

# LONG SHORT-TERM MEMORY

Neural Computation 19(8):1873–1908, 1997

Sepp Hochreiter

Fakultät für Informatik

Technische Universität München

8020 München, Germany

<http://www5.informatik.tu-muenchen.de>

<http://www5.informatik.tu-muenchen.de/~hochrei>

Jürgen Schmidhuber

TUM

Clausius-Drexlin 26

8001 Zugspitze, Switzerland

[jurgen@tik.tum.de](mailto:jurgen@tik.tum.de)

<http://www.tik.tum.de/~jurgen>

## Abstract

Learning to store information over extended time intervals via recurrent backpropagation suffers strong long-term, mostly linear, error flow. We briefly review Hochreiter's 1991 analysis of this problem, then address it by introducing a novel, efficient, gradient-based method called "Long Short-Term Memory" (LSTM). Truncating the gradient where this does not do harm, LSTM can learn to bridge minimal time lags in excess of 1000 discrete time steps by enforcing *constant error reservoirs* within special units. Multiplicative gate units learn to open and close access to the constant error flows. LSTM is local in space and time; its computational complexity per time step and weight is O(1). Our experiments with artificial data involve local, distributed, real-valued, and noisy pattern representations. In comparisons with RTRL, BPTT, Recurrent Cascade-Correlation, Elman nets, and Neural Sequence Chunking, LSTM leads to many more successful runs, and learns much faster. LSTM also solves complex, artificial long time lag tasks that have never been solved by previous recurrent network algorithms.

## 1 INTRODUCTION

Recurrent networks learn in principle via their feed-back connections to store representations of recent input events in form of activations [“short-term memory”, as opposed to “long-term memory” manifested by slowly changing weights]. This is potentially significant for many applications, including speech processing, non-Markovian control, and music composition [e.g., Mozer 1992]. The most widely used algorithms for learning *when* to put in short-term memory, however, take much time or do not work well at all, especially when minimal time lags between inputs and corresponding target signals are long. Although theoretically fascinating, existing methods do not provide clear *practical* advantages over, say, backprop in feed-forward nets with limited time windows. This paper will review an analysis of the problem and suggest a remedy.

The problem, with conventional “Back-Propagation Through Time” (BPTT, e.g., Williams and Zipser 1989, Werbos 1988) or “Real-Time Recurrent Learning” (RTRL, e.g., Bellman and Gallistel 1985), error signals “flowing backwards in time” tend to either [1] blow up or [2] vanish: the temporal evolution of the backpropagated error exponentially depends on the size of the weights [Hochreiter 1991]. Case [1] may lead to oscillating weights, while in case [2] learning to bridge long time lags takes a prohibitively amount of time, or does not work at all [see section 8].

This remedy: This paper presents “Long Short-Term Memory” (LSTM), a novel recurrent network architecture in conjunction with an appropriate gradient-based learning algorithm. LSTM is designed to overcome these error back-flow problems. It can learn to bridge time intervals in excess of 1000 steps even in case of noisy, incompressible input sequences, without loss of short time lag capabilities. This is achieved by an efficient, gradient-based algorithm for an architecture

enforcing *constant* [thus neither neglecting nor vanishing] error flow through internal states of sigmoid units [provided the gradient computation is truncated at certain architecture-specific points — this does not affect long-term error flow though].

**Outline of paper.** Section 2 will briefly review previous work. Section 3 begins with an outline of the detailed analysis of vanishing errors due to Hochreiter [1991]. It will then introduce a new approach to constant error backpropagation for sigmoidal purges, and highlight its problems concerning information storage and retrieval. These problems will lead to the LSTM architecture as described in Section 4. Section 5 will present numerous experiments and comparisons with competing methods. LSTM outperforms them, and also learns to solve complex, artificial tasks no other recurrent net algorithm has solved. Section 6 will discuss LSTM’s limitations and advantages. This appendix contains a detailed description of the algorithm [A.1], and explicit error flow formulas [A.2].

## 2 PREVIOUS WORK

This section will focus on recurrent nets with time-varying inputs [as opposed to nets with stationary inputs and step-function gradient calculations, e.g., Almeida 1987, Elman 1986].

**Gradient-thickening variants.** The approaches of Elman [1986], Elman [1990], Williams [1989], Schmidhuber [1990a], Baumgartner [1993], and many of the related algorithms in Baumgartner’s comprehensive overview [1993] suffer from the same problems as RPTT and RTRL [see Sections 1 and 8].

**Time-lagging.** Other methods that avoid gradient flow for short time lags only are Time-lagging Neural Networks [Lam et al. 1990] and Elman’s method [Elman 1990], which updates unit activations based on a weighted sum of old activations [see also De Weix and Brincker 1991]. Lin et al. [1995] propose variants of time-lag networks called NEON networks.

**Time constants.** To deal with long time lags, Mozer [1992] uses time constants influencing changes of unit activations [DeWeix and Brincker’s above-mentioned approach [1991] may in fact be viewed as a mixture of TBRNN and time constants]. For long time lags, however, the time constants need external fine tuning [Mozer 1992]. Sun et al.’s alternative approach [1993] updates the activation of a recurrent unit  $h_t$  adding the old activation and the [scaled] current net input. The net input, however, tends to forget the stored information, which makes long-term storage impractical.

**Rings’s approach.** Ring [1993] also proposed a method for bridging long time lags. Whenever a unit in his network receives conflicting error signals, he adds a higher-order unit influencing inappropriate connections. Although this approach can sometimes do extremely fast, to bridge a time lag involving 100 steps may require the addition of 100 units. Also, Ring’s net does not generalize to unseen lag durations.

**Fengie et al.’s approaches.** Fengie et al. [1991] investigates methods such as simultaneous annealing, multi-grid random search, time-weighted pseudo-Newton optimization, and discrete error propagation. Their “lattice” and “ $k$ -sequence” gradients are very similar to gradient flow with minimal time lag 100 [see Deppermann 93]. Fengie and Frasconi [1993] also propose an RNN approach for propagating targets. With so-called “state networks”, at a given time, their system can be in one of only  $n$  different states. See also beginning of Section 5. But to solve continuous problems such as the “adding problem” [Section 5.1], their system would require an unacceptable number of states [i.e., state networks].

**Nelmann filters.** Huskerius and Nelmann [1993] use Nelmann filter techniques to improve recurrent net performance. Since they use “a recursive discount factor imposed to reduce exponentially the effects of past dynamics shortcomings,” there is no reason to believe that their Nelmann Filter Trained Recurrent Networks will be useful for very long minimal time lags.

**Second-order nets.** We will see that LSTM uses multiplicative units (MUs) to protect error flow from unwanted perturbations. It is not the first recurrent net method using MUs though. For instance, Watrous and Nullin [1992] use MUs in second-order nets. Some differences to LSTM are: (1) Watrous and Nullin’s architecture does not enforce constant error flow and is not designed

to solve long time lag problems. [2] It has fully connected second-order sigma-pi units, while the LSTM architecture’s MUs are used only to gate access to constant error flow. [3] Waturus and Stollnitz’s algorithm costs  $\mathcal{O}(M^2)$  operations per time step, uses only  $\mathcal{O}(M)$ , where  $M$  is the number of weights. See also Miller and Giles [1998] for additional work on MUs.

**Simple weight guessing.** To avoid long time lag problems of gradient-Descent approaches we may simply randomly initialize all network weights until the resulting net happens to classify all training sequences correctly. In fact, recently we discovered [Stollnitz/Miller and Hochreiter 1993, Hochreiter and Stollnitz/Miller 1993, 1997] that simple weight guessing solves many of the problems in [Hongie 1993], Hongie and Braverman 1993, Miller and Giles 1993, Lin et al. 1993] faster than the algorithms proposed therein. This does not mean that weight guessing is a good algorithm. It just means that the problems are very simple. More realistic tasks require either many free parameters (e.g., input weights) or high weight precision (e.g., for continuous-valued parameters), such that guessing becomes hopelessly infeasible.

**All-optimal sequence chunkers.** Stollnitz/Miller’s hierarchical chunker systems [1993E, 1993F] do have a capability to bridge arbitrary time lags, but only if there is local gradability across the subsequences causing the time lags (see also Mezey 1992). For instance, in his postdoctoral thesis [1993F], Stollnitz/Miller uses hierarchical recurrent nets to rapidly solve certain grammar learning tasks involving minimal time lags in excess of 1000 steps. The performance of chunker systems, however, deteriorates as the noise level increases and the input sequences becomes less compressible. LSTM does not suffer from this problem.

## 3 CONSTANT ERROR BACKPROP

### 3.1 EXPONENTIALITY DECAYING ERROR

Conventional RNN [e.g. Williams and Zipser 1992]. Output unit  $k$ ’s target at time  $t$  is denoted by  $y_k[t]$ . Using mean squared error,  $k$ ’s error signal is

$$e_k[t] = \mathbb{E}_k[\text{error}_k[t]](y_k[t] - y^k[t]),$$

where

$$y^k[t] = \mathbb{E}_k[\text{out}_k[t]]$$

is the activation of a non-input unit  $k$  with differentiable activation function  $\mathbb{E}_k$ ,

$$\text{out}_k[t] = \sum_i w_{ki} y^i[t] - 1$$

is unit  $k$ ’s current net input, and  $w_{ki}$  is the weight on the connection from unit  $i$  to  $k$ . Some non-output unit  $k$ ’s backpropagated error signal is

$$e_k[t] = \mathbb{E}_k[\text{out}_k[t]] \left( \sum_i w_{ki} e_i[t] + 1 \right).$$

The corresponding contribution to  $w_{ki}$ ’s total weight update is  $m\alpha_k[t]y^i[t](-1)$ , where  $m$  is the learning rate, and  $i$  stands for an arbitrary unit connected to unit  $k$ .

