

Refined Gate: A Simple and Effective Gating Mechanism for Recurrent Units

Zhanzhan Cheng^{12*}, Yunlu Xu^{2*}, Mingjian Chen², Yu Qiao³, Shiliang Pu², Yi Niu², Fei Wu¹

¹Zhejiang University, Hangzhou, China; ²Hikvision Research Institute, China

³Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, China

ABSTRACT

Recurrent neural network (RNN) has been widely studied in sequence learning tasks, while the mainstream models (e.g., LSTM and GRU) rely on the gating mechanism (in control of how information flows between hidden states). However, the *vanilla* gates in RNN (e.g., the input gate in LSTM) suffer from the problem of *gate undertraining*, which can be caused by various factors, such as the saturating activation functions, the gate layouts (e.g., the gate number and gating functions), or even the suboptimal memory state *etc.*. Those may result in failures of learning gating switch roles and thus the weak performance. In this paper, we propose a new gating mechanism within general gated recurrent neural networks to handle this issue. Specifically, the proposed gates directly short connect the extracted input features to the outputs of *vanilla* gates, denoted as refined gates. The refining mechanism allows enhancing gradient back-propagation as well as extending the gating activation scope, which can guide RNN to reach possibly deeper minima. We verify the proposed gating mechanism on three popular types of gated RNNs including LSTM, GRU and MGU. Extensive experiments on 3 synthetic tasks, 3 language modeling tasks and 5 scene text recognition benchmarks demonstrate the effectiveness of our method.

1 INTRODUCTION

Recurrent neural networks (RNN) receive extensive research interests because of their powerful ability to handle sequential data in various applications such as action recognition[10], image captioning [41], text recognition [6], language translation [3], and speech recognition [14], etc. Specifically, the gated recurrent neural networks (GRNN) including Long Short Term Memory (LSTM) [21], Gated Recurrent Unit (GRU) [7] and Minimal Gated Unit (MGU) [44] are the prevailing variants of RNN, which can successfully learn the long-term dependency attributing to their *gates* on input x and recurrent hidden h (i.e., the *forget/input/output* gates in LSTM, the *update/reset* gates in GRU and the only *forget* gate in MGU).

Then, some recent works [9, 15, 23, 35] were proposed to investigate the *vanilla* gating mechanism. For example, [15] and [23] verified the effectiveness of these gates by analyzing different variants of GRNN such as 8 variants [15] for LSTM and 3 variants [23] for GRU, and expected to find out better variants of GRNN. Unfortunately, all variants didn't obviously outperform the standard GRNN models. [35] designed an Architecture Searcher to search

potential RNN architectures, but the best searched architecture didn't follow human intuition and cannot provide clear suggestions for the design of RNN structures.

Empirically, we find that *vanilla* gates in GRNNs tend to suffer from the problem of insensitive gating activation (e.g., narrowing activation scope and approximately activating values with small standard deviation). It means the gates can't play the switch roles very well and have limited ability to control the information flows between hidden states. We call this issue *gate undertraining*. *Gate undertraining* problem can be derived from three perspectives: (1) The *vanilla* gating functions in popular GRNNs are always represented by the sigmoid function σ (See the *vanilla* gate part of Figure 1 (b)), while the saturation characteristics of σ limits the gate roles as addressed in early works [12, 39]. It leads to the imbalance in training of the gates (slower learning) and the following feedforward neural networks (faster learning). When gates are trapped into the inevitable saturation, it means with parameter perturbed, the change to the gates tend to be small due to the sigmoid operator. Gates become persistently indiscriminate on flows while the remained network part keep tuning in a large speed. Thereby, gates in GRNN themselves might miss the befitting learning stage. (2) Different gate layouts (e.g., the gate number and the gating functions) affects the relation representation between hidden states. For example, MGU with a minimal σ gate may have limited ability to learn the long-dependency relationship between hidden states, because the minimal *vanilla* gate undertake multiple gate roles such as the output and the forget switch roles, which requires the minimal gate more discriminative. (3) The hidden state is responsible for dynamically holding the past and current information. The memory cell should be able to perceive the subtle differences between the current inputs and the past states, which means the gates should be very sensitive to capture flow change.

Above issues motivate us to develop a more discriminative gating mechanism to enhance the gate control ability, which not only avoid the saturation problem, but also perceive the relationship among information flows sensitively. And then the stronger gates can remedy the weak performance caused by undertraining. Inspired by previous researches of residual network [19] and highway network [38], we develop a simple yet effective gating mechanism named refined gate, introducing a path from the input features to the outputs of nonlinear activation, as shown in Figure 1(a). By this way, the refined gate become boundless and then eliminate the saturation problem (See the refined parts of Figure 1). On the other hand, such shortcut paths achieve the identity mapping between the memory state and the input, which plays the differential amplifier role to better depict the switch effects. Note that, different from introducing residual or highway block into RNNs for overcoming the gradient vanishing through multiple RNN layers [25, 43], our

*Both authors contributed equally to this research.

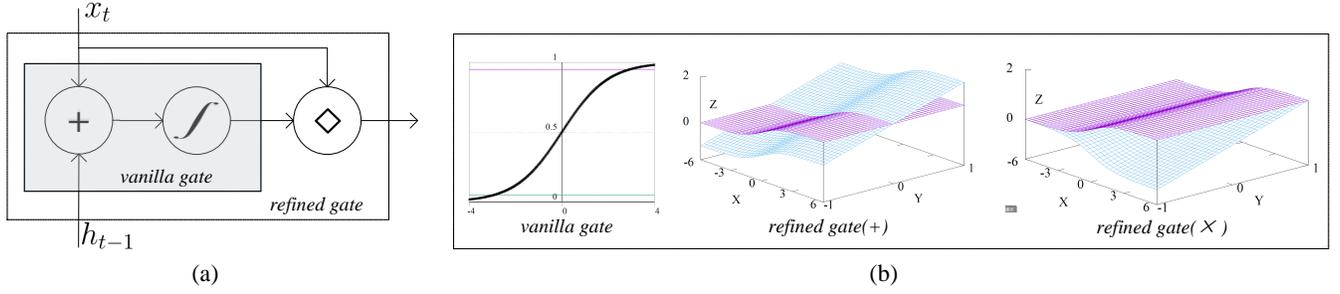


Figure 1: The illustration of refined gate. (a) is the refined gate block where \int and \diamond separately mean the activation function and the refining operation. In (b), the first sub-figure is the *vanilla* gate, the second and the last are the new gates with two refining operations ‘+’ and ‘ \times ’. In refined gates, the X, Y and Z axes correspond to the transformed (x, h) , the x_t and the gate output, respectively. The purple and blue surfaces separately denote the response of *vanilla* gate and refined gates.

focus is that the refined gate can learn to play the pure gating roles sensitively inside the recurrent units.

