

Fast Training and Model Compression of Gated RNNs via Singular Value Decomposition

Rui Dai

Dept. Industrial Engineering
Tsinghua University
Beijing, China
dair17@mails.tsinghua.edu.cn

Lefei Li

Dept. Industrial Engineering
Tsinghua University
Beijing, China
lilefei@tsinghua.edu.cn

Wenjian Yu

BNRist, Dept. Computer Science & Tech.
Tsinghua University
Beijing, China
yu-wj@tsinghua.edu.cn

Abstract—Long Short-Term Memory (LSTM) network and Gated Recurrent Units (GRU) network are two widely-used gated Recurrent Neural Network (RNN) architectures. Both of them usually have a huge model size and require a long time to be trained. In this paper, we first propose a singular value decomposition (SVD) based approach for fast training of LSTM. Then, the factorized model and SVD based training approach are proposed for the GRU network, which adaptively choose the rank parameter for the matrix factorization model and reduce the training time and parameters of the gated RNNs. Experiments are carried out on the image classification and sentiment classification tasks using datasets MNIST and IMDB, respectively. The results show that the proposed LSTM-SVD approach achieves up to 3.9X speedup compared with training the original LSTM model, without loss of accuracy. The approaches for training the GRU network also have about 2X speedup. And, with the factorized models the quantity of RNN cell parameters can be significantly reduced by more than 10X.

Index Terms—Gated recurrent neural networks (RNNs), model compression, singular value decomposition, training acceleration

I. INTRODUCTION

As one of the most powerful deep neural networks (DNNs), the gated recurrent neural network (RNN), like LSTM [1] or GRU network [2], has shown promising results in many machine learning tasks, including language modeling [3], machine translation [4] and speech recognition [5]. The gated RNN usually involves millions of parameters and a large quantity of matrix multiplications both for the training and inference stages, so that its training and deployment require a lot of computation resource. To accelerate the training of large-scale DNNs and to compress the relevant model size, the approaches using GPU-based computing [6] and distributed computation [7] have been proposed. However, such high-performance computing equipments are costly, and are not always available. In fact, the hardware-based acceleration is not feasible for training the models used in intelligent household electrical appliance, mobile phones or other even smaller embedded devices. And, compression of the trained RNN model becomes crucial for the realization of edge computing.

Another line of work aims at reducing the amount of computation associated with the forward-propagation and the back-propagation process. Since floating-point arithmetic costs most time of training the deep-learning model, some work

were focused on replacing floating-point multiplications by binary shifts [8] or integer shifts [9]. Lin et al. converted the multiplications to sign changes during forward propagation and performed quantized back-propagation that converts the remaining multiplication into binary shifts [10]. In [11] and [12], Han et al. proposed techniques including pruning the unimportant connections and then retraining the remaining sparse network, trained quantization and Huffman coding to accelerate the training procedure of the convolution neural network (CNN) and to reduce the storage of model. Their method also improves the energy efficiency and does not degrade the accuracy. Girshick applied truncated SVD to fully connected layers in region-based convolutional network (R-CNN) for faster detection [13]. Spring and Shrivastava proposed a randomized hashing approach to reduce the amount of computation in the training and testing procedures drastically [14].

Besides the approaches mentioned above, low-rank matrix approximation is another key that can be possibly used for the training acceleration and model compression. Denil et al. proposed a matrix-factorization based method to reduce the number of parameters in deep learning models by training only a small number of weights and predicting the rest [15]. Sindhwani et al. used structured matrix transformations with low-rank matrices to accelerate inferences and forward/backward pass [16]. Singular value decomposition (SVD) of a matrix produces the optimal low-rank approximation, and this has also been utilized for the training acceleration. In the fast learning model proposed by Cai et al. [17], SVD is applied to the multi-layer back-propagation (BP) network which has been pre-adjusted by supervised training, and then the analogous back-propagation is used for subsequent training of the factorized network. Finally, the factorized model is converted inversely to the BP network. It should be pointed out that there is few work for the fast training of gated RNNs. Most of existing work is focused on the BP network or CNNs.

For the gated RNNs, recently the training efficiency of LSTM networks was improved by leveraging factorization tricks [18]. One of the tricks is based on matrix factorization and is named F-LSTM (Factorized LSTM). In the approach, two low-rank matrices are initialized randomly to represent the weight matrix in LSTM network. This reduces the number

of parameters and results in the acceleration of the training procedure.

In this paper, a new training acceleration approach is proposed for the LSTM network, which combines the matrix factorization and the SVD trick used in [17]. It adaptively chooses the rank parameter according to the accuracy criterion of matrix approximation, and makes faster convergence of LSTM training than the approach in [18]. The matrix factorization based techniques are also proposed to accelerate the training of GRU networks. Experimental results on MNIST and IMDB datasets show that the proposed technique is superior to F-LSTM, and significantly reduces the training time of the original LSTM model (up to 3.9X). For the GRU network, the training time can also be reduced by about 2X. The proposed techniques also significantly compress the RNN model, with more than 10X parameter reduction shown in experiments.

II. TECHNICAL BACKGROUND

In this section, we briefly introduce the computation carried out for training LSTM network and GRU network, and the F-LSTM approach. $\sigma(\cdot)$ and $\tanh(\cdot)$ denote the sigmoid function and hyperbolic tangent function, respectively. Both work for vector. And, “ \odot ” stands for the operator of element-wise multiplication. The structures of an LSTM cell and a GRU cell are illustrated in Figure 1(a) and 1(b) respectively. An RNN classifier is illustrated in Figure 1(c), where the hidden state of the last time step is connected to a dense layer and softmax is used as the activation function.