Outline of Hochreiter’s analysis [1991, pages 12-13]. Suppose we have a fully connected net where non-input unit indices range from 1 to  $n$ . Let us focus on local error flows from unit  $m$  to unit  $n$  [later we will see that this analysis immediately extends to global error flows]. The error occurring at an arbitrary unit  $m$  at time step  $t$  is propagated “back in time” for  $y$  time steps, to an arbitrary unit  $n$ . This will scale the error by the following factor:

$$\frac{\partial e_n[t-y]}{\partial e_m[t]} = \begin{cases} \frac{\mathbb{E}_n[\text{out}_n[t-y-1]]w_{mn}}{\mathbb{E}_n[\text{out}_n[t-y]](\sum_{i=1}^n \frac{\partial e_i[t-y+1]}{\partial e_m[t]} w_{ni})} & y = 1 \\ \frac{\mathbb{E}_n[\text{out}_n[t-y-1]]w_{mn}}{\mathbb{E}_n[\text{out}_n[t-y]](\sum_{i=1}^n \frac{\partial e_i[t-y+1]}{\partial e_m[t]} w_{ni})} & y > 1 \end{cases}. \quad [1]$$

With  $b_x = n$  and  $b_0 = m$ , we obtain:

$$\frac{\partial E_n[\hat{y} - y]}{\partial \theta_m[\hat{y}]} = \sum_{l_1=1}^m \dots \sum_{l_{q-1}=1}^m \prod_{m=1}^q |\theta'_{m,l_m} [m\theta_{l_m}[\hat{y} - m]]| m_{l_m l_{m-1}} \quad [2]$$

[proof by induction]. The sum of the  $m^{q-1}$  terms  $\prod_{m=1}^q |\theta'_{m,l_m} [m\theta_{l_m}[\hat{y} - m]]| m_{l_m l_{m-1}}$  determines the total error flow [note that since the summation terms may have different signs, increasing the number of units  $m$  does not necessarily increase error flow].

Intuitive complementation of equation [2]. If

$$|\theta'_{m,l_m} [m\theta_{l_m}[\hat{y} - m]]| m_{l_m l_{m-1}} | > 1.0$$

for all  $m$ , then error flow propagates, e.g., with linear  $\theta'_{m,l_m}$  [then the largest product increases exponentially with  $q$ . That is, the error flows up, and conflicting error signals arriving at unit  $n$  can lead to oscillating weights and unstable learning [for error blow-ups or bifurcations see also Hinrichs 1988, Raab et al. 1991, Elaydi 1992]. On the other hand, if

$$|\theta'_{m,l_m} [m\theta_{l_m}[\hat{y} - m]]| m_{l_m l_{m-1}} | < 1.0$$

for all  $m$ , then the largest product decreases exponentially with  $q$ . That is, the error vanishes, and nothing new is learned in acceptable time.

If  $\theta'_{m,l_m}$  is the logistic sigmoid function, then the maximal value of  $|\theta'_{m,l_m}|$  is 0.5. If  $y^{l_{m-1}}$  is constant and not equal to zero, then  $|\theta'_{m,l_m} [m\theta_{l_m}[\hat{y} - m]]|$  takes on maximal values without

$$m_{l_m l_{m-1}} = \frac{1}{y^{l_{m-1}}} \coth \left[ \frac{1}{2} m \theta_{l_m} \right],$$

goes to zero for  $|m_{l_m l_{m-1}}| = \infty$ , and is less than 1.0 for  $|m_{l_m l_{m-1}}| < \infty$  [e.g., if the absolute maximal weight value  $m_{max}$  is smaller than 1.0]. Hence with conventional logistic sigmoid activation functions, the error flow tends to vanish as long as the weights have absolute values below 0.5, especially in the beginning of the training phase. In general the use of larger initial weights will not help although — as soon as  $m$ , for  $|m_{l_m l_{m-1}}| = \infty$  the relevant derivative goes to zero “faster” than the absolute weight can grow [also, some weights will have to change their signs by crossing zero]. Likewise, increasing the learning rate does not help either — it will not change the ratio of long-range error flow and short-range error flow. RPROP is more sensitive to recent distributions. A very similar, more recent analysis was presented by Elaydi et al. 1992].

**Critical error flow.** The local error flow analysis allows immediately看出 that global error flow vanishes, too. To see this, compute

$$\sum_{m: m \text{ output unit}} \frac{\partial E_n[\hat{y} - y]}{\partial \theta_m[\hat{y}]}.$$

**Weak upper bound for scaling factor.** The following, slightly modified vanishing error analysis also takes  $m$ , the number of units, into account. For  $y > 1$ , formula [2] can be rewritten as

$$(\mathbf{M}_{n,T} \mathbf{I}^T \mathbf{I}^T [\hat{y} - 1]) \prod_{m=1}^{q-1} |\theta'_{m,l_m} [\hat{y} - m]| + \mathbf{M}_n \theta'_n [\hat{y} \theta_n [\hat{y} - q]],$$

where the weight matrix  $\mathbf{M}$  is defined by  $[\mathbf{M}]_{ij} := m_{ij}$ ,  $i$ 's outgoing weight vector  $\mathbf{M}_i$  is defined by  $[\mathbf{M}_i]_k := [\mathbf{M}]_{ik} = m_{ik}$ ,  $i$ 's incoming weight vector  $\mathbf{M}_{i,T}$  is defined by  $[\mathbf{M}_{i,T}]_k := [\mathbf{M}]_{ik} = m_{ik}$ , and for  $m = 1, \dots, q$ ,  $[\theta']_{l_m} [\hat{y} - m]$  is the diagonal matrix of first order derivatives defined as:  $[\theta']_{l_m} [\hat{y} - m]_{ij} := 1$  if  $i = l_m$ , and  $[\theta']_{l_m} [\hat{y} - m]_{ij} := |\theta'_{l_m} [\hat{y} - m]|$  otherwise. Here  $\mathbf{I}$  is the transposition operator,  $[\mathbf{M}]_{ij}$  is the element in the  $i$ -th column and  $j$ -th row of matrix  $\mathbf{M}$ , and  $[\mathbf{v}]_k$  is the  $k$ -th component of vector  $\mathbf{v}$ .

Using a matrix norm  $\|\cdot\|_{\square}$  compatible with vector norm  $\|\cdot\|_v$ , we define

$$\mathbf{B}'_{\text{max}} := \max_{m=1,\dots,n} \{\|\mathbf{B}'[y - m]\|_{\square}\}.$$

For  $m \in \{1, \dots, n\}$ ,  $\|\alpha_m\|_v \leq \|\alpha\|_v$  we get  $|\alpha^T y| \leq m \|\alpha\|_v \|y\|_v$ . Since

$$|\mathbf{B}'[m\alpha]y - z| \leq \|\mathbf{B}'[y - z]\|_{\square} \leq \mathbf{B}'_{\text{max}},$$

we obtain the following inequality:

$$\left| \frac{\partial \mathbf{B}_x[y - z]}{\partial \mathbf{B}_x[y]} \right| \leq m \|\mathbf{B}'_{\text{max}}\|^2 \|\mathbf{B}_x\|_v \|\mathbf{B}_{xT}\|_v \|\mathbf{B}\|_E^{-2} \leq m \|\mathbf{B}'_{\text{max}}\| \|\mathbf{B}\|_E^2.$$

This inequality results from

$$\|\mathbf{B}_x\|_v = \|\mathbf{B}\mathbf{B}_x\|_v \leq \|\mathbf{B}\|_{\square} \|\mathbf{B}_x\|_v \leq \|\mathbf{B}\|_{\square}$$

and

$$\|\mathbf{B}_{xT}\|_v = \|\mathbf{x}_T \mathbf{B}_x\|_v \leq \|\mathbf{x}_T\|_{\square} \|\mathbf{B}_x\|_v \leq \|\mathbf{B}_x\|_{\square},$$

whereas  $\mathbf{x}_T$  is the unit vector whose components are 0 except for the  $M$ -th component, which is 1. Note that this is an event, therefore some upper bound — it will be reached only if all  $|\mathbf{B}'[y - m]|$  take maximal values, and if the contributions of all parallel errors will all error flows back from unit  $m$  to unit  $n$  with the same sign. Large  $\|\mathbf{B}\|_{\square}$ , however, typically result in small values of  $\|\mathbf{B}'[y - z]\|_{\square}$ , as confirmed by experiments (see, e.g., Heedlmüller 1991).

For example, with norms

$$\|\mathbf{B}\|_{\square} := \max_n \sum_i |m_{ni}|$$

and

$$\|\alpha\|_v := \max_i |\alpha_i|,$$

we have  $\mathbf{B}'_{\text{max}} = 0.25$  for the logistic sigmoid. We observe that if

$$|m_{ij}| \leq m_{\text{max}} < \frac{1.0}{\mathbf{B}'_{\text{max}}} \approx 4,$$

then  $\|\mathbf{B}\|_{\square} \leq m_{\text{max}} < 1.0$  will result in exponential decay — by setting  $\alpha := \left\lfloor \frac{m_{\text{max}}}{0.1} \right\rfloor < 1.0$ , we obtain

$$\left| \frac{\partial \mathbf{B}_x[y - z]}{\partial \mathbf{B}_x[y]} \right| \leq m \|\mathbf{B}\|_E^2.$$

We refer to Heedlmüller's 1991 thesis for additional results.

## 3.2 CONSTRAIN ERROR FLOW: NAIVE APPROACH

A single unit. To avoid cancelling error signals, there can not exist no constant error flow through a single unit  $y$  with a single connection to itself. According to the rules above, at time  $t$ ,  $y$ 's local error flow is  $\mathbf{B}_y[y] = \mathbf{B}_y[m\alpha_y[y]]\mathbf{B}_y[y] + 1|m_y|$ . To enforce *constant* error flow through  $y$ , we require

$$\mathbf{B}_y[m\alpha_y[y]]m_y = 1.0.$$

Note the similarity to Mozer's dead time constant system [1992] — a time constant of 1.0 is appropriate for potentially infinite time lags<sup>1</sup>.

The constant error approach. Integrating the differential equation above, we obtain  $\mathbf{B}_y[m\alpha_y[y]] = \frac{m\alpha_y[y]}{m_y}$  for arbitrary  $m\alpha_y[y]$ . This means:  $\mathbf{B}_y$  has to be linear, and unit  $y$ 's activation has to remain constant:

$$y_t[y+1] = \mathbf{B}_y[m\alpha_y[y+1]] = \mathbf{B}_y[m\alpha_y[y]]y_t[y] = y_t^*[y].$$

<sup>1</sup>We do not use the designation "time constant" in the differential sense, as, e.g., Heedlmüller [1991].