Concretely, instead of all the *vanilla* gates implemented with a sigmoid function, i.e., $g_{vanilla} = \sigma(x, h)$, we design the refined gate denoted as $g_{refine} = \sigma(x, h) \diamond x$. Here, \diamond is the element-wise refining operation, which combines the activation function σ and x to enhance the traditional gating mechanism. In this paper, we simply illustrate \diamond as two modes: + or \times , and treat the whole refined structure as a new gate for better controlling information flows. We note that the refined gates can provide broader and more dispersive activation scope than *vanilla* gates, as shown in Figure 1(b), which is also boundless but without the saturation problem. We demonstrate that the refined gating structure can be well equipped to general GRNN units boosting performance without introducing any extra parameters as well as any vanishing and exploding gradient problems.

The contributions of this paper are: (1) We provide a deeper understanding of gating mechanism in GRNNs and focus on the widely existing challenging problem: *gate undertraining*. (2) We propose a new gating mechanism enhancing the vanilla gates using simple yet effective refining operations, which is verified to be well adapted in existing GRNN units like LSTM, GRU and MGU. (3) We show intuitive evaluation on gate controlling ability through well-designed sequential tasks, i.e., adding and counting, and offer reasonable illustrations both qualitatively and quantitatively. (4) Experiments on various tasks, including 3 synthetic datasets and multiple real-world datasets (3 language modeling tasks and 5 scene text recognition benchmarks) demonstrate that the proposed gate refinement mechanism can effectively boost GRNN learning.

2 BACKGROUND

2.1 Gated RNNs

The simple recurrent architecture is hard to train properly in practice because of the vanishing and exploding gradient problems[34]. Therefore, various gated RNNs are developed for capturing long-term temporal dependencies, which are usually implemented by introducing various gates to control how information flows in RNN. We here briefly list three prevailing gated RNNs including LSTM,

GRU and MGU in Eq. 1-3.

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f), \quad (1a)$$

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i), \quad (1b)$$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o), \quad (1c)$$

$$\tilde{c}_t = \phi(W_c x_t + U_c h_{t-1} + b_c), \quad (1d)$$

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

$$h_t = o_t \odot \tanh(c_t). \quad (1f)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z), \quad (2a)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r), \quad (2b)$$

$$A = r_t \odot h_{t-1}, \quad (2c)$$

$$\tilde{h}_t = \phi(W_h A + U_h x_t + b_h), \quad (2d)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t. \quad (2e)$$

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f), \quad (3a)$$

$$A = f_t \odot h_{t-1}, \quad (3b)$$

$$\tilde{h}_t = \phi(W_h A + U_h x_t + b_h), \quad (3c)$$

$$h_t = (1 - f_t) \odot h_{t-1} + f_t \odot \tilde{h}_t. \quad (3d)$$

The typical *LSTM* [21] solved the gradient issues by introducing three gates (i.e., the *forget* gate, the *input* gate and the *output* gate) to control how information flows in RNN, formalized as Eq. 1 where ϕ denotes the tanh activation. Different from LSTM, *GRU* [7] discards the output gate as well as memory state, and makes each recurrent unit to adaptively capture dependencies of different time scales by introducing the update gate and the reset gate, which is formalized as Eq. 2. *MGU* [44] is a recent proposed gated RNN with only one forget gate (shown as Eq. 3), which is implemented by coupling the update and reset gate of GRU into a forget gate.

In addition, Greff *et al.* [15] explored 8 variants of the LSTM architecture to evaluate the effects of these gates, and Jozefowicz [23] evaluated 3 variants of GRU, but the results showed that none of variants can obviously outperform the standard models. Chung

et al. [9] also compared the performance of LSTM and GRU on multiple tasks and showed the similar performance.

2.2 Gate Functions

Some nonlinear activation functions like sigmoid and tanh have been extensively used in feedforward neural networks with the inevitable *gate undertraining* problem. Early works [12, 39] attempted to relieve the saturation problem by introducing the adaptive amplitude of activation strategy. When it comes to GRNNs, bounded activating gates still dominate in the literature to avoid exploding gradients while such gates fall into the *gate undertraining* problem, which might result in the poor performance on various sequence data. In common, dominating gated RNNs have the same gating modality: $\text{sigmoid}(W_g h_{t-1} + U_g x_t + b_g)$, where $[W_g, U_g, b_g]$ are the learnable parameters to determine the gating states.

Some methods have been designed for alleviating one of the *gate undertraining* problem: the saturating state. [16] injected noise to *sigmoid* and *tanh* and replaced the soft-saturating with the proposed hard-activations. [27] proposed the rectified linear unit (ReLU) in companion with identity weight matrix initialization, which seemingly sidestepped the saturating sigmoid yet brought new problem of the Dead ReLU. Recently, [5] proposed a new recurrent unit without saturating gates and evaluated with vanilla RNNs on multiple tasks. Existing works except [16] relieved the saturating problem by designing an entirely new recurrent unit forgoing the bounded non-linear gates. While we devote to further excavate the ON/OFF control ability of gates in original structures themselves with only a minor yet universally applicable revision.

2.3 Shortcut Structures

The residual structures [19] and highway networks [38] are influential implementations of the shortcut path for deep neural networks, both of which provide direct path of data flows in order to simplify the model training. As a result, works [25, 29, 43] attempted to focus on the adaption of shortcuts into the recurrent networks. Highway LSTM [43] connects the internal memory cells of the two neighboring layers, where the only change to the conventional LSTM is to add a highway connection of the last-layer cells to the c_t in Eq. 1d. Similar to [43], [29] introduced a simple recurrent unit with highway connections of input x_t on both the memory cells c_t and the output h_t to enhance parallelizable recurrence. Instead of the shortcut path on an internal memory cell c_t , residual LSTM [25] was proposed to add a path to the LSTM output h_t . All the above works were motivated by the shortcut and aimed to combine the structure into a specific RNN, but they focused on building deeper or longer RNNs or enhancing the parallelizability. Conformably, limited by the bound $[0, 1]$ of the traditional gates, these methods bypassed the essential parts in RNNs: the gate activation, which still seriously suffer from the undertraining problem. Different from existing explorations, our concentration is the gates.

