s_t^m]$  and  $o_t = r(W_{MLP} * s_t)$  where  $s_t$  is the concatenated vector,  $W_{MLP}$  is the weight of MLP and  $r$  means the ReLU activation function of MLP.

### Methodology I: Aggregation

Any operation with changing a boundary can cause an emergence of complexity. The natural boundary is the agent itself, and sudden emergence of complexity is possible at this boundary if complexity is transferred from the agent to the system or vice versa from the system to the agent. There are two basic operations, aggregation and specialization, that can be used to transfer complexity between different dimensions (Fromm 2004).

According to Arthur' s second mechanism, internal complexity can be increased by aggregation and composition of sub-agents, which means a number of RNN agents is conglomerated into a single big system. In this way, aggregation and

composition transfer traditional RNN from the outer to the inner dimension, from the system to the agent, for the selected RNNs are accumulated to become a part of a larger group.

For a concrete NOR layer, suppose it is composed of  $n$  subnetworks, and  $i$ -th subnetwork is made up of  $k_i$  RNN neurons. Then, at the time-step  $t$ , given the input  $x_t$ , the operation flow is as follows:

1. **Component I:** Copy the input  $x_t$  into  $n$  duplicates and feed each of them into one subnetwork in component S using equation (3), then we get  $x_t^1, x_t^2, \dots, x_t^n$ .
2. **Component M:** Deliver the memory of each RNN neuron from the last time-step to the current time-step using equation (4), then we get memories  $m_{1,1}^t, \dots, m_{1,k_1}^t$  in first subnetwork and memories  $m_{2,1}^t, \dots, m_{2,k_2}^t$  in second subnetwork, etc.
3. **Component S:** For each subnetwork  $i$ , take advantage of the input  $x_t^i$  and memories  $m_{i,1}^t, \dots, m_{i,k_i}^t$  to get the nonlinear transformation output:  $s_t^i = f([x_t^i, m_{i,1}^t, \dots, m_{i,k_i}^t])$  then, we get  $s_t^1, s_t^2, \dots, s_t^n$ .
4. **Component O:** Concatenate all outputs by equation (6) and use a MLP function to determine how much signals in each subnetwork to flow through the component O by equation (7).

Obviously, the number, the type and the interaction of the aggregated RNNs determine the internal structure or inner complexity of the newly formed layer system. Thus, we propose three kinds of topologies of NOR aggregation method.

**Multi-Agent:** In systems theory, the natural description of complex system is the multi-agent system created by replication and adaptation. “Replication” means to copy and reproduce a new RNN agent, and “adaptation” means they are not totally same and some changes on weights or somewhere else by variation can increase the diversity of the system. As shown in Figure 2(a), there is a NOR layer (called MA-NOR) composed of four parallel RNNs. Figure 3 shows this layer being unrolled into a full network. Each subnetwork of MA-NOR layer is a one-tier RNN, thus at time-step  $t$ , the  $i$ -th subnetwork of component S in MA-NOR is calculated as:  $m_t^i = o_{t-1}^i$  and  $o_t^i = r(W_i x_t + U_i m_t^i)$  where  $r$  means ReLU activation function,  $W_i$  and  $U_i$  are parameters of corresponding RNN neuron,  $o_t^i$  is the nonlinear transformation output and will be delivered to next time-step to be used as  $m_{t+1}^i$ , and  $s_t^i$  is the output of  $i$ -th subnetwork which is equal to  $o_t^i$ .

The nonlinear function in equation (5) of each subnetwork may be more complex. For example, Figure 2(b) shows a NOR layer made up of three two-tier RNNs. At time-step  $t$ , the  $i$ -th subnetwork in component S is calculated as:  $m_t^{i,1} = o_{t-1}^{i,1}$ ,  $o_t^{i,1} = r(W^{i,1} x_t + U^{i,1} m_t^{i,1})$ ,  $m_t^{i,2} = o_{t-1}^{i,2}$ , and  $o_t^{i,2} = r(W^{i,2} o_t^{i,1} + U^{i,2} m_t^{i,2})$ .

