

# Prototypical Recurrent Unit

Dingkun Long<sup>1</sup>, Richong Zhang<sup>1</sup>, Yongyi Mao<sup>2</sup>

<sup>1</sup>School of Computer Science and Engineering, Beihang University, Beijing, China

<sup>2</sup>School of Electrical Engineering and Computer Science, University of Ottawa  
 {longdk, zhangrc}@act.buaa.edu.cn, yymao@eecs.uottawa.ca

## Abstract

Despite the great successes of deep learning, the effectiveness of deep neural networks has not been understood at any theoretical depth. This work is motivated by the thrust of developing a deeper understanding of recurrent neural networks, particularly LSTM/GRU-like networks. As the highly complex structure of the recurrent unit in LSTM and GRU networks makes them difficult to analyze, our methodology in this research theme is to construct an alternative recurrent unit that is as simple as possible and yet also captures the key components of LSTM/GRU recurrent units. Such a unit can then be used for the study of recurrent networks and its structural simplicity may allow easier analysis. Towards that goal, we take a system-theoretic perspective to design a new recurrent unit, which we call the prototypical recurrent unit (PRU). Not only having minimal complexity, PRU is demonstrated experimentally to have comparable performance to GRU and LSTM unit. This establishes PRU networks as a prototype for future study of LSTM/GRU-like recurrent networks. This paper also studies the memorization abilities of LSTM, GRU and PRU networks, motivated by the folk belief that such networks possess long-term memory. For this purpose, we design a simple and controllable task, called “memorization problem”, where the networks are trained to memorize certain targeted information. We show that the memorization performance of all three networks depends on the amount of targeted information, the amount of “interfering” information, and the state space dimension of the recurrent unit. Experiments are also performed for another controllable task, the adding problem, and similar conclusions are obtained.

## Introduction

Deep learning has demonstrated great power in the recent years and appears to have prevailed in a broad spectrum of application domains (see, e.g., [12] [17]). Despite its great successes, the effectiveness of deep neural networks has not been understood at a theoretical depth. Thus developing novel analytic tools and theoretical frameworks for studying deep neural networks is of the greatest importance at the present time, and is anticipated to be a central subject of machine learning research in the years to come.

This work is motivated by the thrust of understanding recurrent neural networks, particularly LSTM/GRU-like networks [13] [8] [9] [4] [23]. These networks have demonstrated to be the state-of-the-art models for time series or sequence data [10] [1] [21]. Recently LSTM/GRU recurrent units have also been successfully adopted for modelling other forms of data (e.g., [3] [22]). Despite these successes, the design of LSTM and GRU recurrent units was in fact heuristical; to date there is little theoretical analysis justifying their effectiveness. A particularly interesting observation regarding these networks is that they appear to possess “long-term memory”, namely, being able to selectively “remember” the information from many time steps ago [7]. As one may naturally expect such memorization capability to have played an important role in the working of these networks, this aspect has not been well studied, analytically or experimentally.

The difficulty in analyzing recurrent networks resides in the complex structure of the recurrent unit, which induces highly complex nonlinear dynamics. To understand LSTM-like recurrent networks, the methodology explored in this theme of research is to maximally simplify the structure of the recurrent unit. That is, we wish to construct an alternative recurrent unit that captures the key components LSTM and GRU but stays as simple as possible. Such a unit can then be used for the study of recurrent networks and its structural simplicity may allow easier analysis in future research.

Towards that goal, the main objective of this present paper is to design such a recurrent unit and verify that this unit performs comparably to LSTM and GRU. To that end, we develop a new recurrent unit, which we call the *Prototypical Recurrent Unit* (PRU). We rationalize our design methodology from a system-theoretic perspective where a recurrent unit is understood as a causal time-invariant system in state-space representations. Insights from previous research suggest that additive evolution appear essential for LSTM-like networks to avoid the “gradient-vanishing” problem under back-propagation [14] [5] [18]. This understanding is also exploited in our design of PRU.

The performance of PRU is verified and compared against LSTM and GRU via extensive experiments. Using these three kinds of recurrent unit, we not only experiment on constructing a standard language model for character prediction [19], but also test the recurrent units for two controlled learning tasks, the Adding Problem [13], and the Memorization Problem. The latter problem is what we propose in this work specifically for studying the memorization capability of the recurrent networks. All experimental results confirm that PRU performs comparably to LSTM and GRU, achieving the purpose of this paper.

As another contribution, our experiments in this work demonstrate that the intrinsic memorization capability of the recurrent units depends critically on the dimension of the state space. The amount of targeted information (for memorization), the duration of memory, and the intensity of the interfering signal also directly impact the memorization performance.

Finally it is perhaps worth noting that although PRU is designed to be a prototype which hopefully allows for easier analysis in future research, our experiments suggest that it can also be used as a practical alternative to LSTM and GRU. A particular advantage of PRU is its time complexity. In this metric, PRU demonstrates to be superior to both LSTM and GRU.

## State-Space Representations

In system theory [15], a (discrete-time) system can be understood as any physical or conceptual device that responds to an *input sequence*  $x_1, x_2, \dots$  and generates an *output sequence*  $y_1, y_2, \dots$ , where the indices of the sequences are discrete time. In general, each  $x_t$  and each  $y_t$  at any time  $t$  may be a vector of arbitrary dimensions. We will then use  $\mathcal{X}$  and  $\mathcal{Y}$  to denote the vector spaces from which  $x_t$  and  $y_t$  take value respectively. We will call  $\mathcal{X}$  the *input space* and  $\mathcal{Y}$  the *output space*. The behaviour of the system is characterized by a function  $J$  that maps the space of all input sequences to the space of all output sequences. Then two systems  $J$  and  $J'$  are *equivalent* if  $J$  and  $J'$  are identical as functions.

The class of systems that are of primary interest are causal systems, namely those in which the output  $y_t$  at each time  $t$  is independent of all future inputs  $x_{t+1}, x_{t+2}, \dots$ . The grand idea in system theory is arguably the introduction of the notion of *state* to causal systems [15]. This makes state-space models the central topic in system theory, resulting in wide and profound impact on system analysis and design. In a nutshell, the state configuration is an quantity internal to the system, serving as a complete summary of the all past inputs so that *given the current state, the current and future outputs are independent of all past inputs*.