3 METHODOLOGY

3.1 Refined Gate

The existing *vanilla* saturating gates limit the model learning and mainly result in the *gate undertraining* problem. It further brings the failures of controlling the information flows between hidden states

as well as the final weak performance. Considering the limitations, we propose a new gating mechanism to exploit the power of GRNNs. Specifically, the refined gate is implemented as a simple shortcut structure by directly short connecting x_t to the outputs of activation function, as shown in Figure 1 (a). The refined gate is described as:

$$\begin{aligned} g_t &= \sigma(\hat{g}_t) \diamond x_t \\ &= \sigma(W_g x_t + U_g h_{t-1} + b_g) \diamond x_t, \end{aligned} \quad (4)$$

in which $\sigma(\hat{g}_t)$ is the traditional gate function, and $[W_g, U_g, b_g]$ are learnable parameters to control information flows among hidden states. \diamond denotes the refining operation which can be formulated in the element-wise addition operation $+$ or the element-wise multiplication operation \times . Here, $+$ and \times correspond to two kinds of different refinement operations. That is, both x and σ will obtain the same gradient in $+$ case (as done in residual network). While in the \times case, x and σ will separately obtain gradients with the scales of σ and x reversely, treated as the mutual refinement.

Notice, instead of previous gating concept that gating value should be limited in $[0, 1]$, the refined gates are **boundless**. It explicitly embraces the decoupled current input features, and differential mapping between the memory state and the input. Therefore, the gates can better depict the switch role, which means a more clear ON/OFF controlling on flows, not ambiguous about current input features. Then the refining mechanism can be directly equipped in GRNNs well as long as it isn't directly applied in hidden state updating such as the Eq. 1e, 2e and 3d due to the gradient explosion [34], which has been demonstrated in Section 3.2.

3.2 GRNNs with Refinement

Refined LSTM It is safe to refine the *input* and *output* gates of LSTM, because only learnable parameters themselves (See Eq. 1b 1c) are directly affected by the refinement operation in back-propagation. Therefore, we here modify the input and output gates of traditional LSTM as follows

$$i'_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \diamond x_t, \quad (5a)$$

$$o'_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \diamond x_t, \quad (5b)$$

where \diamond denotes $+$ or \times operations. The refinement of gates in Eq. 5a and Eq. 5b can be applied independently or jointly.

However, the refinement operation is not suitable for the *forget* gate due to the gradient explosion in memory state learning. That is, if the *forget* gate is refined, $\delta c_t \rightarrow +\infty$ (shown in Eq. 6). Then f_t will be gradient exploding by referring to $\delta f_t = \delta c_t \odot c_{t-1}$, which results in the learning failure.

$$\delta c_t = \frac{\partial c_T}{\partial c_t} = \prod_{T \geq k \geq t} \text{diag}(f_k \diamond x_k). \quad (6)$$

Here, T is the index of the last hidden state.

Refined GRU The refinement on reset gate r_t in GRU (defined in Eq. 2b) is also safe, that is,

$$r'_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \diamond x_t. \quad (7)$$

Similarly, the refinement operation is also not suitable for the *update* gate due to the gradient explosion between hidden state learning, i.e., $\delta h_t \rightarrow +\infty$ because of $\delta h_t = \frac{\partial h_T}{\partial h_t} = \prod_{T \geq k \geq t} \text{diag}(z_k \diamond x_k)$. Then gradient explosion will directly occur in z_t by referring to $\delta z_t = \delta h_t \odot (\tilde{h}_t - h_{t-1})$.

Refined MGU As a variant of GRU, the refinement in MGU (defined in Eq. 3b) can be also revised as

$$f'_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \diamond x_t. \quad (8)$$

Note that, the refinement operation can be just applied into Eq. 3b while not suitable for Eq. 3d due to the same gradient explosion problem.

3.3 Back-propagation of Refining Mechanism

In order to explore the working mechanism of refined gates, we derive the gradients propagation process in gates to analyze the effectiveness of refined gates in both modes.

Referring to the forward process in Eq. 4, given calculated difference δg_t that back-propagates from the following feedforward neural network, we have $\delta \hat{g}_t$, δx_t , δh_{t-1} and $\delta[W_g, U_g, b]$ based on the chain rules, *i.e.*,

$$\delta \hat{g}_t = \delta g_t \odot \frac{\partial(\sigma(\hat{g}_t) \diamond x_t)}{\partial \hat{g}_t} \quad (9a)$$

$$\delta x_t = \delta g_t \odot \frac{\partial(\sigma(\hat{g}_t) \diamond x_t)}{\partial x_t} \quad (9b)$$

$$\delta h_{t-1} = \delta \hat{g}_t U, \quad (9c)$$

$$\delta[W, U, b] = \delta \hat{g}_t [x_t, h_{t-1}, 1]. \quad (9d)$$

Since the refining operation \diamond can be represented two forms $+$ and \times , each of them generally plays individual roles in gate refining process.

Concretely, for operation $+$, we can reformulate Eq. 9 as

$$\delta \hat{g}_t = \delta g_t \odot \hat{g}_t \odot (1 - \hat{g}_t), \quad (10a)$$

$$\delta x_t = \delta \hat{g}_t W + \delta g_t, \quad (10b)$$

$$\delta h_{t-1} = \delta \hat{g}_t U, \quad (10c)$$

$$\delta[W, U, b] = \delta \hat{g}_t [x_t, h_{t-1}, 1]. \quad (10d)$$

We find that GRNN can directly back-propagate δg_t to x_t , which plays a role like identity mapping [19] for x_t , as shown in Eq. 10b. Then the refined gate will relieve the burden on learning x_t and focus on the context relation between hidden states h_t . It is noting that the $+$ refining operation only affects the training of x_t , and makes no difference to the back-propagation towards other parts within a recurrent unit such as Eq. 10a, 10c and 10d.