**Multi-Scale:** The combination of multiple RNNs in a Multi-Agent NOR layer makes it somewhat like an ensemble. And empirically, diversity among the

members of a group of agents is deemed to be a key issue in ensemble structure (Kuncheva and Whitaker 2003). One way to increase the diversity is to use the Multi-Scale topology which introduces new agent type to the system and can learn sequence dependencies in different timescales. Figure 2(c) shows a NOR layer made up of four subnetworks, in which two of them are one-tier RNNs and the others are two-tier RNNs. Two kinds of timescale dependencies are learned in component S, which are formalized as follows:  $m_t^1 = o_{t-1}^1$ ,  $o_t^1 = r(W^1 x_t + U^1 m_t^1)$ ,  $s_t^1 = o_t^1$ ,  $m_t^2 = o_{t-1}^2$ ,  $o_t^2 = r(W^2 x_t + U^2 m_t^2)$ ,  $s_t^2 = o_t^2$ ,  $m_t^{3,1} = o_{t-1}^{3,1}$ ,  $o_t^{3,1} = r(W^{3,1} x_t + U^{3,1} m_t^{3,1})$ ,  $m_t^{3,2} = o_{t-1}^{3,2}$ ,  $o_t^{3,2} = r(W^{3,2} o_t^{3,1} + U^{3,2} m_t^{3,2})$ ,  $s_t^3 = o_t^{3,2}$ ,  $m_t^{4,1} = o_{t-1}^{4,1}$ ,  $o_t^{4,1} = r(W^{4,1} x_t + U^{4,1} m_t^{4,1})$ ,  $m_t^{4,2} = o_{t-1}^{4,2}$ , and  $o_t^{4,2} = r(W^{4,2} o_t^{4,1} + U^{4,2} m_t^{4,2})$ ,  $s_t^4 = o_t^{4,2}$ .

**Self-Similarity:** The above mentioned aggregation and composition operation lead to big RNN agent-groups. While in turn, they can also be combined to form even bigger groups. Such repeated aggregation and high accumulation makes the fractal and self-similar structure come into being. As shown in Figure 2(d), we also use three paths. But after each path first learns its own intermediate representation, the second layers gather all intermediate representations of three paths to learn high-level abstract features. In this way, different paths do not learn and train independently. The connections among each other help the model easily share information. Thus, it becomes possible that the whole model learns and trains to be an organic rather than parallel independent structure. We formalize the cooperation of component S as follows:  $m_t^{1,1} = o_{t-1}^{1,1}$ ,  $o_t^{1,1} = r(W^{1,1} x_t + U^{1,1} m_t^{1,1})$ ,  $m_t^{2,1} = o_{t-1}^{2,1}$ ,  $o_t^{2,1} = r(W^{2,1} x_t + U^{2,1} m_t^{2,1})$ ,  $m_t^{3,1} = o_{t-1}^{3,1}$ ,  $o_t^{3,1} = r(W^{3,1} x_t + U^{3,1} m_t^{3,1})$ ,  $m_t^{1,2} = o_{t-1}^{1,2}$ ,  $o_t^{1,2} = r(W^{1,2} [o_t^{1,1}, o_t^{2,1}, o_t^{3,1}] + U^{1,2} m_t^{1,2})$ ,  $s_t^1 = o_t^{1,2}$ ,  $m_t^{2,2} = o_{t-1}^{2,2}$ ,  $o_t^{2,2} = r(W^{2,2} [o_t^{1,1}, o_t^{2,1}, o_t^{3,1}] + U^{2,2} m_t^{2,2})$ ,  $s_t^2 = o_t^{2,2}$ ,  $m_t^{3,2} = o_{t-1}^{3,2}$ , and  $o_t^{3,2} = r(W^{3,2} [o_t^{1,1}, o_t^{2,1}, o_t^{3,1}] + U^{3,2} m_t^{3,2})$ ,  $s_t^3 = o_t^{3,2}$ .