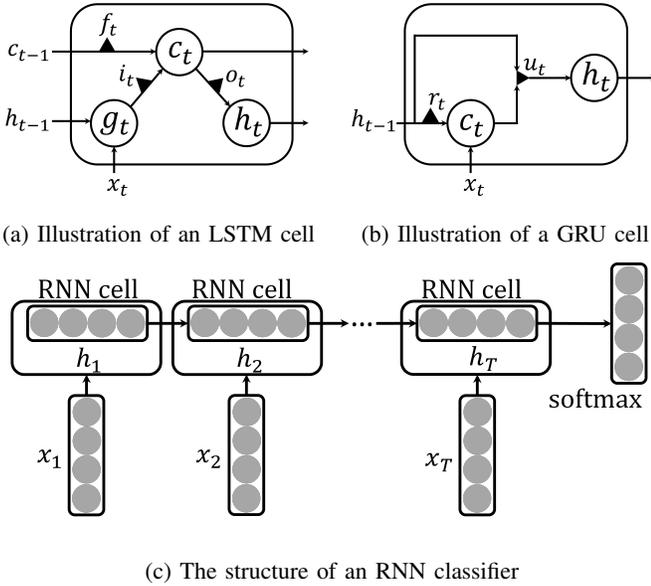


Fig. 1: Illustration of RNN cells and an RNN classifier.

A. Long Short-Term Memory Network

The computations associated with the LSTM cell can be written in the following recurrent form:

$$LSTM : \mathbf{h}_{t-1}, \mathbf{c}_{t-1}, \mathbf{x}_t \rightarrow \mathbf{h}_t, \mathbf{c}_t, \quad (1)$$

where $\mathbf{x}_t \in \mathbb{R}^{n_i}$ is the input at time step t ($t = 1, 2, \dots, T$), $\mathbf{h}_t \in \mathbb{R}^{n_h}$ is the corresponding hidden state and $\mathbf{c}_t \in \mathbb{R}^{n_h}$ is the corresponding memory. The states of the cell's gates before non-linearity activation at time step t are computed:

$$\mathbf{T}_t = \begin{bmatrix} \bar{\mathbf{i}}_t \\ \bar{\mathbf{f}}_t \\ \bar{\mathbf{o}}_t \\ \bar{\mathbf{g}}_t \end{bmatrix} = \mathbf{W} \mathbf{x}'_t + \mathbf{b} = \mathbf{W} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} + \mathbf{b}, \quad (2)$$

where $\bar{\mathbf{i}}_t, \bar{\mathbf{f}}_t, \bar{\mathbf{o}}_t, \bar{\mathbf{g}}_t \in \mathbb{R}^{n_h}$, $\mathbf{W} \in \mathbb{R}^{4n_h \times (n_i + n_h)}$, and \mathbf{b} , $\mathbf{T}_t \in \mathbb{R}^{4n_h}$. Here n_i and n_h denote the dimensions of input sequences and hidden states, respectively. Then, at time step t the states of the input gate, forget gate, output gate and the new memory (denoted by \mathbf{i}_t , \mathbf{f}_t , \mathbf{o}_t and \mathbf{g}_t respectively) are computed:

$$\mathbf{i}_t = \sigma(\bar{\mathbf{i}}_t), \mathbf{f}_t = \sigma(\bar{\mathbf{f}}_t), \mathbf{o}_t = \sigma(\bar{\mathbf{o}}_t), \mathbf{g}_t = \sigma(\bar{\mathbf{g}}_t). \quad (3)$$

The final memory \mathbf{c}_t and the hidden state \mathbf{h}_t are computed as:

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t, \quad (4)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \quad (5)$$

B. Gated Recurrent Unit Network

Similar to the LSTM cell, the GRU cell has the following recurrent form:

$$GRU : \mathbf{h}_{t-1}, \mathbf{x}_t \rightarrow \mathbf{h}_t. \quad (6)$$

The cell's gates before nonlinearity activation at time step t are computed:

$$\mathbf{T}_t^g = \begin{bmatrix} \bar{\mathbf{r}}_t \\ \bar{\mathbf{u}}_t \end{bmatrix} = \mathbf{W}_g \mathbf{x}'_t + \mathbf{b}_g = \mathbf{W}_g \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} + \mathbf{b}_g, \quad (7)$$

where $\bar{\mathbf{r}}_t, \bar{\mathbf{u}}_t \in \mathbb{R}^{n_h}$, $\mathbf{x}'_t \in \mathbb{R}^{n_i + n_h}$, $\mathbf{b}_g \in \mathbb{R}^{2n_h}$, and $\mathbf{W}_g \in \mathbb{R}^{2n_h \times (n_i + n_h)}$. The state of reset gate $\mathbf{r}_t = \sigma(\bar{\mathbf{r}}_t)$ and the state of update gate $\mathbf{u}_t = \sigma(\bar{\mathbf{u}}_t)$. The new memory is computed as:

$$\mathbf{T}_t^c = \mathbf{W}_c \mathbf{x}'_t + \mathbf{b}_c = \mathbf{W}_c \begin{bmatrix} \mathbf{x}_t \\ \mathbf{r} \odot \mathbf{h}_{t-1} \end{bmatrix} + \mathbf{b}_c, \quad (8)$$

$$\mathbf{c}_t = \tanh(\mathbf{T}_t^c), \quad (9)$$

where $\mathbf{x}'_t \in \mathbb{R}^{n_i + n_h}$, $\mathbf{W}_c \in \mathbb{R}^{n_h \times (n_i + n_h)}$, and $\mathbf{b}_c \in \mathbb{R}^{n_h}$. Finally, the hidden state is computed as

$$\mathbf{h}_t = (\mathbf{1} - \mathbf{u}_t) \odot \mathbf{c}_t + \mathbf{u}_t \odot \mathbf{h}_{t-1}, \quad (10)$$

where $\mathbf{1}$ denotes the vector with all 1 elements.

C. The F-LSTM Approach

To speed up the training of LSTM, Kuchaiev and Ginsburg proposed F-LSTM [18]. Its idea is to approximate matrix \mathbf{W} in Eq. (2) with the product of two smaller matrices, i.e., $\mathbf{W} = \mathbf{W}_2 \cdot \mathbf{W}_1$. This makes the number of parameters reduced from $4n_h \times (n_i + n_h)$ to $(5n_h + n_i)r$, where r is the rank parameter. In F-LSTM, \mathbf{W}_1 and \mathbf{W}_2 are initialized randomly before training.