In this experiments, this will be ensured by using the identity function  $\mathbf{I}_k : \mathbf{I}_k[a] = a$ ,  $\forall a$ , and by setting  $m_{\text{gg}} = 1.0$ . We refer to this as the constant error approach [CEA]. CIDE will be LSTM's central feature [see Section 3].

Of course unit  $\mathbb{I}$  will not only be connected to itself but also to other units. This involves two obvious, related problems [also inherent in all other gradient-based approaches]:

1. **Input weight conflict:** for simplicity, let us focus on a single additional input weight  $m_{\text{ig}}$ . Assume that the total error sum be reduced by switching on unit  $\mathbb{I}$  in response to a certain input, and keeping it active for a long time [until it helps to compute a desired output]. Obviously  $\mathbb{I}$  is non-zero, since the same incoming weight has to be used for both storing certain inputs and ignoring others,  $m_{\text{ig}}$  will often receive conflicting weight update signals during this time [recall that  $\mathbb{I}$  is linear]: these signals will attempt to make  $m_{\text{ig}}$  participate in [1] storing the input  $\mathbb{I}$  by switching on  $\mathbb{I}$  and [2] protecting the input  $\mathbb{I}$  by preventing  $\mathbb{I}$  from being switched off by irrelevant later inputs]. This conflict makes learning difficult, and calls for a more context-sensitivity mechanism for controlling “write operations” through input weights.

2. **Output weight conflict:** assume  $\mathbb{I}$  is switched on and currently stores some previous input. For simplicity, let us focus on a single additional outgoing weight  $m_{\text{og}}$ . The same  $m_{\text{og}}$  has to be used for both retrieving  $\mathbb{I}$ 's content at certain times and preventing  $\mathbb{I}$  from disturbing  $\mathbb{I}$  at other times. As long as unit  $\mathbb{I}$  is non-zero,  $m_{\text{og}}$  will attract conflicting weight update signals generated during sequence processing: these signals will attempt to make  $m_{\text{og}}$  participate in [1] accessing the information stored in  $\mathbb{I}$  and — at different times — [2] protecting unit  $\mathbb{I}$  from being perturbed by  $\mathbb{I}$ . For instance, with many tasks there are certain “short time lag errors” that can be reduced in early training stages. However, at later training stages  $\mathbb{I}$  may suddenly start to cause needless errors in situations that already seemed under control by attempting to participants in reducing more difficult “long time lag errors”. Again, this conflict makes learning difficult, and calls for a more context-sensitivity mechanism for controlling “read operations” through output weights.

Of course, input and output weight conflicts are not specific for long time lags, but occur for short time lags as well. Their effects, however, become particularly pronounced in the long time lag case: as the time lag increases, [1] stored information must be protected against perturbation for longer and longer periods, and — especially in advanced stages of learning — [2] more and more already correct outputs also require protection against perturbation.

Thus to this problem allows this naive approach does not work well except in cases of certain simple problems involving local input/output representations and non-repeating input patterns [see Hochreiter 1991 and Silen et al. 1993]. This next section shows how to do it right.

## 4. USING SHORT-TERM MEMORY

Memory cells and gate units. To construct an architecture that allows for constant error flow through special, self-connected units without the disadvantages of the naive approach, we extend the constant error approach [CEA] provided by the self-connected, linear unit  $\mathbb{I}$  from Section 3.2 by introducing additional features. A multiplicative *output gate*  $m_{\text{og}}$  is introduced to protect the memory contents stored in  $\mathbb{I}$  from perturbation by irrelevant inputs. Likewise, a multiplicative *input gate*  $m_{\text{ig}}$  is introduced which protects other units from perturbation by currently irrelevant memory contents stored in  $\mathbb{I}$ .

The resulting, more complex unit is called a *memory cell* [see Figure 1]. This  $j$ -th memory cell is denoted  $\mathbb{m}_j$ . Each memory cell is built around a central linear unit with a fixed self-connection [the CIDE]. In addition to  $m_{\text{ig}}_{ij}$ ,  $\mathbb{m}_j$  gets input from a multiplicative unit  $m_{\text{og}}^{\text{out},j}$  [the “output gate”], and from another multiplicative unit  $m_{\text{ig}}^{\text{in},j}$  [the “input gate”].  $m_{\text{ig}}^{\text{in},j}$ 's activation at time  $N$  is denoted by  $y^{\text{in},j}[N]$ , and  $y^{\text{out},j}[N]$ . We have

$$y^{\text{out},j}[N] = \mathbf{b}_{\text{out},j}[\text{act}(b_{\text{out},j}[N])]; y^{\text{in},j}[N] = \mathbf{b}_{\text{in},j}[\text{act}(b_{\text{in},j}[N])];$$

whereas

$$m_{out,j}[n] = \sum_m m_{out,j,m} y^m[n-1],$$

and

$$m_{in,j}[n] = \sum_m m_{in,j,m} y^m[n-1].$$

At time  $n$

$$m_{s,j}[n] = \sum_m m_{s,j,m} y^m[n-1].$$

The summation indices  $m$  may stand for input units, gates units, memory cells, or even conventional hidden units if there are any [see also paragraph on “network topology” below]. All these different types of units may contain useful information about the current state of the net. For instance, an input gate [output gate] may use inputs from other memory cells to decide whether to store [process] certain information in its memory cell. There can be recurrent self-connections like  $m_{s,j,s_j}$ . It is up to the user to define the network topology. See Figure 8 for an example.

At time  $n$ ,  $s_j$ ’s output  $y^{s,j}[n]$  is computed as

$$y^{s,j}[n] = y^{out,j}[n]/m_{s,j}[n],$$

where the “internal state”  $s_{s,j}[n]$  is

$$s_{s,j}[0] = 1, s_{s,j}[n] = s_{s,j}[n-1] + g^{bs,j}[n] y^{s,j}[n], \text{ for } n > 0.$$

The differentiable function  $g$  squares  $m_{s,j}$ ; the differentiable function  $b$  scales memory cell outputs computed from the internal state  $s_{s,j}$ .

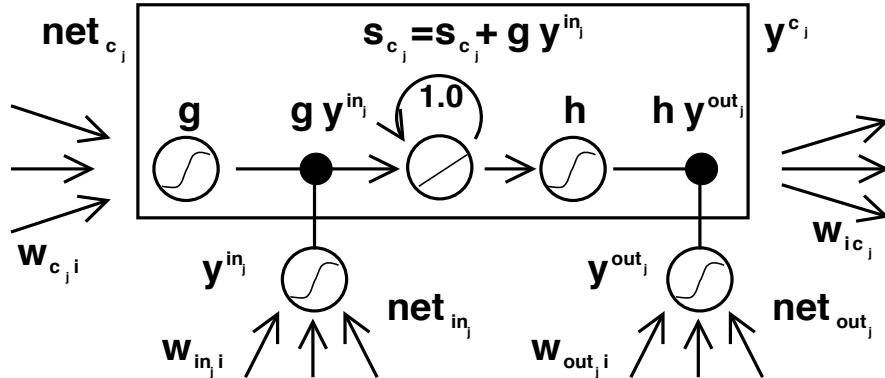


Figure 1: Illustrations of memory cell  $s_j$  [its body] and its gate units  $g_j$  and  $h_j$ . The self-connection connection [with weight 1.0] facilitates feedback while it is being off it does stop. It holds the body of the “membrane access command” GAC. The gate units open with close access to GAC. See Unit with appendix H.1 for details.

**Gates units:** The initial input weight coefficients,  $w_{g,j}$  controls the error flow to memory cell  $s_j$ ’s input connections  $m_{s,j,i}$ . The circumstant  $s_j$ ’s output weight coefficients,  $w_{h,j}$  controls the error flow from unit  $j$ ’s output connections. In other words, this net can use  $b_{g,j}$  to decide whether to keep or override information in memory cell  $s_j$ , and  $w_{h,j}$  to decide whether to access memory cell  $s_j$  and whether to prevent other units from using information in  $s_j$  [see Figure 1].

Diver signals trapped within a memory cell’s GAC control settings – But different error signals flowing into this cell [at different times] via its output gate may get superimposed. This output gate will have to learn which errors to trap in its GAC, by appropriately scaling them. This input

gates will have to learn when to release errors, again by appropriately scaling them. Essentially, the multiplicative gate units open and close access to constant error flow through GRU.

Distributed output representations typically do require output gates. Not always are both gate types necessary, though — one may be sufficient. For instance, in Experiments 2a and 2b in Section 5, it will be possible to use input gates only. In fact, output gates are not required in case of local output smoothing — preventing memory cells from perturbing already learned outputs can be done by simply setting the corresponding weights to zero. Even in this case, however, output gates can be beneficial: this prevents the net’s attempts at storing long time lag memories [which are usually hard to learn] from perturbing activations representing easily learnable short time lag memories. (This will prove quite useful in Experiment 1, for instance.)

**Network topology.** We use networks with one input layer, one hidden layer, and one output layer. The (fully) self-connected hidden layer contains memory cells and corresponding gate units [for convenience, we refer to both memory cells and gate units as being located in the hidden layer]. The hidden layer may also contain “conventional” hidden units processing inputs to gate units and memory cells. All units [except for gate units] in all layers have directed connections [some are inputs] to all units in the layer above [or to all higher layers — Experiments 2a and 2b].

**Memory cell blocks.**  $\mathcal{M}$  memory cells sharing the same input gate and the same output gate form a structure called an “memory cell block of size  $\mathcal{M}$ ”. Memory cell blocks facilitate information storage — as with conventional neural nets, it is not necessary to make a distributed input within a single cell. Since each memory cell block has its own gate units as a single memory cell [remember that], this block architecture can be even slightly more efficient [see paragraph “computational complexity”]. A memory cell block of size 1 is just a simple memory cell. In the experiments [Section 5], we will use memory cell blocks of various sizes.

