

# Energy-Efficient LSTM Networks for Online Learning

Tolga Ergen<sup>✉</sup>, Ali H. Mirza, and Suleyman Serdar Kozat, *Senior Member, IEEE*

**Abstract**—We investigate variable-length data regression in an online setting and introduce an energy-efficient regression structure build on long short-term memory (LSTM) networks. For this structure, we also introduce highly effective online training algorithms. We first provide a generic LSTM-based regression structure for variable-length input sequences. To reduce the complexity of this structure, we then replace the regular multiplication operations with an energy-efficient operator, i.e., the ef-operator. To further reduce the complexity, we apply factorizations to the weight matrices in the LSTM network so that the total number of parameters to be trained is significantly reduced. We then introduce online training algorithms based on the stochastic gradient descent (SGD) and exponentiated gradient (EG) algorithms to learn the parameters of the introduced network. Thus, we obtain highly efficient and effective online learning algorithms based on the LSTM network. Thanks to our generic approach, we also provide and simulate an energy-efficient gated recurrent unit (GRU) network in our experiments. Through an extensive set of experiments, we illustrate significant performance gains and complexity reductions achieved by the introduced algorithms with respect to the conventional methods.

**Index Terms**—ef-operator, exponentiated gradient (EG), gradient descent, long short-term memory (LSTM), matrix factorization.

## I. INTRODUCTION

### A. Preliminaries

NEURAL networks are extensively studied in the literature thanks to their highly strong modeling capabilities [1]–[4]. Especially, recurrent neural networks (RNNs) are the main source of interest in these studies due to their inherent memory unit that can store time (or state) information, which boosts their capability to model time series data [5]. However, due to lacking control structures, basic RNNs might suffer from exponential growth or decay in the norm of the gradient of their parameters during training, which are also known as the exploding and vanishing gradient problems [6], [7]. Thus, basic RNNs are not usually able to capture long- and short-term dependences present in the data [6]. To address these issues, an advanced

RNN architecture, i.e., the long short-term memory (LSTM) network, is introduced, which uses several information gates to regulate the information flow [8]. However, the LSTM networks have several additional nonlinear control structures (gates) and parameters, which result in complexity and training problems [8].

To this end, in this article, we investigate the efficient training of LSTM networks for data regression. In the literature, the LSTM networks are usually trained in a batch setting, where all data sequences are available and processed together for training [9], [10]. However, in big data applications, such approaches might cause storage problems due to the need to store all data sequences at one place [5], [10], [11]. Furthermore, in certain scenarios, we sequentially receive data instances, which prevents training in a batch setting [10]. Hence, we investigate efficient training of the LSTM network in an online setting, where we sequentially receive a data sequence with its label to train the parameters of the LSTM network and forget the data sequence after using it.

In the current literature, there exist several online training methods for the LSTM network [5], [11]–[14]. Among these methods, the first-order gradient-based algorithms are generally employed due to their computational efficiency in training the LSTM network [11], [12], [15], [16]. The first-order gradient-based training algorithms usually perform additive updates, i.e., each parameter is updated through an addition operation, e.g., the stochastic gradient descent (SGD) algorithm [5], [11], [12]. However, such algorithms suffer from slow convergence rate and poor performance, especially when only a few components of the input data are related to the desired label [17]. To circumvent these issues, the first-order training method with multiplicative updates, i.e., the exponentiated gradient (EG) algorithm, is introduced [17], [18]. However, since the EG algorithm employs multiplicative updates, it requires more computational resources compared with additive updates, which restricts its usage in real-life applications [17], [19], [20].

Recently, many applications require LSTMs to be implemented on embedded systems [21], [22]. However, since embedded devices have constraints, e.g., power, resource, and budget, it becomes either highly costly or impossible to employ LSTMs in real-life applications [21]–[23]. In particular, LSTMs have several parameters to train and require performing various arithmetic operations, which trigger its complexity issues. Among arithmetic operations, since the multiplication operation consumes more energy (or resources),

Manuscript received May 2, 2018; revised April 8, 2019; accepted August 12, 2019. Date of publication September 13, 2019; date of current version August 4, 2020. This work was supported in part by the Tubitak Project under Grant 117E153. (Corresponding author: Tolga Ergen.)

T. Ergen is with the Department of Electrical Engineering, Stanford University, Stanford, CA 94305 USA (e-mail: ergen@stanford.edu).

A. H. Mirza and S. S. Kozat are with the Department of Electrical and Electronics Engineering, Bilkent University, 06800 Ankara, Turkey (e-mail: mirza@ee.bilkent.edu.tr; kozat@ee.bilkent.edu.tr).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2019.2935796

it is the decisive factor to determine the computational complexity of training LSTMs. To be more precise, [22] shows that a multiplication operation consumes more than four times the energy required by an addition operation. Thus, especially, matrix-vector multiplications in LSTMs prohibit their implementation in real-life applications such as embedded devices.

In addition to the complexity and energy issues, the performance of LSTMs might be degraded by multiplication operations. Particularly, multiplicative terms in LSTMs might cause the exploding and vanishing gradient problems [6], [7] so that the gradient-based training methods, e.g., SGD, can provide less than adequate training performance. Moreover, training algorithms that employ multiplicative updates, e.g., EG [17], [18], exacerbate such problems even further.

In order to address these issues, we introduce a novel energy-efficient LSTM network and the training methods based on the EG [17] and SGD [11] algorithms. Particularly, we first introduce an LSTM-based regression structure to process variable-length input sequences. We then introduce an energy-efficient LSTM network, which has a significantly smaller number of multiplication operations (only required for certain scaling operations) compared with the classical LSTM network. In order to further reduce the complexity, we also apply a matrix factorization method [24] to the LSTM parameters, such that the number of parameters that needs to be learned is significantly reduced. For this structure, we also introduce online training algorithms based on the EG and SGD algorithms. Thus, unlike the methods in the literature [11], [25], we not only enjoy high performance provided by the LSTM network but also achieve low computational complexity in training. Here, thanks to our generic approach, we also apply this approach to the gated recurrent unit (GRU) network [26] in our experiments. Through an extensive set of simulations, we illustrate significant performance gains and complexity reductions with respect to the conventional methods [11].

### B. Prior Art and Comparisons

Various first-order gradient-based training algorithms have been introduced to train the RNN architectures in an online manner [11], [12], [25], [27]. These first-order gradient-based training algorithms usually employ additive updates in order not to exacerbate complexity issues, e.g., the SGD algorithm [11], [17]. However, the training algorithms with additive updates suffer from slow convergence and inadequate performance, specifically when the input data contains sparse information [17]. In addition, even the first-order algorithms with additive updates might suffer from high complexity while training certain complex RNN architectures, e.g., the LSTM network [5], [12]. To mitigate the complexity issues, [24] applies a matrix factorization method to the parameters of the LSTM network. Thus, they significantly reduce the total number of parameters to be learned. On the other hand, [20], [23], and [28] replace the regular multiplication operation in neural networks with an energy-efficient operator, i.e., the ef-operator. Unlike the regular multiplication operation, the ef-operator only requires sign multiplication and

addition, and thus, it significantly reduces the complexity of neural networks [23], [28]. However, since both approaches employ additive updates, they provide restricted performance in certain tasks [17], [18].

To remedy the performance issues relevant to additive updates, the first-order gradient-based algorithms with multiplicative updates, e.g., the EG algorithm, are introduced [17]. Although such algorithms provide higher performance and faster convergence rate than the algorithms with additive updates, they are highly complex due to the multiplicative structure [19], [20]. As an example, Srinivasan *et al.* [25] derived the backpropagation algorithm for a simple neural network with one hidden layer using both the GD and EG algorithms. Their calculations clearly illustrate high computational complexity in the application of the EG algorithm. In order to achieve high performance provided by multiplicative updates while enjoying low computational complexity, in this article, we introduce an energy-efficient LSTM network, where we replace the regular multiplication operation with the ef-operator [19]. To further reduce the computational complexity, we also apply a matrix factorization method to the matrices in the classical LSTM architecture [24]. Thus, we significantly diminish the number of parameters to be trained in our LSTM network. We then train the introduced network with a training method based on the EG algorithm, as well as a training method based on the SGD algorithm.

### C. Contributions

Our main contributions are as follows.

- 1) As the first time in the literature, we introduce an energy-efficient LSTM network, where we apply a matrix factorization method to reduce the computational complexity of our network. Here, we also replace each regular multiplication operation with an energy-efficient operator that only requires sign multiplication and addition to further reduce the computational complexity.
- 2) We introduce online training methods based on the EG and SGD algorithms to train our energy-efficient LSTM architecture, where we derive online updates for each parameter. Here, the energy-efficient LSTM network trained with our algorithms achieves substantial performance gains with respect to the classical LSTM architecture [8] trained with the conventional training methods [11].
- 3) We achieve these substantial performance gains with a computational complexity that is significantly less than the conventional methods in the literature [11].
- 4) Through an extensive set of simulations, we demonstrate significant performance improvements achieved by the introduced methods with respect to the conventional methods [11]. Moreover, since our approach is generic, we also introduce an energy-efficient GRU network in Section IV.

### D. Organization of This Article

The organization of this article is as follows. We describe the variable-length online regression problem and provide our

LSTM-based structure in Section II. In Section III, we first introduce the basic energy-efficient RNN and LSTM networks using the ef-operator and then apply the matrix factorization method to these networks, where we also introduce our training methods based on the EG and SGD algorithms. In Section IV, we demonstrate the merits of the introduced energy-efficient networks and training algorithms through several experiments, where we also provide an energy-efficient GRU network. Finally, we present concluding remarks in Section V.

## II. MODEL AND PROBLEM DESCRIPTION

In this article, all vectors are column vectors and denoted by boldface lowercase letters. Matrices are represented by boldface uppercase letters. For a matrix  $\mathbf{U}$  (or a vector  $\mathbf{u}$ ),  $\mathbf{U}^T$  (or  $\mathbf{u}^T$ ) is its ordinary transpose. The time index is given as subscript, e.g.,  $\mathbf{u}_t$  is the vector at time  $t$ . For a vector  $\mathbf{u}$ ,  $\|\mathbf{u}\|$  is the  $\ell_1$ -norm. For a vector  $\mathbf{u}_t$ ,  $u_{t,i}$  is the  $i$ th element of that vector. Similarly, for a matrix  $\mathbf{U}$ ,  $u_{ij}$  is the entry at the  $i$ th row and  $j$ th column of  $\mathbf{U}$ . Given a vector  $\mathbf{u}$ ,  $\mathbf{D}^{(u)}$  is the diagonal matrix with the entries of  $\mathbf{u}$  at its diagonal.

