

# Reproducing and Regularizing the SCRN Model

## First Author

Affiliation / Address line 1  
Affiliation / Address line 2  
Affiliation / Address line 3  
email@domain

## Second Author

Affiliation / Address line 1  
Affiliation / Address line 2  
Affiliation / Address line 3  
email@domain

## Abstract

We reproduce the Structurally Constrained Recurrent Network (SCRN) model, and then regularize it using the existing widespread techniques, such as naïve dropout, variational dropout, and weight tying. We show that when regularized and optimized appropriately the SCRN model can achieve performance comparable with the ubiquitous LSTM model in language modeling task on English data, while outperforming it on non-English data.

## 1 Introduction

Recurrent neural networks (RNN) have demonstrated tremendous success in sequence modeling in general and in language modeling in particular. The most basic RNN (Elman, 1990) suffers from the problem of vanishing and exploding gradients (Bengio et al., 1994) and is hard to train efficiently. One of the most widespread and efficient alternatives to the basic RNN is the Long-Short Term Memory (LSTM) model (Hochreiter and Schmidhuber, 1997), which effectively addresses the problem of vanishing gradients. However, LSTM is a fairly complex model with excessive number of parameters and its inner functionality is not obvious. This complexity has motivated some of the researchers to find more apparent and less complex alternatives. One of such alternative models is a Structurally Constrained Recurrent Network (SCRN) proposed by Mikolov et al. (2015). They encouraged some of the hidden units to change their state slowly by making part of the recurrent weight matrix close to identity, thus forming a kind of longer term memory and showed that their SCRN model can outperform the simple RNN and achieve the performance comparable with the LSTM under no regularization and small parameter budget. It is natural to try to regularize the SCRN model under larger budgets: *Will it approach the performance of LSTM?* Our experiments show that (1) under naïve dropout the SCRN demonstrates performance close to that of the LSTM, but (2) under variational dropout and weight tying the LSTM demonstrates better performance.

Recently Ororbia II et al. (2017) introduced Delta-RNN architecture for which SCRN serves as a predecessor. They showed that, when regularized using naïve dropout (Zaremba et al., 2014), Delta-RNN performs comparably to LSTM and GRU (Chung et al., 2014). However, they did not compare to regularized SCRN, and they did not consider more recent regularization techniques, such as variational dropout (Gal and Ghahramani, 2016) and weight tying (Inan et al., 2017; Press and Wolf, 2017).

## 2 Baseline SCRN Model

Let  $\mathcal{W}$  be a finite vocabulary of words. We assume that words have already been converted into indices. Based on one-hot word embeddings  $\mathbf{x}_{1:k} = \mathbf{x}_1, \dots, \mathbf{x}_k$  for a sequence of words  $w_{1:k}$ , the baseline SCRN model (Mikolov et al., 2015) produces two sequences of states,  $\mathbf{s}_{1:k}$  and  $\mathbf{h}_{1:k}$ , according to<sup>1</sup>

$$\mathbf{s}_t = (1 - \alpha)\mathbf{x}_t\mathbf{B} + \alpha\mathbf{s}_{t-1}, \quad (1)$$

$$\mathbf{h}_t = \sigma(\mathbf{x}_t\mathbf{P} + \mathbf{s}_t\mathbf{A} + \mathbf{h}_{t-1}\mathbf{R}), \quad (2)$$

---

<sup>1</sup>Vectors are assumed to be row vectors, which are right multiplied by matrices ( $\mathbf{x}\mathbf{W} + \mathbf{b}$ ). This choice is somewhat non-standard but it maps better to the way networks are implemented in code using matrix libraries such as TensorFlow.

where  $\mathbf{B} \in \mathbb{R}^{|\mathcal{W}| \times d_s}$ ,  $\mathbf{P} \in \mathbb{R}^{|\mathcal{W}| \times d_h}$ ,  $\mathbf{A} \in \mathbb{R}^{d_s \times d_h}$ ,  $\mathbf{R} \in \mathbb{R}^{d_h \times d_h}$ ,  $d_s$  and  $d_h$  are dimensions of  $\mathbf{s}_t$  and  $\mathbf{h}_t$ ,  $\sigma(\cdot)$  is the logistic sigmoid function. Mikolov et al. (2015) refer to  $\mathbf{s}_t$  as a slowly changing *context state*, and to  $\mathbf{h}_t$  as a quickly changing *hidden state*. The last couple of states ( $\mathbf{s}_k, \mathbf{h}_k$ ) is assumed to contain information on the whole sequence  $w_{1:k}$  and is further used for predicting the next word  $w_{k+1}$  of a sequence according to the probability distribution