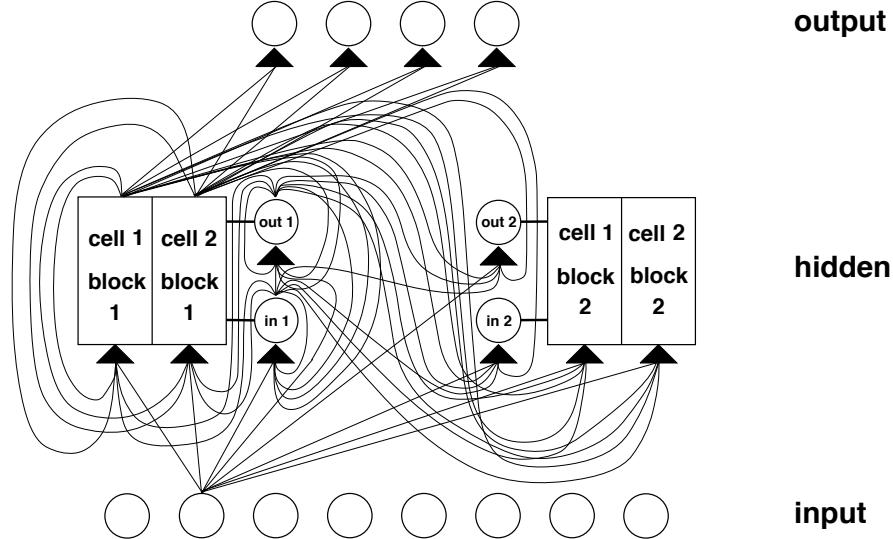
**Learning.** We use a variant of RTRL [e.g., Hollnagel and Hallinan 1998] which properly takes into account the altered, multiplicative dynamics caused by input and output gates. However, to ensure non-diverging error backprop through internal states of memory cells, as with truncated BPTT [e.g., Williams and Peng 1990], errors arriving at “memory cell not inputs” [for cell  $m_j$ , this includes  $m_{j+1}, m_{j+2}, \dots, m_{j+\text{out}_j}$ ] do not get propagated back further in time [although this does not change the incoming weights]. Only within<sup>1</sup> memory cells, errors are propagated back through previous internal states  $s_{i,j}$ . To visualize this: once an error signal arrives at a memory cell output, it gets scaled by output gate activation and  $\eta'$ . Then it is within the memory cell’s GRU, where it can flow back indefinitely without ever being scaled. Only when it leaves the memory cell through the input gate and  $y$ , it is scaled once more by input gate activation and  $\eta'$ . It then comes to change the incoming weights [before it is truncated [as applicable for explicit formulas]].

**Computational complexity.** As with Hochreiter’s famous recurrent backprop algorithm [Hochreiter 1991], only the derivatives  $\frac{\partial e_i}{\partial w_{kl}}$  need to be stored and updated. Hence the LSTM algorithm is very efficient, with an excellent update complexity of  $O(1)$ , where  $\mathcal{M}$  the number of weights [see details in appendix A.1]. Hence, LSTM and BPTT for fully recurrent nets have the same update complexity per time step [while RTRL’s is much worse]. Unlike full BPTT, however, LSTM is local to space and time<sup>2</sup>: there is no need to store activation values collected during sequence processing in a stack with potentially unlimited size.

**Alluse problem and solutions.** In the beginning of the learning phase, error reduction may be possible without storing information over time. The network will thus tend to alluse memory cells, e.g., as bias cells [i.e., it might make their activations constant and use the outgoing connections as adaptive thresholds for other units]. The potential difficulty is: it may take a long time to reduce alluse memory cells and make them available for further learning. A similar “alluse problem” appears if two memory cells store the same [inconsistent] information. There are at least two solutions to the alluse problem: **[1] Separated softmax constraint** [e.g., Hallinan 1991]: a memory cell and the corresponding gate units are added to the network whenever the

<sup>1</sup>For intra-cellular backprop in a quite different context see also Oja and Srihari 1982.

<sup>2</sup>Following Sahai and Culurciello 2009, we say that a recurrent net algorithm is *local to space* if the update complexity per time step and weight does not depend on network size. We say that a method is *local to time* if its storage requirements do not depend on input sequence length. For instance, RTRL is local in time but not in space. BPTT is local in space but not in time.



**Figure 8:** *Diagram of a cell with 8 input units, 2 output units, and 2 memory cell blocks of size 2. Each unit has 2 input units, each unit has 2 output units, each unit has 2 blocks. The first memory cell of Block 1, rabbit/unit1, connects to the two for Figure 1, while goat units have both unit 1 and unit 2 by multiplying Figure 1 by 2. Inputs are cell-stateless, i.e. with regard to self-connections, all units connect with their corresponding parts of Figure 1. This example assumes linear connectivity: each unit with each cell memory cell has all non-self-connections. For simplicity, however, only the first, unconnected update rule, shown here, uses through connections to output units, and through local self-connections within cell blocks (local storage here — see Figure 1). When there is no feedback across all “units” the terms memory cells are split units. Therefore, one connection stored stores the previous term local to this unit from within the connection unambiguously (except from connections to output units), although the connections themselves are ambiguous. That’s why this unrolled LSTM algorithm is so efficient, despite its ability to update using long time lags. See itself with opposite 1, 1 for details. Figure 8 only shows the modifications made by Dapperhank to — only the ones of the non-hyper units be modified.*

error stops decreasing (see Experiment 8 in Section 5). **2.2. Output gate bias:** each output gate gets a negative initial bias, to push initial memory cell activations towards zero. Memory cells with more negative bias automatically get “allocated” later (see Experiments 1, 2, 3, 4 in Section 5).

**Internal state drift and round-off.** If incoming cell  $s_t$ ’s inputs are mostly positive or mostly negative, then its internal state  $s_t$  will tend to drift away over time. This is potentially dangerous, for the  $W[s_t]$  will then adopt very small values, and the gradient will vanish. One way to circumvent this problem is to choose an appropriate function  $W$ . But  $W[s_t] = 0$ , for instance, has the disadvantage of unrestricted memory cell output range. Our simple but effective way of solving drift problems at the beginning of learning is to initially bias the input gate  $b_{\text{out}}$  towards zero. Although there is a tradeoff between the magnitudes of  $W[s_t]$  on the one hand and of  $b_{\text{out}}$  and  $b_{\text{mem}}$  on the other, the potential negative effect of input gate bias is negligible compared to the ones of the shifting effect. With logistic sigmoid activation functions, there appears to be no need for fine-tuning the initial bias, as confirmed by Experiments 1 and 2 in Section 5.5.

## 5. EXPERIMENTS

**Introduction.** All the tasks are appropriate to demonstrate the quality of a model long time lag

algorithm<sup>9</sup>. First of all, minimal time lags between relevant input signals and corresponding teacher signals must be long for all training sequences. In fact, many previous recurrent net algorithms sometimes manage to generalize from very short training sequences to very long test sequences. See, e.g., Holland [1991]. But a real long time lag problem does not have very short time lag exemplars in the training set. For instance, Elman's training procedure, BPTT, offline RTRL, online RTRL, etc., fail miserably on real long time lag problems. See, e.g., Hochreiter [1991] and Aboroz [1992]. A second important requirement is that the tasks should be complex enough such that they cannot be solved quickly by simple-minded strategies such as random weight guessing.

**Guessing can outperform many long time lag algorithms.** Recently we discovered [Sohni-Hukar and Hochreiter 1993; Hochreiter and Sohni-Hukar 1993, 1997] that many long time lag tasks used in previous work can be solved more quickly by simple random weight guessing than by the proposed algorithms. For instance, guessing solved a variant of Bengio and Frasconi's "parity problem" [1992] much faster<sup>10</sup> than the seven methods tested by Bengio et al. [1991] and Bengio and Frasconi [1992]. Similarly for some of Miller and Giles' problems [1993]. Of course, this does not mean that guessing is a good algorithm. It just means that some previously used problems are not extremely appropriate to demonstrate the quality of previously proposed algorithms.

**What's common to Experiments 1–3?** All our experiments [except for Experiment 1] involve long minimal time lags — there are no short time lag training exemplars facilitating learning. Solutions to most of our tasks are sparse in weight space. This requires either many parameters/inputs or high weight precision, such that random weight guessing becomes infeasible.

We always use on-line learning [as opposed to batch learning], and logistic sigmoidal activation functions. For Experiments 1 and 2, initial weights are chosen in the range  $[-0.5, 0.5]$ , for the other experiments in  $[-0.1, 0.1]$ . Training sequences are generated randomly according to the various task descriptions. In slight deviation from the notation in Appendix A1, recall discrete time step of each input sequence involves three processing stages: [1] use current input to set the input units, [2] Compute activations of hidden units [including input gates, output gates, memory cells], [3] Compute output unit activations. Except for Experiments 1, 2a, and 3b, sequence elements are randomly generated on-line, and error signals are generated only at sequence ends. Net activations are reset after each processed input sequence.

For comparisons with recurrent nets taught by gradient descent, we give results only for RTRL, except for comparison 2a, which also includes BPTT. Note, however, that untruncated BPTT [see, e.g., Williams and Peng 1990] computes exactly the same gradient as offline RTRL. With long time lag problems, offline RTRL [or BPTT] and the online version of RTRL [as activation routes, online weight changes] lead to almost identical, negative results [as confirmed by additional simulations in Hochreiter 1991; see also Aboroz 1992]. This is because offline RTRL, online RTRL, and full BPTT all suffer badly from exponential error cloning.

Our LSTM architectures are selected quite arbitrarily. If nothing is known about the complexity of a given problem, a more systematic approach would be: start with a very small net consisting of one memory cell. If this does not work, try two cells, etc. Alternatively, use sequential network construction [e.g., Thallam 1991].