The F-LSTM approach has the following shortages. Firstly, as \mathbf{W}_1 and \mathbf{W}_2 are initialized randomly, the training procedure may converge slowly compared with other approach

which assigns better initial values to \mathbf{W}_1 and \mathbf{W}_2 . Secondly, it is not clear how to choose the parameter r . In this work, we propose a new fast training approach to overcome the above shortages. We will also extend the work to accelerate the training of GRU networks.

III. METHODOLOGY

Three training approaches for RNNs are proposed here: LSTM-SVD, Factorized GRU (F-GRU) and GRU-SVD. All of them can achieve a remarkable speedup and model-size compression without loss of performance.

A. The SVD Based Approach for Training LSTM

The SVD method [19] is one of the most powerful and widely-used tools for matrix factorization. For a weight matrix $\mathbf{W}_{(m \times n)}$ in an RNN model, the rank- r truncated SVD can be obtained by performing the full SVD and then keeping the r largest singular values and corresponding singular vectors.

$$\begin{aligned} \mathbf{W} &= \mathbf{U}_{(m \times m)} \cdot \boldsymbol{\Sigma}_{(m \times n)} \cdot \mathbf{V}_{(n \times n)}^T \\ &\approx \mathbf{U}_{(m \times r)} \cdot \left(\boldsymbol{\Sigma}_{(r \times r)} \cdot \mathbf{V}_{(r \times n)}^T \right) \\ &= \mathbf{W}_2 \cdot \mathbf{W}_1 = \mathbf{W}_f. \end{aligned} \quad (11)$$

From Eq. (11), we see that the truncated SVD can be expressed as the matrix factorization in form of $\mathbf{W}_2 \cdot \mathbf{W}_1$. And, the result matrix \mathbf{W}_f best approximates \mathbf{W} among all matrices with rank r [20]. Suppose the singular values of \mathbf{W} are $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. The approximation error is

$$\frac{\|\mathbf{W} - \mathbf{W}_f\|_2}{\|\mathbf{W}\|_2} = \frac{\sigma_{r+1}}{\sigma_1}. \quad (12)$$

So, if we want to constrain the relative error of the approximation within ε , i.e.,

$$e_r = \frac{\|\mathbf{W} - \mathbf{W}_f\|_2}{\|\mathbf{W}\|_2} \leq \varepsilon, \quad (13)$$

the rank parameter r should be:

$$r^* = \min_{\sigma_{r+1} \leq \sigma_1 \varepsilon} r. \quad (14)$$

This gives an adaptive strategy for choosing r value according to accuracy criterion.

Replacing \mathbf{W} with $\mathbf{W}_2 \cdot \mathbf{W}_1$, we obtain the factorized model of LSTM. To overcome the drawback brought by randomly initializing \mathbf{W}_1 and \mathbf{W}_2 , we instead pre-train the original model for several iterations, and then make truncated SVD on \mathbf{W} to obtain the initial values of \mathbf{W}_1 and \mathbf{W}_2 . After that the factorized model is trained with BPTT algorithm. While doing the truncated SVD, Eq. (14) is used to choose the parameter r for the matrix factorization. This generates an SVD based approach for training LSTM (named LSTM-SVD), which includes the following steps.

- Step 1: Pre-train the LSTM network for several iterations, e.g., for 1 epoch;
- Step 2: Apply SVD to matrix \mathbf{W} obtained from the pre-trained network, and then choose parameter r according to an accuracy ε and Eq. (14);

- Step 3: Train the factorized model with \mathbf{W}_1 and \mathbf{W}_2 , and finally obtain the converged factorized LSTM model.

For the factorized models, BPTT algorithm is still used for training. For LSTMs, the following can be inferred according to the chain rule:

$$\begin{aligned} \nabla_{\mathbf{T}_t} L &= \begin{bmatrix} \nabla_{\mathbf{i}_t} L \\ \nabla_{\mathbf{f}_t} L \\ \nabla_{\mathbf{o}_t} L \\ \nabla_{\mathbf{g}_t} L \end{bmatrix} \\ &= \begin{bmatrix} \nabla_{\mathbf{c}_t} L \odot \mathbf{g}_t \odot (\mathbf{1} - \mathbf{i}_t) \odot \mathbf{i}_t \\ \nabla_{\mathbf{c}_t} L \odot \mathbf{c}_{t-1} \odot (\mathbf{1} - \mathbf{f}_t) \odot \mathbf{f}_t \\ \nabla_{\mathbf{h}_t} L \odot \tanh(\mathbf{c}_t) \odot (\mathbf{1} - \mathbf{o}_t) \odot \mathbf{o}_t \\ \nabla_{\mathbf{c}_t} L \odot \mathbf{i}_t \odot (\mathbf{1} - \mathbf{g}_t^2) \end{bmatrix}, \end{aligned} \quad (15)$$

where L is the cost function, and $\nabla_{\mathbf{v}} L$ stands for the gradient of L with respect to \mathbf{v} , i.e., $\nabla_{\mathbf{v}} L = \left[\frac{\partial L}{\partial v_1}, \frac{\partial L}{\partial v_2}, \dots, \frac{\partial L}{\partial v_n} \right]^T$.

For convenience we denote:

$$\mathbf{L}'_t = \begin{cases} \sum_{s=t}^T \nabla_{\mathbf{h}_s} L \odot \nabla_{\mathbf{c}_t} \mathbf{h}_s & , 1 \leq t \leq T \\ \mathbf{0} & , t = T + 1 \end{cases}. \quad (16)$$

Here $\mathbf{1}$ and $\mathbf{0}$ stand for the vectors with all 1 and 0 elements, respectively. The BPTT algorithm for training the factorized LSTM model is illustrated in Algorithm 1. With the converged \mathbf{W}_1 and \mathbf{W}_2 , we do not convert them back to \mathbf{W} . And, they are used in the testing stage with a slight modification on the forward-propagation algorithm.