$$\Pr(w_{k+1}|w_{1:k}) = \text{softmax}(\mathbf{s}_k \mathbf{U} + \mathbf{h}_k \mathbf{V}), \quad (3)$$

where  $\mathbf{U} \in \mathbb{R}^{d_s \times |\mathcal{W}|}$  and  $\mathbf{V} \in \mathbb{R}^{d_h \times |\mathcal{W}|}$  are output embedding matrices. For the sake of simplicity we omit bias terms in (2) and (3).

Training the model involves minimizing the negative log-likelihood over the corpus  $w_{1:K}$ :

$$-\sum_{k=1}^K \log \Pr(w_k|w_{1:k-1}) \rightarrow \min_{\Theta}, \quad (4)$$

which is usually done by truncated backpropagation through time (Werbos, 1990). Here  $\Theta$  denotes the set of all model parameters.

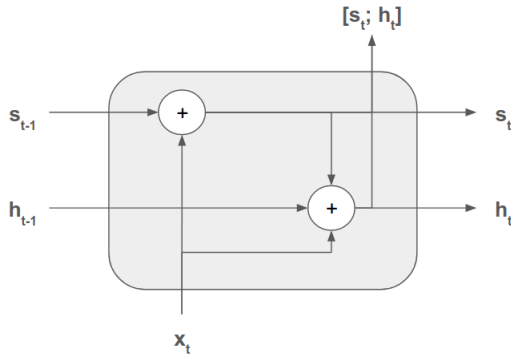


Figure 1: SCR cell.

$\mathbf{P} \in \mathbb{R}^{d_h \times d_h}$ . In this case, the model spends  $|\mathcal{W}| \cdot (2d_h + d_s)$  parameters on input/output embeddings, which is  $d_s \cdot |\mathcal{W}|$  parameters less than (5). E.g., on the Penn Tree Bank (PTB) dataset (Marcus et al., 1993), where  $|\mathcal{W}| = 10,000$ , the reduction is 1M parameters for the SCR model with  $d_s = 100$ .

Notice that SCR is a slight modification of the vanilla RNN model (Elman, 1990), and its simplicity is in stark contrast with the complexity of the widespread LSTM model.

**Dense embeddings:** The original model spends

$$2 \cdot |\mathcal{W}| \cdot (d_s + d_h) \quad (5)$$

parameters to embed words into dense vectors at input and at output. We believe that it is more beneficial to first embed words  $w_t$  into dense vectors  $\mathbf{w}_t = \mathbf{x}_t \mathbf{E} \in \mathbb{R}^{d_h}$  using only one embedding matrix  $\mathbf{E} \in \mathbb{R}^{|\mathcal{W}| \times d_h}$  and then use  $\mathbf{w}_t$  instead of  $\mathbf{x}_t$  in (1) and (2) with the appropriate change of shapes for matrices:  $\mathbf{B} \in \mathbb{R}^{d_h \times d_s}$  and

### 3 Stacking and Regularizing the SCR

#### 3.1 Stacking recurrent cells

It is well known that stacking at least two layers in recurrent neural networks is beneficial, but between two and three layers the results are mixed (Karpathy et al., 2016; Laurent and von Brecht, 2017). To make our results comparable to the previous works on LSTM language modeling (Zaremba et al., 2014; Gal and Ghahramani, 2016; Inan et al., 2017), we experiment with two-layered architectures: the output of the first layer (1, 2) is concatenated  $[\mathbf{s}_t; \mathbf{h}_t]$  and is fed as input into the second layer.