We sequentially receive  $\{d_t\}_{t \geq 1}$ ,  $d_t \in \mathbb{R}$ , and matrices  $\{X_t\}_{t \geq 1}$ , i.e., defined as  $X_t = [\mathbf{x}_{t,1}, \mathbf{x}_{t,2}, \dots, \mathbf{x}_{t,n_t}]$ , where  $\mathbf{x}_{t,j} \in \mathbb{R}^p, \forall j \in \{1, 2, \dots, n_t\}$ , and  $n_t \in \mathbb{Z}^+$  is the number of columns in  $X_t$  that may vary with respect to time  $t$ . Here, we aim to find a relation between the desired label  $d_t$  and the corresponding input vector sequence  $X_t$ . To find this relation, after receiving each  $X_t$ , we generate an estimate  $\hat{d}_t$  based on the current and past observations. We then receive the desired value  $d_t$  and suffer the loss  $L(\hat{d}_t, d_t)$  based on our estimate. This framework can be encountered in several machine learning and signal processing applications [29]. As an example, in sequential prediction under the square loss, at each time  $t$ , we receive a set of features, i.e.,  $X_t$  in our case, related to the desired label  $d_t$ . We then generate the estimate through a function, i.e.,  $\hat{d}_t = \kappa(X_t)$ . After the desired label  $d_t$  is observed, we suffer the square loss  $L(\hat{d}_t, d_t) = (d_t - \hat{d}_t)^2$ .

In this article, we use RNNs to obtain  $\hat{d}_t$ . Since we have variable-length data sequences, we use a structure as shown in Fig. 1 to obtain fixed-length sequences. The basic RNN architecture is defined by the following equations [10]:

$$\mathbf{h}_{t,j} = f(\mathbf{W}\mathbf{x}_{t,j} + \mathbf{R}\mathbf{h}_{t,j-1}) \quad (1)$$

$$\mathbf{z}_{t,j} = g(\mathbf{U}\mathbf{h}_{t,j}) \quad (2)$$

where  $\mathbf{x}_{t,j} \in \mathbb{R}^p$  is the input vector,  $\mathbf{h}_{t,j} \in \mathbb{R}^m$  is the state vector, and  $\mathbf{z}_{t,j} \in \mathbb{R}^m$  is the output vector for the  $j$ th RNN unit. Here,  $f(\cdot)$  and  $g(\cdot)$  usually set to the hyperbolic tangent function, and they apply to vectors pointwise. Moreover,  $\mathbf{W}$ ,  $\mathbf{R}$ , and  $\mathbf{U}$  are the parameters of the basic RNN architecture, where the sizes are chosen according to the size of the input and output vectors.

We use the LSTM network as a special variant of RNNs to obtain  $\hat{d}_t$ . Among different implementations of LSTM, we choose the most widely used one, i.e., the LSTM network without peephole connections [12]. Since we receive variable-length data sequence, we apply the LSTM network

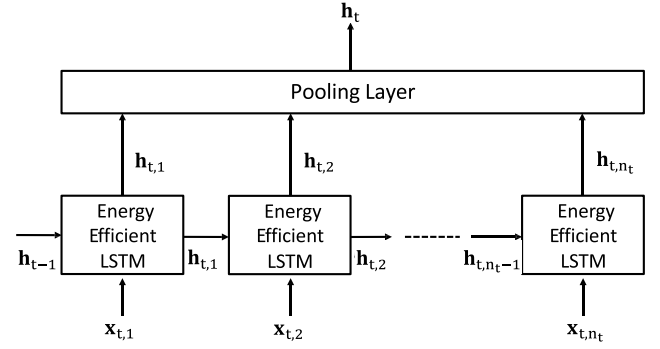


Fig. 1. Detailed schematic of energy-efficient LSTM-based architecture.

to each column of  $X_t$  as shown in Fig. 1, where the internal LSTM equations for the  $j$ th unit are as follows [8], [12]:

$$\tilde{\mathbf{c}}_{t,j} = g(\mathbf{W}^{(c)}\mathbf{x}_{t,j} + \mathbf{R}^{(c)}\mathbf{h}_{t,j-1} + \mathbf{b}^{(c)}) \quad (3)$$

$$\mathbf{i}_{t,j} = \sigma(\mathbf{W}^{(i)}\mathbf{x}_{t,j} + \mathbf{R}^{(i)}\mathbf{h}_{t,j-1} + \mathbf{b}^{(i)}) \quad (4)$$

$$\mathbf{f}_{t,j} = \sigma(\mathbf{W}^{(f)}\mathbf{x}_{t,j} + \mathbf{R}^{(f)}\mathbf{h}_{t,j-1} + \mathbf{b}^{(f)}) \quad (5)$$

$$\mathbf{c}_{t,j} = \mathbf{D}_{t,j}^{(i)}\tilde{\mathbf{c}}_{t,j} + \mathbf{D}_{t,j}^{(f)}\mathbf{c}_{t,j-1} \quad (6)$$

$$\mathbf{o}_{t,j} = \sigma(\mathbf{W}^{(o)}\mathbf{x}_{t,j} + \mathbf{R}^{(o)}\mathbf{h}_{t,j-1} + \mathbf{b}^{(o)}) \quad (7)$$

$$\mathbf{h}_{t,j} = \mathbf{D}_{t,j}^{(o)}g(\mathbf{c}_{t,j}) \quad (8)$$

where  $\mathbf{c}_{t,j} \in \mathbb{R}^m$  is the state vector,  $\mathbf{x}_{t,j} \in \mathbb{R}^p$  is the input vector,  $\mathbf{h}_{t,j} \in \mathbb{R}^m$  is the output vector. Here,  $\mathbf{i}_{t,j}$ ,  $\mathbf{f}_{t,j}$ , and  $\mathbf{o}_{t,j}$  are the input, forget, and output gates, respectively. The function  $g(\cdot)$  applies to vectors pointwise and commonly set to  $\tanh(\cdot)$ . Similarly, the sigmoid function  $\sigma(\cdot)$  applies to vectors pointwise. The sizes of the other matrices and vectors are determined according to the size of the input and output vectors. After the consecutive applications of the LSTM network to each column as shown in Fig. 1, we take the average of the outputs of the LSTM networks, i.e., the mean pooling method, in order to obtain a fixed-length representation, i.e., denoted as  $\mathbf{h}_t \in \mathbb{R}^m$  at time  $t$ . Using the fixed-length vectors, we generate the final estimate as

$$\hat{d}_t = \mathbf{w}_t^T \mathbf{h}_t \quad (9)$$

where  $\mathbf{w}_t \in \mathbb{R}^m$  represents the regression coefficients at time  $t$ . In this framework, we aim to train the system parameters, such that the total loss at time  $t$ , i.e.,  $\sum_{i=1}^t L(\hat{d}_i, d_i)$ , is minimized.

For the pooling operation in Fig. 1, we use the mean pooling method to obtain the fixed-length output vectors as  $\mathbf{h}_t = (1/n_t) \sum_{j=1}^{n_t} \mathbf{h}_{t,j}$ . However, there are certain other pooling methods in the literature, and we can also employ them in our approach. As an example, we can apply the max and last pooling methods in our case by using  $\mathbf{h}_t = \max_j \mathbf{h}_{t,j}$  and  $\mathbf{h}_t = \mathbf{h}_{t,n_t}$ , respectively. With such changes, our derivations can be extended to the other pooling methods.

## III. ONLINE LEARNING WITH ENERGY-EFFICIENT RNN ARCHITECTURES

In this section, we first apply the ef-operator to the basic RNN and LSTM architectures. We then introduce our energy-efficient RNN and LSTM architectures using the

matrix factorization method along with the ef-operator. Finally, we introduce online training algorithms based on the SGD and EG algorithms, where we provide the required updates for each parameter.

#### A. RNN With EF-Operator

In this section, we study a modified version of the basic RNN architecture, where we replace the regular multiplication operations with the ef-operator.

Let  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ , and the ef-operator [19] on  $\mathbf{a}$  and  $\mathbf{b}$  is defined as

$$\mathbf{a} \diamond \mathbf{b} := \sum_{i=1}^p \text{sign}(a_i \times b_i)(|a_i| + |b_i|) \quad (10)$$

where the  $\text{sign}(\cdot)$  function returns the sign of its input. Equation (10) can also be written as

$$\mathbf{a} \diamond \mathbf{b} := \sum_{i=1}^p \text{sign}(a_i)b_i + \text{sign}(b_i)a_i.$$

From the above-mentioned definition, it is obvious that the ef-operator only uses addition and sign multiplications, which are all energy-efficient operators.

In a similar manner, we define the ef-operator for matrix multiplications as follows:

$$(\mathbf{A} \diamond \mathbf{B})_{ij} = \mathbf{a}_i \diamond \mathbf{b}_j$$

where  $\mathbf{a}_i$  and  $\mathbf{b}_j$  are the  $i$ th row of  $\mathbf{A}$  and  $j$ th column of  $\mathbf{B}$ , respectively.

By applying the ef-operator, (1) can be written as

$$\mathbf{h}_{t,j} = f(\mathbf{a}_h \odot (\mathbf{W} \diamond \mathbf{x}_{t,j}) + \mathbf{b}_h \odot (\mathbf{R} \diamond \mathbf{h}_{t,j-1})) \quad (11)$$

where  $\mathbf{a}_h \in \mathbb{R}^m$  and  $\mathbf{b}_h \in \mathbb{R}^m$  are the scaling coefficients introduced in [23] and [28] and  $\odot$  is the element-by-element multiplication of two vectors of the same size. Here, the scaling coefficients are the crucial factors to match the performance of the classical multiplicative networks. In particular, even though we eliminate several matrix-vector multiplications, which are one of the main pillars of RNNs, these coefficients keep the modeling capabilities of the network at the same level by introducing an additional multiplicative term before the nonlinear function is applied. In addition, note that the weight matrices, i.e.,  $\mathbf{W}$  and  $\mathbf{R}$ , in (1) and (11) are not necessarily the same, however, their function is the same, i.e., both are weight matrices that multiply the input vector. In (11),  $\mathbf{W} \diamond \mathbf{x}_{t,j}$  and  $\mathbf{R} \diamond \mathbf{h}_{t,j-1}$  are given as follows:

$$\mathbf{W} \diamond \mathbf{x}_{t,j} = [\mathbf{w}_1 \diamond \mathbf{x}_{t,j} \quad \mathbf{w}_2 \diamond \mathbf{x}_{t,j} \quad \dots \quad \mathbf{w}_m \diamond \mathbf{x}_{t,j}]^T$$

where  $\mathbf{w}_i$  represents the transpose of the  $i$ th row of  $\mathbf{W}$  and  $\mathbf{w}_i \diamond \mathbf{x}_{t,j}$  is given as

$$\mathbf{w}_i \diamond \mathbf{x}_{t,j} = \sum_{k=1}^p \text{sign}(x_{t,jk})w_{ik} + \text{sign}(w_{ik})x_{t,jk}$$

where  $x_{t,jk}$  is the  $k$ th element of  $\mathbf{x}_{t,j}$ . Similarly

$$\mathbf{R} \diamond \mathbf{h}_{t,j-1} = [\mathbf{r}_1 \diamond \mathbf{h}_{t,j-1} \quad \mathbf{r}_2 \diamond \mathbf{h}_{t,j-1} \quad \dots \quad \mathbf{r}_m \diamond \mathbf{h}_{t,j-1}]^T$$

where  $\mathbf{r}_i$  represents the transpose of the  $i$ th row of  $\mathbf{R}$  and  $\mathbf{r}_i \diamond \mathbf{h}_{t,j}$  is given as

$$\mathbf{r}_i \diamond \mathbf{h}_{t,j} = \sum_{k=1}^m \text{sign}(h_{t,jk})r_{ik} + \text{sign}(r_{ik})h_{t,jk}$$

where  $h_{t,jk}$  is the  $k$ th element of  $\mathbf{h}_{t,j}$ . Likewise, (2) is modified as follows:

$$\mathbf{z}_{t,j} = g(\mathbf{b}_z \odot (\mathbf{U} \diamond \mathbf{h}_{t,j})) \quad (12)$$

where  $\mathbf{b}_z$  is the scaling coefficient.