Algorithm 1 The BPTT Algorithm for the factorized model of LSTM

Input: $\mathbf{x}_t \in \mathbb{R}^{n_i}, \mathbf{i}_t, \mathbf{f}_t, \mathbf{o}_t, \mathbf{g}_t, \mathbf{c}_t, \mathbf{h}_t \in \mathbb{R}^{n_h}, (t = 1, \dots, T);$
 $\mathbf{W}_1 \in \mathbb{R}^{r \times (n_i + n_h)}, \mathbf{W}_2 \in \mathbb{R}^{4n_h \times r}, \mathbf{b} \in \mathbb{R}^{4n_h}; \eta \in \mathbb{R}$

Output: $\Delta \mathbf{W}_1 \in \mathbb{R}^{r \times (n_i + n_h)}, \Delta \mathbf{W}_2 \in \mathbb{R}^{4n_h \times r}, \Delta \mathbf{b} \in \mathbb{R}^{4n_h}$

- 1: $lw1 \leftarrow \mathbf{0}^{r \times (n_i + n_h)}, lw2 \leftarrow \mathbf{0}^{4n_h \times r}, lb \leftarrow \mathbf{0}^{4n_h}$ $\{lw1$ denotes $\nabla_{\mathbf{W}_1} L$, $lw2$ denotes $\nabla_{\mathbf{W}_2} L$, and lb denotes $\nabla_{\mathbf{b}} L\}$
 - 2: **for** $t \leftarrow T$ to 1 **do**
 - 3: $\nabla_{\mathbf{c}_t} L \leftarrow \nabla_{\mathbf{h}_t} L \odot \mathbf{o}_t \odot (\mathbf{1} - \tanh^2(\mathbf{c}_t)) + \mathbf{L}'_{t+1}$
 - 4: calculate $\nabla_{\mathbf{T}_t} L$ using Eq. (15)
 - 5: $lw1 += \mathbf{W}_2^T \cdot \nabla_{\mathbf{T}_t} L \cdot \mathbf{x}'_t{}^T$
 - 6: $lw2 += \nabla_{\mathbf{T}_t} L \cdot (\mathbf{W}_1 \mathbf{x}'_t{}^T)$
 - 7: $lb += \nabla_{\mathbf{T}_t} L$
 - 8: $\mathbf{L}'_t \leftarrow \nabla_{\mathbf{c}_t} L \odot \mathbf{f}_t$
 - 9: $\nabla_{\mathbf{h}_{t-1}} L \leftarrow ((\mathbf{W}_2 \mathbf{W}_1)^T \cdot \nabla_{\mathbf{T}_t} L) [-n_h :]$ $\{“[-n :]”$ means the last n elements of a vector. $\}$
 - 10: **end for**
 - 11: $\Delta \mathbf{W}_1 \leftarrow \eta \cdot lw1$
 - 12: $\Delta \mathbf{W}_2 \leftarrow \eta \cdot lw2$
 - 13: $\Delta \mathbf{b} \leftarrow \eta \cdot lb$
 - 14: **return** $\Delta \mathbf{W}_1, \Delta \mathbf{W}_2, \Delta \mathbf{b}$.
-

B. The Factorization-Based Approaches for GRU

For the GRU network, both the factorization approach and the SVD based approach in last subsection can also be applied. They reduce the time for training the GRU network. The

difference to training LSTM is that there are two parameter matrices in GRU network.

The first matrix is \mathbf{W}_g including $2n_h(n_i + n_h)$ parameters. So, if it is factorized as: $\mathbf{W}_g \approx \mathbf{W}_{gf} = \mathbf{W}_{g2} \cdot \mathbf{W}_{g1}$, where $\mathbf{W}_{g1} \in \mathbb{R}^{r_g \times (n_i + n_h)}$ and $\mathbf{W}_{g2} \in \mathbb{R}^{2n_h \times r_g}$. The number of parameters is cut down to $(n_i + 3n_h)r_g$, where r_g is the rank parameter of \mathbf{W}_g . Hence, for the approach with factorized GRU model (denoted by F-GRU), the calculation of \mathbf{T}_t^g in Eq. (7) changes to

$$\mathbf{T}_t^g = \mathbf{W}_{g2} \cdot (\mathbf{W}_{g1} \cdot \mathbf{x}_t^g) + \mathbf{b}_g. \quad (17)$$

The second weight matrix \mathbf{W}_c can also be approximated by a rank- r_c matrix $\mathbf{W}_{cf} = \mathbf{W}_{c2} \cdot \mathbf{W}_{c1}$, where $\mathbf{W}_{c1} \in \mathbb{R}^{r_c \times (n_i + n_h)}$ and $\mathbf{W}_{c2} \in \mathbb{R}^{n_h \times r_c}$. Then, \mathbf{T}_t^c in Eq. (8) changes to

$$\mathbf{T}_t^c = \mathbf{W}_{c2} \cdot (\mathbf{W}_{c1} \cdot \mathbf{x}_t^c) + \mathbf{b}_c. \quad (18)$$

Notice that \mathbf{W}_c has $n_h(n_i + n_h)$ parameters and the factorization can reduce this number down to $(n_i + 2n_h)r_c$. Compared with \mathbf{W}_g , this reduction seems not prominent. In practice, we can factorize both \mathbf{W}_g and \mathbf{W}_c , but we prefer only do factorization for the larger matrix \mathbf{W}_g . Parameters of the F-GRU model are randomly initialized before training and the model is also trained with the BPTT method.

An alternative approach is named GRU-SVD, which has the same idea of LSTM-SVD. We first pre-train a conventional GRU model for several iterations. Then, we apply SVD to \mathbf{W}_g or both \mathbf{W}_g and \mathbf{W}_c , and each matrix is substitute by the product of two low-rank matrices. Thus the subsequent work is the same as the F-GRU approach. Finally, we can get the trained GRU-SVD model. Our experiment shows that the large singular values of \mathbf{W}_g are more distinctive than those of \mathbf{W}_c . Therefore, a preferable strategy is applying SVD to \mathbf{W}_g only while keeping the original \mathbf{W}_c matrix in the GRU-SVD approach.

