

A Parameter Initialization Derivation

Following the derivation of Glorot and Kaiming initialization (Glorot and Bengio, 2010; He et al., 2015), we assume the values of each input vector \mathbf{x}_t are i.i.d with zero mean and a small variance:

$$E[x_{t,i}] = 0, \quad \text{Var}[x_{t,i}] \ll 1 \quad \forall 1 \leq i \leq d .$$

We initialize each weight matrix with zero mean and a variance of $1/d$. After a matrix multiplication $\mathbf{y} = \mathbf{W}\mathbf{x}_t$, each value y_i would have

$$E[y_i] = E\left[\sum_j w_{i,j}x_{t,j}\right] = 0$$

$$\text{Var}[y_i] = \sum_j \text{Var}[w_{i,j}] \cdot \text{Var}[x_{t,j}] = \text{Var}[x] ,$$

which means the scale of the values after matrix multiplication remains the same.

A.1 Computing $\text{Var}[\mathbf{c}_t]$

Let $f_{t,i}$ be the i -th entry of the forget gate \mathbf{f}_t :

$$f_{t,i} = \sigma(\mathbf{w}_{f,i}^\top \mathbf{x}_t + v_{f,i} c_{t-1,i} + b_{f,i}) .$$

The pre-activation value will be sufficiently close to 0 because the parameters are initialized with zero mean and small variance and the bias value is initially 0. As a result,

$$E[f_{t,i}] = \sigma(0) = 0.5 .$$

The state value $c_{t,i}$ is computed according to

$$c_{t,i} = f_{t,i} \cdot c_{t-1,i} + (1 - f_{t,i}) \cdot (\mathbf{w}_i^\top \mathbf{x}_t) .$$

Substituting the expectation of $f_{t,i}$ in, we get:⁶

$$c_{t,i} = \mathbf{w}_i^\top \left(\frac{\mathbf{x}_t}{2} + \frac{\mathbf{x}_{t-1}}{4} + \frac{\mathbf{x}_{t-2}}{8} + \dots \right) .$$

Therefore, $E[c_{t,i}] = 0$ as $E[\mathbf{w}_i^\top \mathbf{x}] = 0$. The variance of $c_{t,i}$ however depends on the correlation between input vectors. When the input vectors are independent:

$$\text{Var}[c_{t,i}] = \text{Var}[\mathbf{w}_i^\top \mathbf{x}] \left(\frac{1}{2^2} + \frac{1}{4^2} + \frac{1}{8^2} + \dots \right)$$

$$\approx \text{Var}[\mathbf{w}_i^\top \mathbf{x}] \cdot \frac{1}{3} = \text{Var}[x]/3 .$$

However, the two vectors in the input sequence, for instance \mathbf{x}_t and \mathbf{x}'_t , are not necessarily independent, for example because two words in an input

⁶We are ignoring the correlation between $\mathbf{f}_{t,i}$ and $\mathbf{f}_{t,i'}$ here because their variance is small.

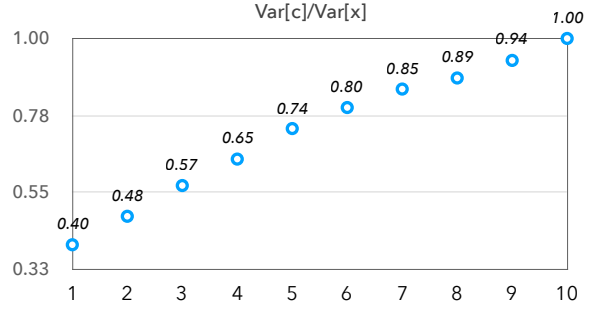


Figure 6: Empirical estimation of the variance ratio $\text{Var}[\mathbf{c}_t]/\text{Var}[\mathbf{x}_t]$ at each layer in a randomly initialized SRU model. We use the pre-trained word2vec embeddings as input, resulting an initial ratio slightly higher than $1/3$. As expected, the ratio increases to 1 in deep layers.

sentence are often correlated. When the input vectors are perfectly correlated $\mathbf{x}_t = \mathbf{x}_{t-1} = \dots = \mathbf{x}$, on the other hand,

$$\text{Var}[c_{t,i}] = \text{Var}[\mathbf{w}_i^\top \mathbf{x}] = \text{Var}[x] .$$

In practice, multiple SRU layers are stacked to construct a deep network. The internal state \mathbf{c}_t and \mathbf{h}_t would be a weighted combination of inputs $\{\mathbf{x}_1 \dots \mathbf{x}_t\}$, which will increase the correlation of the state vectors at different steps. These state vectors are again fed into the next layer, and keep increasing the correlation. As a result, we expect the actual ratio between the variance of \mathbf{c}_t and that of the input of the current layer \mathbf{x}_t lies between the two derived values,

$$\frac{1}{3} \leq \frac{\text{Var}[c]}{\text{Var}[x]} \leq 1 , \quad (5)$$

and would finally converge to the upper bound value of 1. Figure 6 confirms our expectation by computing the empirical value of $\text{Var}[c]/\text{Var}[x]$ in deep SRU networks.