#### B. LSTM With EF-Operator

In this section, we replace the regular multiplication operators in the classical LSTM architecture with the ef-operator, as illustrated in Fig. 2. Based on this modification, (3)–(8) are written as follows:

$$\tilde{\mathbf{c}}_{t,j} = g(\mathbf{a}_{\tilde{c}} \odot (\mathbf{W}^{(\tilde{c})} \diamond \mathbf{x}_{t,j}) + \mathbf{b}_{\tilde{c}} \odot (\mathbf{R}^{(\tilde{c})} \diamond \mathbf{h}_{t,j-1}) + \mathbf{b}^{(\tilde{c})}) \quad (13)$$

$$\mathbf{i}_{t,j} = \sigma(\mathbf{a}_i \odot (\mathbf{W}^{(i)} \diamond \mathbf{x}_{t,j}) + \mathbf{b}_i \odot (\mathbf{R}^{(i)} \diamond \mathbf{h}_{t,j-1}) + \mathbf{b}^{(i)}) \quad (14)$$

$$\mathbf{f}_{t,j} = \sigma(\mathbf{a}_f \odot (\mathbf{W}^{(f)} \diamond \mathbf{x}_{t,j}) + \mathbf{b}_f \odot (\mathbf{R}^{(f)} \diamond \mathbf{h}_{t,j-1}) + \mathbf{b}^{(f)}) \quad (15)$$

$$\mathbf{c}_{t,j} = \mathbf{i}_{t,j} \blacklozenge \tilde{\mathbf{c}}_{t,j} + \mathbf{f}_{t,j} \blacklozenge \mathbf{c}_{t,j-1} \quad (16)$$

$$\mathbf{o}_{t,j} = \sigma(\mathbf{a}_o \odot (\mathbf{W}^{(o)} \diamond \mathbf{x}_{t,j}) + \mathbf{b}_o \odot (\mathbf{R}^{(o)} \diamond \mathbf{h}_{t,j-1}) + \mathbf{b}^{(o)}) \quad (17)$$

$$\mathbf{h}_{t,j} = \mathbf{o}_{t,j} \blacklozenge g(\mathbf{c}_{t,j}) \quad (18)$$

where  $\mathbf{a}_{(\cdot)}, \mathbf{b}_{(\cdot)} \in \mathbb{R}^m$  are the scaling coefficients and the  $\blacklozenge$  operation is defined as follows:

$$\begin{aligned} \mathbf{a} \blacklozenge \mathbf{b} &:= [\mathbf{a}_1 \diamond \mathbf{b}_1 \quad \mathbf{a}_2 \diamond \mathbf{b}_2 \quad \dots \quad \mathbf{a}_p \diamond \mathbf{b}_p]^T \\ &:= \text{sign}(\mathbf{a}) \odot \mathbf{b} + \text{sign}(\mathbf{b}) \odot \mathbf{a}. \end{aligned}$$

In (13),  $\mathbf{W}^{(\tilde{c})} \diamond \mathbf{x}_{t,j}$  is written as

$$\mathbf{W}^{(\tilde{c})} \diamond \mathbf{x}_{t,j} = [\mathbf{w}_1^{(\tilde{c})} \diamond \mathbf{x}_{t,j} \quad \mathbf{w}_2^{(\tilde{c})} \diamond \mathbf{x}_{t,j} \quad \dots \quad \mathbf{w}_m^{(\tilde{c})} \diamond \mathbf{x}_{t,j}]^T \quad (19)$$

where  $\mathbf{w}_i^{(\tilde{c})} \diamond \mathbf{x}_{t,j}$  is given as

$$\mathbf{w}_i^{(\tilde{c})} \diamond \mathbf{x}_{t,j} = \sum_{k=1}^p \text{sign}(x_{t,jk})w_{ik}^{(\tilde{c})} + \text{sign}(w_{ik}^{(\tilde{c})})x_{t,jk}$$

and

$$\mathbf{R}^{(\tilde{c})} \diamond \mathbf{h}_{t,j-1} = [\mathbf{r}_1^{(\tilde{c})} \diamond \mathbf{h}_{t,j-1} \quad \mathbf{r}_2^{(\tilde{c})} \diamond \mathbf{h}_{t,j-1} \quad \dots \quad \mathbf{r}_m^{(\tilde{c})} \diamond \mathbf{h}_{t,j-1}]^T \quad (20)$$

where  $\mathbf{r}_i^{(\tilde{c})} \diamond \mathbf{h}_{t,j}$  is given as

$$\mathbf{r}_i^{(\tilde{c})} \diamond \mathbf{h}_{t,j} = \sum_{k=1}^m \text{sign}(h_{t,jk})r_{ik}^{(\tilde{c})} + \text{sign}(r_{ik}^{(\tilde{c})})h_{t,jk}$$

where  $\mathbf{w}_i^{(\tilde{c})}$  and  $\mathbf{r}_i^{(\tilde{c})}$  are the  $i$ th row of  $\mathbf{W}^{(\tilde{c})}$  and  $\mathbf{R}^{(\tilde{c})}$ , respectively.



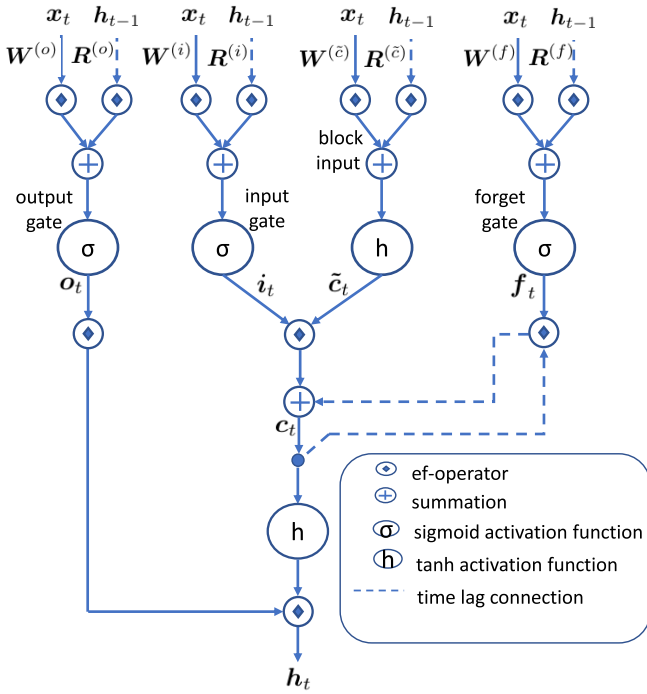


Fig. 2. Detailed schematic of energy-efficient LSTM block at time  $t$ . Note that the solid lines represent the direct connections, while the dotted lines represent the time-lagged connections. For the sake of simplicity, bias terms are not shown in the figure.

For the other multiplications, we change the parameters in either (19) or (20) according to the chosen coefficient matrix. Other than that, we follow the same procedures in (19) and (20).

*Remark 1:* Compared with the original LSTM network in (3)–(8), we convert  $4m(m+p) + 3m$  regular multiplication operations into sign multiplication and addition operations thanks to the ef-operator. However, due to the scaling factors introduced in (13)–(15) and (17), we have  $8m$  additional regular multiplications. Overall, since for large  $m$  and  $p$  values  $8m \ll 4m(m+p) + 3m$ , we significantly reduce the number of regular multiplications. Thus, we provide a substantial decrease in the computational complexity and energy consumption compared with the classical LSTM network.

### C. Energy-Efficient RNN With Weight Matrix Factorization

In this section, we apply the matrix factorization method [24] to the weight matrices of the basic RNN architecture in order to reduce the number of parameters to be trained.

We first factorize the weight matrices in (11) as  $W \approx MN$  and  $R \approx PQ$ , where  $W \in \mathbb{R}^{m \times p}$ ,  $M \in \mathbb{R}^{m \times d}$ ,  $N \in \mathbb{R}^{d \times p}$ ,  $R \in \mathbb{R}^{m \times m}$ ,  $P \in \mathbb{R}^{m \times f}$ , and  $Q \in \mathbb{R}^{f \times m}$ .

*Remark 2:* We factorize the RNN weight matrices into two smaller matrices. The rank of these two smaller matrices is selected, such that  $d, f \ll \min(m, p)$ . Thus, we significantly reduce the number of parameters that needs to be learned, e.g.,  $W$  has  $mp$  entries, while  $M$  and  $N$  have  $d(m+p) \ll mp$ .

The energy-efficient RNN with weight matrix factorization can be written as

$$h_{t,j} = f(a_h \odot (M \diamond N \diamond x_{t,j}) + b_h \odot (P \diamond Q \diamond h_{t,j-1})). \quad (21)$$

In (21),  $M \diamond N \diamond x_{t,j}$  is given as follows:

$$M \diamond N \diamond x_{t,j} = [\mu_1 \diamond x_{t,j} \quad \mu_2 \diamond x_{t,j} \dots \mu_m \diamond x_{t,j}]^T$$

where  $\mu_i \in \mathbb{R}^p$  is the  $i$ th row of  $M \diamond N$  and  $\mu_i \diamond x_{t,j}$  is given as

$$\mu_i \diamond x_{t,j} = \sum_{k=1}^p \text{sign}(x_{t,jk}) \mu_{ik} + \text{sign}(\mu_{ik}) x_{t,jk} \quad (22)$$

and

$$P \diamond Q \diamond h_{t,j-1} = [v_1 \diamond h_{t,j-1} \quad v_2 \diamond h_{t,j-1} \dots v_m \diamond h_{t,j-1}]^T$$

where  $v_i \in \mathbb{R}^m$  is the  $i$ th row of  $P \diamond Q$  and  $v_i \diamond h_{t,j}$  is given as

$$v_i \diamond h_{t,j} = \sum_{k=1}^m \text{sign}(h_{t,jk}) v_{ik} + \text{sign}(v_{ik}) h_{t,jk}. \quad (23)$$

In a similar manner, (12) is modified as follows:

$$z_{t,j} = g(b_z \odot (S \diamond T \diamond h_{t,j})) \quad (24)$$

where we factorize the  $U$  matrix as  $U \approx ST$  so that the number of columns in  $S$  (or the number of rows in  $T$ ) is significantly smaller than the number of rows in  $S$  (or the number of columns in  $T$ ).