In this perspective, a recurrent unit can be regarded precisely as a causal time-invariant system in a state-space representation. We now formalize such a state-space representations.

At each time instant  $t$ , in addition to the input variable  $x_t$  and output variable  $y_t$ , the representation of a recurrent unit also contains a state variable  $s_t$ , taking values in a vector space  $\mathcal{S}$ , which will be referred to as the *state space*. Before the system is excited by the input, or at time  $t = 0$ , it is assumed that the state variable  $s_0$  takes certain initial configuration, which is assumed customarily to be the origin  $0 \in \mathcal{S}$ .

The behavior of the recurrent unit is governed by two functions  $F : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{S}$  and  $G : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{Y}$  as follows. At each time instant  $t$ , function  $F$  maps the current input  $x_t$  and the previous state  $s_{t-1}$  to the current state  $s_t$ , namely, via

$$s_t = F(x_t, s_{t-1}), \quad (1)$$

and function  $G$  maps the current input  $x_t$  and the current state  $s_t$  to the current output  $y_t$ , namely, via

$$y_t = G(x_t, s_t). \quad (2)$$

That is, in general a recurrent unit can be specified by the tuple  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)$  according to (1) and (2). We call such specification of the recurrent unit *Type-I state-space representation* of the unit, and denote it by  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_I$ .

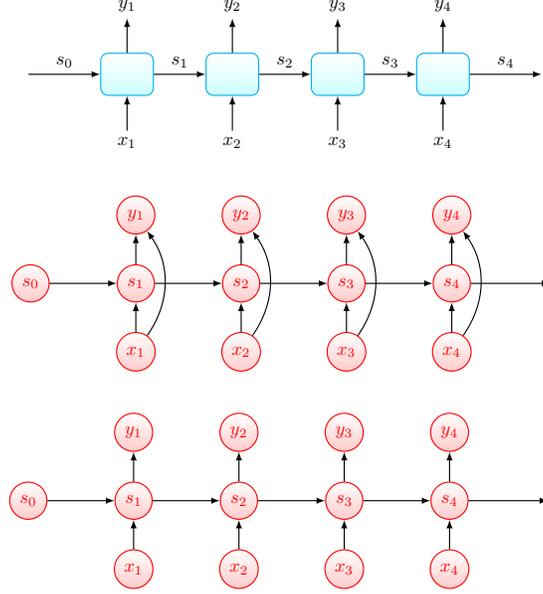


Figure 1: A recurrent network(top) and the dependency structure of variables in Type-I representation (middle) and Type-II representation (bottom).

It is remarkable that Type-I state-space representation is generic for any causal time-invariant system and hence generic for any recurrent unit. To illustrate this, we take the LSTM network as an example. The standard formulation of the LSTM network is given by the following equations:

$$i_t = \sigma(W_i[c_{t-1}, h_{t-1}, x_t] + b_i) \quad (3)$$

$$f_t = \sigma(W_f[c_{t-1}, h_{t-1}, x_t] + b_f) \quad (4)$$

$$o_t = \sigma(W_o[c_{t-1}, h_{t-1}, x_t] + b_o) \quad (5)$$

$$\tilde{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_g) \quad (6)$$

$$c_t = i_t \odot \tilde{c}_t + f_t \odot c_{t-1} \quad (7)$$

$$h_t = o_t \odot \tanh(c_t) \quad (8)$$

where  $\odot$  is the element-wise product. In these equations, if we take  $(c_t, h_t)$  as state  $s_t$ , and  $h_t$  as  $y_t$ , Equations(3-7) can be expressed as Equation (1), and Equations (5) and (8) can be expressed as Equation (2). We then arrive at a Type-I representation. It is also easy to verify that the recurrent unit in RNN [6] and GRU networks can all be expressed this way.

As a clarification which might be necessary for the remainder of this paper, we pause to remark that in this paper (and under a system-theoretic perspective), the notion of a recurrent unit and that of a recurrent (neural) network are synonyms. In particular, a recurrent unit that operates over  $n$  time instances may be viewed as  $n$  copies of the same recurrent unit connected in a chain-structured network as shown in Figure 1 (top). In this “time-unfolded” view, the dependency structure between the variables in Type-I representation is shown in Figure 1 (middle).

Since we aim at designing a *simpler* recurrent unit, we now introduce another simpler representation, which we call *Type-II state-space representation*. This representation is identical to the Type-I representation except that the function  $G$  is made to have domain  $\mathcal{S}$ , or alternatively put, the current output  $y_t$  at each time  $t$  is made dependent only of the current state  $s_t$ . That is,  $G$  acts only on  $s_t$  and generates  $y_t$  via

$$y_t = G(s_t). \quad (9)$$

Under this representation, the recurrent unit is specified again by the tuple  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)$ , but according to (1) and (9). We denote this representation by  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_{\text{II}}$ . A diagram exhibiting the dependency structure of the variables in this representation is shown in Figure 1 (bottom).

The following lemma suggests that Type-II representation has precisely the same expressive power as Type-I representation.

**Lemma 1** *Given its input and output spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , a recurrent unit can be represented by  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_{\text{I}}$  for some choice of  $\mathcal{S}$ ,  $F$  and  $G$  if and only if it can be represented by  $(\mathcal{X}, \mathcal{Y}, \tilde{\mathcal{S}}, f, g)_{\text{II}}$  for some choice of  $\tilde{\mathcal{S}}$ ,  $f$  and  $g$ .*

As the proof of this lemma contains certain insights into state-space representations, we sketch it here.

The “if” part of the proof is trivial, since function  $G$  in Equation (9) is a special case of function  $G$  in Equation (2). The “only if” part can be proved by construction, proceeded as follows. Let  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_{\text{I}}$  be given. Define  $\tilde{\mathcal{S}} := \mathcal{X} \times \mathcal{S}$ . Let function  $f : \mathcal{X} \times \tilde{\mathcal{S}} \rightarrow \tilde{\mathcal{S}}$  be defined as follows: for each  $(x, x', s) \in \mathcal{X} \times \mathcal{X} \times \mathcal{S} = \mathcal{X} \times \tilde{\mathcal{S}}$ ,  $f(x, x', s) = (x^*, s^*) \in \mathcal{X} \times \mathcal{S} = \tilde{\mathcal{S}}$ , where  $x^* = 0 \in \mathcal{X}$  and  $s^* = F(x, s) \in \mathcal{S}$ . Define function  $g : \tilde{\mathcal{S}} \rightarrow \mathcal{Y}$  as follows: for each  $(x, s) \in \mathcal{X} \times \mathcal{S} = \tilde{\mathcal{S}}$ ,  $g(x, s) = G(x, s)$ . Now the lemma can be proved by identifying that systems  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_{\text{I}}$  and  $(\mathcal{X}, \mathcal{Y}, \tilde{\mathcal{S}}, f, g)_{\text{II}}$  are equivalent. This latter fact can be easily established using proof by induction.