## Methodology II: Specialization

We have mentioned that the emergence of complexity is usually connected to a transfer of complexity at the boundary of the system. Aggregation and composition transfer complexity from the system to the agent, and from the outer dimension to the inner dimension. Another way to cross the agent boundary is specialization or inheritance, which transfers complexity from the agent to the system, and from the inner dimension to the outer dimension (Fromm 2004). Specialization is related to Arthur's first mechanism. It increases structural sophistication outside of the agent by adding new agent forms. Through inheritance and specialization, objects become objects of a certain class and agents become agents of a certain type, and the more such an agent becomes a particular class or type, the more it needs to delegate special tasks that it cannot handle alone to other agents (Fromm 2004).

The effect of specialization is the emergence of delegation and division of labor in the newly formed groups. Thus, the formalization

of  $k$ -th output in Component S can be rewritten as the following:  $s_t^k = g(f_1([X, M]), f_2([X, M]), \dots, f_L([X, M]))$  where  $f_l$  is the  $l$ -th specialized agent function,  $g$  means the cooperation of all specialized agents, and  $L$  is the number of specialized agents. This equation denotes the function  $f$  in equation (5) is implemented by the separated operations  $f_1, f_2, \dots, f_L$  and  $g$ .

**Gate-Specialization:** We see gate mechanism as one of the specialization methods. As shown in Figure 6, a general RNN agent is separated into two specialized RNN agents, one for gate duty and the other for generalization duty. A concrete Gate-NOR is shown in Figure 7. In the original Multi-Agent NOR layer, each RNN agent is specialized as one generalization-specific RNN and one gate-specific RNN. We formalize it as:  $o_t^1 = \sigma(W^1x_t + U^1m_t^1)$ ,  $o_t^2 = r(W^2x_t + U^2m_t^2)$ ,  $s_t^1 = o_t^1 \odot o_t^2$ ,  $o_t^3 = \sigma(W^3x_t + U^3m_t^3)$ ,  $o_t^4 = r(W^4x_t + U^4m_t^4)$ , and  $s_t^2 = o_t^3 \odot o_t^4$ , where  $\sigma$  denotes the sigmoid activation and  $\odot$  denotes element-wise multiplication.

**Relationship with LSTM and GRU:** We see Long Short-term Memory (LSTM) (Hochreiter and Schmidhuber 1997) and Gated Recurrent Unit (GRU) (Chung et al. 2014) as two special cases of Network of Recurrent Neural Networks. Take LSTM for example, at time-step  $t$ , given input  $x_t$  and previous memory cell  $c_{t-1}$  and hidden state  $h_{t-1}$ , the transition equations of standard LSTM can be expressed as follows:  $i = \sigma(W_ix_t + U_ih_{t-1})$ ,  $f = \sigma(W_fx_t + U_fh_{t-1})$ ,  $o = \sigma(W_ox_t + U_oh_{t-1})$ ,  $g = \tanh(W_gx_t + U_g h_{t-1})$ ,  $c_t = c_{t-1} \odot f + g \odot i$ , and  $s_t = \tanh(c_t) \odot o$ . From the perspective of NOR, LSTM is made up of four RNNs, in which three of four (i.e.,  $i, f, o$  RNNs) are specialized for gate tasks to control how much information to let through in different parts. Moreover, there is only a shared memory  $h_{t-1}$  which can be accessed by each RNN cell in LSTM. While in turn, LSTM and GRU can also be combined to form even bigger groups.

## Experiments

In order to evaluate the performance of the presented model structures, we design experiments on the following tasks: sentiment classification, question classification, and named entity recognition. We compare all models under comparable parameter numbers to validate the capacity of better utilizing the parametric space. In order to verify the effectiveness and universality of the experiments, we conduct three comparative tests under total parameters of different orders of magnitude, see Table 1. Every experiment is repeated 20 times with different random initializations and then we report the mean results. It's worth noting that our aim here is to compare the model performance under the same hyper-parameter settings, not to achieve best performance for one single model.