In what follows the superscript index in round brackets denotes the layer index,  $\xi(p) = [\xi_1, \dots, \xi_d]$  is a random vector (dropout mask) with  $\xi_i \sim \text{Bernoulli}(1-p)$ ,  $d$  is the dimensionality of the corresponding layer,  $p$  is a dropout rate, and  $\odot$  is the element-wise (Hadamard) product. One time-step of a stacked SCR model is fully specified by the following equations:

- Input embedding layer:

$$\mathbf{y}_t^{(0)} = \mathbf{x}_t \mathbf{E}. \quad (6)$$

- Two SCR layers: for  $l = 1, 2$

$$\begin{aligned} \mathbf{s}_t^{(l)} &= (1 - \alpha) \mathbf{y}_t^{(l-1)} \mathbf{B}^{(l)} + \alpha \mathbf{s}_{t-1}^{(l)}, \\ \mathbf{h}_t^{(l)} &= \sigma \left( \mathbf{y}_t^{(l-1)} \mathbf{P}^{(l)} + \mathbf{s}_t^{(l)} \mathbf{A}^{(l)} + \mathbf{h}_{t-1}^{(l)} \mathbf{R}^{(l)} \right), \\ \mathbf{y}_t^{(l)} &= [\mathbf{s}_t^{(l)}; \mathbf{h}_t^{(l)}]. \end{aligned} \quad (7)$$

- Softmax prediction:

$$\Pr(w_{t+1}|w_{1:t}) = \text{softmax} \left( \mathbf{y}_t^{(2)} \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix} \right). \quad (8)$$

In the equations above, the matrices  $\mathbf{U}$  and  $\mathbf{V}$  are concatenated along the first dimension to produce a  $(d_s + d_h) \times |\mathcal{W}|$  matrix.

### 3.2 Naïve dropout

Due to high complexity of deep neural networks the regularization techniques are crucial for good generalization performance. Zaremba et al. (2014) proposed one of the ways how dropout (Srivastava et al., 2014) can be used to regularize recurrent neural networks. They apply dropout only to non-recurrent connections, while keeping the recurrent connections without change. This method, usually referred to as *naïve dropout*, improves the baseline results of the LSTM model without any other modifications. Application of the same dropout technique to the SCRNN model (6, 7, 8) is specified by

$$\begin{aligned} \hat{\mathbf{y}}_t^{(0)} &= \mathbf{y}_t^{(0)} \odot \boldsymbol{\xi}_t^{(0)}(p_i), \\ \hat{\mathbf{y}}_t^{(l)} &= \mathbf{y}_t^{(l)} \odot \boldsymbol{\xi}_t^{(l)}(p_o), \end{aligned}$$

correspondingly, where  $p_i$  and  $p_o$  are input and output dropout rates.

### 3.3 Variational dropout

The approach of Zaremba et al. (2014) have led many to believe that dropout cannot be extended to recurrent connections, leaving them with no regularization. However, Gal and Ghahramani (2016) showed that it is possible to derive a variant of dropout which successfully regularizes recurrent connections. In their dropout variant, which is usually referred to as *variational dropout*, they repeat the same dropout mask at each time step for inputs, recurrent layers, and outputs:

$$\begin{aligned} \bar{\mathbf{y}}_t^{(0)} &= \mathbf{y}_t^{(0)} \odot \boldsymbol{\xi}^{(0)}(p_i), \\ \bar{\mathbf{h}}_t^{(l)} &= \mathbf{h}_t^{(l)} \odot \boldsymbol{\xi}^{(l)}(p_h), \\ \bar{\mathbf{y}}_t^{(l)} &= \mathbf{y}_t^{(l)} \odot \boldsymbol{\xi}^{(l)}(p_o), \end{aligned} \quad (9)$$

where  $p_h$  is a recurrent dropout rate. This is in contrast to the naïve dropout where different masks are sampled at each time step for the inputs and outputs alone and no dropout is used with the recurrent connections. Notice that we do not regularize the context state  $\mathbf{s}_t$  in horizontal (recurrent) direction.

### 3.4 Tying word embeddings

Tying input and output word embeddings in word-level RNNLM is a regularization technique, which was introduced earlier (Bengio et al., 2001; Mnih and Hinton, 2007) but has been widely used relatively recently, and there is empirical evidence (Press and Wolf, 2017) as well as theoretical justification (Inan et al., 2017) that such a simple trick improves language modeling quality while decreasing the total number of trainable parameters almost two-fold, since most of the parameters are due to embedding matrices. In case of the SCRNN model, reusing input embeddings at output can be done by setting

$$\mathbf{V} = \mathbf{E}^\top$$

in the softmax layer (8).