The significance of this lemma is that every recurrent unit can be represented using Type-II representation, in which the current output is made only dependent of the current state. In the proof of this result, we see that to convert a Type-I representation to a Type-II representation, it may require *increasing the dimension of the state space*. In the worst case, although often unnecessary in practice, one can make the state space  $\tilde{\mathcal{S}}$  equal to the cartesian product  $\mathcal{X} \times \mathcal{S}$  of the input space  $\mathcal{X}$  and the state space  $\mathcal{S}$  in the Type-I representation.

## Prototypical Recurrent Unit (PRU)

Given that there is no loss of expressive power in Type-II representation, to arrive at a simplified recurrent unit, we will stay within this representation. That is, for some given choices of vector spaces  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathcal{S}$ , we will design two functions  $F : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{S}$  and  $G : \mathcal{S} \rightarrow \mathcal{S}$  for  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_{\text{II}}$ . It is our hope that the designed recurrent unit captures the essence of recurrent unit in LSTM and GRU networks, but stays as simple as possible.

From the previous literature [20], the following properties of LSTM and GRU appear crucial for their effectiveness.

1. The recurrent unit behaves according to a nonlinear system, where the nonlinearity is induced by the use of nonlinear activation functions such as sigmoid, tanh, or ReLU functions.
2. The evolution from state  $s_t$  to state  $s_{t+1}$  is additive. It has been understood that such a property is critical for eliminating the problem of vanishing or blowing-up gradient in backpropagation.

Based on this understanding, our design philosophy is to impose these two properties *minimally* on the recurrent unit. Our hypothesis is that if these two properties are indeed essential, the resulting recurrent unit will behave in a way similar to GRU and LSTM recurrent units and can be used as a prototypical example for in-depth understanding of LSTM/GRU-like recurrent networks.

Such a design philosophy naturally results in the following new recurrent unit, which we call the Prototypical Recurrent Unit (PRU) and now describe.

We begin with some notations. We consider  $\mathcal{X} = \mathbb{R}^m$ ,  $\mathcal{Y} = \mathbb{R}^l$  and  $\mathcal{S} = \mathbb{R}^k$ ; all vectors are taken as column vectors; the sigmoid (logistic or soft-max) function will be denoted by  $\sigma$ ; when an activation function ( $\sigma$ , tanh, or any other function  $h : \mathbb{R} \rightarrow \mathbb{R}$ ) applies to a vector, it acts on the vector element-wise and outputs a vector of the same length.

With these notations, we describe the functions  $F$  and  $G$  in  $(\mathcal{X}, \mathcal{Y}, \mathcal{S}, F, G)_{\text{II}}$  that defines PRU.

**Function  $F$ :** The function  $F$  is defined by the following sequence of function compositions involving two other variables  $u_t \in \mathcal{S}$  and  $c_t \in \mathbb{R}^k$  (we note that although here  $c_t$  is a  $k$ -dimensional vector, it should not be interpreted as a state configuration in  $\mathcal{S}$  due to its physical meaning).

$$u_t = \tanh(U_s s_{t-1} + U_x x_t + b_u) \quad (10)$$

where  $U_s$  is a  $k \times k$  matrix,  $U_x$  is a  $k \times m$  matrix, and  $b_u$  is a  $k$ -dimensional vector.

$$c_t = \sigma(C_s^T s_{t-1} + C_x^T x_t + b_c) \quad (11)$$

where  $C_s$  is a  $k \times k$  matrix,  $C_x$  is a  $k \times m$  matrix, and  $b_c$  is a  $k$ -dimensional vector.

$$s_t = c_t \odot s_{t-1} + (1 - c_t) \odot u_t \quad (12)$$

where  $\odot$  is the element-wise product.

**Function  $G$ :** The function  $G$  is defined as follows.

$$y_t = h(Ws_t + b) \quad (13)$$

where  $W$  is an  $l \times k$  matrix,  $b$  is a length  $l$  vector, and  $h$  is an activation function. Depending on the applications and the physical meaning of output  $y_t$ ,  $h$  can be chosen as  $\sigma$ ,  $\tanh$ ,  $\text{ReLU}$ , or even the identity function.

At this point, we have completely defined PRU, which is parameterized by  $\theta := (C_x, C_s, b_c, U_x, U_s, b_u, W, b)$ .

## Experimental Study

Our experimental study serves two purposes.

First, we wish to verify that the designed PRU behaves similarly as LSTM and GRU. For this purpose, experiments need to be performed not only for real-world applications, in which one has no control over the datasets, but also for certain meaningful tasks where we have full control over the data. Such controllable tasks will allow a comparison of these recurrent units over arbitrary ranges of data parameter settings, so as to fully demonstrate the performances of the compared recurrent units and reduce the risk of being biased by the statistics of a particular dataset.

Second, we wish to take the opportunity to investigate a fundamental aspect of recurrent networks, namely, their memorization capabilities. It has been experimentally observed and intuitively justified that LSTM/GRU-like recurrent unit has “long-term memory” [11]. Motivate by such observations, we are interested in thoroughly studying the memorization capability of these recurrent units and understand what factors may influence their memorization performance.

As such, we consider four different learning tasks, where the recurrent networks are trained to solve four different problems: the Memorization Problem, the Adding Problem, the Character Prediction Problem, and the MNIST Image Classification Problem. The Character Prediction Problem and the MNIST Image Classification Problem are two well-known problems in the real-world application domain [19] [16]. The Adding Problem is a controllable task, first introduced in [13]. The Memorization Problem is also a controllable task that we introduce in this work, inspired by the idea of a similar task presented in [2].