### Outline of experiments.

- Experiment 1 focuses on a standard benchmark test for recurrent nets: the undirected Reber grammar. Since it allows for training sequences with short time lags, it is not a long time lag problem. We include it because [1] it provides a nice seeming failure LSTM's output gates are truly beneficial, and [2] it is a popular benchmark for recurrent nets that has been used by many authors — we want to include at least one experiment where conventional BPTT and RTRL do not fail completely [LSTM, however, clearly outperforms them]. The undirected Reber grammar's minimal time lags represent a harder case in the sense that it is still possible to learn to bridge them with conventional algorithms. Only slightly longer

<sup>9</sup>It should be mentioned, however, that different input representations and different types of noise may lead to worse guessing performance [Mishra, Bengio, personal communication, 1997].

minimal time lags would make this almost impossible. The more interesting tasks in our paper, however, are those that RDRRL, RSTT, etc. cannot solve at all.

- Experiment 2 focuses on noise-free and noisy sequences involving numerous input symbols distracting from the few important ones. This most difficult task [Task 2] involves hundreds of distracter symbols at random positions, and minimal time lags of 1000 steps. LSTM solves it, while RSTT and RDRRL already fail in case of 10-step minimal time lags [see also, e.g., Hochreiter 1991 and Abawajy 1998]. For this reason RDRRL and RSTT are omitted in the remaining, more complex experiments, all of which involve much longer time lags.
- Experiment 3 addresses long time lag problems with noise and signal on the same input line. Experiments 8n/8l focus on Wong et al.'s 1991 “E-sequence problem”. Because this problem actually can be solved quickly by random weight guessing, we also include a far more difficult E-sequence problem [8l] which requires to learn real-valued, conditional expectations of noisy targets, given the inputs.
- Experiments 4 and 5 involve distributed, continuous-valued input representations and require learning to store precise, real values for very long time periods. Relevant input signals can occur at quite different positions in input sequences. Again minimal time lags involve hundreds of steps. Similar tasks have been solved by other recurrent net algorithms.
- Experiment 6 involves tasks of an different complete type that also has not been solved by other recurrent net algorithms. Again, relevant input signals can occur at quite different positions in input sequences. This experiment shows that LSTM can extract information concealed by the temporal order of widely separated inputs.

Subsection 5.7 will provide a detailed summary of experimental conditions in two tables for reference.

### 5.1 EXPERIMENT 1: REBER GRAMMAR

**Task.** Our first task is to learn the “symbolical Reber grammar”, e.g. Smith and Zipser [1989], Choromanski et al. [1992], and Ballmann [1991]. Since it allows for training sequences with short time lags [of as few as 3 steps], it is not a long time lag problem. We include it for two reasons: [1] it is an popular recurrent net benchmark used by many authors — we wanted to have at least one experiment where RDRRL and RSTT do not fail completely, and [2] it shows nicely how output gates can be beneficial.

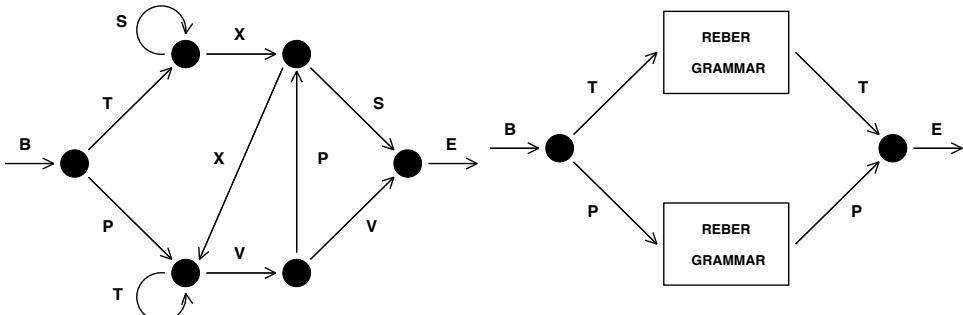


Figure 1: Transition Diagram from the symbolical Reber grammar. Note the equivalence of any of the grammars (see Figure 2).

Starting at the leftmost node of the directed graph in Figure 1, symbol strings are generated sequentially (beginning with the empty string) by following edges — and appending the associated

symbols to the current string — until the rightmost node is reached. The net’s task is to read strings, one symbol at a time, and to grammatically predict the next symbol [error signals occur at every time step]. To correctly predict the symbol before last, the net has to remember the second symbol.

**Comparison.** We compare LSTM to “Dinam nets trained by Dineam’s training procedure” [DLN] [results taken from Choromans et al. 1989], Hallinan’s “Recurrent Network-Correlation” [RNC] [results taken from Hallinan 1991], and RTRL [results taken from Smith and Zipser 1989], [where only the few successful trials are listed]. It should be mentioned that Smith and Zipser actually make the task easier by increasing the probability of short time lag exemplars. We didn’t do this for LSTM.

**Training/Testing.** We use a local input/output representation [ $\overline{I}$  input units,  $\overline{O}$  output units]. Following Hallinan, we use 25 training strings and 25 test strings. The training set is generated randomly; training exemplars are picked randomly from the training set. Test sequences are generated randomly, too, but sequences already used in the training set are not used for testing. After string presentation, all activations are reinitialized with zeros. A trial is considered successful if all string symbols of all sequences in both test set and training set are predicted correctly — that is, if the output unit(s) corresponding to the possible next symbol(s) is/are always the most active ones.

**Architectures.** Architectures for RTRL, DLN, RNC and reported in the references listed above. For LSTM, we use  $\overline{I} = 5$  incoming cell blocks. Each block has  $\overline{C} = 11$  incoming cells. The output layer’s only incoming connections originate at incoming cells. Each incoming cell and each gate unit receives incoming connections from all incoming cells and gate units [the hidden layer is fully connected — less connectivity may work as well]. The input layer has forward connections to all units in the hidden layer. The gate units are binary. These architectures parameters make it easy to store at least  $\overline{I}$  input signals [architectures RNC and  $\overline{I}=1$  are employed to obtain comparable numbers of weights for both architectures: 255 for  $\overline{I}=1$  and 270 for  $\overline{I}=5$ ]. Other parameters may be appropriate as well. However, all sigmoid functions are logistic with output range [0, 1], except for  $H$ , whose range is  $[-1, 1]$ , and  $g$ , whose range is  $[-\frac{1}{2}, \frac{1}{2}]$ . All weights are initialized in  $[-0.5, 0.5]$ , except for the output gate biases, which are initialized to -1, -4, and -8, respectively [see values given below, equation (8) of Section II]. The trial learning rates of 0.1, 0.2 and 0.5.

**Results.** We use 5 different, randomly generated pairs of training and test sets. With each such pair we run 10 trials with different initial weights. See Table 1 for results [mean of 50 trials]. Unlike the other methods, LSTM always learns to solve the task. Even when we ignore the unsuccessful trials of the other approaches, LSTM learns much faster.

**Importance of output gates.** This experiment proves to a nice example whether the output gates is truly beneficial. Learning to store the first T or G should not perturb activations representing the more easily learned transitions of the original Polder grammar. This is the job of the output gates. Without output gates, we did not achieve fast learning.

## 5.2 EXPERIMENT 2: NOISE-FREE SEQUENCES WITH LONG TIME LAGS

**Task 2a:** noise-free sequences with long time lags. There are  $p+1$  possible input symbols denoted  $v_1, \dots, v_{p-1}, v_p = u, v_{p+1} = y$ .  $v_k$  is “locally” represented by the  $p+1$ -dimensional vector where k-th component is 1 [all other components are 0]. A net with  $\overline{I}+1$  input units and  $\overline{O}+1$  output units sequentially observes input symbol sequences, one at a time, grammatically trying to predict the next symbol — error signals occur at every single time step. To emphasize the “long time lag problem”, we use a training set consisting of only two very similar sequences:  $[y, v_1, v_2, \dots, v_{p-1}, u]$  and  $[u, v_1, v_2, \dots, v_{p-1}, y]$ . Both is selected with probability 0.5. To predict the final element, the net has to learn to store a representation of the first element for  $p$  time steps.

We compare “Real-Time Recurrent Learning” for fully recurrent nets [RTRL], “End-to-Propagation Through Time” [ETT], the sometimes very successful 2-net “Neural Sequence Blocker” [SH, Salimbeni/Huller 1992], and our new method [LSTM]. In all cases, weights are initialized in  $[-0.5, 0.5]$ . Due to limited computation time, training is stopped after 5 million sequences pro-

method	hidden units	# weights	learning rate	% of success	success after
IRIDRL	8	≈ 170	0.05	“some fraction”	170,000
IRIDRL	12	≈ 190	0.1	“some fraction”	25,000
IDLM	12	≈ 180		0	>200,000
IRCR	7-8	≈ 110-120		0%	180,000
LSTM	8 blocks, size 1	≈ 200	0.1	100%	30,000
LSTM	8 blocks, size 2	≈ 270	0.1	100%	31,000
LSTM	8 blocks, size 3	≈ 270	0.5	0%	10,000
LSTM	8 blocks, size 1	≈ 200	0.5	0%	3,000
LSTM	8 blocks, size 2	≈ 270	0.5	100%	8,000

Table 1: *Diamond-Delon sequence generation: percentage of successful trials and number of experiments needed until success for IRIDRL (results taken from Schmidhuber and Ziegler 2020), “Diamond and Delon by Diamond’s grammar” (results taken from Choromanski et al. 2020), “Diamond-Delonza-Diamondization” (results taken from Huttermann 2020) and our own approach (LSTM). “Diamond” denotes the first 8 rows and “Delon” — the corresponding rows 9-16 generated all the indicated lengths. Only LSTM always learns to solve this task (only two failures out of 100 trials). Please notice we ignore the unsuccessful trials of the other approaches, LSTM learns much faster (the number of required training examples is the bottom row numbers Leibniz 8,000 and 20,000).*

tations. A successful run is one that fulfills the following criterion: after training, during 10,000 successive, randomly chosen input sequences, the measured absolute error of all output units is always below 0.05.

**Architectures.** IRIDRL: one self-recurrent hidden unit,  $p+1$  non-recurrent output units. Recell layer has connections from all layers below. All units use the logistic activation function sigmoid in [0,1].

IRDTT: same architecture as the one trained by IRIDRL.

OH: both not architectures like IRIDRL’s, but one has an additional output for predicting the hidden unit of the other one (see Schmidhuber 1992b for details).