For operation \times , we reformulate Eq. 9 as

$$\delta \hat{g}_t = \delta g_t \odot \hat{g}_t \odot (1 - \hat{g}_t) \times x_t, \quad (11a)$$

$$\delta x_t = \delta \hat{g}_t W + \delta g_t \odot \sigma(\hat{g}_t), \quad (11b)$$

$$\delta h_{t-1} = \delta \hat{g}_t U, \quad (11c)$$

$$\delta[W, U, b] = \delta \hat{g}_t [x_t, h_{t-1}, 1]. \quad (11d)$$

Here x_t can be treated as a scaling factor to dynamically adjust amplitudes of activation outputs, as shown in Eq. 11a. In fact, the tunable amplitude plays important roles to eliminate the saturation problem of activation functions, as demonstrated in [39]. Differently, we just use the learnt x_t as the amplitude factor for each neural unit without any additional parameters. Apart from the amplitude gains on x_t , the residual item $\delta g_t \odot \sigma(\hat{g}_t)$ in Eq. 11b also help gates learn their core contents, similar to the $+$ mode. Besides, it is a slightly larger change than the $+$ refining mode, where the \times operation will

affect the back-propagation of the whole recurrent model such as Eq. 11b, 11c and 11d.

4 EXPERIMENTS

In the section, we first explore the proposed architecture in three variants of GRNNs on sequential MNIST [28] task. For further illustrating gating functions, we design two independent numerical tasks with qualitative and quantitative analysis on gate controlling ability. Then we evaluate our refined gate on two general tasks: Language Modeling (LM) and the Scene Text Recognition (STR).

4.1 Ablation on Sequential Digits

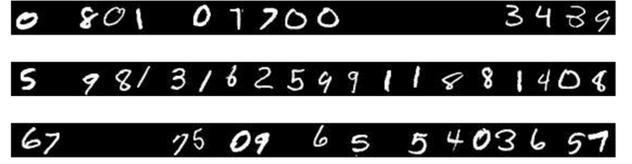


Figure 2: Samples of the sequence recognition task using the MNIST-based dataset.

Table 1: Ablation results of refined gates on multiple GRNN architectures. We report the average accuracy and standard deviation of 20 repetitive trails.

| Model | Refined Gate | Accuracy (+) | Accuracy (\times) |
|-------|----------------|------------------|------------------------------------|
| LSTM | – | 86.17 \pm 0.05 | 86.17 \pm 0.05 |
| | Input | 86.73 \pm 0.21 | 87.57 \pm 0.20 |
| | Output | 87.75 \pm 0.18 | 88.42 \pm 0.20 |
| | Input & Output | 87.85 \pm 0.21 | 88.62 \pm 0.19 |
| GRU | – | 85.22 \pm 0.04 | 85.22 \pm 0.04 |
| | Reset | 87.70 \pm 0.10 | 88.09 \pm 0.28 |
| MGU | – | 85.52 \pm 0.10 | 85.52 \pm 0.10 |
| | Forget | 88.34 \pm 0.19 | 88.67 \pm 0.15 |

We construct a more challenging dataset from MNIST for evaluation. The sequence recognition images with resolution of 560 x 32 pixels are composed of 12 to 20 digitals randomly selected from digitals of original 28 x 28 pixels [28] with jittering and non-overlapping over all-zero background. The training set contains 50,000 images and the test set has 10,000 images. Figure 2 shows three sequential samples.

We train the networks composed of a CNN encoding part with 7-layer convolution layers [36], a single GRNN layer and a connectionist temporal classification [13] as decoder. Our optimization is using AdaDelta with learning rate 1 for straightforward and convenient evaluation.

Table 1 shows the test results of our refined variants on LSTM, GRU and MGU structures. It can be seen that in all three GRNN models, both $+$ and \times refining operations significantly improve the recognition performance, and the \times operation achieves slightly better results than $+$ which might be attributed to the broader activation space of \times as shown in Figure 3. More details in LSTM,

the refined *output* gate (short in RO) outperforms the refined *input* gate (short in RI), because the RO affects input gates in the chain of back-propagation, while the RI can't affect other gates. Besides, the joint training of RO and RI (short in RIO) can achieve similar results with the RO, or may be slightly improving. In what follows, we use + as the refining operation by default.

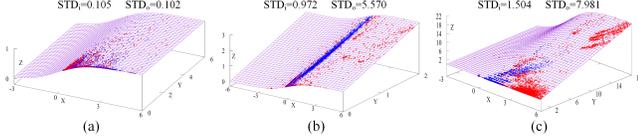


Figure 3: The illustration of activation distribution in gates. (a), (b) and (c) separately denote the distribution of vanilla gates, + and \times refined gates in LSTM. The blue and red points represent the response of *input* and *output* gates.

4.2 Ability of Gate Controlling

4.2.1 ON/OFF Indication on Adding Task. Adding problems [21] were once designed for evaluating long-term memorizing ability of RNN variants. While we focus on the gate controlling power, we make a minor revision of the task. The inputs have two binary sequences with a pre-set length $L \in \{10, 20, 50, 500, 5000, 50000\}$. For example, given the sample '011000000' and '010010000' as two inputs, the corresponding annotation is denoted as '000110000'. Note that the addition number are binary sequences in a reverse order. Since in a unidirectional RNN state learning setting, the current output depends on the information of the last state (i.e., the carry of each single bit addition is added to the next bit).

For each input length, we randomly generate a training set of 10,000 binary sequences and a test set of 5,000 samples, respectively. Thus, we expect the gates to turn ON/OFF for indicating the addition information, i.e., 0/1 carry bit. We evaluate effects of gates by feeding feature sequence into a single-direction RNN layer with 4 hidden states, after that following a fully connected layer to transform the hidden state vector as a sequential 2-class predictor.

As the adding problem is the simple addition operation, once the model converged, the predict accuracy tends to be 100%. Therefore, we compare the number of the completely converge epoches instead of the prediction accuracy to evaluate GRNNs as shown in Table 2. Table 2 shows all listed GRNN models can converge

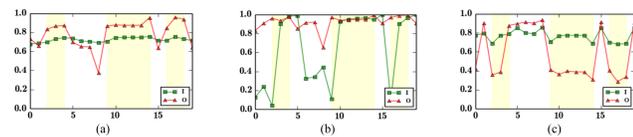


Figure 4: The adding process in LSTM. (a), (b) and (c) separately denote the activation values of σ in LSTM, LSTM-RI and LSTM-RO in + refining mode along the binary sequence. The green and red points represent the response of *input* and *output* gating functions. The yellow regions mean carry bit regions.