All models in these experiments have the architecture shown in the top diagram of Figure 1 with single layer. In Memorization problem and Adding problem, we use single layer structure, in the other two experiments we will use a two-layer stacked structure. In the description of the experiments, when we speak of “state space dimension”, for both PRU and GRU, it refers to the length of the vector passed between two consecutive recurrent units in the diagram. In LSTM networks, there are two vectors of the same length passed between two consecutive recurrent units. Although from a system-theoretic perspective, two times this length should be regarded as the state space dimension, this choice would put LSTM in disadvantage. This is because the output of the unit depends only on one of the vectors. For this reason, for LSTM networks, the term “state space dimension” refers to half of the true state-space dimension.

Experiments on Memorization Problem and Adding Problem are performed on the computer(Intel(R) Core(TM) i5-4570 CPU @3.20Hz), whereas experiment on Character Prediction Problem and MNIST Image Classification Problem are performed on a GeForce GTX 970 GPU. Time cost is evaluated in unit of second.

### Memorization Problem

To describe this problem, let us first imagine a “memorization machine”  $\mathcal{M}_{\text{mem}}$  that behaves as follows. For any given non-negative integers  $I$  and  $N$ , an input sequence  $x_1, x_2, \dots, x_{I+N}$  of scalar values are fed to the machine, where

for  $t = 1, 2, \dots, I$ ,  $x_t$  takes on value in  $\{+1, -1\}$  each with probability  $1/2$ , and for  $t = I + 1, I + 2, \dots, I + N$ ,  $x_t$  is drawn independently from a Gaussian distribution with zero mean and variance  $\delta^2$ . After processing the input sequence, the machine generates an output vector  $(x_1, x_2, \dots, x_I)^T$  of dimension  $I$ . That is, as a function, the machine  $\mathcal{M}_{\text{mem}}$  behaves according to

$$\mathcal{M}_{\text{mem}}(x_1, x_2, \dots, x_{I+N}) = (x_1, x_2, \dots, x_I)^T.$$

Then in the Memorization Problem, the objective is to train a model that simulates the behaviour of  $\mathcal{M}_{\text{mem}}$ , namely, capable of “memorizing” the “ $I$  bit” “targeted information” in the beginning of the input sequence, after  $N$  symbols of “noise” or “interfering signal” enter the model. Obviously, the Memorization Problem is configured by three parameters:  $I$ ,  $N$ , and  $\delta$ , where  $I$  represents the amount of targeted information,  $N$  represents the duration of memory, and  $\delta$  represents the intensity of noise that might interfere with the memorization behaviour of the model.

**Modelling:** Under a recurrent network model, it is natural to regard the input space  $\mathcal{X}$  as  $\mathbb{R}$  and the output space  $\mathcal{Y}$  as  $\mathbb{R}^I$ , and one may freely configure the dimension  $k$  of the state space  $\mathcal{S}$ . Except at the final time  $t = I + N$ , the output  $y_t$  is discarded, and final output  $y_{I+N}$  is used to simulate the output  $\mathcal{M}_{\text{mem}}(x_1, x_2, \dots, x_{I+N})$  of the memorization machine.

**Datasets:** For each problem setting  $(I, N, \delta^2)$ , we generate 50000 training examples and 1000 testing examples according to the specification of the problem.

**Training:** The training of each model is performed by optimizing the Mean Square Error (MSE) defined as

$$\mathcal{E}_{\text{MSE}}(\theta) := \mathbb{E} \|\mathcal{M}_{\text{mem}}(x_1, x_2, \dots, x_{I+N}) - y_{I+N}\|^2 \quad (14)$$

where the expectation operation  $\mathbb{E}$  is taken as averaging over the training examples. Mini-batched Stochastic Gradient Descent (SGD, in fact more precisely, mini-batched Back-Propagation Through Time) is used for this optimization. The batch size is chosen as 100, the learning rate as  $10^{-3}$ , and the number of epochs as 1000. Each component of the model parameters is initialized to random values drawn independently from the zero-mean unit-variance Gaussian distribution.

**Evaluation Metrics:** A trained model is evaluated using MSE defined in (14), where the expectation operation  $\mathbb{E}$  is taken as averaging over the testing examples. For experiment setting  $(I, N, \delta^2, k)$ , each studied model is trained 50 times with different random initializations, and the average MSE is taken the performance metric for the experiment setting. Time complexity for the three models are also evaluated.

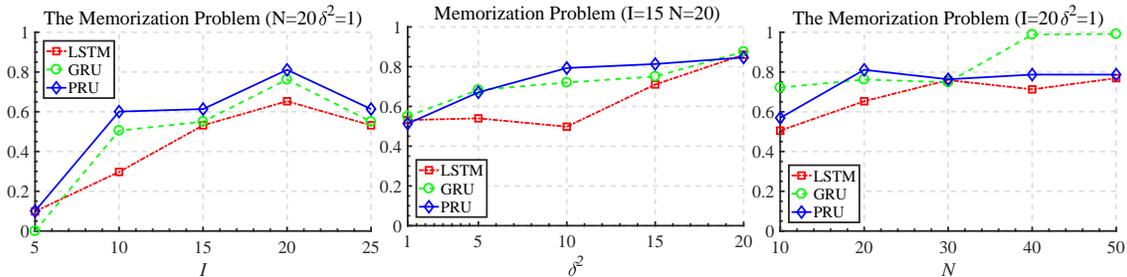


Figure 2: MSE comparison of LSTM, GRU and PRU networks in Memorization Problem with same state space dimension

**Results:** Results are obtained for LSTM, GRU and PRU under various problems settings  $(I, N, \delta^2)$  and model state-space dimensions  $k$ .

Figure 2 shows the performance comparison of the three recurrent units with same state space dimension. It can be seen that the three units perform similarly, among which LSTM’s performance is superior to the other two with same state space dimension. However, in this comparison, LSTM uses the most parameters. Figure 3 demonstrates the performances under the same number of parameters. It can be observed that PRU outperforms GRU, even catches up the performance of LSTM to a certain extent. In addition, with respect to any given parameter, the performance trends of the three units are identical.

Figure 4 shows how the performance of each unit is related to the problem parameters  $I$ ,  $N$ , and  $\delta^2$ . For every unit and a fixed state space dimension  $k$ , the following performance trend can be observed.