## 4 Experimental setup

We use perplexity (PPL) to evaluate the performance of the language models. Perplexity of a model over a sequence  $[w_1, \dots, w_T]$  is given by

$$\text{PPL} = \exp \left( -\frac{1}{T} \sum_{k=1}^T \log \Pr(w_k | w_{1:k-1}) \right).$$

**Data sets:** The baseline model is trained and evaluated on the PTB (Marcus et al., 1993), while all regularized and stacked configurations are trained and evaluated on the PTB and the WikiText-2 (Merity et al., 2017) data sets. For the PTB we utilize the standard training (0-20), validation (21-22), and test (23-24) splits along with pre-processing per Mikolov et al. (2010). WikiText-2 is an alternative to PTB, which is approximately two times as large in size and three times as large in vocabulary.

**Baseline Model:** To reproduce the results of the baseline (single-layer and non-regularized) SCRNN model we use the original `Torch` implementation<sup>2</sup> released by Mikolov et al. (2015). We have spent fair amount of time and effort to make their script run, as several of its dependencies have not been updated for few years, and are not compatible with the up-to-date versions of the others. To simplify the path for other researchers we release a script<sup>3</sup>, which installs necessary versions of the dependencies. We use exactly the same set of hyperparameters which were reported in the original paper (Table 1). We also

Hyperparameter	Mikolov et al. (2015)	Our implementation
$\alpha$ in (1)	0.05	0.05
batch size	32	20
initial LR	0.05	0.8
LR decay	1/1.5	0.5 (0.75)
LR decayed if	valid PPL doesn't improve	valid PPL doesn't improve (after the first 6 epochs)
BPTT steps	50	35 (70)
BPTT frequency	5	35 (70)
gradients	renormalized	norms clipped at 5
weights initialized over	$[-0.05, 0.05]$	$[-0.3, 0.3]$ ( $[-0.05, 0.05]$ )

Table 1: Hyperparameters of the baseline SCRNN model. Abbreviations: LR — learning rate, BPTT — backpropagation through time. Values in brackets correspond to the  $(d_h, d_s) = (300, 40)$  configuration when they differ from others.

implement the baseline SCRNN model ourselves<sup>4</sup> using `TensorFlow` (Abadi et al., 2016) which, unlike `Torch`, uses static computational graphs, and thus has different style of truncated backpropagation through time<sup>5</sup> (BPTT). The choice of hyperparameters (Table 1) is motivated by the previous work on word-level language modeling (Zaremba et al., 2014).

**Stacked and Regularized Models:** In the previous works on regularizing the LSTM, small-sized models usually had  $d_h = 200$  and medium-sized models had  $d_h = 650$ . The inner simplicity of the SCRNN cell allows us slightly larger hidden sizes: we use  $d_h = 240$  for small models and  $d_h = 750$  for medium models. Context state sizes  $d_s$  are chosen to be 40 (small) and 120 (medium), so that total number of parameters does not exceed the budget, which is 5M parameters for small models, and 20M parameters for medium models. We find empirically, that a good ratio between context size and hidden size in the SCRNN model is around 1/6. We optimize hyperparameters separately under naïve dropout and under variational dropout. Some of the hyperparameters are tuned using random search according to the marginal distributions:

- $p_i \sim U[0.01, 0.5]$ ,
- $p_h \sim U[0.01, 0.5]$ ,
- $p_o \sim U[0.01, 0.5]$ ,
- initial LR  $\sim U[0.5, 0.99]$ ,
- LR decay  $\sim U[0.5, 0.89]$ ,
- initialization scale  $\sim U[0.05, 0.3]$ .

where  $U[a, b]$  means continuous uniform distribution over the interval  $[a, b]$ , and initialization scale is a number  $r$  such that all model weights are initialized uniformly over  $[-r, r]$ . Other hyperparameters

<sup>2</sup><https://github.com/facebookarchive/SCRNNs>

<sup>3</sup><http://masked>

<sup>4</sup>Our implementation is available at <http://MASKED>

<sup>5</sup><https://r2rt.com/styles-of-truncated-backpropagation.html>

are tuned manually through trial-and-error. When performing random search we first choose ranges mentioned above. After 100 runs the initial ranges are shrunk to the neighborhoods of the values that give best performances, and the random search is performed again. We repeat this procedure until hyperparameters converge to their (sub)optimal values. To prevent exploding gradients we clip the norm of the gradients (normalized by minibatch size) at 5. For training (4) we use stochastic gradient descent.