Table 2: The epochs of complete converge according to L . ∞ means non-convergence. 'K' means thousand. Lower values represents faster and better training.

| Architecture | Mode | L=10 | L=20 | L=50 | L=100 | L=500 |
|--------------|----------|------|------|----------|----------|----------|
| LSTM | None | 15 | 32 | 92 | ∞ | ∞ |
| LSTM-RI | + | 6 | 6 | 6 | 12 | ∞ |
| LSTM-RI | \times | 13 | 18 | 26 | 57 | ∞ |
| LSTM-RO | + | 6 | 6 | 6 | 12 | 12 |
| LSTM-RO | \times | 16 | 16 | 16 | 16 | 25 |
| GRU | None | 23 | 60 | ∞ | ∞ | ∞ |
| GRU-RR | + | 22 | 22 | 68 | ∞ | ∞ |
| GRU-RR | \times | 17 | 29 | 140 | ∞ | ∞ |
| MGU | None | 23 | 92 | ∞ | ∞ | ∞ |
| MGU-RF | + | 21 | 23 | 66 | ∞ | ∞ |
| MGU-RF | \times | 21 | 29 | 91 | ∞ | ∞ |

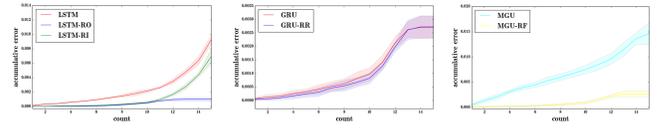


Figure 5: Distribution of accumulative error rates on different GRNNs with + refined gates. The horizontal and vertical axis represent the count value and the accumulative error rate, respectively.

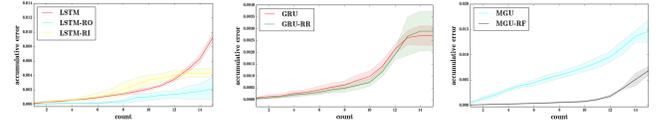


Figure 6: Distribution of accumulative error rate on different GRNNs with \times refined gates. The horizontal and vertical axis represent the count value and the accumulative error rate, respectively.

when L equals to 10 or 20, and LSTM can hold longer. Results can be quantitatively evaluated in two aspects: **converging speed** and **input length**. It costs more epoches for converging with the input length increasing, and the refined RNN models obviously obtain better performance than their vanilla structures. It worth noting that the LSTM-RO remarkably converges rapidly even when $L=50,000$ (not listed in the table) while vanilla models can only hold less than $L=50$, which illustrates the prominent superior performance of LSTM-RO comparing to the others (holding *1000 times longer* input length with less training epoches). We also speculate that the faster converge of the + than \times mode is due to the intrinsic context in this task, where the mathematical carry bit can be more easily obtained through residual operations. Note that we mention input length here only for evaluating the ON/OFF gate roles in different cases. In other words, even with much longer sequences, the goal of the networks is always the same, to learn the *carry bit*, which requires only the previous one memory state and the current input addition numbers. As for the controlling of long-dependency flows, we would go further to discuss in next section.

From the qualitative perspective of the **gate controlling ability**, we attribute the essence of the sequential learning in the adding task to indicating the carry bit. Conformably, we find that the refining operation makes gates more sensitive to the carry bit, when noticing the different amplitude of activation values in the gates, as shown in Figure 4. For detail, the *input* gate curve in LSTM is quite steady as a saboteur (Figure 4 (a)), while the input curve of LSTM-RI has drastic change clearly indicating each carry bit (Figure 4 (b)). Similarly, LSTM-RO has a very sensitive activation in its output gate (Figure 4 (c)). The adding task verifies the refined gates are more likely to perceive the gating roles of the adding problem, which shows that the refining operation can help eliminate the *gate undertraining*.

4.2.2 Dependency Controlling in Counting Task. We further evaluate the dependency controlling of gate in the counting task. The input is a binary sequence with length of 20. Each binary digit (0 or 1) in input sequences is encoded into a 2-dimensional one-hot encoding. We evaluate effects of gates by feeding feature sequences into an LSTM with 2 hidden states, after that following a fully connected layer to predict the repetition number of the last appearing digit. For instance, given the input sequence ended with '11110010110100100000', the predict value should be 5. In this set-

Table 3: Perplexity results of different recurrent models on the word-level PTB task. Lower is better.

| Model | Validation | Test |
|-------------------------|-------------|-------------|
| SoA Reported: | | |
| Unregularized LSTM [42] | 119.4 | 115.6 |
| Regularized LSTM [42] | 86.2 | 82.7 |
| Noisy LSTM+ NAH [16] | 111.7 | 108.0 |
| Variational LSTM [11] | 81.9 | 79.7 |
| Zoneout [26] | 84.4 | 80.6 |
| QRNN [4] | 82.9 | 79.9 |
| Our Setup: | | |
| LSTM [42] | 85.9 | 82.1 |
| LSTM-RI (Ours) | 82.6 | 80.1 |
| LSTM-RO (Ours) | 82.3 | 79.6 |
| LSTM-RIO (Ours) | 82.1 | 79.5 |

ting, gates are responsible for controlling information flow related to the repetition number of the last appearing digit. Figure 5 and 6 shows the accumulative error rates along different count ranges, in which the refined gates effectively decrease the error rates, especially in predicting larger counting number. We attribute it to improvement on the gating power of controlling the long-dependency flows. Similarly, as mentioned in adding task above (See Table 2), the refined RNN shows obvious improvement on the long-dependence sequential flows (holding over 1000 times longer input length with LSTM-RO compared with vanilla LSTM structure).