(Le, Jaitly, and Hinton 2015) showed that when initializing the recurrent weight matrix to be the identity matrix and biases to be zero, simple RNN composed of ReLU activation function (named as IRNN) can be comparable with or even

outperform LSTM. In our experiments, all basic RNN neurons are simple RNNs applied with ReLU function. We also keep the number of the hidden units same over all RNN neurons in a NOR model. Obviously, our baseline model is a single giant simple RNN (Elman 1990) applied with ReLU activation. At the same time, two improved RNNs (GRU (Chung et al. 2014) and LSTM (Hochreiter and Schmidhuber 1997)) have been widely and successfully used in NLP in recent years, so we also choose them as our baseline models.

The pre-trained 300-D Glove 6B vectors and 300-D Google News vectors were obtained for the word embeddings. During training we fix all word embeddings and learn only the other parameters in all models. The embeddings for out-of-vocabulary words are set to zero vectors. We pad or crop the input sentences to a fixed length. The trainings are done through stochastic gradient descent over shuffled mini-batches with the optimizer Adam (Kingma and Ba 2014). All models are regularized by using dropout (Srivastava et al. 2014) method. At the same time, in order to avoid overfitting, early stopping is applied to prevent unnecessary computation when training. More details on hyper-parameters setting can be found in our codes, which are publicly available.

### Sentiment Classification

We evaluate our models on the task of sentiment classification on the popular Stanford Sentiment Treebank (SST) benchmark (Socher et al. 2013), which consists of 11855 movie reviews and is split into train (8544), dev (1101) and test (2210). SST provides detailed phrase-level annotation and all sentences along with the phrases are annotated with 5 labels: very positive, positive, neutral, negative, and very negative. In our experiments, we only use the sentence-level annotation. One of our goals is to avoid expensive phrase-level annotation, like (Qian, Huang, and Zhu 2016). Another is that, in practice, phrase-level annotation is hard to provide.

All models use the same architecture: embedding layer  $\rightarrow$  dropout layer  $\rightarrow$  RNN/NOR layer  $\rightarrow$  RNN/NOR layer  $\rightarrow$  max-pooling layer  $\rightarrow$  dropout layer  $\rightarrow$  softmax layer. The first layer is the word embedding layer, next are two-layer RNN/NOR layers as the non-linear feature transformation layer. Then a max-pooling layer max-pools all transformed feature vectors by selecting the max value in each position to get sentence representation. Finally, a softmax layer is used as output layer to get the final result. To benefit from the regularization, two dropout layers with rate of 0.5 are added after embedding layer and before softmax layer. The initial learning rates of all models are set to 0.0002. We use publicly available 300-D Glove 840B vectors to initialize word embeddings. Three different network sizes are tested for each architecture, such that the number of parameters are roughly 200k, 400k and 800k (see Table 1). We set the minibatch size as 20. Finally, we use the Cross-Entropy criterion as loss function.

Accuracy (%) comparison over different experiments on SST corpus.

The results of the experiments are shown in Table 2. It is obvious that NOR models get superior performances compared with IRNN baseline, especially when the network size is big enough. All models improve as network size grows. Among all NOR models, Gate-NOR gets the best results. However, we find that LSTM and GRU get much better results in three comparative tests.

### Question Type Classification

Question classification is an important step in a question answering system which classifies a question into a specific type. For this task, we use TREC (Li and Roth 2002) benchmark, which divides all questions into 6 categories: location, human, entity, abbreviation, description and numeric. TREC provides 5452 labeled questions in the training set and 500 questions in the test set. We randomly select 10% of the training data as the validation set.

All network types use the same architecture: embedding layer  $\rightarrow$  dropout layer  $\rightarrow$  RNN/NOR layer  $\rightarrow$  max-pooling layer  $\rightarrow$  dropout layer  $\rightarrow$  softmax layer. Dropout rates are set to 0.5. Three hidden layer sizes are chosen such that the total number of parameters for the whole model is roughly 100k, 200k, 400k, see Table 1. All networks use a learning rate of 0.0005 and are trained to minimize the Cross Entropy Error.