For GRUs, it's not hard to derive the following expressions:

$$\nabla_{\mathbf{u}_t} L = \nabla_{\mathbf{h}_t} L \odot (\mathbf{h}_{t-1} - \mathbf{c}_t), \quad (19)$$

$$\nabla_{\mathbf{r}_t} L = ((\mathbf{W}_{c2} \mathbf{W}_{c1})^T \cdot (\mathbf{1} - \mathbf{c}_t^2) \odot \nabla_{\mathbf{c}_t} L) [-n_h :] \odot \mathbf{h}_{t-1}, \quad (20)$$

$$\nabla_{\mathbf{T}_t^g} L = \begin{bmatrix} \nabla_{\tilde{\mathbf{u}}_t} L \\ \nabla_{\tilde{\mathbf{r}}_t} L \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{u}_t} L \odot \mathbf{u}_t \odot (\mathbf{1} - \mathbf{u}_t) \\ \nabla_{\mathbf{r}_t} L \odot \mathbf{r}_t \odot (\mathbf{1} - \mathbf{r}_t) \end{bmatrix}. \quad (21)$$

Here, “[$-n_h$:]” following Python notation, means the last n elements of a vector. The BPTT algorithm for a factorized GRU cell is illustrated in Algorithm 2.

IV. EXPERIMENTS

The proposed approaches have been tested on two tasks - MNIST classification task and IMDB sentiment classification task. Each of the two tasks corresponds to a publicly available dataset, whose statistics are summarized in Table I. The training approaches for LSTM and GRU networks are implemented based on TensorFlow [21]. In the SVD based approach, the

Algorithm 2 The BPTT Algorithm for the factorized model of GRU

Input: $\mathbf{x}_t \in \mathbb{R}^{n_i}, \mathbf{r}_t, \mathbf{u}_t, \mathbf{c}_t, \mathbf{h}_t \in \mathbb{R}^{n_h} (t = 1, \dots, T);$
 $\mathbf{W}_{g1} \in \mathbb{R}^{r_g \times (n_i + n_h)}, \mathbf{W}_{g2} \in \mathbb{R}^{2n_h \times r_g}, \mathbf{b}_g \in \mathbb{R}^{2n_h},$
 $\mathbf{W}_{c1} \in \mathbb{R}^{r_c \times (n_i + n_h)}, \mathbf{W}_{c2} \in \mathbb{R}^{n_h \times r_c}, \mathbf{b}_c \in \mathbb{R}^{n_h}; \eta \in \mathbb{R}$
Output: $\Delta \mathbf{W}_{g1} \in \mathbb{R}^{r_g \times (n_i + n_h)}, \Delta \mathbf{W}_{g2} \in \mathbb{R}^{2n_h \times r_g}, \Delta \mathbf{b}_g \in \mathbb{R}^{2n_h},$
 $\Delta \mathbf{W}_{c1} \in \mathbb{R}^{r_c \times (n_i + n_h)}, \Delta \mathbf{W}_{c2} \in \mathbb{R}^{n_h \times r_c}, \Delta \mathbf{b}_c \in \mathbb{R}^{n_h}$

- 1: $lwg1 \leftarrow \mathbf{0}^{r_g \times (n_i + n_h)}, lwg2 \leftarrow \mathbf{0}^{2n_h \times r_g}, lbg \leftarrow \mathbf{0}^{2n_h}$
- 2: $lwc1 \leftarrow \mathbf{0}^{r_c \times (n_i + n_h)}, lwc2 \leftarrow \mathbf{0}^{n_h \times r_c}, lbc \leftarrow \mathbf{0}^{n_h}$ { $lwg1, lwg2, lbg, lwc1, lwc2$ and lbc denote $\nabla_{\mathbf{W}_{g1}} L, \nabla_{\mathbf{W}_{g2}} L, \nabla_{\mathbf{b}_g} L, \nabla_{\mathbf{W}_{c1}} L, \nabla_{\mathbf{W}_{c2}} L,$ and $\nabla_{\mathbf{b}_c} L$ respectively}
- 3: **for** $t \leftarrow T$ to 1 **do**
- 4: $\nabla_{\mathbf{c}_t} L \leftarrow \nabla_{\mathbf{h}_t} L \odot (\mathbf{1} - \mathbf{u}_t)$
- 5: calculate $\nabla_{\mathbf{T}_t^g} L$ using Eq. (21)
- 6: $lwg1 += \mathbf{W}_{g2}^T \cdot \nabla_{\mathbf{T}_t^g} L \cdot (\mathbf{x}_t^g)^T$
- 7: $lwg2 += \nabla_{\mathbf{T}_t^g} L \cdot (\mathbf{W}_{g1} \mathbf{x}_t^g)^T$
- 8: $lbg += \nabla_{\mathbf{T}_t^g} L$
- 9: $lwc1 += \mathbf{W}_{c2}^T \cdot \nabla_{\mathbf{c}_t} L \odot (\mathbf{1} - \mathbf{c}_t^2) \cdot (\mathbf{x}_t^c)^T$
- 10: $lwc2 += \nabla_{\mathbf{c}_t} L \odot (\mathbf{1} - \mathbf{c}_t^2) \cdot (\mathbf{W}_{c1}^T \mathbf{x}_t^c)^T$
- 11: $lbc += \nabla_{\mathbf{c}_t} L \odot (\mathbf{1} - \mathbf{c}_t^2)$
- 12: $\nabla_{\mathbf{h}_{t-1}} L \leftarrow ((\mathbf{W}_{g2} \mathbf{W}_{g1})^T \cdot \nabla_{\mathbf{T}_t^g} L) [-n_h :]$
- 13: **end for**
- 14: $\Delta \mathbf{W}_{g1} \leftarrow \eta \cdot lwg1, \Delta \mathbf{W}_{g2} \leftarrow \eta \cdot lwg2, \Delta \mathbf{b}_g \leftarrow \eta \cdot lbg$
- 15: $\Delta \mathbf{W}_{c1} \leftarrow \eta \cdot lwc1, \Delta \mathbf{W}_{c2} \leftarrow \eta \cdot lwc2, \Delta \mathbf{b}_c \leftarrow \eta \cdot lbc$
- 16: **return** $\Delta \mathbf{W}_{g1}, \Delta \mathbf{W}_{g2}, \Delta \mathbf{b}_g, \Delta \mathbf{W}_{c1}, \Delta \mathbf{W}_{c2}, \Delta \mathbf{b}_c.$