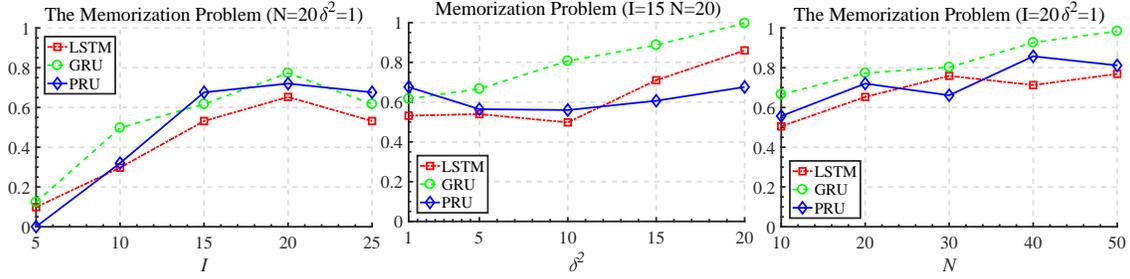


Figure 3: MSE comparison of LSTM, GRU and PRU networks in Memorization Problem with same number of parameters.

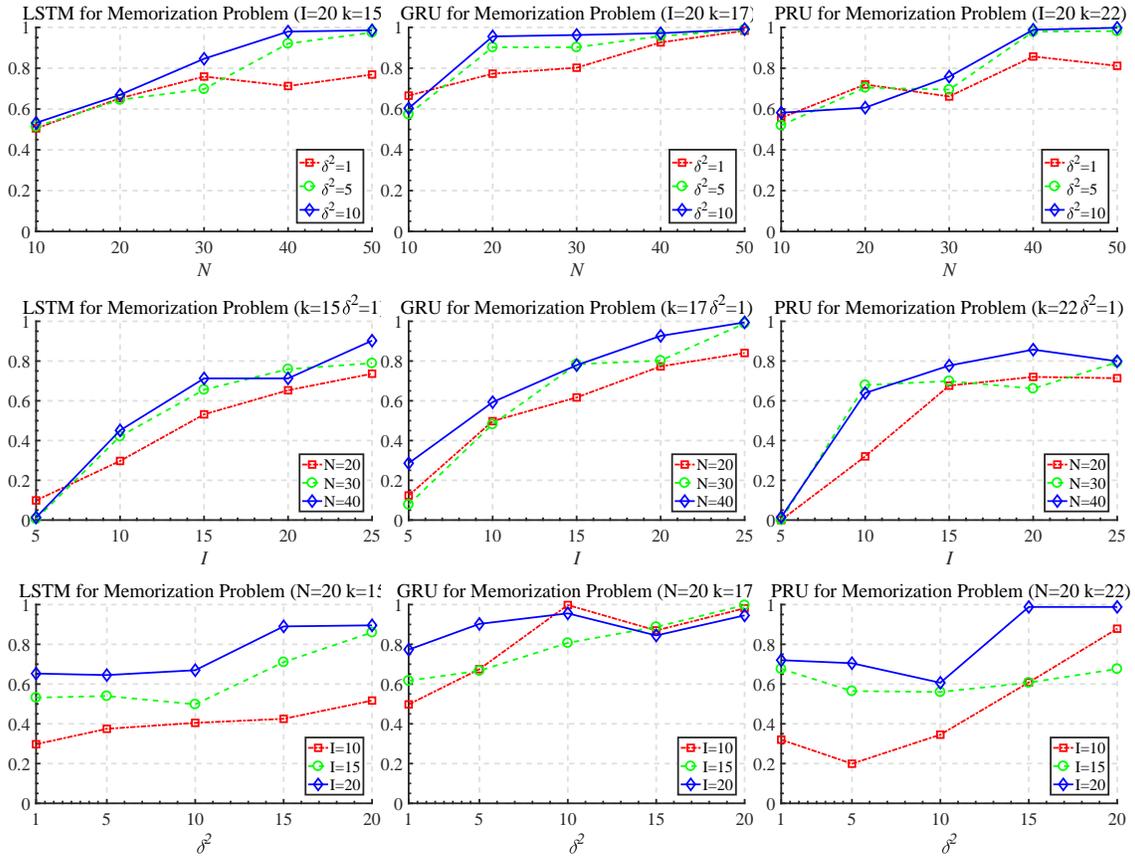


Figure 4: MSE comparison of LSTM, GRU and PRU networks with varying parameters in Memorization Problem

- The performance degrades with increasing  $I$ . That is, when the amount of targeted information increases, it becomes more difficult for the unit to memorize this information.
- The performance degrades with increasing  $N$ . That is, over a long period of time, the units tend to forget the targeted information.
- The performance degrades with increasing  $\delta^2$ . That is, when the interfering signal become stronger, it is more difficult to memorize the targeted information.

Figure 5 shows how the performance of each unit varies with the state space dimension  $k$ . It is apparent from the

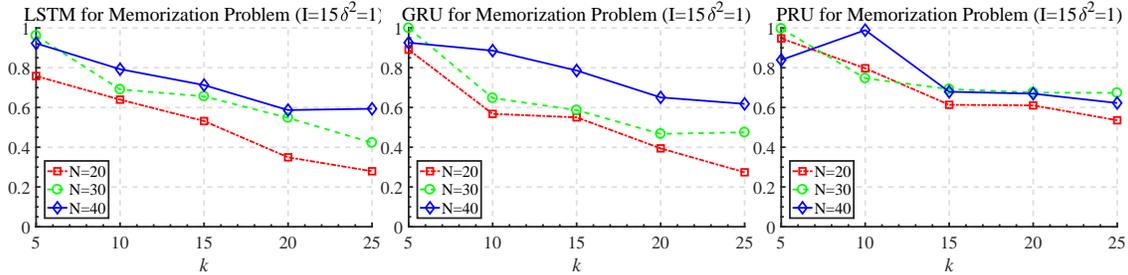


Figure 5: MSE comparison of LSTM, GRU and PRU networks with varying state-space dimension  $k$ .

figure that as the dimension of state space increases, the performance of each unit improves. This behaviour is sensible, since the role of the state variable in a recurrent unit may be intuitively understood as the “container” for “storing” information, and large state space would result in larger “storage capacity”.

From Table 1 (measured at  $k = 15$  and  $I = 15$ ), it can be observed that PRU has lowest time complexity, significantly below GRU and LSTM. This is a direct consequence of PRU’s structural simplicity.