LSTM: like with IRIDRL, but the hidden unit is replaced by a memory cell and an input gate (no output gate required).  $\sigma$  is the logistic sigmoid, and  $I$  is the identity function  $I : \mathbb{R}[\theta] = \theta$ . The memory cell and input gate are added once the error has stopped decreasing (see diffuse problem: solution [1] in Section 2).

**Results.** Using IRIDRL and a short 8 times step delay ( $\tau = 8$ ),  $\frac{1}{2}$  of all trials were successful. We had one success with  $p = 10$ . With long time lags, only the natural sequence shunker and LSTM achieved successful trials, while IRDTT and OH failed. With  $p = 100$ , the 8-unit sequence shunker solved the task in only  $\frac{1}{8}$  of all trials. LSTM, however, always learned to solve this task. Comparing successful trials only, LSTM learned much faster. See Table 2 for details. It should be mentioned, however, that a Diamond/Delon shunker can also always quickly solve this task (Schmidhuber 1992a, 1992b).

**Task 2B:** no local regularities. With this task alone, the shunker sometimes learns to correctly predict the final element, but only because of predictably local regularities in the input stream that allows for compressing the sequence. In an additional, more difficult task (involving many more different possible sequences), we remove compressibility by replacing the deterministic suffix sequence  $[v_1, v_2, \dots, v_{p-1}]$  by a random subsequence (of length  $p-1$ ) over the alphabet  $v_1, v_2, \dots, v_{p-1}$ . We obtain 2 classes (two sets of sequences)  $\{(v, v_1, v_2, \dots, v_{p-1}, y) \mid 1 \leq b_1, b_2, \dots, b_{p-1} \leq p-1\}$  and  $\{(v, v_1, v_2, \dots, v_{p-1}, y) \mid 1 \leq b_1, b_2, \dots, b_{p-1} \leq n-1\}$ . Again, every next sequence element has to be predicted. This only totally predictable targets. However, now  $v$  and  $y$ , which occur at sequence ends. Training examples are chosen randomly from the 2 classes. Architectures and parameters are the same as in Experiment 1a. A successful run is one that fulfills the following criterion: after training, during 10,000 successive, randomly chosen input

Method	Fixing $\eta$	Learning rate	# weights	# Successful trials	Success after
RTRL	0	1.0	88	68	1,000,000
RTRL	0	1.0	88	53	888,000
RTRL	0	10.0	88	22	888,000
RTRL	10	1.0-10.0	100	0	> 5,000,000
RTRL	100	1.0-10.0	1000	0	> 5,000,000
RETT	100	1.0-10.0	1000	0	> 5,000,000
CH	100	1.0	1000	88	88,000
LSTM	100	1.0	1000	100	5,000

Table 6: *Task E: Memorizing of successful trials with execution of branching sequences with success, from “Word-Block Memorization” [18], “Word-Hopscotch Memorization” [20], branch sequences branching [21], and the own method [22, 23]. Trials continue until the success of 100 trials. While RTRL does step learning, only CH and LSTM perform successful trials. When failure the algorithm has unsuccessful trials of the other approaches, LSTM learns much faster.*

sequences, the maximal absolute error of all output units is below 0.05 at sequence end.

**Results.** As expected, the chunker failed to solve this task (see Fig. RETT and RTRL, of course). LSTM, however, was always successful. On average (mean of 18 trials), success for  $\eta = 100$  was achieved after 5,000 sequence presentations. This demonstrates that LSTM does not require sequence regularities to work well.

**Task E:** every long times lags — no local regularities. This is the most difficult task in this subsection. To our knowledge no other recurrent net algorithm can solve it. Now there are  $y+1$  possible input symbols denoted  $v_1, \dots, v_{y-1}, v_y, v_{y+1} = u, v_{y+2} = x, v_{y+3} = y, v_{y+4} = z, \dots, v_k$  and also called “*transition symbols*”. Again,  $v_i$  is locally represented by the  $y+1$ -dimensional vector whose 0-th component is 1 [all other components are 0]. A net with  $y+1$  input units and  $k$  output units sequentially observes input symbol sequences, one at a time. Training sequences are randomly chosen from the union of two long similar subsets of sequences:  $\{[v, y, v_1, v_2, \dots, v_{y+3}, x, y] \mid 1 \leq b_1, b_2, \dots, b_{y+3} \leq z\}$  and  $\{[v, x, v_1, v_2, \dots, v_{y+3}, x, z] \mid 1 \leq b_1, b_2, \dots, b_{y+3} \leq z\}$ . To produce a training sequence, we [1] randomly generate a sequence profile of length  $y+2$ , [2] randomly generates a sequence suffix of additional elements  $\{z, y, x, z, y\}$  with probability  $\frac{3}{10}$  or, alternatively, an  $x$  with probability  $\frac{1}{10}$ . In the latter case, we [3] conclude the sequence with  $x$  or  $y$ , depending on the second element. For a given  $k$ , this leads to a uniform distribution on the possible sequences with length  $y+2+k$ . The minimal sequence length is  $y+2$ ; the expected length is

$$2 + \sum_{k=0}^{\infty} \frac{1}{10} \left( \frac{3}{10} \right)^k (y+2) = y+10.$$

The expected number of occurrences of element  $v_i$ ,  $1 \leq i \leq y$ , in a sequence is  $\frac{y+10}{10} \approx \frac{y}{2}$ . The goal is to predict the last symbol, which always occurs after the “trigger symbol”  $x$ . Error signals are generated only at sequence ends. To predict the final element, the net has to learn to store a representation of the second element for at least  $y+1$  time steps [until it sees the trigger symbol  $x$ ]. Success is defined as “prediction error [for final sequence element] of both output units always below 0.5, for 10,000 successive, randomly chosen input sequences”.

**Architecture/Learning.** The net has  $y+1$  input units and  $k$  output units. Weights are initialized in  $[-0.5, 0.5]$ . To avoid too much learning time variation due to different weight initializations, the hidden layer gets two memory cells [two cell blocks of size 1 — although one would be sufficient]. There are no other hidden units. The output layer receives connections only from memory cells. Memory cells and gate units receive connections from input units, memory cells and gate units [i.e., the hidden layer is fully connected]. No bias weights are used.  $M$  and  $g$  are logistic sigmoid with output ranges  $[-1, 1]$  and  $[-2, 2]$ , respectively. The learning rate is 0.01.

$y$ [times lag -1]	$n$ [# random inputs]	$\frac{z}{y}$	# weights	Success after
50	50	1	3200	30,000
100	100	1	3200	31,000
200	200	1	12800	33,000
500	500	1	32000	35,000
1,000	1,000	1	32000	36,000
		2	32000	36,000
1,000	500	2	32000	36,000
1,000	200	5	12800	75,000
1,000	100	10	3200	135,000
1,000	50	20	3200	200,000

Table 8: *Task B:* LSTM and long unlabelled time lags  $y + 1$  and  $n$  full of zeros.  $y$  is the number of available hidden symbols ( $y + 1$  is the number of input units).  $\frac{z}{y}$  is the expected number of occurrences of a given hidden symbol in a sequence. The neighbour column lists the number of training sequences required by LSTM (green), RNN (red) and the other competitors (blue) for success of solving this task. If  $n$  is less than the number of hidden symbols (purple numbers), however, no progression in the time lag learning process may slowly. The lower blue numbers are the expected slow-down due to bottleneck frequency of hidden symbols.

Note that the *unlabelled* time lag is  $y + 1$  — this not never sees short training sequences facilitating the classification of long test sequences.

**Results.** 20 trials were made for all tested pairs  $[y, n]$ . Table 8 lists the mean of the number of training sequences required by LSTM to achieve success (RNN and RNNL have no chance of solving non-trivial tasks with minimal time lags of 1000 steps).

**Scaling.** Table 8 shows that if we let the number of input symbols (#weights) increase in proportion to the time lag, learning time increases very slowly. This is a similar remarkable property of LSTM not shared by any other method we are aware of. Indeed, RNNL and RNNTT are far from scaling reasonably — instead, they appear to scale exponentially, and appear quite useless when the time lags exceed us from 10 steps.

**Distractor influence.** In Table 8, the column labelled  $\frac{z}{y}$  gives the expected frequency of distractor symbols. Increasing this frequency decreases learning speed, an effect due to weight oscillations caused by frequently selected input symbols.

### 5.5 EXPERIMENT 3: NOISE AND SIGNAL ON SAME CHANNEL

This experiment serves to illustrate that LSTM does not encounter fundamental problems if noise and signal are mixed on the same input line. We initially focus on Bengio et al.’s simple 1993 “2-sequence problem”; in Experiment 3 we will then pose a more challenging 3-sequence problem.

**Task 8a (“2-sequence problem”).** The task is to observe and then classifying input sequences. There are two classes, each occurring with probability 0.5. There is only one input line. Only the first  $N$  real-valued sequence elements carry relevant information about the class. Sequence elements at positions  $t > N$  are generated by a Gaussian with mean zero and variance 0.2. Once  $N = 1$ : the first sequence element is 1.0 for class 1, and -1.0 for class 2. Once  $N = 2$ : the first three elements are 1.0 for class 1 and -1.0 for class 2. The target at the sequence end is 1.0 for class 1 and 0.0 for class 2. Correct classification is defined as “absolute output error at sequence end below 0.2”. Given a constant  $T$ , the sequence length is randomly selected between  $T$  and  $T + T/10$  [in difference to Bengio et al.’s problem is that they also permit shorter sequences of length  $T/2$ ].

**Guessing.** Bengio et al. [1993] and Bengio and Frasconi [1993] tested 6 different methods on the 2-sequence problem. As discussed, however, their random weight guessing easily outper-

$M$	$N$	stop: ST1	stop: ST2	# weights	ST2: fraction misclassified
100	8	27,820	32,820	102	0.000135
100	1	28,820	31,820	102	0.000117
1000	8	116,820	122,820	102	0.000078

**Table 8: Task 8a: Misclassification of sequences problem.**  $M$  is unlabeled sequence length.  $N$  is the number of information-conveying elements of sequence length. The column labeled by ST1 [ST2] gives the number of sequences misclassified by unlabeled stopping criterion ST1 [ST2]. The rightmost column lists the fraction of misclassified past-including sequences (with absolute error  $> 0.5$ ) from a test set consisting of 2000 sequences (selected after ST2 was satisfied). This column can exceed 1.0. The Gaussian, Unconstrained, and LSTM problem is as simple. Most problems involve processing values  $\in \{0, 1\}$  from 0.0001 to 0.000078 from misclassified to most classified results.

forms the one all sequences the problem is as simple<sup>9</sup>. See Schmidhuber and Hochreiter [1993] and Hochreiter and Schmidhuber [1993, 1997] for additional results in this vein.