## 5 Results

**Baseline:** To assure that our implementation of the baseline SCRNN is adequate, we evaluate it against the original SCRNN code by Mikolov et al. (2015) (Table 2). As one can see, the original SCRNN code does

Hidden size	Context size	Mikolov et al. (2015)		Our implementation	
		Valid PPL	Test PPL	Valid PPL	Test PPL
40	10	133.6 (133)	127.5 (127)	134.5	128.0
90	10	125.4 (124)	120.3 (119)	124.9	118.6
100	40	127.6 (120)	122.9 (115)	124.9	118.7
300	40	130.1 (120)	124.4 (115)	120.4	116.2

Table 2: Reproducing the baseline model. For the original implementation (columns 3 and 4), values outside brackets were obtained when running the script from <https://github.com/facebookarchive/SCRNNs>, and values in brackets were reported in the paper of Mikolov et al. (2015).

not fully reproduce the results reported in the original paper. Their hyperparameters (Table 1) work well for the case when  $(d_s, d_h) \in \{(40, 10), (90, 10)\}$ , but are not optimal for the other two configurations. Our implementation together with our set of hyperparameters (Table 1) brings the validation and test perplexities of *all* the configurations close to that which are reported in the paper.

**Stacked and Regularized Models:** Tuning the stack of two SCRNNs results in hyperparameters in Table 3. The results of evaluating these models against regularized and stacked LSTMs on PTB and

Hyperparameter	SCRNN + ND		SCRNN + VD	
	Small	Medium	Small	Medium
$p_i$	0.2 (0.15)	0.55 (0.45)	0.15 (0.1)	0.4 (0.3)
$p_h$	—	—	0.15 (0.1)	0.4 (0.3)
$p_o$	0.2 (0.15)	0.55 (0.45)	0.15 (0.1)	0.4 (0.3)
$\alpha$	0.9	0.9	0.9	0.9
initial LR	0.8	0.8	0.8	0.6
LR decay	0.5	0.65	0.87	0.9
LR decayed	when valid PPL does not improve		after 10 epochs	after 25 epochs
initialization scale	0.3	0.3	0.3	0.3

Table 3: Tuned hyperparameters: ND – naïve dropout, VD – variational dropout. Values in brackets correspond to the WikiText-2 in cases when they differ from those used for the PTB.

WikiText-2 are provided in Table 4. Regularization *does* benefit the simple and intuitive SCRNN model, which demonstrates performance comparable to the sophisticated LSTM model under the same regularization and the same parameter budget. This shows that at some point less complex models become less competitive, even when regularized and optimized appropriately. It is important to mention that no architectural modifications were applied to the original SCRNN model except stacking.

Our feeling is that complexity *is* needed for the recent regularization techniques to work well. For example, consider the variational dropout of the recurrent connections (9) in the SCRNN: dropping the

Word-level model	PTB		WikiText-2	
	Small	Medium	Small	Medium
LSTM + ND (Jozefowicz et al., 2015)	—	<b>79.8</b>	—	—
LSTM + ND (Kim et al., 2016)	97.6	85.4	116.8 <sup>†</sup>	—
LSTM + ND (Zaremba et al., 2014)	—	82.7	—	96.2 <sup>‡</sup>
SCRN + ND	95.8	84.9	115.0	86.2
SCRN + ND + WT	<b>94.1</b>	85.7	<b>112.6</b>	<b>85.2</b>
LSTM + VD (Gal and Ghahramani, 2016)	—	78.6	—	—
LSTM + VD (Inan et al., 2017)	87.3	77.7	105.9	95.3
LSTM + VD + WT (Inan et al., 2017)	<b>85.1</b>	<b>73.9</b>	100.5	<b>87.7</b>
SCRN + VD	97.2	90.7	100.0	90.6
SCRN + VD + WT	96.8	90.7	<b>97.7</b>	88.4