Table 1: Time cost per epoch in Memorization Problem,  $k = 3$ ,  $I = 2$

$N$	20	30	40	50	60
LSTM	40.445	60.575	80.653	102.435	125.245
GRU	30.355	45.863	60.524	76.328	93.558
PRU	20.472	31.791	43.434	55.131	70.244

## Adding Problem

To describe the Adding Problem, let  $\mathcal{M}_{\text{add}}$  be an “adding machine”, which is a function mapping a length- $N$  sequence  $(x_1, x_2, \dots, x_N)$  to a real number. In particular, each  $x_t$ ,  $t = 1, 2, \dots, N$ , is a vector in  $\mathbb{R}^2$ , and we may write  $x_t$  as  $(x_t[1], x_t[2])^T$ . At each  $t$ ,  $x_t[1]$  is a random value drawn independently from the zero-mean Gaussian distribution with variance  $\delta^2$ ; and in the sequence  $(x_1[2], x_2[2], \dots, x_N[2])$ , there are exactly two 1’s, the locations of which are randomly assigned; the remaining values of the sequence all are equal to 0. The behaviour of the adding machine is given by

$$\mathcal{M}_{\text{add}}(x_1, x_2, \dots, x_N) := \sum_{t=1}^N x_t[1] \cdot x_t[2].$$

The objective of the Adding Problem is then to train a model that simulates the behaviour of  $\mathcal{M}_{\text{add}}$ . Obviously, the Adding Problem is parametrized by the pair  $(N, \delta^2)$ . Intuitively, the Adding Problem demands higher “memorization capacity” than the Memorization Problem, since only counting the locations of the two 1’s in the second component the input sequence, there are  $\binom{N}{2}$  possibilities.

**Modelling:** Under a recurrent network model, it is natural to take input space  $\mathcal{X} = \mathbb{R}^2$  and output space  $\mathcal{Y} = \mathbb{R}$ . Except at the final time  $t = N$ , the output  $y_t$  is discarded, and final output  $y_N$  is used to simulate the output  $\mathcal{M}_{\text{add}}(x_1, x_2, \dots, x_N)$  of the adding machine.

**Datasets:** For each problem setting  $(N, \delta^2)$ , we generate 2000 training examples and 400 testing examples according to the specification of the problem.

**Training:** The training of each model is performed by optimizing the MSE between the  $y_N$  and  $\mathcal{M}_{\text{add}}(x_1, x_2, \dots, x_N)$ . A mini-batched SGD method is used for optimization, where we use the same set of training parameters as those in the Memorization Problem, except that the batch size is chosen as 50.

**Evaluation Metrics:** MSE is used as the evaluation metric, and the same averaging process as that for the Memorization Problem is applied.

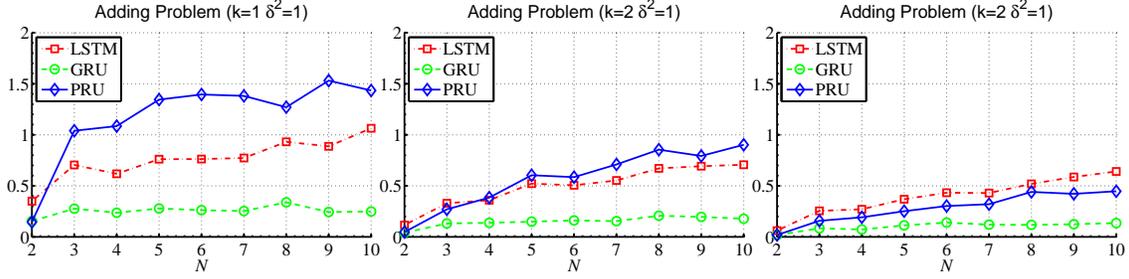


Figure 6: MSE comparison of LSTM, GRU and PRU in Adding Problem.

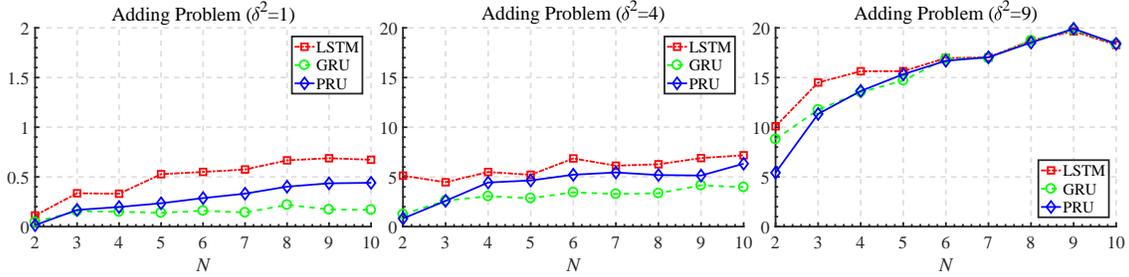


Figure 7: MSE comparison of LSTM, GRU and PRU in Adding Problem with same number of parameters.

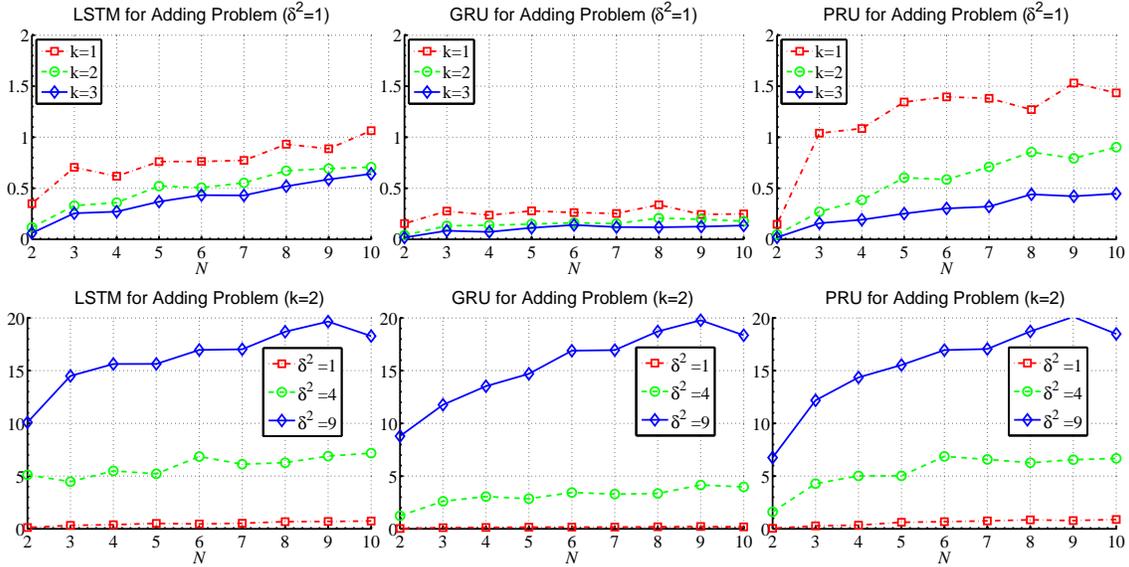


Figure 8: MSE comparison of LSTM, GRU, and PRU in Adding Problem under varying experiment settings. The MSE value has been normalized over  $\delta^2$ .