relative error constraint ε is set to 0.2, except explicitly stated. All experiments were run on a machine equipped with a 16-core Dual-Core AMD Opteron™ Processor 8224 SE and 64 GB memory.

TABLE I: Sizes of the datasets

Dataset	Training size	Validation size	Test size
MNIST	55,000	5,000	10,000
IMDB	30,000	10,000	10,000

A. MNIST Classification Task

MNIST is a large dataset including images of handwritten digits, and each image is represented by a 28×28 pixel matrix of gray-scale values. MNIST has been studied in many image classification tasks [22]. Considering each image as a sequence of pixel rows, each image can be handled as a sequence of vectors with 28 steps. We use the cross-entropy loss as the cost function and Adam optimizer [23] with a learning rate decay is used for training. Specifically, the initial learning rate is 0.0009 which decays by a factor of 0.95 per 100 steps.

We first test the approaches using LSTM network. Each image can be handled as a sequence of vectors with $T = 28$ steps and the input size $n_i = 28$. The size of hidden state is $n_h = 768$. With the SVD based training approach (LSTM-SVD), we first pre-train the LSTM network for one epoch.

Some leading singular values of weight matrix \mathbf{W} after the pre-training are shown in Figure 2. Compared with the initial random weight matrix, there is distinct decay on the singular values of the pre-adjusted \mathbf{W} . Thus, the pre-adjusted \mathbf{W} can be well approximated by a low-rank matrix. With the preset ε and (14), we can calculate the value of rank parameter r^* . In this experiment, $r^* = 48$. Then, with the factorized LSTM model, the number of parameters in the LSTM cell is reduced from 2,448,384 to 185,664 (13X reduction), and the total number of parameters is reduced from 2,456,074 to 193,354 (**13X**).

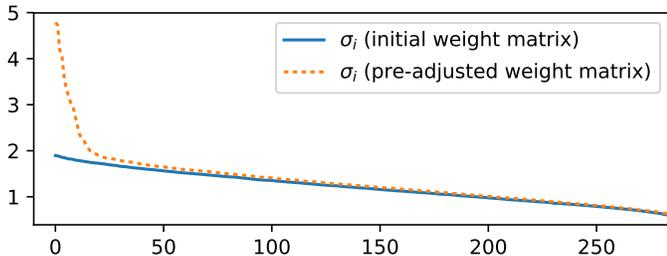


Fig. 2: Distributions of leading singular values of the initial random weight matrix \mathbf{W} and the pre-adjusted \mathbf{W} in the MNIST classification task.

With the same low-rank factorization, we then train the network with the F-LSTM approach. In Figure 3, we show the training loss of LSTM-SVD, F-LSTM and the conventional LSTM. From it we can see that at the same step count LSTM-SVD gains lower loss than F-LSTM even though they both have same number of parameters. More results of the three approaches are listed in Table II. For the training stage, we use the early stopping strategy for regularization, so that the training is stopped if the validation accuracy is not improved for 4 epochs. From the table, we see that the approaches converge with different speed. LSTM-SVD needs less epoch than F-LSTM for reaching convergence, which makes the latter's costs 13% more training time than the former. Compared with the conventional LSTM model, the LSTM-SVD approach reduces the training time for **3.9X**, while keeping same accuracy. For the testing stage, the model obtained with LSTM-SVD has better test accuracy than that corresponding to the F-LSTM approach.

TABLE II: Comparisons of LSTM, F-LSTM and LSTM-SVD approaches in the MNIST classification task.

Approach	LSTM	F-LSTM	LSTM-SVD
Training accuracy	100.0%	99.5%	100.0%
Test accuracy	99.0%	98.4%	98.8%
Training time (min.)	229.3	65.6	58.1
Epoch number	15	19	15
Testing time (sec.)	23.9	12.4	12.0

Then, we test the approaches using GRU network. The size of hidden state is still $n_h = 768$. In a GRU cell, for weight matrices \mathbf{W}_c and \mathbf{W}_g , the distributions of the singular values before and after the pre-training step are shown in Figure 4.

For each initial weight matrix, the singular values are relatively equally distributed. After the pre-training, more distinct decay of singular value is exhibited for the both matrices. Between the both, the decay trend of singular value for the pre-adjusted \mathbf{W}_g is obviously more prominent than the pre-adjusted \mathbf{W}_c . This indicates \mathbf{W}_g should be factorized preferentially.

In the experiment, we just factorize the matrix \mathbf{W}_g , and the GRU-SVD approach automatically choose a rank parameter $r_g^* = 103$. With the same low-rank factorization, the F-GRU approach is then performed. Different from LSTM-based approach, with the factorized GRU model, the training loss converges with same speed as the conventional GRU model. The experimental results also indicate that the GRU-SVD approach does not outperform the F-GRU approach on the convergence rate.

In Table III, the results of the GRU-based approaches in the MNIST classification task are listed. From it we see that, the F-GRU approach achieves the same test accuracy as the conventional GRU model, but costing only 69.7% training time. With the factorized model, the testing time can also be reduced to 69.1% of the original model. The performance of the GRU-SVD is close to that of F-GRU, showing 1.3X reduction of the training time of the conventional GRU model. The number of parameters is also reduced if the F-GRU or GRU-SVD approaches are used, as listed in Table III. It should be pointed out more speedup and model compression can be expected providing that both \mathbf{W}_g and \mathbf{W}_c are factorized with lower ranks. For example, the F-GRU approach with rank-60 factorizations for both \mathbf{W}_g and \mathbf{W}_c can achieve **2.1X** speedup while keeping high performance (98.7% test accuracy). This corresponds to a **7.6X** compression rate of the model size.