Table 4: Evaluation of the SCRN against LSTM under different regularization techniques: ND – naïve dropout, VD – variational dropout, WT – weight tying. <sup>†</sup>We reproduced the LSTM-Word-Small model from Kim et al. (2016) on PTB and then evaluated it on WikiText-2. <sup>‡</sup>We ran the open-source implementation from <https://github.com/tensorflow/models/tree/master/tutorials/rnn/ptb> at medium config on WikiText-2 data.

coordinates  $i_1, \dots, i_l$  in the vector  $\mathbf{h}_t$  is equivalent to zeroing out the columns  $i_1, \dots, i_l$  in the matrices  $\mathbf{A}, \mathbf{P}, \mathbf{R}$ . Now consider one layer of the LSTM model:

$$\begin{aligned}
\mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f) \\
\mathbf{i}_t &= \sigma(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i) \\
\mathbf{o}_t &= \sigma(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o) \\
\mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c) \\
\mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t)
\end{aligned}$$

When one drops the coordinates  $i_1, \dots, i_l$  of  $\mathbf{h}_t$  in LSTM, this can be understood as zeroing out the columns  $i_1, \dots, i_l$  in the matrices  $\mathbf{W}_f, \mathbf{W}_i, \mathbf{W}_o, \mathbf{U}_f, \mathbf{U}_i, \mathbf{U}_o$ , i.e. around 3 times more parameters are zeroed out given the same hidden state size  $d_h$ . In other words, the number of recurrent weights is 3 times larger in LSTM than in SCRN. We believe this is one of the reasons why the variational dropout works worse in SCRN. Roughly speaking, there is nothing much to regularize in the horizontal (recurrent) direction in SCRN, as most parameters are in the embedding and softmax layers. However, based on Table 4 one might conclude that SCRN + VD outperforms LSTM + VD on WikiText-2, but we think this is due to insufficient optimization of hyperparameters in the work of Inan et al. (2017).

### 5.1 Ablation analysis

We consider removal of some parts or forms of regularization from one of our well-performing configurations, SCRN + ND + WT, to see whether such removal degrades the performance. It also tells us which parts of the model are most important. The model and dropout variations are listed below.

**Removing regularization:** To understand how much improvement is gained by the use of regularization we completely remove dropout and weight tying:

$$p_i = p_o = 0, \quad \mathbf{V} \neq \mathbf{E}^\top$$

**Removing dropout of the context state:** According to (1), context state  $\mathbf{s}_t$  changes linearly and thus should not suffer from over-fitting. Thus it seems reasonable to try to not regularize it and apply dropout only to the hidden states, i.e. replacing the equation (7) by

$$\mathbf{y}_t^{(l)} = [\mathbf{s}_t^{(l)}; \mathbf{h}_t^{(l)} \odot \boldsymbol{\xi}_t^{(l)}(p_o)], \quad l = 1, 2.$$

**Removing the context state from the softmax layer:** As in the case of the SCRN, the inner state of the LSTM model also consists of two vectors  $\mathbf{c}_t$  and  $\mathbf{h}_t$ , and usually the state  $\mathbf{c}_t$  is not used at softmax. We do the same for the context state  $\mathbf{s}_t$  in our model, i.e. the equation (8) is replaced by

$$\Pr(w_{t+1}|w_{1:t}) = \text{softmax}(\mathbf{h}_t^{(2)}\mathbf{V}).$$

The meaningfulness of removing the context state from the softmax is that, in our opinion, it plays the role of a long-term memory and thus should not be crucial for predicting the next word of a sequence. Moreover, such removal reduces the model size by at least  $d_s \cdot |\mathcal{W}|$  parameters, which can be significant (see Section 2).

The results of the ablation analysis are provided in Table 5. As we can see, without regularization

Model	Small		Medium	
	Valid	Test	Valid	Test
SCRN + ND + WT	98.2	94.1	90.1	85.7
– regularization	123.3	117.6	146.1	140.2
– dropout of context	108.4	103.6	115.0	109.4
– context in softmax	100.8	96.4	91.7	87.6

Table 5: Model ablations for the small and medium SCRN + ND + WT models on PTB set.

our SCRN + ND + WT model fails to generalize well on validation and test sets. Regularizing only the hidden state (and keeping the context state untouched) is less harmful but still degrades the performance of the model. Finally, not using the context state in the output only slightly worsens the performance but at the same time leads to a significant reduction in model size.