Accuracy (%) comparison over different experiments on TREC corpus.

Table 3 shows the accuracy of the different networks on the question type classification task. Here again, NOR models get better results than baseline IRNN model. Among all NOR models, SS-NOR also gets the best result. In this dataset, we find the performances of LSTM and GRU are even not comparable with IRNN, which proves the validity of results in (Le, Jaitly, and Hinton 2015).

### Named Entity Recognition

Named entity recognition (NER) is a classic NLP task which tries to identify the proper names of persons, organizations, locations, or other entities in the given text. We experiment on CoNLL-2003 dataset (Tjong Kim Sang and De Meulder 2003) which consists of 14987 sentences in the training set, 3466 sentences in the validation set and 3684 sentences in the test set.

Recently, popular NER models are based on bidirectional LSTM (Bi-LSTM) combined with conditional random fields (CRF), named as Bi-LSTM-CRF (Lample et al. 2016). The Bi-LSTM-CRF networks can effectively use past and future features via a Bi-LSTM layer and sentence level tag information via a CRF layer. In our experiments, we also adapt this architecture by replacing LSTM with NORs or other variation of RNNs. So the universal architecture of all tested models is: embedding layer  $\rightarrow$  dropout layer  $\rightarrow$  Bi-RNN/Bi-NOR layer  $\rightarrow$  CRF layer. Three hidden layer sizes are chosen such that the total number of parameters for the whole network is roughly 200k, 400k and 800k, see Table 1. We apply 50% dropout after embedding layer. Initial learning

rate is set to 0.005 and every epoch it is reduced by factor 0.95. The size of each minibatch is 20. We train all networks for 25 epochs and early stop the training when there are 5 epochs with no improvement on validation set.

F1 (%) comparison over different experiments on CoNLL-2003 corpus.

Our results are summarized in Table 4. Not surprisingly, all NORs perform much better than giant single RNN-ReLU model. As we can see, GRU performs the worst, followed by IRNN. Compared to GRU and IRNN, LSTM performs very well, especially when network size grows up. At the same time, all NOR models get superior performances than IRNN, GRU and LSTM. Among them, SS-NOR model gets best results.

## Conclusion

In conclusion, we introduced a novel kind of systems theory based neural networks called “Network Of Recurrent neural network”(NOR) which views existing RNNs (for example, simple RNN, GRU, LSTM) as high-level neurons and then utilizes RNN neurons to design higher-level layers. Then we proposed several methodologies to design different NOR topologies according to the evolution of systems theory (Arthur and others 1993). We conducted experiments on three kinds of tasks (including sentiment classification, question type classification and named entity recognition) to evaluate our proposed models. Experimental results demonstrated that NOR models get superior performances compared with single giant RNN models, and sometimes their performances even exceed GRU and LSTM.

## References

- [Anderson 1972] Anderson, P. W. 1972. More is different. *Science* 177(4047):393-396.
- [Arthur and others 1993] Arthur, W. B., et al. 1993. On the evolution of complexity. Technical report.
- [Bertschinger and Natschläger 2004] Bertschinger, N., and Natschläger, T. 2004. Real-time computation at the edge of chaos in recurrent neural networks. *Neural Computation* 16(7):1413-1436.
- [Chung et al. 2014] Chung, J.; Gulcehre, C.; Cho, K.; and Bengio, Y. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555.
- [Elman 1990] Elman, J. L. 1990. Finding structure in time. *Cognitive Science* 14(2):179-211.
- [Fromm 2004] Fromm, J. 2004. *The Emergence of Complexity*. Kassel University Press.