B. IMDB Sentiment Classification Task

Large Movie Review Dataset contains 50,000 highly polar movie reviews from IMDB and is widely used for sentiment classification tasks [24]. Top 20,000 words in this corpus are used to build the vocabulary dictionary. Natural Language Toolkit (NLTK) [25] is used for tokenization and sentences are truncated or padded to the length of 200.

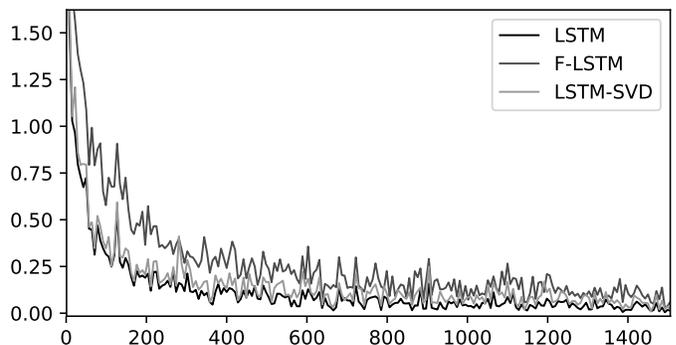


Fig. 3: Training loss versus training steps (mini-batch) for the F-LSTM, LSTM and LSTM-SVD approaches in the MNIST classification task

TABLE III: Comparisons of GRU, F-GRU and GRU-SVD in the MNIST classification task.

Approach	Training accuracy	Test accuracy	Training time (min.)	Testing time (sec.)	num of parameters
GRU	100.0%	98.9%	111.2	16.6	1,843,978
F-GRU	100.0%	98.9%	77.5	11.5	861,518
GRU-SVD	100.0%	99.0%	80.5	11.8	861,518

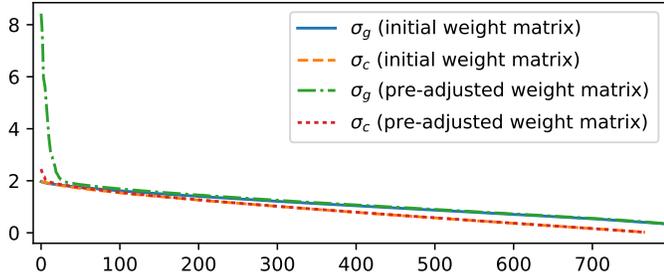


Fig. 4: Distributions of singular values in the GRU-SVD model for MNIST classification task.

In this binary classification task, whose output is either positive or negative, an embedding layer is used before the RNN cell. Dropout is used to avoid overfitting [26]. And, the binary cross-entropy is employed as the cost function. In our experiments, we set hidden size of 768 ($n_h = 768$), embedding size of 32 ($n_i = 32$) and dropout rate 0.5, while the Adam optimizer is used to train the model. The initial learning rate is 0.001, and it decays by a factor of 0.6 each epoch until reaching 0.0001.

We first test the LSTM based approaches. After pre-training the network for 1 epoch, the LSTM-SVD automatically chooses the rank parameter $r^* = 11$ (corresponding to $\varepsilon = 0.2$). Note that the singular value of \mathbf{W} here presents a more sharp drop than that for the MNIST classification task. With the same setting for low-rank factorization, F-LSTM is also used to train the network. And, the early stopping is used for regularization and the training iteration stops when there is no improvement of validation accuracy for 2 epochs. Their results are listed in Table IV, compared with the results from the conventional LSTM.

TABLE IV: Comparison of LSTM, F-LSTM and LSTM-SVD in the IMDB sentiment classification task.

Approach	LSTM	F-LSTM	LSTM-SVD
Training accuracy	97.5%	96.0%	96.0%
Test accuracy	86.2%	86.6%	87.0%
Training time (min.)	577.8	241.5	186.9
Epoch number	11	11	7
Testing time (sec.)	180.0	88.7	89.3
Parameter number	3,102,210	687,202	687,202
RNN parameter number	2,460,672	45,664	45,664

The results show that both F-LSTM and LSTM-SVD achieve better evaluation accuracy than LSTM. The reason might be that overfitting is alleviated by factorization. F-LSTM and LSTM-SVD leads to 2.9X and 3.6X reduction of training time respectively, compared with the LSTM model.

The number of parameters also reduces to 22.2% of the original LSTM model. Note that this proportion is a little larger in comparison with the MNIST task, where the size of the factorized model is only 7.9% of the original LSTM model. The reason is the existence of the embedding layer, which also has a large quantity of parameters. Focusing on the number of parameters within the RNN cell, we can see that the factorized cell has 54X fewer parameters.

Compared with the MNIST task, the pre-adjusted GRU-SVD model in the IMDB task demonstrate different distributions of singular values. Both \mathbf{W}_g and \mathbf{W}_c have a few large singular values that are very distinctive, as illustrated in Figure 5. Hence we factorize both \mathbf{W}_g and \mathbf{W}_c in F-GRU and GRU-SVD approaches. Given $\varepsilon = 0.2$, we have $r_g^* = 99$ and $r_c^* = 315$.

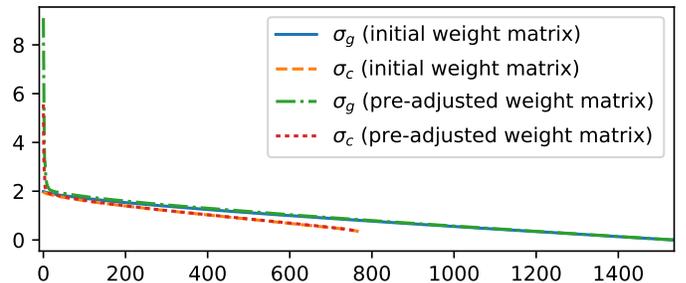


Fig. 5: Distributions of singular values in GRU-SVD model for IMDB sentiment classification task.