4.3 Language Modeling

4.3.1 Word-level PennTree Bank. The Penn TreeBank (PTB) dataset [32] provides data for language modeling. There are 10,000 words in the vocabulary, including 929K words in the training set, 73K in the validation set, and 82K in the test set. Our experimental

settings follow the standard setup [42]: 2 layers of 650 units in each layer with the sequence length of 35 as the encoder. Then a fully connected layer predicts one of the 10,000 words. The minibatch size is 20. We apply 50% dropout on the non-recurrent connections. We train for 39 epoches with a learning rate of 1 and then reduce it by a factor of 1.2 per epoch. Note that state-of-the-art (SoA) methods on the task explore the network design or the training strategies, while we focus on the innermost gating mechanism. As fancy techniques can bring improvements on overall performance but do no help to our fair comparison on the gates themselves, we compare our proposed model in the consistent setting (*i.e.*, 2 layers of 650 or 640 recurrent units), as shown in Table 3. For fair comparison, we implement the Regularized LSTM [42] as baseline in our setups. We also evaluates in Table 4 on the refined model with LSTM, GRU and MGU over various hidden units for comprehensive verification. It clearly shows that all the refined gates achieve improvements.

4.3.2 Char-level PennTree Bank. The character-level PennTree Bank (cPTB) task is to predict the next character in a text corpus at every character position, given all previous text. Our experiments comply with the setting of early work [8]. The model consists of an input embedding layer, an RNN module and an output layer. The RNN module has three layers, and the whole models are trained sequentially on the entire corpus, splitting it into sequences of length 50 for truncated back-propagation through time. The minibatch size is 128. We train models using Adam with an initial learning rate of 0.0002 and drop by a factor of 5 with patience 20.

We compare the performance of different models on the task in Table 5 in terms of test mean bits per character (BPC), where lower value indicates better results. Since reported SoAs are trained under different settings, *e.g.*, recurrent layers can be 1 [24] or 21 [30], we implement the models in the same setup for direct evaluation and focus on the improvement on the same module setting with 1 and 3 layers. It is noting that existing variants [30] motivated to improve gradient propagation among RNN layers, while our method is focused on the roles of gating mechanism, more similar with the work [5] which largely falls behind our method. It shows the improvement of our refined LSTM can be further enlarged in a 3-layer setting.

4.3.3 Wikipedia. We use EnWik8¹, a large public dataset for character-level language modeling. Following standard practice, we use the first 90M characters for training and the remaining 10M split evenly for validation and test. Similar to former work [42], we use a batch size of 128 with the unroll size of 100 for truncated backpropagation during training. We use the Adam optimizer and the same learning rate scheduling following [29] with a maximum of 100 epoches (about 700,000 iterations). The existing models shows close BPCs on tasks even in large variants of network settings (*e.g.*, the number of the recurrent layer, the hidden units in a wide range). Table 6 shows the improvement on the validation and test evaluation under the same settings with existing methods.

4.4 Scene Text Recognition

In addition to language modeling, we evaluate on visual data in text recognition tasks. Following the common settings, we use 2

¹The human knowledge compression contest 2012. <http://prize.hutter1.net/>.

Table 4: Perplexity statistics of three GRNNs (i.e., LSTM, GRU and MGU) and their refined variants on the PTB word-level prediction. Lower is better.

| Param #Hidden | LSTM | | | | GRU | | MGU | |
|------------------|-------|-------------|-------------|--------------|-------|--------------|-------|--------------|
| | LSTM | LSTM-RI | LSTM-RO | LSTM-RIO | GRU | GRU-RR | MGU | MGU-RF |
| 200 | 106.2 | 104.5 | 103.2 | 103.1 | 108.2 | 107.8 | 113.1 | 111.4 |
| 400 | 89.3 | 86.7 | 86.4 | 86.4 | 93.1 | 92.2 | 96.0 | 94.1 |
| 600 | 86.2 | 83.8 | 84.2 | 83.7 | 92.6 | 90.8 | 95.2 | 92.2 |
| 800 | 86.7 | 83.6 | 84.0 | 83.6 | 93.3 | 92.0 | 94.8 | 91.5 |

Table 5: Test BPCs of models on the char-level PTB task. Lower is better.

| Model | #Layer | #Step | Test BPC |
|--------------------------|--------|-------|-------------|
| SoA Reported: | | | |
| RNN [33] | 1 | 150 | 2.89 |
| LSTM [42] | 3 | 150 | 1.48 |
| EURNN [20] | 1 | 150 | 1.69 |
| HM-LSTM [18] | 3 | 100 | 1.30 |
| IndRNN [30] | 6 | 50 | 1.26 |
| IndRNN [30] | 21 | 50 | 1.21 |
| NRU [5] | 1 | 150 | 1.47 |
| nnRNN [24] | 3 | 150 | 1.47 |
| Our Setup: | | | |
| RNN | 1 | 50 | 1.49 |
| Residual LSTM [25] | 3 | 50 | 1.37 |
| IndRNN ¹ [30] | 3 | 50 | 1.36 |
| NRU ² [5] | 3 | 50 | 1.37 |
| LSTM | 1 | 50 | 1.46 |
| LSTM-RI (Ours) | 1 | 50 | 1.41 |
| LSTM-RO (Ours) | 1 | 50 | 1.39 |
| LSTM-RIO (Ours) | 1 | 50 | 1.39 |
| LSTM | 3 | 50 | 1.38 |
| LSTM-RI (Ours) | 3 | 50 | 1.32 |
| LSTM-RO (Ours) | 3 | 50 | 1.31 |
| LSTM-RIO (Ours) | 3 | 50 | 1.29 |

synthetic datasets MJSynth [22] and SynText [17] as our training set, and 5 general text benchmarks (including *III5K*, *SVT*, *IC03*, *IC13* and *IC15* [2]) as evaluation sets without any extra finetuning. We use the ResNet-32 [6] as backbone, sequence modeling (Bidirectional LSTM) for a contextual feature, and an attention-based decoder (in which one LSTM layer is applied for sequential decoding). Here in the sequence modeling, we choose the LSTM with refined output gate for evaluation. We train the model as: (1) backbone without sequence modeling as baseline, (2) traditional LSTM without refined gates, (3) LSTM-RO in its encoder and (4) LSTM-RO in both encoder and decoder.