- [Hannebauer 2002] Hannebauer, M. 2002. *Autonomous Dynamic Reconfiguration in Multi-Agent Systems: Improving the Quality and Efficiency of Collaborative Problem Solving*. Springer-Verlag.
- [Hochreiter and Schmidhuber 1997] Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural Computation* 9(8):1735-1780.
- [Kingma and Ba 2014] Kingma, D., and Ba, J. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- [Kuncheva and Whitaker 2003] Kuncheva, L. I., and Whitaker, C. J. 2003. Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine Learning* 51(2):181-207.
- [Lample et al. 2016] Lample, G.; Ballesteros, M.; Subramanian, S.; Kawakami, K.; and Dyer, C. 2016. Neural architectures for named entity recognition. arXiv preprint arXiv:1603.01360.
- [Le, Jaitly, and Hinton 2015] Le, Q. V.; Jaitly, N.; and Hinton, G. E. 2015. A simple way to initialize recurrent networks of rectified linear units. arXiv preprint arXiv:1504.00941.
- [Lehn 2002] Lehn, J.-M. 2002. Toward complex matter: Supramolecular chemistry and self-organization. *Proceedings of the National Academy of Sciences* 99(8):4763-4768.
- [Li and Roth 2002] Li, X., and Roth, D. 2002. Learning question classifiers. In *Proceedings of the 19th International Conference on Computational Linguistics*, 1-7. Association for Computational Linguistics.
- [Mele, Pels, and Polese 2010] Mele, C.; Pels, J.; and Polese, F. 2010. A brief review of systems theories and their managerial applications. *Service Science* 2(1-2):126-135.
- [Pascanu, Mikolov, and Bengio 2013] Pascanu, R.; Mikolov, T.; and Bengio, Y. 2013. On the difficulty of training recurrent neural networks. In *ICML*.
- [Qian, Huang, and Zhu 2016] Qian, Q.; Huang, M.; and Zhu, X. 2016. Linguistically regularized LSTMs for sentiment classification. arXiv preprint arXiv:1611.03949.
- [Socher et al. 2013] Socher, R.; Perelygin, A.; Wu, J. Y.; Chuang, J.; Manning, C. D.; Ng, A. Y.; Potts, C.; et al. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1631-1642.
- [Sontag 1997] Sontag, E. 1997. Recurrent neural networks: Some systems-theoretic aspects. In *Dealing with Complexity: A Neural Network Approach*. Citeseer.
- [Srivastava et al. 2014] Srivastava, N.; Hinton, G. E.; Krizhevsky, A.; Sutskever, I.; and Salakhutdinov, R. 2014. Dropout: A simple way to prevent neural

networks from overfitting. *Journal of Machine Learning Research* 15(1):1929–1958.

[Thompson and Stewart 2002] Thompson, J. M. T., and Stewart, H. B. 2002. *Nonlinear Dynamics and Chaos*. John Wiley & Sons.

[Tjong Kim Sang and De Meulder 2003] Tjong Kim Sang, E. F., and De Meulder, F. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, 142–147. Association for Computational Linguistics.

[Upton, Janeka, and Ferraro 2014] Upton, J.; Janeka, I.; and Ferraro, N. 2014. The whole is more than the sum of its parts: Aristotle, metaphysical. *Journal of Craniofacial Surgery* 25(1):59–63.

[Von Bertalanffy 1968] Von Bertalanffy, L. 1968. General system theory. *New York* 41973(1968):40.

[Von Bertalanffy 1972] Von Bertalanffy, L. 1972. The history and status of general systems theory. *Academy of Management Journal* 15(4):407–426.

[Walonick 1993] Walonick, D. S. 1993. General information on systems theory. <http://www.statpac.org/systems-theory.htm>.

[Werbos 1988] Werbos, P. J. 1988. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks* 1:339–356.

[Weston, Chopra, and Bordes 2014] Weston, J.; Chopra, S.; and Bordes, A. 2014. Memory networks. arXiv preprint arXiv:1410.3916.

[Wierzbicki 2015] Wierzbicki, A. P. 2015. Systems theory, theory of chaos, emergence. In *Technen: Elements of Recent History of Information Technologies with Epistemological Conclusions*. Springer. 175–188.

*Note: Figure translations are in progress. See original paper for figures.*

*Source: ChinaXiv – Machine translation. Verify with original.*