The computational results of the GRU-based approaches are listed in Table V. Similar to the results on the MNIST classification task, the results indicate that F-GRU performs a little bit better than GRU-SVD. The F-GRU model achieves the same test accuracy as the original GRU model, but only costs 72.0% training time. The GRU-SVD approach has comparable performance as F-GRU. If manually setting the rank parameters, we find out that F-GRU could perform even better. For example, if setting the rank parameters for \mathbf{W}_g and \mathbf{W}_c to 15 and 30 respectively, the training time can be reduced to 247 minutes (2.1X reduction if compared with the original GRU). And, the total number of parameters can be reduced to 29.2% of that in the original GRU.

V. CONCLUSIONS

The training acceleration and model compression of gated RNNs is explored in this work. With the ideas of matrix factorization and SVD, several approaches are proposed for reducing the training time of LSTM and GRU networks. Among them, LSTM-SVD automatically chooses the rank parameter and results in up to 3.9X speedup for the training of LSTM

TABLE V: Comparison of different GRU-based approaches for the IMDB sentiment classification task.

Approach	Training accuracy	Test accuracy	Training time (min.)	Testing time (sec.)	Total number of parameters	Parameters in GRU cell
GRU	99.5%	86.5%	521.5	123.5	2,487,042	1,845,504
F-GRU	98.0%	86.5%	375.6	86.3	1,369,026	727,488
GRU-SVD	97.5%	86.8%	390.4	85.6	1,369,026	727,488

network in the MINST classification and IMDB sentiment classification problems. It outperforms the existing F-LSTM approach in terms of runtime and accuracy. F-GRU and GRU-SVD approaches are for the training of GRU network and bring about 2X reduction of the training time, as shown in the experiments. With some setting of rank parameters, F-GRU performs better and achieves up to 2.1X speedup. The developed approaches induce less error, and can also compress the RNN model by up to **13X**.

An interesting phenomenon is that the obtained factorized models often show better performance on the testing set. A possible reason is that the factorization has the effect of regularization and thus alleviates overfitting. This remains a topic of future research. Besides, how to further improve the proposed technique and make it adapted to actual edge-computing scenarios will be explored in the future.

ACKNOWLEDGMENT

This work was partially supported by the National Key Research and Development Program of China (grant No. 2016YFF0204100). W. Yu is the corresponding author.

REFERENCES

- [1] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [2] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv:1412.3555*, 2014.
- [3] M. Sundermeyer, R. Schlüter, and H. Ney, "LSTM neural networks for language modeling," in *Proceedings of Thirteenth Annual Conference of the International Speech Communication Association*, 2012.
- [4] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv:1409.0473*, 2014.
- [5] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2013)*, IEEE, 2013, pp. 6645–6649.
- [6] Z. Luo, H. Liu, and X. Wu, "Artificial neural network computation on graphic process unit," in *Proceedings of IEEE International Joint Conference on Neural Networks (IJCNN'05)*, vol. 1. IEEE, 2005, pp. 622–626.
- [7] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le *et al.*, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1223–1231.
- [8] H. K. Kwan and C. Tang, "Multiplierless multilayer feedforward neural network design suitable for continuous input-output mapping," *Electronics Letters*, vol. 29, no. 14, pp. 1259–1260, 1993.
- [9] E. L. Machado, C. J. Miosso, R. von Borries, M. Coutinho, P. d. A. Berger, T. Marques, and R. P. Jacobi, "Computational cost reduction in learned transform classifications," *arXiv:1504.06779*, 2015.
- [10] Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio, "Neural networks with few multiplications," *arXiv:1510.03009*, in *Proc. ICLR 2016*, 2015.
- [11] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *Advances in Neural Information Processing Systems*, 2015, pp. 1135–1143.
- [12] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv:1510.00149*, in *Proc. ICLR 2016*, 2015.
- [13] R. Girshick, "Fast R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2015, pp. 1440–1448.
- [14] R. Spring and A. Shrivastava, "Scalable and sustainable deep learning via randomized hashing," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2017, pp. 445–454.
- [15] M. Denil, B. Shakibi, L. Dinh, N. de Freitas *et al.*, "Predicting parameters in deep learning," in *Advances in Neural Information Processing Systems*, 2013, pp. 2148–2156.
- [16] V. Sindhwani, T. Sainath, and S. Kumar, "Structured transforms for small-footprint deep learning," in *Advances in Neural Information Processing Systems*, 2015, pp. 3088–3096.
- [17] C. Cai, D. Ke, Y. Xu, and K. Su, "Fast learning of deep neural networks via singular value decomposition," in *Proceedings of Pacific Rim International Conference on Artificial Intelligence*. Springer, 2014, pp. 820–826.
- [18] O. Kuchaiev and B. Ginsburg, "Factorization tricks for LSTM networks," *arXiv:1703.10722*, in *Proc. ICLR 2017*, 2017.
- [19] V. Klema and A. Laub, "The singular value decomposition: Its computation and some applications," *IEEE Transactions on Automatic Control*, vol. 25, no. 2, pp. 164–176, 1980.
- [20] C. Eckart and G. Young, "The approximation of one matrix by another of lower rank," *Psychometrika*, vol. 1, no. 3, pp. 211–218, 1936.
- [21] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv:1603.04467*, 2016.
- [22] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [23] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv:1412.6980*, 2014.
- [24] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Portland, Oregon, USA: Association for Computational Linguistics, June 2011, pp. 142–150. [Online]. Available: <http://www.aclweb.org/anthology/P11-1015>
- [25] S. Bird, "NLTK: the natural language toolkit," in *Proceedings of the COLING/ACL on Interactive presentation sessions*. Association for Computational Linguistics, 2006, pp. 69–72.
- [26] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.