## 6 Performance on non-English data

It is interesting to see whether the results extend to non-English languages. For this purpose we conduct evaluation of the medium-sized versions of LSTM + VD + WT against SCRN + ND + WT on non-English data which comes from the 2013 ACL Workshop on Machine translation<sup>6</sup> with pre-processing per Botha and Blunsom (2014). Hyperparameters tuned on Wikitext-2 (Table 3) were used for all languages and we did not perform any language-specific tuning. Corpora statistics and the results of evaluation are provided in Table 6. Surprisingly, the SCRN model performs very well (compared to LSTM)

	French	Spanish	German	Czech	Russian
Number of tokens	1M	1M	1M	1M	1M
Vocabulary size	25K	27K	37K	46K	62K
LSTM + ND (Kim et al., 2016)	222	200	286	493	357
LSTM + VD + WT	205	193	277	488	351
SCRN + ND + WT	<b>199</b>	<b>179</b>	<b>258</b>	<b>420</b>	<b>306</b>

Table 6: Evaluation of medium-sized models on non-English data.

when it comes to modeling morphologically rich languages. Also, it is noteworthy that the higher the type-token ratio, the bigger the advantage of SCRN over LSTM. We defer the detailed study of this phenomenon to our future work.

## 7 Conclusion

Being originally implemented in `Torch`, the SCRN model is fully reproducible in `Tensorflow` despite the difference in styles of truncated BPTT in these two libraries. Being conceptually much simpler,