### D. Energy-Efficient LSTM With Weight Matrix Factorization

In this section, we apply the matrix factorization method [24] to the weight matrices of the LSTM network to diminish the number of parameters in the network.

We factorize the LSTM neural network weight matrices into two sub-matrices of lower rank as  $W^{(\cdot)} \approx M^{(\cdot)}N^{(\cdot)}$  and  $R^{(\cdot)} \approx P^{(\cdot)}Q^{(\cdot)}$ , where  $W^{(\cdot)} \in \mathbb{R}^{m \times p}$ ,  $M^{(\cdot)} \in \mathbb{R}^{m \times d}$ ,  $N^{(\cdot)} \in \mathbb{R}^{d \times p}$ ,  $R^{(\cdot)} \in \mathbb{R}^{m \times m}$ ,  $P^{(\cdot)} \in \mathbb{R}^{m \times f}$ , and  $Q^{(\cdot)} \in \mathbb{R}^{f \times m}$ , such that  $d, f \ll \min(p, m)$ . We then apply this factorization to the LSTM neural network in (13)–(18) by replacing the weight matrices with their factorized forms. The modifications for the  $j$ th LSTM unit in Fig. 1 are as follows.

Here,  $M^{(\tilde{c})} \diamond N^{(\tilde{c})} \diamond x_{t,j}$  is given as follows:

$$M^{(\tilde{c})} \diamond N^{(\tilde{c})} \diamond x_{t,j} = [\mu_1^{(\tilde{c})} \diamond x_{t,j} \quad \mu_2^{(\tilde{c})} \diamond x_{t,j} \dots \mu_m^{(\tilde{c})} \diamond x_{t,j}]^T \quad (25)$$

where  $\mu_i^{(\tilde{c})} \in \mathbb{R}^p$  is the  $i$ th row of  $M^{(\tilde{c})} \diamond N^{(\tilde{c})}$  and  $\mu_i^{(\tilde{c})} \diamond x_{t,j}$  is given as

$$\mu_i^{(\tilde{c})} \diamond x_{t,j} = \sum_{k=1}^p \text{sign}(x_{t,jk}) \mu_{ik}^{(\tilde{c})} + \text{sign}(\mu_{ik}^{(\tilde{c})}) x_{t,jk}$$

and

$$P^{(\tilde{c})} \diamond Q^{(\tilde{c})} \diamond h_{t,j-1} = [v_1^{(\tilde{c})} \diamond h_{t,j-1} \quad v_2^{(\tilde{c})} \diamond h_{t,j-1} \dots v_m^{(\tilde{c})} \diamond h_{t,j-1}]^T \quad (26)$$

where  $v_i^{(\tilde{c})} \in \mathbb{R}^m$  is the  $i$ th row of  $P^{(\tilde{c})} \diamond Q^{(\tilde{c})}$  and  $v_i^{(\tilde{c})} \diamond h_{t,j}$  is given as

$$v_i^{(\tilde{c})} \diamond h_{t,j} = \sum_{k=1}^m \text{sign}(h_{t,jk}) v_{ik}^{(\tilde{c})} + \text{sign}(v_{ik}^{(\tilde{c})}) h_{t,jk}.$$

For the other weight matrices, we replace the factorized form of the chosen weight matrix in (25) and (26). Then, we follow the same operations in (25) and (26).

*Remark 3:* We reduce the total number of LSTM network parameters by applying weight matrix factorization. In the original LSTM equations in (3)–(8), we have  $4m(m+p)$  scalar parameters in the weight matrices, i.e.,  $\mathbf{W}^{(\cdot)}$  and  $\mathbf{R}^{(\cdot)}$ . However, in our energy-efficient LSTM network, we have  $4d(m+p) + 8mf$ , which is significantly less than  $4m(m+p)$ , provided that  $d, f \ll \min(m, p)$ .

### E. Online Training Algorithms

In this section, we derive the online updates to train the parameters of the introduced energy-efficient networks. We first derive the online updates based on the SGD algorithm. We then derive the online updates based on the EG algorithm.

We first employ the SGD algorithm [11] to obtain the online updates for each parameter. For  $\mathbf{w}_t$ , the SGD update is computed as follows:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}_t} L \quad (27)$$

where  $\nabla_{\mathbf{w}_t}$  represents the gradient of a certain function with respect to  $\mathbf{w}_t$  and  $\eta_t$  is the learning rate. Here,  $L$  is the instantaneous loss, i.e., the squared error  $(d_t - \hat{d}_t)^2$ , and we denote it as  $L$  rather than  $L(\hat{d}_t, d_t)$  for notational simplicity. Note that SGD updates, e.g., (27), are additive updates since the gradient information is being added at each time step.

On the other hand, for  $\mathbf{w}_t$ , the EG update [17] is computed as follows:

$$w_{t+1,i} = \frac{w_{t,i} r_{t,i}}{\sum_{j=1}^m w_{t,j} r_{t,j}} \quad (28)$$

where

$$r_{t,i} = \exp(-\eta_t L'_{w_{t,i}})$$

and  $w_{t,i}$  is the  $i$ th component of  $\mathbf{w}_t$  and  $L'_{w_{t,i}}$  is the partial derivative of the instantaneous loss function with respect to  $w_{t,i}$ . As seen in (28), EG updates are multiplicative updates since the gradient information is encapsulated in the exponent part and multiplied at each time step. In order to eliminate the multiplication and exponentiation in (28), we use the first-order Taylor series expansion along with the ef-operator as follows:

$$w_{t+1,i} = \frac{w_{t,i} \diamond \hat{r}_{t,i}}{\sum_{j=1}^m (w_{t,j} \diamond \hat{r}_{t,j})} \quad (29)$$

where

$$\hat{r}_{t,i} = 1 - \eta_t L'_{w_{t,i}}.$$

Note that since the division in (29) is the same for all possible  $i$  values, it is just a scaling factor in the implementation of the algorithm. Thus, this operation does not require high amount of energy and computational resources unlike the regular multiplication operation.

*Remark 4:* Since the weight vector  $\mathbf{w}_t$  might contain negative components, we use the slightly modified version of the original EG algorithms, i.e., the  $\text{EG}_+^+$  algorithm [17], which

uses a weight vector  $\mathbf{w}_t^+ - \mathbf{w}_t^-$ . In this algorithm, the weight vector is updated as follows:

$$w_{t+1,i}^+ = \frac{w_{t,i}^+ \diamond \hat{r}_{t,i}^+}{\sum_{j=1}^m (w_{t,j}^+ \diamond \hat{r}_{t,j}^+ + w_{t,j}^- \diamond \hat{r}_{t,j}^-)}$$

$$w_{t+1,i}^- = \frac{w_{t,i}^- \diamond \hat{r}_{t,i}^-}{\sum_{j=1}^m (w_{t,j}^+ \diamond \hat{r}_{t,j}^+ + w_{t,j}^- \diamond \hat{r}_{t,j}^-)}$$

where

$$\hat{r}_{t,i}^+ = 1 - \eta_t L'_{w_{t,i}^+}$$

$$\hat{r}_{t,i}^- = 1 + \eta_t L'_{w_{t,i}^-}.$$

In Sections III-E1 and III-E2, we derive the updates for the parameters of the proposed energy-efficient RNN and LSTM networks.

1) *Online Training of Energy-Efficient RNN:* We compute the first-order gradient of the loss function with respect to each parameter in order to perform SGD and EG updates.

In the basic RNN architecture, we have  $\mathbf{z}_t = \sum_{j=1}^{n_t} \mathbf{z}_{t,j}/n_t$  as the output at time  $t$ . Although our structure in Fig. 1 is generic in the sense that it can process variable-length data sequences, here, we only derive the equations for  $n_t = 1$  for notational and presentation simplicity. However, at the end of this section, we also provide the required extensions to obtain the equations for generic  $n_t$  values. With this modification, we have  $\mathbf{z}_t = \mathbf{z}_{t,1}$ , and thus, we generate the estimate as  $\hat{d}_t = \mathbf{w}_t^T \mathbf{z}_{t,1}$ .

Under the square loss, we compute the first-order derivative of the loss function with respect to  $b_{zi}$ , i.e., the  $i$ th element of  $\mathbf{b}_z$ , as follows:

$$\frac{\partial L}{\partial b_{zi}} = \frac{\partial L}{\partial \hat{d}_t} \frac{\partial \hat{d}_t}{\partial \mathbf{z}_{t,1}} \frac{\partial \mathbf{z}_{t,1}}{\partial b_{zi}} = -2(d_t - \hat{d}_t) \mathbf{w}_t^T$$

$$[g'(\boldsymbol{\varphi}_t) \odot ((\mathbf{u}_i \diamond \mathbf{h}_{t,1}) \mathbf{e}_i + \mathbf{b}_z \odot \boldsymbol{\lambda}_t^{(Uh)})] \quad (30)$$

where  $g'$  is the derivative of  $g(\cdot)$  with respect to its argument,  $\mathbf{e}_i$  is a vector of zeros except a 1 at the  $i$ th index, and

$$\boldsymbol{\varphi}_t = \mathbf{b}_z \odot (\mathbf{U} \diamond \mathbf{h}_{t,1}).$$

Moreover, for the  $i$ th element of  $\boldsymbol{\lambda}_t^{(Uh)} = \partial(\mathbf{U} \diamond \mathbf{h}_{t,1})/\partial b_{zi}$ , we use the following formula in (30):

$$\lambda_{t,i}^{(Uh)} = \sum_{j=1}^m (u_{ij} 2\delta(h_{t,1j}) \gamma_{t,ij}^{(b_z)} + \text{sign}(u_{ij}) \gamma_{t,ij}^{(b_z)}) \quad (31)$$

where  $\delta(\cdot)$  is the dirac delta function, we compute the derivative of the sign function as  $d(\text{sign}(x))/dx = 2\delta(x)$  [30] and

$$\gamma_{t,ij}^{(b_z)} = \frac{\partial h_{t,1j}}{\partial b_{zi}} = f'_j(\boldsymbol{\theta}_t) b_{hj} \lambda_{t-1,j}^{(Rh)} \quad (32)$$

where  $f'_i$  is the derivative of the  $i$ th element of  $f(\cdot)$  with respect to its argument and

$$\boldsymbol{\theta}_t = \mathbf{a}_h \odot (\mathbf{W} \diamond \mathbf{x}_{t,1}) + \mathbf{b}_h \odot (\mathbf{R} \diamond \mathbf{h}_{t-1,1}).$$