**Results:** Figure 6 shows the performance comparison of LSTM, GRU and PRU in the Adding Problem with the same state space dimension, Figure 7 shows the performance comparison of the three recurrent units with the same number of parameters and Figure 8 shows the performance trend of each of the three units with respect varying parameters. In Adding Problem, overall the three units perform comparably, GRU superior to the other two units. It is worth noting in Figure 6, with low state-space dimension ( $k = 1$ ), PRU appears under-perform LSTM. But as the state-space

dimension increases, PRU catches up (at  $k = 2$ ) and even out-performs LSTM (at  $k = 3$ ). This may be explained as follows. First the Adding Problem demands higher “memorization capacity”. But as we discussed earlier, PRU uses the Type-II representations, which may need larger state-space for the same representation power. In Figure 7, it can be seen that the performance difference between PRU and GRU is narrowing with same number of parameters, additionally, under certain parameter settings, the performance between PRU and GRU is almost with no difference even the same.

These results also suggest that the three studied units all have identical performance trends with respect to state-space dimension or any given problem parameter. Conclusions similar to those in the Memorization Problem may be obtained. The time complexity of PRU is also the lowest among the three for the Adding Problem (table below, measured at  $k = 3$ ).

Table 2: Time cost per epoch for Adding Problem,  $k = 3$

$N$	2	4	6	8	10
LSTM	0.4678	0.8097	1.1826	1.5450	2.0387
GRU	0.4346	0.7270	0.9954	1.3106	1.5756
PRU	0.3191	0.5822	0.7367	0.9380	1.2037

## Character Prediction Problem

Let  $\mathcal{M}_{\text{char}}$  be a “character-prediction machine”, which takes an input sequence  $(x_1, x_2, \dots, x_N)$  of arbitrary length  $N$  and produces an output sequence of the same length. The input sequence is fed to the machine one symbol per time unit, and at each time  $t$ , the machine is characterized by a function  $\mathcal{M}_{\text{char}}^t$  defined by

$$\mathcal{M}_{\text{char}}^t(x_1, \dots, x_t) := x_{t+1}.$$

That is, for every input sequence, the output of the machine is the input sequence shifted in time. Here each symbol  $x^t$  is a character in a  $K$ -character alphabet. Each character in the alphabet is represented by a length- $K$  *one-hot* vector. The objective of the Character Prediction Problem is then to train a model that simulates the behaviour of  $\mathcal{M}_{\text{char}}$ .

**Modelling:** Naturally, both  $\mathcal{X}$  and  $\mathcal{Y}$  are taken as  $\mathbb{R}^K$  in the models. The output  $y_t$  is computed by a soft-max classifier.

**Dataset:** A Shakespeare drama dataset<sup>1</sup> is used in this experiment, where each sentence is taken as an input sequence. The dataset consists of 1,115,393 occurrences of characters from an alphabet of size 64, where 90% of the sentences are used for training set and the rest is held out for testing.

**Training and Evaluation:** The objective of this problem is to minimize the (expected) cross entropy loss (CEL)

$$\mathcal{E}_{\text{CEL}}(\theta) = -\mathbb{E} \left( \frac{1}{N} \sum_{t=1}^N \sum_{i=1}^K \mathcal{M}_{\text{char}}^t(x_1, \dots, x_t)[i] \log y_t[i] \right)$$

where we have used  $v[i]$  to denote the  $i^{\text{th}}$  component of vector  $v$ , and used  $N$  to denote the length of the input sequence. Mini-batched SGD with Adadelta dynamic learning rate [24] is used for optimization. All parameters are randomly initialized in  $[-0.1, 0.1]$ , and the base learning rate is set to 0.8. CEL is used to evaluate the models.

**Results:** Figure 9 plots the performances of the three units as functions of SDG epoch number. The three units show very close performances for the chosen three settings of state space dimension. The table below lists the CEL performance of the three recurrent units at the end of SDG iterations, where PRU appears slightly outperform the other two.

State Space Dimension	LSTM	GRU	PRU
64	1.2752	1.3015	1.2245
96	1.2211	1.2132	1.1894
128	1.1584	1.1968	1.1410

The average training time for PRU, GRU, and LSTM per epoch are respectively 83.65, 104.65 and 188.86 seconds respectively, with PRU leading by a significant margin.

<sup>1</sup><https://github.com/karpathy/char-rnn>

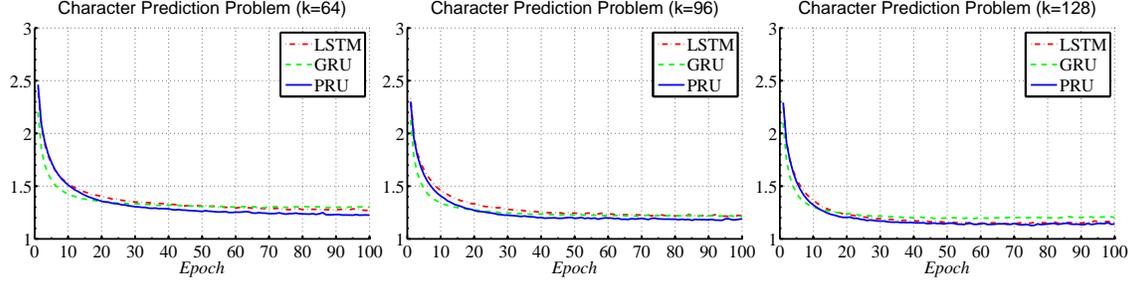


Figure 9: CEL comparison of LSTM, GRU and PRU in Character Prediction Problem

## MNIST Image Classification Problem

The MNIST dataset contains images of handwritten digits ('0'-'9'), all the images are of size  $28 \times 28$ , we treat each row of the image (28 pixels) as a single input in the input sequence. Let  $\mathcal{M}_{\text{image}}$  be a "image-classification machine", which takes an input sequence  $(x_1, x_2, \dots, x_N)$  with  $N = 28$  and predict the label  $y$  of the input sequence. the machine is asked to predict the category of the image after seeing all the pixels and  $\mathcal{M}_{\text{image}}$  can be defined by

$$\mathcal{M}_{\text{image}}(x_1, \dots, x_N) := y$$

The objective of the MNIST Image Classification Problem is then to train a model that simulates the behaviour of  $\mathcal{M}_{\text{char}}$ .