<sup>6</sup><http://www.statmt.org/wmt13/translation-task.html>

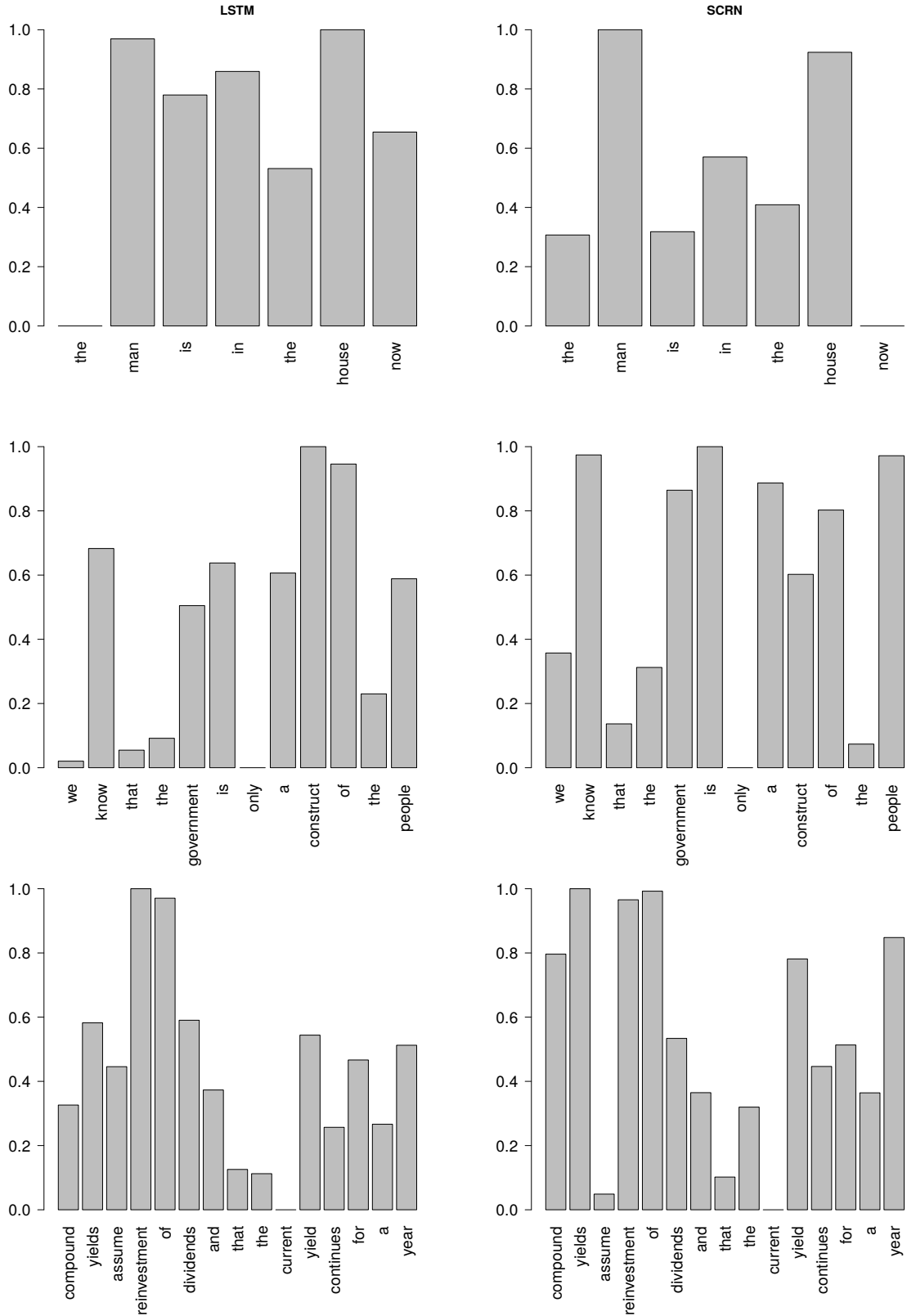


Figure 2:  $L_1$  norm of deltas between consecutive states of LSTM (left) and SCRN (right) trained on Penn Treebank plotted over words of example sentences. The main observation is that the norm is in general lower for low-information content words, such as the article *the*, and higher for informative words, such as *government*, and this difference is more pronounced in SCRN.



SCRN architecture demonstrates performance comparable to the widely used LSTM model in language modeling task under the existing regularization techniques. However, it is interesting to see if the same holds true under the most recently proposed regularization and optimization techniques for recurrent neural network language models, such as averaged stochastic gradient descent, DropConnect (Merity et al., 2018), mixture of softmaxes (Yang et al., 2018). We defer such study to our future work, which will also include larger-scale language modeling experiments for different languages as well as subword-level modeling. In addition, we intend to explore how useful the SCRNN might be in other tasks that the architectures such as the LSTM currently hold state-of-the-art performance in.

## References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283.
- Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2001. A neural probabilistic language model.
- Jan Botha and Phil Blunsom. 2014. Compositional morphology for word representations and language modelling. In *Proc. of ICML*.
- 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*.
- Jeffrey L Elman. 1990. Finding structure in time. *Cognitive science*, 14(2):179–211.
- Yarin Gal and Zoubin Ghahramani. 2016. A theoretically grounded application of dropout in recurrent neural networks. In *Proc. of NIPS*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Hakan Inan, Khashayar Khosravi, and Richard Socher. 2017. Tying word vectors and word classifiers: A loss framework for language modeling. In *Proc. of ICLR*.
- Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. 2015. An empirical exploration of recurrent network architectures. In *Proc. of ICML*, pages 2342–2350.
- Andrej Karpathy, Justin Johnson, and Li Fei-Fei. 2016. Visualizing and understanding recurrent networks.
- Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. 2016. Character-aware neural language models. In *Proc. of AAAI*, pages 2741–2749.
- Thomas Laurent and James von Brecht. 2017. A recurrent neural network without chaos. In *Proc. of ICLR*.
- Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. In *Proc. of ICLR*.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018. Regularizing and optimizing lstm language models. In *Proc. of ICLR*.
- Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Proc. of INTERSPEECH*.
- Tomas Mikolov, Armand Joulin, Sumit Chopra, Michael Mathieu, and Marc’Aurelio Ranzato. 2015. Learning longer memory in recurrent neural networks. In *Proc. of ICLR Workshop Track*.
- Andriy Mnih and Geoffrey Hinton. 2007. Three new graphical models for statistical language modelling. In *Proc. of ICML*.

- Alexander G Ororbia II, Tomas Mikolov, and David Reitter. 2017. Learning simpler language models with the differential state framework. *Neural computation*, 29(12):3327–3352.
- Ofir Press and Lior Wolf. 2017. Using the output embedding to improve language models. In *Proc. of EACL*.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958.
- Paul J Werbos. 1990. Backpropagation through time: what it does and how to do it. *Proc. of the IEEE*, 78(10):1550–1560.
- Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W Cohen. 2018. Breaking the softmax bottleneck: a high-rank rnn language model. In *Proc. of ICLR*.
- Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. 2014. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*.