Similarly, we have the following derivative for  $u_{ik}$ :

$$\begin{aligned} \frac{\partial L}{\partial u_{ik}} &= \frac{\partial L}{\partial \hat{d}_t} \frac{\partial \hat{d}_t}{\partial \mathbf{z}_{t,1}} \frac{\partial \mathbf{z}_{t,1}}{\partial u_{ik}} \\ &= -2(d_t - \hat{d}_t) \mathbf{w}_t^T [g'(\varphi_t) \odot (b_{zi}(\text{sign}(h_{t,1k}) \\ &\quad + 2\delta(u_{ik})h_{t,1k})\mathbf{e}_i + \mathbf{b}_z \odot \boldsymbol{\alpha}_i^{(U_h)})] \end{aligned} \quad (33)$$

where

$$\alpha_{t,i}^{(U_h)} = \sum_{j=1}^m (u_{ij} 2\delta(h_{t,1j}) \gamma_{t,ij}^{(u_{ik})} + \text{sign}(u_{ij}) \gamma_{t,ij}^{(u_{ik})}) \quad (34)$$

and

$$\gamma_{t,ij}^{(u_{ik})} = \frac{\partial h_{t,1j}}{\partial u_{ik}} = f'_j(\boldsymbol{\theta}_t) b_{hj} \alpha_{t-1,j}^{(R_h)}.$$

*Remark 5:* When we take the derivative of  $L$  with respect to the other parameters, the position of the term with  $\mathbf{e}_i$  changes. As an example, for the derivative of  $L$  with respect to  $b_{hi}$ , the  $\mathbf{r}_i \diamond \mathbf{h}_{t-1,1}$  term appears in (32) when  $j = i$ ; otherwise, (32) does not change. If we write (32) in a vector form, the contribution of  $\mathbf{r}_i \diamond \mathbf{h}_{t-1,1}$  can be written as  $(\mathbf{R} \diamond \mathbf{h}_{t-1,1})\mathbf{e}_i$ . As seen in this case, for the other derivatives, the position and the form of the term with  $\mathbf{e}_i$  slightly change; other than that, we follow the same procedure in (30)–(34).

*Remark 6:* For the other  $n_t$  values, i.e.,  $n_t \neq 1$ , the recursion in (32) is performed through the outputs of the different RNN blocks at a certain time  $t$ , as shown in Fig. 1. Thus, rather than having  $t - 1$  in (32), we have multiple recursions based on another index at time  $t$ . Besides this slight change, all of our derivations hold for generic  $n_t$  values.

*Remark 7:* For the energy-efficient RNN architecture with weight matrix factorization, we take the derivative of the loss function with respect to the parameters of each factorized matrix. As an example, in (33), instead of only taking the derivative with respect to the parameters of  $\mathbf{U}$ , we compute the derivatives of the loss with respect to the entries of both  $\mathbf{S}$  and  $\mathbf{T}$ , i.e., the factorized versions of  $\mathbf{U}$ . Other than such changes, we follow the same procedures in (30)–(34).

With the derived gradients, we can update each parameter of the basic RNN architecture as in (27) and (29).

**2) Online Training of Energy-Efficient LSTM:** Here, we derive the first-order gradient of the loss function with respect to each LSTM parameter to obtain the online updates based on the SGD and EG algorithms. We again derive the derivatives for the  $n_t = 1$  case for notational and presentation simplicity. However, at the end of this section, we also provide the required extensions to obtain the equations for generic  $n_t$  values. With this modification, we have  $\mathbf{h}_t = \mathbf{h}_{t,1}$ , and hence, we generate the estimate  $\hat{d}_t = \mathbf{w}_t^T \mathbf{h}_{t,1}$ .

We first compute the derivative of  $L$  with respect to  $w_{ij}^{(\bar{c})}$ , i.e., the element at the  $i$ th row and the  $j$ th column of  $\mathbf{W}^{(\bar{c})}$ , as follows:

$$\frac{\partial L}{\partial w_{ij}^{(\bar{c})}} = \frac{\partial L}{\partial \hat{d}_t} \frac{\partial \hat{d}_t}{\partial \mathbf{h}_{t,1}} \frac{\partial \mathbf{h}_{t,1}}{\partial w_{ij}^{(\bar{c})}} = -2(d_t - \hat{d}_t) \mathbf{w}_t^T \frac{\partial(\mathbf{o}_{t,1} \diamond g(\mathbf{c}_{t,1}))}{\partial w_{ij}^{(\bar{c})}}. \quad (35)$$

In (35), we calculate the partial derivative as

$$\begin{aligned} \frac{\partial(\mathbf{o}_{t,1} \diamond g(\mathbf{c}_{t,1}))}{\partial w_{ij}^{(\bar{c})}} &= \frac{\partial \mathbf{o}_{t,1}}{\partial w_{ij}^{(\bar{c})}} \odot \text{sign}(g(\mathbf{c}_{t,1})) \\ &\quad + \mathbf{o}_{t,1} \odot 2\delta(g(\mathbf{c}_{t,1})) \odot g'(\mathbf{c}_{t,1}) \odot \frac{\partial \mathbf{c}_{t,1}}{\partial w_{ij}^{(\bar{c})}} \\ &\quad + 2\delta(\mathbf{o}_{t,1}) \odot \frac{\partial \mathbf{o}_{t,1}}{\partial w_{ij}^{(\bar{c})}} \odot g(\mathbf{c}_{t,1}) \\ &\quad + \text{sign}(\mathbf{o}_{t,1}) \odot g'(\mathbf{c}_{t,1}) \odot \frac{\partial \mathbf{c}_{t,1}}{\partial w_{ij}^{(\bar{c})}}. \end{aligned} \quad (36)$$

For (36), we now compute the derivatives of  $\mathbf{o}_{t,1}$  and  $\mathbf{c}_{t,1}$  with respect to  $w_{ij}^{(\bar{c})}$ . With  $\boldsymbol{\lambda}_{t-1}^{(R^{(o)h})} = \partial(\mathbf{R}^{(o)} \diamond \mathbf{h}_{t-1,1}) / \partial w_{ij}^{(\bar{c})}$  as in (31), the derivative of (17) is as follows:

$$\frac{\partial \mathbf{o}_{t,1}}{\partial w_{ij}^{(\bar{c})}} = \mathbf{D}_{t,1}^{(\sigma'(\zeta^{(o)}))} (\mathbf{b}_o \odot \boldsymbol{\lambda}_{t-1}^{(R^{(o)h})}) \quad (37)$$

where

$$\boldsymbol{\zeta}_{t,1}^{(o)} = \mathbf{a}_o \odot (\mathbf{W}^{(o)} \diamond \mathbf{x}_{t,1}) + \mathbf{b}_o \odot (\mathbf{R}^{(o)} \diamond \mathbf{h}_{t-1,1}) + \mathbf{b}^{(o)}. \quad (38)$$

In order to calculate (36), we also compute the derivative of  $\mathbf{c}_{t,1}$  with respect to  $w_{ij}^{(\bar{c})}$ . For this derivative, we obtain the following recursive relation from (16):

$$\begin{aligned} \frac{\partial \mathbf{c}_{t,1}}{\partial w_{ij}^{(\bar{c})}} &= \text{sign}(\tilde{\mathbf{c}}_{t,1}) \odot \frac{\partial \tilde{\mathbf{c}}_{t,1}}{\partial w_{ij}^{(\bar{c})}} + 2\delta(\tilde{\mathbf{c}}_{t,1}) \odot \frac{\partial \tilde{\mathbf{c}}_{t,1}}{\partial w_{ij}^{(\bar{c})}} \odot \mathbf{i}_{t,1} \\ &\quad + \text{sign}(\mathbf{i}_{t,1}) \odot \frac{\partial \tilde{\mathbf{c}}_{t,1}}{\partial w_{ij}^{(\bar{c})}} + 2\delta(\mathbf{i}_{t,1}) \odot \frac{\partial \tilde{\mathbf{c}}_{t,1}}{\partial w_{ij}^{(\bar{c})}} \odot \tilde{\mathbf{c}}_{t,1} \\ &\quad + \text{sign}(\mathbf{c}_{t-1,1}) \odot \frac{\partial \mathbf{f}_{t,1}}{\partial w_{ij}^{(\bar{c})}} + 2\delta(\mathbf{c}_{t-1,1}) \odot \frac{\partial \mathbf{c}_{t-1,1}}{\partial w_{ij}^{(\bar{c})}} \odot \mathbf{f}_{t,1} \\ &\quad + \text{sign}(\mathbf{f}_{t,1}) \odot \frac{\partial \mathbf{c}_{t-1,1}}{\partial w_{ij}^{(\bar{c})}} + 2\delta(\mathbf{f}_{t,1}) \odot \frac{\partial \mathbf{f}_{t,1}}{\partial w_{ij}^{(\bar{c})}} \odot \mathbf{c}_{t-1,1}. \end{aligned} \quad (39)$$

For (39), we compute the derivatives of (13)–(15) with respect to  $w_{ij}^{(\bar{c})}$  as follows:

$$\frac{\partial \mathbf{i}_{t,1}}{\partial w_{ij}^{(\bar{c})}} = \mathbf{D}_{t,1}^{(\sigma'(\zeta^{(i)}))} (\mathbf{b}_i \odot \boldsymbol{\lambda}_{t-1}^{(R^{(i)h})}) \quad (40)$$

$$\frac{\partial \mathbf{f}_{t,1}}{\partial w_{ij}^{(\bar{c})}} = \mathbf{D}_{t,1}^{(\sigma'(\zeta^{(f)}))} (\mathbf{b}_f \odot \boldsymbol{\lambda}_{t-1}^{(R^{(f)h})}) \quad (41)$$

$$\begin{aligned} \frac{\partial \tilde{\mathbf{c}}_{t,1}}{\partial w_{ij}^{(\bar{c})}} &= \mathbf{D}_{t,1}^{(\sigma'(\zeta^{(\bar{c})}))} ((\text{sign}(x_{t,1j}) + 2\delta(w_{ij}^{(\bar{c})})x_{t,1j})\mathbf{e}_i \\ &\quad + \mathbf{b}_{\bar{c}} \odot \boldsymbol{\lambda}_{t-1}^{(R^{(\bar{c})h})}). \end{aligned} \quad (42)$$

Using (40)–(42), we compute (39). Then, we compute (36) using (39) and (37) in order to calculate (35). After obtaining (35), we update the parameter using the SGD- and EG-based algorithms as in (27) and (29).