Table 7 shows the recognition accuracies on the 5 benchmarks. For more convincing evaluation, we show the results of 1 layer and 2 layers of BiLSTMs in the top of its encoder. Results show LSTM-RO in encoder can outperform the baseline, and LSTM-RO in both the encoder and decoder get the best performance. We also compare our results with three strong baseline [6, 31, 37] and achieve competitive results, but only falls behind [37] on *III5K* and *SVT* because [37]

Table 6: Validation and test BPCs of different recurrent models on EnWik8 dataset. Lower is better.

| Model | #Layer | #Step | Valid | Test |
|----------------------------|--------|-------|-------------|-------------|
| SoA Reported: | | | | |
| LSTM [42] | 3 | 100 | - | 1.67 |
| Grid-LSTM [42] | 3 | 100 | - | 1.58 |
| MI-LSTM [40] | 3 | 100 | - | 1.44 |
| LN LSTM [1] | 3 | 100 | - | 1.46 |
| HM-LSTM [18] | 3 | 100 | - | 1.40 |
| HyperLSTM [18] | 3 | 100 | - | 1.39 |
| QRNN [4] | 4 | 200 | - | 1.33 |
| SRU [29] | 6 | 100 | 1.29 | 1.30 |
| Our Setup: | | | | |
| Residual LSTM [25] | 3 | 100 | 1.39 | 1.40 |
| QRNN(k=1) ³ [4] | 3 | 100 | 1.39 | 1.40 |
| SRU ⁴ [29] | 3 | 100 | 1.36 | 1.38 |
| LSTM | 3 | 100 | 1.38 | 1.40 |
| LSTM-RI(Ours) | 3 | 100 | 1.37 | 1.39 |
| LSTM-RO(Ours) | 3 | 100 | 1.36 | 1.38 |
| LSTM-RIO(Ours) | 3 | 100 | 1.36 | 1.37 |

used the deeper 45-layers ResNet as feature encoder. Intrinsically, the 2-layer LSTM-RO further enlarges the improvement above the baseline results from 1.3% to 1.7%, largely surpassing the traditional LSTM gain of 0.4%.

Note that the comprehensive analysis of modules in recent work [2] shows the 2-layer BiLSTM improve the baseline model (without the whole sequence modeling modules) by only slight gain in average (only 0.5% improvement in the model with vanilla LSTM layers in our setup), the **1.7% improvement with only minor revision** above the dominated sequence modeling method (LSTM) without extra cost is considerably significant in the STR tasks.

5 CONCLUSION

In this paper, we propose the a new gating mechanism for handling the *gate undertraining* problem. The refining operation is implemented by directly short connecting the inputs to the outputs

¹The SoA models are trained in different settings, and can not be directly comparable in the unified setups. We incorporated the IndRNN module with its released code at https://github.com/Sunnydreamrain/IndRNN_Theano_Lasagne for evaluation in the same setting. And NRU, QRNN and SRU are similar.

²We use the code at <https://github.com/apsarath/NRU>.

³We use the code at <https://github.com/salesforce/pytorch-qrnn>.

⁴We use the code at <https://github.com/taolei87/sru>.

Table 7: Recognition accuracies on the 5 general recognition datasets in a lexicon-free setting. For more straightforward evaluation, we compute the average (AVG) accuracies of the 5 benchmark results.

| Method | #Layer | #Step | Type | III5K | SVT | IC03 | IC13 | IC15 | AVG |
|----------------------|--------|-------|--------------|-------------|-------------|-------------|-------------|-------------|-------------------|
| SoA Reported: | | | | | | | | | |
| FAN [6] | 1 | 65 | – | 83.7 | 82.2 | 91.5 | 89.4 | 63.3 | 82.0 |
| ASTER [37] | 2 | 26 | – | 91.9 | 88.8 | 93.5 | 89.8 | – | – |
| MORAN [31] | 1 | 26 | – | 84.2 | 82.2 | 91.0 | 90.1 | 65.6 | 82.6 |
| Our Setup: | | | | | | | | | |
| Baseline | 0 | 65 | – | 86.1 | 84.7 | 92.6 | 90.3 | 67.4 | 84.2 |
| LSTM | 1 | 65 | – | 86.2 | 85.5 | 93.8 | 90.8 | 67.4 | 84.7(+0.5) |
| LSTM-RO (Ours) | 1 | 65 | Only Encoder | 86.2 | 84.9 | 93.9 | 91.5 | 68.9 | 85.1(+0.9) |
| LSTM-RO (Ours) | 1 | 65 | Both | 86.6 | 85.5 | 94.3 | 91.9 | 69.0 | 85.5(+1.3) |
| LSTM | 2 | 65 | – | 86.8 | 84.1 | 93.1 | 91.5 | 67.8 | 84.6(+0.4) |
| LSTM-RO (Ours) | 2 | 65 | Only Encoder | 87.2 | 85.3 | 93.5 | 92.0 | 69.2 | 85.4(+1.2) |
| LSTM-RO (Ours) | 2 | 65 | Both | 87.7 | 86.1 | 93.8 | 92.1 | 69.7 | 85.9(+1.7) |

of activating function, which can effectively enhance the controlling ability of gates. We illustrate the refined gates, offer a deeper understanding of the gating mechanism, and then evaluate their performance on several popular GRNNs. Extensive experiments on various sequence data including the scene text recognition, language modeling and arithmetic tasks, demonstrate the effectiveness of the refining mechanism. In future, we will further explore the gating mechanism in GRNNs.