A.2 Computing $\text{Var}[\mathbf{h}_t]$

Given the result in Equation (5), we proceed to compute $\text{Var}[\mathbf{h}_t]$. The i -th entry of \mathbf{h}_t is similarly computed as

$$h_{t,i} = r_{t,i} \cdot c_{t,i} + (1 - r_{t,i}) \cdot x_{t,i}$$

$$\text{where } r_{t,i} = \sigma(\mathbf{w}_{r,i}^\top \mathbf{x}_t + v_{r,i} c_{t-1,i} + b_{r,i}) .$$

The highway reset gate is not necessarily initialized with a zero bias. Let the initial bias be b and $u = \mathbf{w}_{r,i}^\top \mathbf{x}_t + v_{r,i} c_{t-1,i}$ denote the rest of terms

in the sigmoid function. We have $E[u] = 0$ and $\text{Var}[u] \ll 1$ because \mathbf{x}_t and \mathbf{c}_{t-1} have small variance.

We approximate the value of $r_{t,i}$ using its Taylor expansion at $u = 0$:

$$\begin{aligned} r_{t,i} &= \sigma(u + b) \\ &\approx \frac{e^b}{e^b + 1} + \frac{e^b \cdot u}{(e^b + 1)^2} \\ E[r_{t,i}^2] &\approx \frac{e^{2b}}{(e^b + 1)^2} + \frac{e^{2b} \cdot u^2}{(e^b + 1)^4} . \end{aligned}$$

We can ignore the term with u^2 since $\text{Var}[u] \ll 1$, which gives us

$$E[r_{t,i}^2] \approx \frac{e^{2b}}{(e^b + 1)^2} .$$

Substituting this result in $\text{Var}[h_{t,i}]$,

$$\begin{aligned} \text{Var}[h_{t,i}] &= E[r_{t,i}^2 \mathbf{c}_{t,i}^2 + (1 - r_{t,i})^2 x_{t,i}^2] \\ &= \frac{e^{2b} \cdot \text{Var}[c]}{(e^b + 1)^2} + \frac{\text{Var}[x]}{(e^b + 1)^2} \quad (6) \end{aligned}$$

Since from (5) we have $\text{Var}[x]/3 \leq \text{Var}[c] \leq \text{Var}[x]$, we get the bound of $\text{Var}[h_{t,i}]$

$$\frac{e^{2b} + 3}{3(e^b + 1)^2} \leq \frac{\text{Var}[h]}{\text{Var}[x]} \leq \frac{e^{2b} + 1}{(e^b + 1)^2}$$

which is equivalent to

$$\frac{1}{3} \leq \frac{\text{Var}[h]}{\text{Var}[x]} \leq \frac{1}{2}$$

when $b = 0$.

A.3 Computing the Scaling Constant α

Finally, we compute the scaling constant α (Section 3.2). Using the result in Equation (6), when α is introduced we get:

$$\begin{aligned} \text{Var}[h_{t,i}] &= \frac{e^{2b} \cdot \text{Var}[c]}{(e^b + 1)^2} + \frac{\alpha^2 \cdot \text{Var}[x]}{(e^b + 1)^2} \\ &\approx \frac{e^{2b} + \alpha^2}{(e^b + 1)^2} \cdot \text{Var}[x] , \end{aligned}$$

as $\text{Var}[c] \rightarrow \text{Var}[x]$ according to Equation (5) and the empirical evaluation (Figure 6). This implies $e^{2b} + \alpha = (1 + e^b)^2$ if we want $\text{Var}[h] \approx \text{Var}[x]$. By solving for α we have

$$\alpha = \sqrt{1 + 2 \cdot e^b} ,$$

and $\alpha = \sqrt{3}$ when $b = 0$.

B Experimental Details

We include additional experimental setup and results in this section.

B.1 Classification

The data and pre-processing code are obtained from the code repository of Harvard NLP.⁷

We use a batch size of 32 and a dropout probability of 0.5 for all models. In addition, we increment the dropout to 0.55 or 0.6 for the 8-layer SRU model. Following the implementation of (Kim, 2014), out-of-vocabulary words that are not in the pre-trained embeddings are initialized with random vectors with values from $[-0.25, 0.25]$.

B.2 Question Answering

We use a word embedding dropout of 0.5 and a recurrent dropout of 0.2. In the setup of Chen et al. (2017a), the bi-LSTM models concatenates the output of each layer and feed it to subsequent layers. This helps the gradient propagation and improves the final performance. With highway connection, this is no longer necessary. In SRU and Q-RNN (with highway), only the output of the last layer is given to subsequent layers.