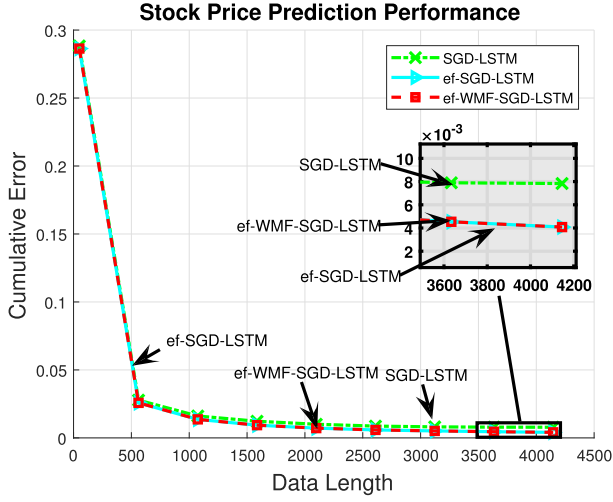


Fig. 3. Daily stock price prediction performances of the algorithms with the SGD updates on the Alcoa Corporation stock price data set.

As in Remark 5, when we take the derivative with respect to the other parameters, only the location of the term with  $e_i$  changes. Similar to the RNN case, when  $n_t \neq 1$ , the recursion in (39) is performed through the outputs of the different LSTM blocks at a certain time  $t$ , as shown in Fig. 1. Moreover, for the factorized LSTM network, we compute the derivatives of the loss function with respect to each factorized matrix parameter as in Remark 7.

With the derived gradients, we can update each parameter of the energy-efficient LSTM architecture as in (27) and (29).

#### IV. SIMULATIONS

In this section, we illustrate the performances of our algorithms on various data sets under different scenarios. We first compare the regression performances of our algorithms on a financial data set, i.e., the Alcoa Corporation stock price rate data set [31]. We then evaluate the regression performances on the real-life data sets, i.e., the kinematic [32] and elevators [33] data sets. Since our approach is generic, we also compare the performances of our training algorithms on two different RNNs, i.e., the LSTM and GRU neural networks. We then compare the structural complexity of our energy-efficient algorithms with the conventional structures.

Throughout this section, “Model 1” represents the conventional LSTM network (LSTM). Similarly, “Model 2” represents the introduced LSTM network with the ef-operator (ef-LSTM), and “Model 3” represents the introduced LSTM network with the ef-operator and weight matrix factorization (ef-WMF-LSTM).

##### A. Financial Data Set

In this section, we compare the performances of our algorithms on a financial data set. We consider the Alcoa Corporation stock price data set [31], for which we have the daily stock price values. In this case, our aim is to predict future stock prices based on the past prices, where we examine the past five days for prediction. Here, we evaluate the regression

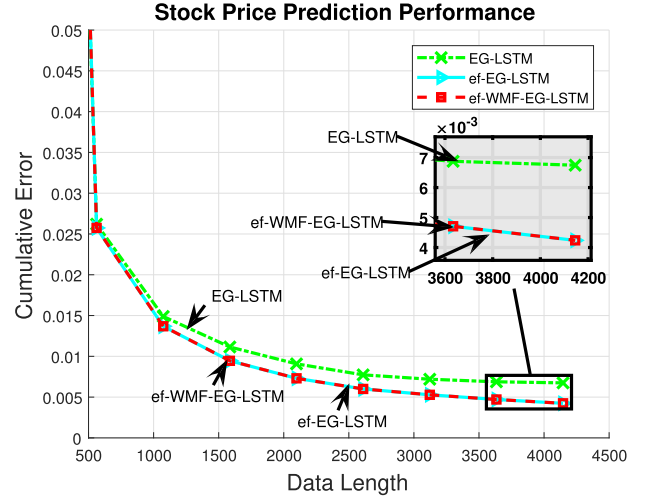


Fig. 4. Daily stock price prediction performances of the algorithms with the EG updates on the Alcoa Corporation stock price data set.

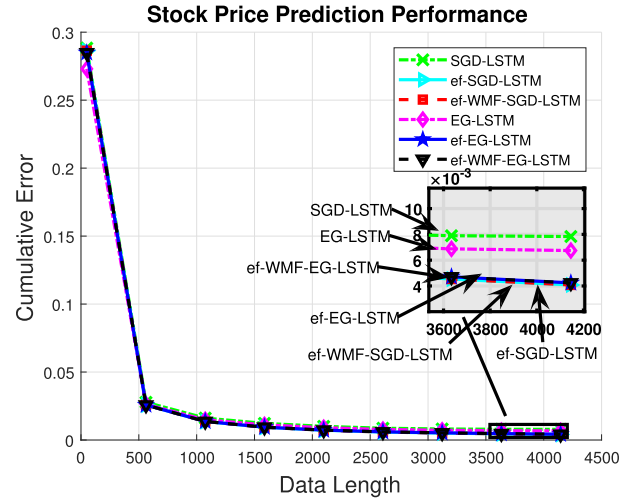


Fig. 5. Comparison of all the algorithms with the SGD and EG updates on the Alcoa Corporation stock price data set.

performance of Model 1 and consider this performance to be a benchmark for our proposed models, i.e., Model 2 and Model 3. To provide a fair setup, we select the same values for the common parameters of all the models. In addition, for Model 3, we set the rank of factorized LSTM weight matrices as 2 based on our observations in Section IV-C. For all these experiments, we perform 100 trials and plot the averaged curves. Moreover, we set the learning rate as  $\eta = 0.1$ ,  $\mathbf{X}_t \in \mathbb{R}^5$ , and the output dimensionality as  $m = 5$ .

Since Model 2 and Model 3 have an additive structure, the gradient of each parameter becomes more robust to the vanishing and exploding gradient problems. Thus, in Fig. 3, Model 2 and Model 3 outperform Model 1 in terms of the error performance. Although both the models, i.e., Model 2 and Model 3, perform similarly, we consider Model 3 as superior due to having a smaller number of network parameters to be trained (see details in Section IV-E). Likewise, in Fig. 4, Model 2 and Model 3 have smaller error than Model 1. In Fig. 5, we evaluate the combined results of all the models with both



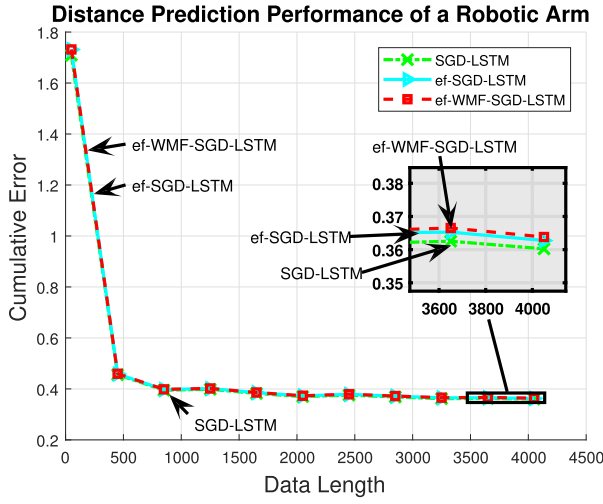


Fig. 6. Distance prediction performances of the algorithms with the SGD updates on the kinematic data set.

the SGD and EG updates. Model 2 and Model 3 with the SGD updates provide slightly smaller steady-state errors compared with all other models. Overall, Model 3 with the SGD updates outperforms its competitors in terms of both error performance and complexity issues.

### B. Real-Life Data Sets

In this section, we compare the performances of our algorithms using two real-life data sets, i.e., the kinematic [32] and elevators [33] data sets. We first evaluate the performances of the models on the kinematic data set [32], which contains the data related to a realistic simulation of the forward dynamics of an eight-link all-revolute robot arm. Our aim is to predict the distance of the end-effector from a target. In order to provide a fair experimental environment, we select the same common parameters for all the models. For this data set, the input vector is  $X_t \in \mathbb{R}^8$ ,  $m = 8$ , and the learning rate for all the models is  $\eta = 0.1$ . For Model 3, we select the rank of network matrices as 2. As shown in Fig. 6, all the models perform similarly. However, in terms of computational complexity and the total number of network parameters, Model 3 has the lowest complexity and the total number of parameters. In Fig. 7, Model 3 outperforms all other models thanks to having a smaller number of parameters and less complicated optimization problem for the parameters. In Fig. 8, we compare the models with both the SGD and EG updates. We observe that Model 1 with the SGD updates achieves a slightly smaller error compared with Model 2 and Model 3.

In addition to the kinematic data set, we also evaluate the performances on the elevators data set [33], which is obtained from the movements of an F16 aircraft, and we aim to predict the variable that expresses the movements of the aircraft. In this case, we have  $X_t \in \mathbb{R}^{18}$ ,  $m = 18$ , and select the learning rate as  $\eta = 0.1$ . The rank for Model 3 is 2. In Fig. 9, all the models using the EG updates outperform the models using the SGD updates, which arises from the

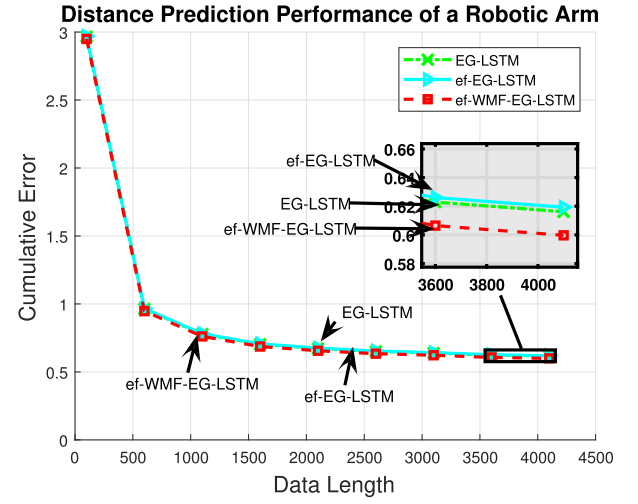


Fig. 7. Distance prediction performances of the algorithms with the EG updates on the kinematic data set.

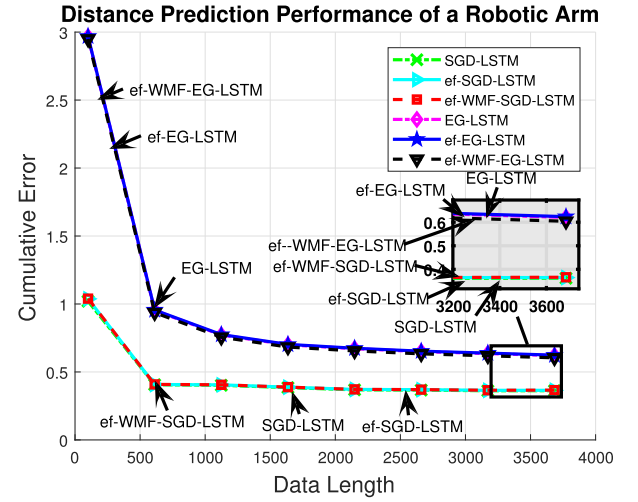


Fig. 8. Comparison of sequential prediction performances of the algorithms on the kinematic data set.

sparseness of  $X_t$  unlike the previous experiments. Among the models with the SGD updates, Model 3 has the smallest error, and for the EG updates, all the models provide comparable performances. Although all the models perform similarly, Model 3 is the ideal choice because of having less complexity and a smaller number of network parameters compared with the other models.