**Modelling:** It is natural to take input space  $\mathcal{X} = \mathbb{R}^{28}$  and  $\mathcal{Y} = \mathbb{R}^{10}$ , since the number of image category is 10. in the models. The output  $y$  is computed by a soft-max classifier.

**Dataset:** The MNIST dataset<sup>2</sup> contains 60,000 images in the training set, and 10,000 in the test set.

**Training and Evaluation:** The objective of this problem is to maximize the prediction accuracy (PA)

$$\mathcal{E}_{\text{PA}}(\theta) = \frac{1}{K} \sum_{i=1}^K \mathcal{G}(\mathcal{M}_{\text{image}}(x_1^i, \dots, x_N^i), y_i)$$

where we have used  $K$  to denote the total number of examples,  $(x_1^i, \dots, x_N^i)$  and  $y_i$  denote the input and label of an example respectively, if  $\mathcal{M}_{\text{image}}(x_1^i, \dots, x_N^i)$  matches  $y_i$ , the value of  $\mathcal{G}(\mathcal{M}_{\text{image}}(x_1^i, \dots, x_N^i), y_i)$  is 1, otherwise 0. Mini-batched SGD with Adadelta dynamic learning rate [24] is used for optimization. All parameters are randomly initialized in  $[-0.1, 0.1]$ , and the base learning rate is set to 0.1. PA is used to evaluate the models.

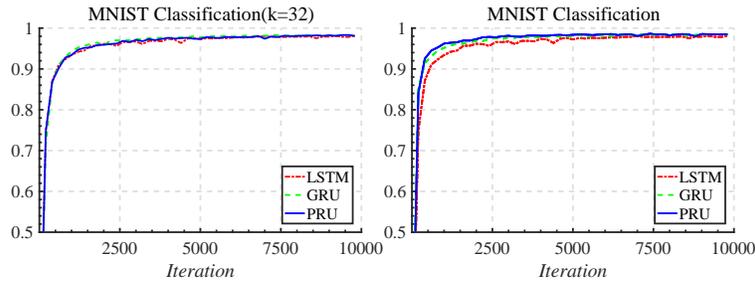


Figure 10: Test prediction accuracy comparison of LSTM, GRU and PRU in MNIST Image Classification Problem: Under same state space dimension(left), With same number of parameters(right)

**Results:** Results are obtained for LSTM, GRU and PRU in two conditions: under of model state-space dimensions  $k$  and the same number of parameters  $p$ . Figure 10 plots the performances of the three units as functions of iteration

<sup>2</sup> <http://yann.lecun.com/exdb/mnist/>

Experiment Setting	LSTM	GRU	PRU
Same state space dimension	0.9815	0.9821	0.9852
Same number of parameters	0.9815	0.9843	0.9856

Table 3: Prediction Accuracy for MNIST Image Classification

number. The three units show very close performances for the chosen three settings of state space dimension. The table below lists the PA performance of the three recurrent units at the end of SDG iterations, the performance on MNIST is similar to that on the character prediction problem, where PRU appears slightly outperform the other two.

The average training time for PRU, GRU, and LSTM per epoch are respectively 110.73, 154.94 and 210.04 seconds respectively, with PRU leading by a significant margin.

## Concluding Remarks

This paper presents a new recurrent unit, PRU. Having very simple structure, PRU is shown to perform similarly to LSTM and GRU. This potentially allows the use of PRU as a prototypical example for analytic study of LSTM-like recurrent networks. Its complexity advantage may also make it a practical alternative to LSTM and GRU.

This work is only the beginning of a journey towards understanding recurrent networks. It is our hope that PRU may provide some convenience to this important endeavor.

## References

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *Computer Science*, 2014.
- [2] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.
- [3] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. *arXiv preprint arXiv:1601.06733*, 2016.
- [4] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. pages 1724–1735, 2014.
- [5] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [6] Jeffrey L Elman. Distributed representations, simple recurrent networks, and grammatical structure. *Machine learning*, 7(2-3):195–225, 1991.
- [7] Felix Gers. *Long short-term memory in recurrent neural networks*. PhD thesis, Universität Hannover, 2001.
- [8] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction with lstm. *Neural computation*, 12(10):2451–2471, 2000.
- [9] Felix A Gers, Nicol N Schraudolph, and Jürgen Schmidhuber. Learning precise timing with lstm recurrent networks. *Journal of machine learning research*, 3(Aug):115–143, 2002.
- [10] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6645–6649. IEEE, 2013.
- [11] K Greff, R. K. Srivastava, J Koutnik, B. R. Steunebrink, and J Schmidhuber. Lstm: A search space odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 2015.
- [12] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [14] Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent network architectures. In *Proceedings of The 32nd International Conference on Machine Learning*, pages 2342–2350, 2015.
- [15] H. K. Khalil. *Nonlinear Systems*. Prentice-Hall, Englewood Cliffs, NJ, 3rd edition, 2002.
- [16] Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. Character-aware neural language models. 2016.
- [17] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [18] Tomas Mikolov, Armand Joulin, Sumit Chopra, Michael Mathieu, and Marc’Aurelio Ranzato. Learning longer memory in recurrent neural networks. *arXiv preprint arXiv:1412.7753*, 2014.
- [19] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Interspeech*, volume 2, page 3, 2010.

- [20] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. *ICML (3)*, 28:1310–1318, 2013.
- [21] Iulian V Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI-16)*, 2016.
- [22] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of video representations using lstms. *CoRR*, *abs/1502.04681*, 2, 2015.
- [23] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. *arXiv preprint arXiv:1601.06759*, 2016.
- [24] Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.