**LSTM architecture.** We use a R-layer net with 1 input unit, 1 output unit, and 8 cell blocks of size 1. The output layer receives connections only from memory cells. Memory cells and gate units receive inputs from input units, memory cells and gate units, and have bias weights. Gate units and output unit are logistic sigmoid in  $[0, 1]$ ,  $b$  in  $[-1, 1]$ , and  $y$  in  $[-1, 1]$ .

**Training/Testing.** All weights (except the bias weights to gate units) are randomly initialized in the range  $[-0.1, 0.1]$ . The first input gate bias is initialized with  $-1.0$ , the second with  $-0.1$ , and the third with  $-0.0$ . The first output gate bias is initialized with  $-0.0$ , the second with  $-0.0$  and the third with  $-0.0$ . This precise initialization reduces learning matter though, as confirmed by additional experiments. The learning rate is 1.0. All activations are reset to zero at the beginning of a new sequence.

We stop training (and judge the task as being solved) according to the following criterion: ST1: none of 200 sequences from a randomly chosen test set is misclassified. ST2: ST1 is satisfied, and mean absolute test set error is below 0.01. In case of ST2, an additional test set consisting of 2000 randomly chosen sequences is used to determine the fraction of misclassified sequences.

**Results.** See Table 8. The results are means of 10 trials with different weight initializations in the range  $[-0.1, 0.1]$ . LSTM is able to solve this problem, though by far not as fast as random weight guessing (see paragraph “Guessing” above). Clearly, this trivial problem shows not precisely a very good testbed to compare performances of various non-trivial algorithms. Still, it demonstrates that LSTM does not encounter fundamental problems when faced with signal and noise on the same platform.

**Task 8b.** Architecture, parameters, etc. like in Task 8a, but now with Gaussian noise [mean 0 and variance 0.0] added to the information-conveying elements ( $0 \leq M$ ). We stop training (and judge the task as being solved) according to the following, slightly modified criteria: ST1: less than 8 out of 200 sequences from a randomly chosen test set are misclassified. ST2: ST1 is satisfied, and mean absolute test set error is below 0.001. In case of ST2, an additional test set consisting of 2000 randomly chosen sequences is used to determine the fraction of misclassified sequences.

**Results.** See Table 9. The results represent means of 10 trials with different weight initializations. LSTM easily solves this problem.

**Task 8c.** Architecture, parameters, etc. like in Task 8a, but with a few essential changes that make the task non-trivial: the targets are 0.5 and 0.8 for class 1 and class 2, respectively, and there is Gaussian noise on the targets [mean 0 and variance 0.1; std dev. 0.31]. To minimize mean squared error, the system has to learn the *conditional correlations* of the targets given the inputs. Misclassification is defined as “absolute difference between output and noise-free target  $> 0.2$  for

<sup>9</sup>It should be mentioned, however, that different input representations and different types of noise may lead to worse guessing performance [Wirsching, 1999].

$H$	$N$	stop: ST1	stop: ST2	# weights	ST2: fraction misclassified
100	8	11,700	18,200	102	0.00028
100	1	22,200	28,700	102	0.01600
1000	1	181,000	188,000	102	0.01207

Table 3: *Task B1: modified E-sequence problem.* Same as in Table 2, but now the background-noising elements are also generated by noise.

$H$	$N$	stop	# weights	fraction misclassified	ave. Differences to mean
100	8	248,000	102	0.00028	0.010
100	1	252,500	102	0.00001	0.008

Table 3: *Task B2: modified, noisy shuffling E-sequence problem.* Same as in Table 2, but with noisy and-noisy targets. The system has to learn the conditional probabilities of the targets given the inputs. The rightmost column provides the average differences between network output and expected target. Unlike in task B1, the task cannot be solved globally by random weightassing.

class 1 and 0.8 for class 2)  $> 0.1$ .” This network output is considered acceptable if the mean absolute difference between noise-free target and output is below 0.015. Since this requires high weight precision, *Task B2* [unlike *Task B1*] *cannot be solved globally by random guessing*.

**Training/Testing.** The learning rate is 0.1. We stop training according to the following criterion: none of 500 sequences from a randomly chosen test set is misclassified, and mean absolute difference between noise-free target and output is below 0.015. An additional test set consisting of 2500 randomly chosen sequences is used to determine the fraction of misclassified sequences.

**Results.** See Table 3. The results represent means of 10 trials with different weight initializations. Despite the noisy targets, LSTM still can solve the problem by learning the requested target values.

### 5.2. INDIVIDUALITY 2: ADDING PROPORTION

The difficult task in this section is of a type that has never been solved by other recurrent net algorithms. It shows that LSTM can solve long time lag problems involving distributed, continuous-valued representations.

**Task 3.** Each element of each input sequence is a pair of components. The first component is a real value randomly chosen from the interval  $[-1, 1]$ ; the second is either 1.0, 0.0, or -1.0, and is used as a marker: at the end of each sequence, the task is to output the sum of the first components of those pairs that are marked by second components equal to 1.0. Sequences have random lengths between the minimal sequence length  $H$  and  $H + \frac{H}{10}$ . In a given sequence exactly two pairs are marked as follows: one first randomly selected and marked one of the first ten pairs (whose first component was call  $\Delta_1$ ). Then one randomly selected and marked one of the first  $\frac{H}{10} - 1$  still unmarked pairs (whose first component was call  $\Delta_2$ ). The second components of all remaining pairs are zero except for the first and final pair, whose second components are -1. [In this case sense without the first pair of the sequence gets marked, we set  $\Delta_1$  to zero.] An error signal is generated only at the sequence end: the target is  $0.5 + \frac{\Delta_1 + \Delta_2}{H}$  [the sum  $\Delta_1 + \Delta_2$  scaled to the interval  $[0, 1]$ ]. A sequence is processed correctly if the absolute error at the sequence end is below 0.01.

**Architecture.** We use a 5-layer net with 5 input units, 1 output unit, and 5 cell blocks of size 5. The output layer receives connections only from memory cells. Memory cells and gate units receive inputs from memory cells and gate units [i.e., the hidden layer is fully connected — less connectivity may work as well]. The input layer has forward connections to all units in this hidden

$M$	minimized lag	# weights	# wrong predictions	Success after
100	50	98	1 out of 2000	70,000
200	50	98	0 out of 2000	200,000
1000	50	98	1 out of 2000	800,000

**Table 5: ADDITIONAL 4: Results from the Adding Problem.**  $M$  is the minimal sequence length,  $M/5$  the minimized time lag. “# wrong predictions” is the number of incorrectly processed sequences (errors > 0.02) from a test set containing 2000 sequences. The rightmost column gives the number of including sequences required to achieve the stopping criterion. All values are ratios of 10 labels. When  $M = 1000$  the number of required including samples varies between 200,000 and 8,000,000, according 200,000 by only 8 cases.

layer. All non-input units have bias weights. This architecture guarantees make it easy to store at least  $K$  input signals [in cell block size of 1 words each, too]. All activation functions are logistic with output range  $[0, 1]$ , except for  $b$ , whose range is  $[-1, 1]$ , and  $y$ , whose range is  $[-2, 2]$ .

**State drift versus initial bias.** Note that the task requires storing the precise values of real numbers for long durations — the system must learn to protect memory cell contents against even minor internal state drift [see Section 4]. To study the significance of the drift problem, we make the task even more difficult by biasing all non-input units, thus artificially inducing internal state drift. All weights [including the bias weights] are randomly initialized in the range  $[-0.1, 0.1]$ . Following Section 4’s naming for state drifts, the first input gate bias is initialized with  $-0.1$ , the second with  $-0.0$  [through this precise values manually monitor, as confirmed by additional experiments].

**Training/Testing.** The learning rate is 0.5. Training is stopped once the average training error is below 0.01, and the 2000 most recent sequences were processed correctly.

**Results.** With a test set consisting of 2000 randomly chosen sequences, the average test set error were always below 0.01, and there were never more than 8 incorrectly processed sequences. Table 6 shows details.

This experiment demonstrates: [1] LSTM is able to work well with distributed representations. [2] LSTM is able to learn to perform calculations involving continuous values. [3] Since the system manages to store continuous values without deterioration for minimal lags of  $\frac{M}{5}$  time steps, there is no significant, harmful internal state drift.

## 5.5 EXPERIMENT 5: MULTIPLICATION PROPERTY

One may argue that LSTM is at best biased towards tasks such as the Adding Problem from the previous subsection. Solutions to the Adding Problem may exploit the GPU’s built-in integration capabilities. Although this GPU property may be viewed as a feature rather than a disadvantage [integration seems to be a natural subset of many tasks occurring in the real world], this question arises whether LSTM can solve tasks with inherently non-integrative solutions. To test this, we change the problem by requiring the final target to equal the product [instead of the sum] of familiar marked inputs.

**Task.** Like the task in Section 5.1, except that the first component of each pair is a real value randomly drawn from the interval  $[0, 1]$ . In this rare case where the first pair of the input sequence gets marked, we set  $M_1$  to 1.0. The target at sequence end is the product  $M_1 \times M_2$ .

**Architecture.** Like in Section 5.1. All weights [including the bias weights] are randomly initialized in the range  $[-0.1, 0.1]$ .