B.3 Machine Translation

We use the OpenNMT PyTorch implementation for the translation experiments. Table 6 shows the list of configuration options used for training. For evaluation, we use beam size 5 and length penalty 0.6.

-layers	4 to 6	-share_embedding
-rnn_size	512	-position_encoding
-word_vec_size	512	-param_init 0
-batch_type	tokens	-max_grad_norm 0
-normalization	tokens	-dropout 0.1
-batch_size	5120	-label_smoothing 0.1
-accum_count	5	-epoch 40
-optim	adam	-param_init_glorot
-learning_rate	2	
-adam_beta2	0.998	
-decay_method	noam	
-warmup_steps	16000	

Table 6: Translation training configuration.

⁷<https://github.com/harvardnlp/sent-conv-torch>

Epoch	Transformer base		w/ SRU (4 layer)		w/ SRU (5 layer)	
	Valid	Test	Valid	Test	Valid	Test
20	26.1	27.3	26.2	27.6	26.6	27.9
21	26.2	27.3	26.3	27.7	26.6	28.1
22	26.1	27.4	26.3	27.8	26.7	28.0
23	26.2	27.4	26.4	27.7	26.8	28.1
24	26.2	27.4	26.4	27.8	26.7	28.0
25	26.3	27.4	26.4	27.7	26.6	28.1
26	26.5	27.5	26.5	27.7	26.7	28.1
27	26.4	27.6	26.4	27.6	26.8	28.1
28	26.4	27.6	26.4	27.7	26.7	28.2
29	26.4	27.5	26.4	27.8	26.8	28.2
30	26.5	27.7	26.4	27.8	26.9	28.1
31	26.4	27.6	26.6	27.7	26.9	28.3
32	26.5	27.5	26.5	27.8	26.9	28.3
33	26.5	27.5	26.5	27.8	27.1	28.3
34	26.4	27.6	26.5	27.9	26.9	28.2
35	26.4	27.6	26.5	27.9	26.9	28.2
36	26.5	27.6	26.5	27.8	26.9	28.3
37	26.5	27.5	26.5	27.8	26.9	28.2
38	26.5	27.6	26.5	28.0	27.0	28.2
39	26.5	27.6	26.7	27.8	27.0	28.2
40	26.6	27.6	26.6	27.9	27.0	28.2

Table 7: Average BLEU scores after each epoch.

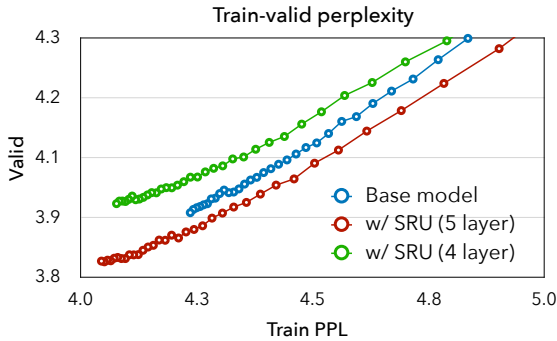


Figure 7: Training and validation perplexity curves of the base model and two SRU models.

Table 7 shows the averaged BLEU score of each model from 20th to 40th epoch. The improvement over the Transformer base model is consistent across different epochs.

Figure 7 plots the training and validation perplexity of three models. With a higher dropout (0.2) used for the SRU, the 5-layer model gets consistent lower validation perplexity over the base model and the 4-layer model. We also see that models with SRU exhibit much faster training progress with much lower training perplexity, suggesting the models could be tuned better with further training regularization.

B.4 Character-level Language Modeling

We train all models using a weight decay of 10^{-7} and a gradient clipping of 0.3. We set the learning rate factor of Noam scheduling to 3 and the warmup steps to 32,000 (8,000 for models without projection). We tune the dropout probability from $\{0.2, 0.3\}$.

The projection (bottleneck) trick is implemented as follows. Recall that the batched multiplication of SRU is computed as

$$\begin{pmatrix} \mathbf{W} \\ \mathbf{W}_f \\ \mathbf{W}_r \end{pmatrix} [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L] .$$

The stacked parameter matrices on the left is reparameterized by a low-rank factorization,

$$\begin{pmatrix} \mathbf{W} \\ \mathbf{W}_f \\ \mathbf{W}_r \end{pmatrix} = \mathbf{P}^\top \mathbf{Q} ,$$

where $\mathbf{Q} \in \mathbb{R}^{d_{in} \times d'}$ and $\mathbf{P} \in \mathbb{R}^{3d_{out} \times d'}$ are two new parameter matrices to be learned, and d' is the projection dimension that is much smaller than the input and output dimension of the SRU.