REFERENCES

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer Normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [2] Jeonghun Baek, Geewook Kim, Junyeop Lee, Sungrae Park, Dongyoon Han, Sangdoon Yun, Seong Joon Oh, and Hwalsuk Lee. 2019. What is Wrong with Scene Text Recognition Model Comparisons? Dataset and Model Analysis. In *ICCV*. to appear.
- [3] Dzmitry Bahdanau, Jan Chorowski, Dmitriy Serdyuk, Yoshua Bengio, et al. 2016. End-to-end attention-based large vocabulary speech recognition. In *ICASSP*. IEEE, 4945–4949.
- [4] James Bradbury, Stephen Merity, Caiming Xiong, and Richard Socher. 2017. Quasi-Recurrent Neural Networks. In *ICLR*.
- [5] Sarath Chandar, Chinnadhurai Sankar, Eugene Vorontsov, Samira Ebrahimi Kahou, and Yoshua Bengio. 2019. Towards Non-saturating Recurrent Units for Modelling Long-term Dependencies. *arXiv preprint arXiv:1902.06704* (2019).
- [6] Zhanzhan Cheng, Fan Bai, Yunlu Xu, Gang Zheng, Shiliang Pu, and Shuigeng Zhou. 2017. Focusing Attention: Towards Accurate Text Recognition in Natural Images. In *ICCV*. 5076–5084.
- [7] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In *EMNLP*. 1724–1735.
- [8] Junyoung Chung, Sungjin Ahn, and Yoshua Bengio. 2017. Hierarchical Multiscale Recurrent Neural Networks. In *ICLR*.
- [9] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. *arXiv preprint arXiv:1412.3555* (2014).
- [10] Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Trevor Darrell, and Kate Saenko. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*. 2625–2634.
- [11] Yarín Gal and Zoubin Ghahramani. 2016. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. In *NIPS*. 1019–1027.
- [12] Su Lee Goh and Danilo P Mandic. 2003. Recurrent neural networks with trainable amplitude of activation functions. *Neural Networks* 16, 8 (2003), 1095–1100.
- [13] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks. In *ICML*. 369–376.
- [14] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *ICASSP*. 6645–6649.
- [15] Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. 2017. LSTM: A Search Space Odyssey. *IEEE TNNLS* 28, 10 (2017), 2222–2232.
- [16] Caglar Gulcehre, Marcin Moczulski, Misha Denil, and Yoshua Bengio. 2016. Noisy Activation Functions. In *ICML*. 3059–3068.
- [17] A. Gupta, A. Vedaldi, and A. Zisserman. 2016. Synthetic Data for Text Localisation in Natural Images. In *CVPR*. 2315–2324.
- [18] David Ha, Andrew Dai, and Quoc V. Le. 2016. HyperNetworks. *arXiv preprint arXiv:1609.09106* (2016).
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *CVPR*. 770–778.
- [20] Mikael Henaff, Arthur Szlam, and Yann Lecun. 2016. Recurrent Orthogonal Networks and Long-Memory Tasks. In *ICML*. 2034–2042.
- [21] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9, 8 (1997), 1735–1780.
- [22] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Synthetic Data and Artificial Neural Networks for Natural Scene Text Recognition. *arXiv preprint arXiv:1406.2227* (2014).
- [23] Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. 2015. An Empirical Exploration of Recurrent Network Architectures. In *ICML*. 2342–2350.
- [24] Giancarlo Kerg, Kyle Goyette, Maximilian Puelma Touzel, Gauthier Gidel, Eugene Vorontsov, Yoshua Bengio, and Guillaume Lajoie. 2019. Non-normal Recurrent Neural Network (nnRNN): learning long time dependencies while improving expressivity with transient dynamics. *arXiv preprint arXiv:1905.12080* (2019).
- [25] Jaeyoung Kim, Mostafa Elkhamy, and Jungwon Lee. 2017. Residual LSTM: Design of a Deep Recurrent Architecture for Distant Speech Recognition. In *ISCA*. 1591–1595.
- [26] David Krueger, Tegan Maharaj, Janos Kramer, Mohammad Pezeshki, Nicolas Ballas, Nan Rosemary Ke, Anirudh Goyal, Yoshua Bengio, Hugo Larochelle, Aaron C Courville, et al. 2016. Zoneout: Regularizing RNNs by Randomly Preserving Hidden Activations. *arXiv preprint arXiv:1606.01305* (2016).
- [27] Quoc V Le, Navdeep Jaitly, and Geoffrey E Hinton. 2015. A simple way to initialize recurrent networks of rectified linear units. *arXiv preprint arXiv:1504.00941* (2015).
- [28] Yann Lecun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [29] Tao Lei, Yu Zhang, Sida I Wang, Hui Dai, and Yoav Artzi. 2018. Simple Recurrent Units for Highly Parallelizable Recurrence. In *EMNLP*. 4470–4481.
- [30] Shuai Li, Wanqing Li, Christopher David Cook, Ce Zhu, and Yanbo Gao. 2018. Independently Recurrent Neural Network (IndRNN): Building A Longer and Deeper RNN. In *CVPR*. 5457–5466.
- [31] Canjie Luo, Lianwen Jin, and Zenghui Sun. 2019. MORAN: A Multi-Object Rectified Attention Network for Scene Text Recognition. *PR* 90 (2019), 109–118.
- [32] Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics* 19, 2 (1993), 313–330.
- [33] Behnam Neyshabur, Yuhuai Wu, Ruslan Salakhutdinov, and Nathan Srebro. 2016. Path-Normalized Optimization of Recurrent Neural Networks with ReLU Activations. In *NIPS*. 3477–3485.
- [34] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In *ICML*. 1310–1318.

- [35] Martin Schrimpf, Stephen Merity, James Bradbury, and Richard Socher. 2017. A Flexible Approach to Automated RNN Architecture Generation. *arXiv preprint arXiv:1712.07316* (2017).
- [36] Baoguang Shi, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. 2016. Robust Scene Text Recognition with Automatic Rectification. In *CVPR*. 4168–4176.
- [37] Baoguang Shi, Mingkun Yang, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. 2018. ASTER: An Attentional Scene Text Recognizer with Flexible Rectification. *IEEE TPAMI* 41, 9 (2018), 2035–2048.
- [38] Rupesh Kumar Srivastava, Klaus Greff, and Jurgen Schmidhuber. 2015. Training Very Deep Networks. In *NIPS*. 2377–2385.
- [39] Edmondo Trentin. 2001. Networks with trainable amplitude of activation functions. *Neural Networks* 14, 4-5 (2001), 471–493.
- [40] Yuhuai Wu, Saizheng Zhang, Ying Zhang, Yoshua Bengio, and Ruslan Salakhutdinov. 2016. On Multiplicative Integration with Recurrent Neural Networks. In *NIPS*. 2856–2864.
- [41] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Rich Zemel, and Yoshua Bengio. 2015. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *ICML*. 2048–2057.
- [42] Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. 2014. Recurrent Neural Network Regularization. *arXiv preprint arXiv:1409.2329* (2014).
- [43] Yu Zhang, Guoguo Chen, Dong Yu, Kaisheng Yao, Sanjeev Khudanpur, and James R Glass. 2016. Highway long short-term memory RNNs for distant speech recognition. In *ICASSP*. 5755–5759.
- [44] Guo-Bing Zhou, Jianxin Wu, Chen-Lin Zhang, and Zhi-Hua Zhou. 2016. Minimal Gated Unit for Recurrent Neural Networks. *IJAC* 13, 3 (2016), 226–234.