### C. Rank Effect on WMF

In this section, we illustrate the effects of the rank on the performances of the introduced models. For this purpose, we use the Alcoa Corporation stock price data set. In Table I, we observe that as the rank decreases, the training time also decreases for the LSTM-based WMF models with both the SGD and EG updates, and the errors stay approximately the same. From this fact, we conclude that we can use lower rank weight matrices for our proposed models and still get the same performance in less amount of time. Thus, with our approach,

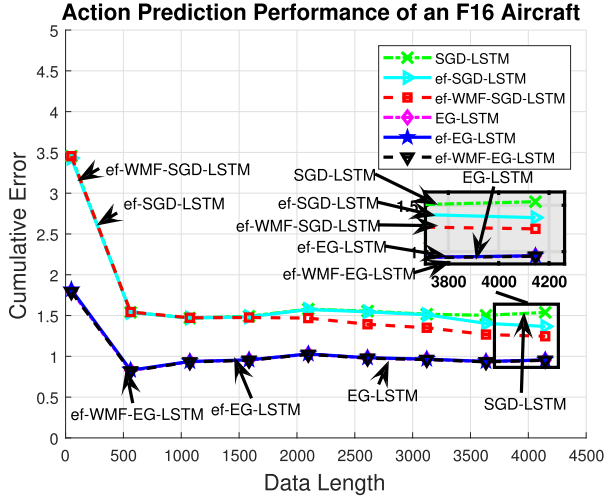


Fig. 9. Comparison of Movement prediction performances of the algorithms on the elevators data set.

TABLE I

TRAINING TIMES (IN s) FOR ONE TRIAL AND TIME ACCUMULATED ERRORS FOR THE WMF ALGORITHMS USING DIFFERENT RANK WEIGHT MATRICES ON THE ALCOA CORPORATION STOCK PRICE DATA SET. NOTE THAT THIS EXPERIMENT IS PERFORMED WITH A COMPUTER THAT HAS I5-6400 PROCESSOR, 2.7-GHz CPU, AND 16-GB RAM

	Rank=2		Rank=3	
	Error	Time	Error	Time
WMF-SGD-LSTM	$4.1 \times 10^{-3}$	103.15	$4 \times 10^{-3}$	120.31
WMF-EG-LSTM	$4.2 \times 10^{-3}$	115.38	$4.1 \times 10^{-3}$	138.2

one can significantly reduce the number of parameters to be trained in an LSTM network while enjoying high performance. Based on these observations, in all experiments, we select the rank as 2. Note that we do not reduce the rank to 1 since WMF significantly degrades the performance in that case.

#### D. LSTM and GRU Neural Networks

In this section, we evaluate the performances of the algorithms on the real-life and financial data sets. Since our approach is generic, in the sense, that it can be applied to any RNN structure, we also include the GRU-based algorithms to provide comparative analysis. The GRU network is defined by the following equations [26]:

$$\tilde{z}_{t,j} = \sigma(\mathbf{W}^{(\tilde{z})}\mathbf{x}_{t,j} + \mathbf{R}^{(\tilde{z})}\mathbf{h}_{t,j-1}) \quad (43)$$

$$\mathbf{r}_{t,j} = \sigma(\mathbf{W}^{(r)}\mathbf{x}_{t,j} + \mathbf{R}^{(r)}\mathbf{h}_{t,j-1}) \quad (44)$$

$$\tilde{y}_{t,j} = g(\mathbf{W}^{(y)}\mathbf{x}_{t,j} + \mathbf{r}_{t,j} \odot (\mathbf{R}^{(y)}\mathbf{h}_{t,j-1})) \quad (45)$$

$$\mathbf{y}_{t,j} = \tilde{y}_{t,j} \odot \tilde{z}_{t,j} + \mathbf{y}_{t,j-1} \odot (1 - \tilde{z}_{t,j}) \quad (46)$$

where  $\mathbf{x}_{t,j} \in \mathbb{R}^p$  is the input vector and  $\mathbf{y}_{t,j} \in \mathbb{R}^m$  is the output vector. Here,  $\tilde{z}_{t,j}$  and  $\mathbf{r}_{t,j}$  are the update and reset gates, respectively. The functions  $g(\cdot)$  and  $\sigma(\cdot)$  apply to vectors pointwise and commonly set to the  $\tanh(\cdot)$  and sigmoid functions, respectively. In order to obtain an energy-efficient version of the GRU network, we apply the matrix factorization method and the ef-operator as in the LSTM case.

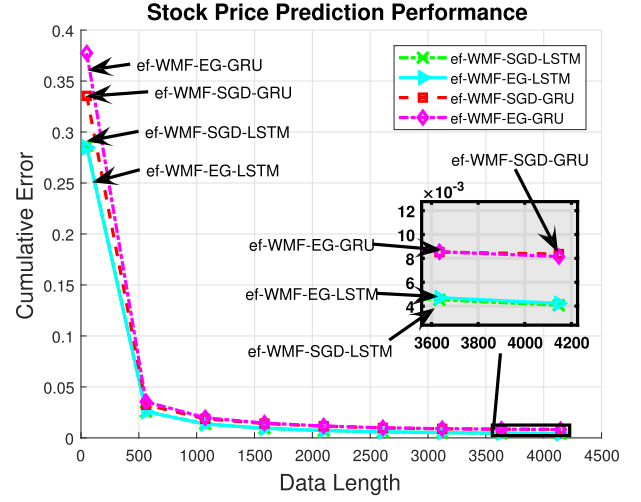


Fig. 10. Comparison of energy-efficient LSTM and GRU networks on the Alcoa Corporation data set.

TABLE II

RELATIVE ENERGY CONSUMPTIONS (IN pJ) OF THE INTRODUCED RNN NETWORKS AT EACH TIME STEP. HERE, WE USE THE ENERGY CONSUMPTION DATA OF ARITHMETIC OPERATIONS FOR A 45-nm CMOS PROCESS [34]

Networks Datasets	LSTM		GRU	
	LSTM	ef-LSTM	GRU	ef-GRU
Alcoa	980	526	741	390
Kinematic	2451.2	1187.2	1848	883.2
Elevators	12139	5263.2	9126	3931.2

Here, we select the same parameters with Sections IV-A–IV-C. Since Model 3 provides the best performance in Sections IV-A–IV-C, we compare the LSTM and GRU Networks on three data sets using Model 3. For the Alcoa Corporation stock price data set, the LSTM-based algorithm with the SGD updates achieves the smallest steady-state error, as shown in Fig. 10. In Fig. 11, the GRU-based algorithm with the SGD updates outperforms the other network models. For the elevators data set, the LSTM-based algorithm with the EG updates achieves the smallest steady-state error among all the network models, as shown in Fig. 12. Overall, since the LSTM architecture has an output gate to control its memory content unlike the GRU network, it generally outperforms the GRU network on various real-life scenarios.

Moreover, in order to illustrate the energy efficiency of the introduced architectures, we provide energy consumption data for each data set in Table II. We observe that our approach almost halves the energy consumption for each case, and as the dimensionality of the data set increases, the provided energy efficiency even further increases.

#### E. Training Times and Structural Complexity

In this section, we first provide the training times (in s) of all the LSTM network-based models with both the SGD

TABLE III  
TRAINING TIMES (IN s FOR ONE TRIAL) OF THE INTRODUCED ENERGY-EFFICIENT LSTM NETWORKS

Algorithms Datasets	SGD-LSTM	ef-SGD-LSTM	ef-WMF-SGD-LSTM	EG-LSTM	ef-EG-LSTM	ef-WMF-EG-LSTM
Alcoa	55.71	144.98	103.15	66.43	155.07	115.38
Kinematic	99.84	181.10	167.50	117.01	200.61	178.17
Elevators	545.97	1052.59	384.57	676.73	1185.97	404.44

TABLE IV  
TIME ACCUMULATED ERRORS OF THE INTRODUCED ALGORITHMS

Algorithms Datasets	SGD-LSTM	ef-SGD-LSTM	ef-WMF-SGD-LSTM	EG-LSTM	ef-EG-LSTM	ef-WMF-EG-LSTM
Alcoa	0.0084	0.0041	0.0041	0.0069	0.0042	0.0042
Kinematic	0.3595	0.3621	0.3629	0.6139	0.6171	0.5975
Elevators	1.5607	1.367	1.2478	0.9587	0.9611	0.9564

TABLE V  
TOTAL NUMBER OF PARAMETERS TO BE LEARNED FOR THE INTRODUCED NETWORKS

Algorithms Datasets	LSTM	ef-LSTM	ef-WMF-LSTM	GRU	ef-GRU	ef-WMF-GRU
Alcoa	220	260	220	150	180	150
Kinematic	544	608	352	384	432	240
Elevators	2664	2808	792	1944	2052	540

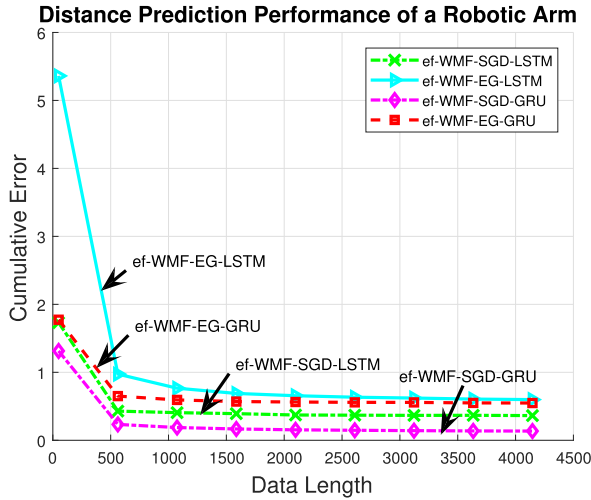


Fig. 11. Comparison of energy-efficient LSTM and GRU networks on the kinematic data set.

and EG updates. We then give the total number of network parameters for each model, i.e., the structural complexity of the corresponding model. Finally, we compare all the models based on the training times, the number of network parameters, and error performance.

In Table III, we provide the training times of all the network models for each data set. Note that all the experiments are performed with a computer that has i5-6400 processor, 2.7-GHz CPU, and 16-GB RAM. Among all the network models, Model 1 has the fastest training performance when the data size is small, and however, it does not have the smallest

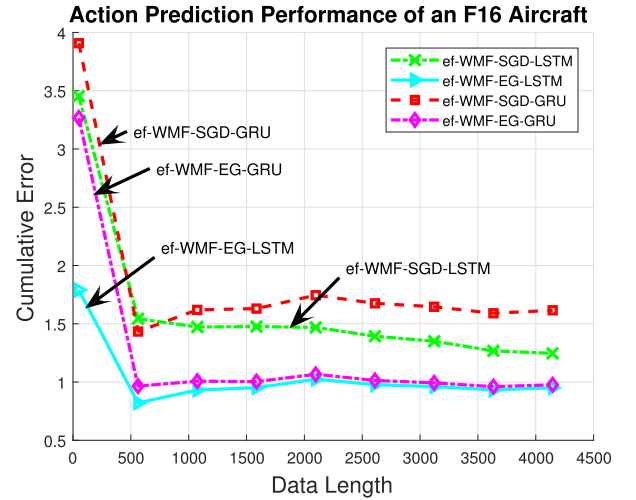


Fig. 12. Comparison of energy-efficient LSTM and GRU networks on the elevators data set.

time accumulated errors, i.e., defined as

$$\sum_{t=1}^T (d_t - \hat{d}_t)^2$$

for a data set with  $T$  samples, as stated in Table IV. Model 3 achieves intermediate training times (and the fastest training when the data size is large as in the elevators data set) and the smallest cumulative errors, as shown in Tables III and IV, respectively. In Table V, we provide the total number of network parameters for each model. We observe that Model 3 has the smallest number of network parameters among all

other models thanks to the weight matrix factorization. Overall, based on training times, error performances, and the number of network parameters, Model 3 is the best choice among all the models.

## V. CONCLUSION

In this article, we have studied the variable-length data regression in an online framework and introduced an energy-efficient regression structure based on the LSTM network. In particular, we have introduced a generic LSTM-based regression structure to obtain fixed-length representations from variable-length data sequences. In order to reduce the complexity of this structure, we first eliminate the regular multiplications by replacing them with an energy-efficient operator, i.e., the ef-operator. We then apply a factorization method to all the matrices in the classical LSTM network in order to diminish the total number of parameters. For this energy-efficient and factorized LSTM network, we have introduced the online training algorithms based on the SGD [11] and EG [17] algorithms. Hence, we obtain highly efficient and effective online learning algorithms based on the LSTM network. Thanks to the generic structure of our approach, we have also introduced an energy-efficient GRU network in our simulations. Through several experiments involving real and financial data, we demonstrate significant performance improvements and complexity reductions achieved by the introduced algorithms with respect to the conventional methods.

## REFERENCES

- [1] D. F. Specht, "A general regression neural network," *IEEE Trans. Neural Netw.*, vol. 2, no. 6, pp. 568–576, Nov. 1991.
- [2] N. Wang, M. J. Er, and M. Han, "Generalized single-hidden layer feedforward networks for regression problems," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 6, pp. 1161–1176, Jun. 2015.
- [3] T. Ding and A. Hirose, "Fading channel prediction based on combination of complex-valued neural networks and chirp Z-transform," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 9, pp. 1686–1695, Sep. 2014.
- [4] J.-T. Chien and Y.-C. Ku, "Bayesian recurrent neural network for language modeling," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 27, no. 2, pp. 361–374, Feb. 2016.
- [5] A. C. Tsoi and A. C. Tsoi, "Gradient based learning methods," in *Adaptive Processing of Sequences and Data Structures: International Summer School on Neural Networks 'E. R. Caianiello' Vietri sul Mare, Salerno, Italy September 6–13, 1997 Tutorial Lectures*, G. L. Giles and M. Gori, Eds. Berlin, Germany: Springer, 1998, pp. 27–92. doi: 10.1007/BFb0053994.
- [6] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Trans. Neural Netw.*, vol. 5, no. 2, pp. 157–166, Mar. 1994.
- [7] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *Proc. ICML*, vol. 28, 2013, pp. 1310–1318.
- [8] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [9] J. Mazumdar and R. G. Harley, "Recurrent neural networks trained with backpropagation through time algorithm to estimate nonlinear load harmonic currents," *IEEE Trans. Ind. Electron.*, vol. 55, no. 9, pp. 3484–3491, Sep. 2008.
- [10] H. Jaeger, "Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the echo state network approach," *GMD-Forschungszentrum Informationstechnik*, vol. 5, Jan. 2002.
- [11] A. W. Smith and D. Zipser, "Learning sequential structure with the real-time recurrent learning algorithm," *Int. J. Neural Syst.*, vol. 1, no. 2, pp. 125–131, 1989.
- [12] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, and J. Schmidhuber, "LSTM: A search space odyssey," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 10, pp. 2222–2232, Oct. 2016.
- [13] F. A. Gers, J. A. Pérez-Ortiz, D. Eck, and J. Schmidhuber, "DEKF-LSTM," in *Proc. ESANN*, 2002, pp. 369–376.
- [14] J. A. Pérez-Ortiz, F. A. Gers, D. Eck, and J. Schmidhuber, "Kalman filters improve LSTM network performance in problems unsolvable by traditional recurrent nets," *Neural Netw.*, vol. 16, no. 2, pp. 241–250, 2003.
- [15] D. Monner and J. A. Reggia, "A generalized LSTM-like training algorithm for second-order recurrent neural networks," *Neural Netw.*, vol. 25, pp. 70–83, Jan. 2012. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0893608011002036>
- [16] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional LSTM and other neural network architectures," *Neural Netw.*, vol. 18, no. 5, pp. 602–610, 2005. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0893608005001206>
- [17] J. Kivinen and M. K. Warmuth, "Exponentiated gradient versus gradient descent for linear predictors," *Inf. Comput.*, vol. 132, no. 1, pp. 1–63, 1997.
- [18] S. I. Hill and R. C. Williamson, "Convergence of exponentiated gradient algorithms," *IEEE Trans. Signal Process.*, vol. 49, no. 6, pp. 1208–1215, Jun. 2001.
- [19] H. Tuna, I. Onaran, and A. E. Cetin, "Image description using a multiplier-less operator," *IEEE Signal Process. Lett.*, vol. 16, no. 9, pp. 751–753, Sep. 2009.
- [20] C. E. Akbaş, A. Bozkurt, A. E. Çetin, R. Çetin-Atalay, and A. Üner, "Multiplication-free neural networks," in *Proc. 23rd Signal Process. Commun. Appl. Conf. (SIU)*, May 2015, pp. 2416–2418.
- [21] M. S. Razlighi, M. Imani, F. Koushanfar, and T. Rosing, "LookNN: Neural network with no multiplication," in *Proc. Design Autom. Test Eur. Conf. Exhib. (DATE)*, Mar. 2017, pp. 1775–1780.
- [22] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 1135–1143.
- [23] A. Afrasiyabi, B. Nasir, O. Yildiz, F. T. Y. Vural, and A. E. Cetin, "An energy efficient additive neural network," in *Proc. 25th Signal Process. Commun. Appl. Conf. (SIU)*, May 2017, pp. 1–4.
- [24] O. Kuchaiev and B. Ginsburg, "Factorization tricks for LSTM networks," *CoRR*, vol. abs/1703.10722, 2017. [Online]. Available: <http://arxiv.org/abs/1703.10722>
- [25] N. Srinivasan, V. Ravichandran, K. L. Chan, J. R. Vidhya, S. Ramakrishnan, and S. M. Krishnan, "Exponentiated backpropagation algorithm for multilayer feedforward neural networks," in *Proc. 9th Int. Conf. Neural Inf. Process. (ICONIP)*, vol. 1, Nov. 2002, pp. 327–331.
- [26] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," 2014, *arXiv:1412.3555*. [Online]. Available: <https://arxiv.org/abs/1412.3555>
- [27] E. Levin, "A recurrent neural network: Limitations and training," *Neural Netw.*, vol. 3, no. 6, pp. 641–650, 1990. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/0893608090900540>
- [28] A. Afrasiyabi, O. Yildiz, B. Nasir, F. T. Yarman-Vural, and A. E. Çetin, "Energy saving additive neural network," *CoRR*, 2017. [Online]. Available: <https://arxiv.org/abs/1702.02676>
- [29] N. D. Vanli and S. Kozat, "A comprehensive approach to universal piecewise nonlinear regression based on trees," *IEEE Trans. Signal Process.*, vol. 62, no. 20, pp. 5471–5486, Oct. 2014.
- [30] R. N. Bracewell and R. N. Bracewell, *The Fourier Transform & Its Applications*, vol. 31999. New York, NY, USA: McGraw-Hill, 1986.
- [31] *Summary for Alcoa Inc. Common Stock*. Accessed: Apr. 1, 2018. [Online]. Available: <http://finance.yahoo.com/quote/AA?litr=1>
- [32] C. E. Rasmussen *et al.* *Delve Data Sets*. Accessed: Apr. 1, 2018. [Online]. Available: <http://www.cs.toronto.edu/~delve/data/datasets.html>
- [33] J. Alcalá-Fdez, A. Fernández, J. Luengo, J. Derrac, and S. García, "KEEL data-mining software tool: Data set repository, integration of algorithms and experimental analysis framework," *J. Multiple-Valued Logic Soft Comput.*, vol. 17, nos. 2–3, pp. 255–287, 2011.
- [34] M. Horowitz, "Energy table for 45nm process," Stanford VLSI wiki.



**Tolga Ergen** received the B.S. and M.S. degrees in electrical and electronics engineering from Bilkent University, Ankara, Turkey, in 2016 and 2018, respectively. He is currently pursuing the Ph.D. degree with the Electrical Engineering Department, Stanford University, Stanford, CA, USA. His current research interests include machine learning, optimization, and neural networks.





**Ali H. Mirza** received the B.Sc. degree (Hons.) in electrical engineering from the University of Engineering and Technology, Lahore, Pakistan, in 2014. He is currently pursuing the Ph.D. degree with the Department of Electrical and Electronics Engineering, Bilkent University, Ankara, Turkey.

His current research interests include machine learning, big data signal processing, online learning, and deep neural networks.



**Suleyman Serdar Kozat** (A'10–M'11–SM'11) received the B.S. degree (Hons.) from Bilkent University, Ankara, Turkey, and the M.S. and Ph.D. degrees in electrical and computer engineering from the University of Illinois at Urbana–Champaign, Urbana, IL, USA.

He joined the IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA, as a Research Staff Member and later became a Project Leader with the Pervasive Speech Technologies Group, where he focused on problems related to statistical signal processing and machine learning. He was a Research Associate with the Cryptography and Anti-Piracy Group, Microsoft Research, Redmond, WA, USA. He is currently a Professor with the Electrical and Electronics Engineering Department, Bilkent University. He has coauthored over 200 papers in refereed high-impact journals and conference proceedings and holds several patent inventions (used in several different Microsoft and IBM products). He holds several patent inventions due to his research accomplishments with the IBM Thomas J. Watson Research Center and Microsoft Research. His current research interests include cybersecurity, anomaly detection, big data, data intelligence, adaptive filtering, and machine learning algorithms for signal processing.

Dr. Kozat received many international and national awards. He is the Elected President of the IEEE Signal Processing Society, Turkey Chapter.