**Training/Testing.** The learning rate is 0.1. At test performances twice as soon as less than  $w_{\text{avg}}$  of the 2000 most recent training sequences had to achieve errors exceeding 0.01, unless  $w_{\text{avg}} = 100$ , and  $w_{\text{avg}} = 18$ . Why these values?  $w_{\text{avg}} = 100$  is sufficient to learn storage of the relevant inputs. It is not enough though to fine-tune the precise final outputs.  $w_{\text{avg}} = 18$ , however,

$m$	minimum lag	# weights	$m_{\text{avg}}$	# wrong predictions	MSE	Success after
100	50	28	100	169 out of 2500	0.0228	100,000
100	50	28	100	15 out of 2500	0.0189	1,000,000

Table S: **Task 2: Multihead LSTM**: Results from the Multihead LSTM.  $m$  is the minimum sequence length,  $M$  is the minimum time lag. The first row of best self-supervised LSTM sequences is shown as base lines.  $m_{\text{avg}}$  is the LSTM model recall (including sequences both its error  $> 1.00$ ). “# wrong predictions” is the number of test sequences with error  $> 1.00$ . MSE is the mean squared error on the test set. The rightmost column lists numbers of including sequences required to achieve the stopping criterion. All values are means of 10 trials.

leads to quite satisfying results.

**Results.** For  $m_{\text{avg}} = 100$  [ $m_{\text{avg}} = 100$ ] with a test set consisting of 2500 randomly chosen sequences, the average test set error was always below 0.028 [0.018], and there were never more than 169 [15] incorrectly predicted sequences. Table S shows details. [A net with additional standard hidden units or with a hidden layer allows the memory cells may learn this fine-tuning part more quickly.]

The experiment demonstrates: LSTM can solve tasks involving both continuous-valued representations and non-integrative information processing.

### 5.3 EXPERIMENT 3: INDIVISORIAL CRITERIA

In this subsection, LSTM solves rather difficult [but artificial] tasks that have never been solved by previous recurrent net algorithms. The experiment shows that LSTM is able to abstract information controlled by the temporal order of widely separated inputs.

**Task 3a:** three relevant, widely separated symbols. The goal is to classifying sequences. Elements and targets are represented locally [input vectors with only one non-zero bit]. The sequence starts with an  $\mathbb{B}$ , ends with a  $\mathbb{C}$  [the “trigger symbol”], and otherwise consists of randomly chosen symbols from the set  $\{\mathbb{A}, \mathbb{B}, \mathbb{C}, \mathbb{D}\}$  except for two elements at positions  $\mathbb{U}_1$  and  $\mathbb{U}_2$  that are either  $\mathbb{A}$  or  $\mathbb{B}$ . The sequence length is randomly chosen between 100 and 110,  $\mathbb{U}_1$  is randomly chosen between 10 and 20, and  $\mathbb{U}_2$  is randomly chosen between 50 and 60. There are 8 sequences classes  $\mathbb{A}, \mathbb{B}, \mathbb{C}, \mathbb{D}$  with respect to the temporal order of  $\mathbb{A}$  and  $\mathbb{B}$ . The rules are:  $\mathbb{A}, \mathbb{A} = \mathbb{A}; \mathbb{A}, \mathbb{B} = \mathbb{B}; \mathbb{B}, \mathbb{A} = \mathbb{B}; \mathbb{B}, \mathbb{B} = \mathbb{A}$ .

**Task 3b:** three relevant, widely separated symbols. Again, the goal is to classifying sequences. Elements/targets are represented locally. The sequence starts with an  $\mathbb{B}$ , ends with a  $\mathbb{C}$  [the “trigger symbol”], and otherwise consists of randomly chosen symbols from the set  $\{\mathbb{A}, \mathbb{B}, \mathbb{C}, \mathbb{D}\}$  except for three elements at positions  $\mathbb{U}_1$ ,  $\mathbb{U}_2$  and  $\mathbb{U}_3$  that are either  $\mathbb{A}$  or  $\mathbb{B}$ . The sequence length is randomly chosen between 100 and 110,  $\mathbb{U}_1$  is randomly chosen between 10 and 20,  $\mathbb{U}_2$  is randomly chosen between 50 and 60, and  $\mathbb{U}_3$  is randomly chosen between 80 and 90. There are 8 sequences classes  $\mathbb{A}, \mathbb{B}, \mathbb{C}, \mathbb{D}, \mathbb{A}, \mathbb{B}, \mathbb{C}, \mathbb{D}, \mathbb{A}, \mathbb{B}, \mathbb{C}$  with respect to the temporal order of this  $\mathbb{A}$ s and  $\mathbb{B}$ s. The rules are:  $\mathbb{A}, \mathbb{A}, \mathbb{A} = \mathbb{A}; \mathbb{A}, \mathbb{A}, \mathbb{B} = \mathbb{B}; \mathbb{A}, \mathbb{B}, \mathbb{A} = \mathbb{B}; \mathbb{A}, \mathbb{B}, \mathbb{B} = \mathbb{A}; \mathbb{B}, \mathbb{A}, \mathbb{A} = \mathbb{B}; \mathbb{B}, \mathbb{A}, \mathbb{B} = \mathbb{A}; \mathbb{B}, \mathbb{B}, \mathbb{A} = \mathbb{B}; \mathbb{B}, \mathbb{B}, \mathbb{B} = \mathbb{B}$ .

There are no many output units as there are classes. One class is locally represented by a binary target vector with one non-zero component. With both tasks, error signals occur only at the end of a sequence. The sequence is classified correctly if the final absolute error of all output units is below 0.8.

**Architecture.** We use a 3-layer net with 8 input units, 8 [38] and blocks of size 8 and 9 [38] output units for Task 3a [38]. Again all non-input units have their weights, and the output layer receives connections from memory cells only. Memory cells and gate units receive inputs from input units, memory cells and gate units [i.e., the hidden layer is fully connected — less connectivity may work as well]. The architecture parameters for Task 3a [38] make it easy to

stores at least 8 [0] input signals. All activation functions are logistic with output range [0, 1], except for  $h$ , whose range is  $[-1, 1]$ , and  $g$ , whose range is  $[-2, 2]$ .

**Training/Testing.** The learning rate is 0.5 [0.1] for Experiment 3a [3b]. Training is stopped once the average training error falls below 0.1 and the 2000 most recent sequences were classified correctly. All weights are initialized in the range  $[-0.1, 0.1]$ . The first input gate bias is initialized with  $-0.1$ , the second with  $-0.2$ , and [for Experiment 3b] the third with  $-0.3$  [again, not confirmed by additional experiments that the precise values hardly matter].

**Results.** With a test set consisting of 2500 randomly chosen sequences, the average test set error was always below 0.1, and there were never more than 8 incorrectly classified sequences. Table 3 shows details.

The experiment shows that LSTM is able to extract information contained by the temporal order of unlikely sequential inputs. In Task 3a, for instance, the delays between first and second relevant input and between second relevant input and sequence end are at least 20 time steps.

task	# weights	# wrong predictions	Success after
Task 3a	15k	1 out of 2500	81,900
Task 3b	20k	8 out of 2500	571,100

Table 3: **EXPERIMENT 3:** Results from the Unlikely Order Problem. “# wrong predictions” is the number of incorrectly classified sequences [ $\text{error} > 0.5$  from all test set input units] from a test set containing 2500 sequences. The rightmost column gives the number of training sequences required to achieve the stopping criterion. The results from Task 3a are means of 20 trials; those from Task 3b of 10 trials.

**Typical solutions.** In Experiment 3a, those three LSTM distinguish between temporal orders  $[0, 1]$  and  $[1, 0]$ . One of many possible solutions is to store the first  $0$  or  $1$  in cell Block 1, and the second  $1/0$  in cell Block 2. Before the first  $0/1$  occurs, Block 1 can see that it is still empty by means of its recurrent connections. After the first  $0/1$ , Block 1 can close its input gate. Once Block 1 is filled and closed, this fact will become visible to Block 2 [recall that all gate units and all memory cells receive connections from all non-output units].

Typical solutions, however, require only one memory cell Block. This Block stores the first  $0$  or  $1$ ; once the second  $0/1$  occurs, it changes its state depending on the first stored symbol. Solution type 1 exploits this connection between memory cell output and input gate unit — this following sequence census: Different input gate activations: “ $0$  occurs in conjunction with a filled Block”; “ $0$  occurs in conjunction with an empty Block”. Solution type 2 is based on a strong positive connection between memory cell output and memory cell input. The previous occurrence of  $0$  [0] is represented by a positive [negative] internal state. Since this input gate opens for the second time, so does the output gate, and the memory cell output is fed back to its own input. This census  $[0, 0]$  is represented by a positive internal state, because  $0$  contributes to the new internal state twice [as current internal state and cell output feedback]. Similarly,  $[1, 1]$  gate represented by a negative internal state.

## 5.7 SUMMARY OF EXPERIMENTAL CONDITIONS

The two tables in this subsection provide an overview of the most important LSTM parameters and architectural details for Experiments 1–3. The conditions of the simple experiments 1a and 2a differ slightly from those of the other, more systematic experiments, due to historical reasons.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Task	in	leg	R	s	in	out	m	r	neg	neg	igR	igS	h1	g1	m
1-1	0	0	0	1	0	0	250	0	-1,-2,-3,-4	0	r	gs	h1	g2	0.1
1-2	0	0	0	2	0	0	250	0	-1,-2,-3	0	r	gs	h1	g2	0.1

For the conditions not used grayed

Test	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
	in	leg	in	s	in	out	in	s	out	leg	leg	leg	leg	leg	leg	leg
1-E	X	X	X	X	X	X	X	X	X	-1,-E,-S	r	pa	h1	pa	X.2	
1-E!	X	X	X	1	X	X	X	X	X	-1,-E,-E,-E!	r	pa	h1	pa	X.2	
1-E-	X	X	X	X	X	X	X	X	X	-1,-E,-S	r	pa	h1	pa	X.2	
2a	1XX	1XX	1	1	1XX	1XX	1XX	1XX	1XX	no cap	none	none	i1	pl	1.X	
2b	1XX	1XX	1	1	1XX	1XX	1XX	1XX	1XX	no cap	none	none	i1	pl	1.X	
2c-1	ST	ST	S	1	ST	S	ST	S	ST	none	none	none	h1	pa	X.11	
2c-2	1XX	1XX	S	1	1XX	S	1XX	S	1XX	none	none	none	h1	pa	X.11	
2c-3	2XX	2XX	S	1	2XX	S	2XX	S	2XX	none	none	none	h1	pa	X.11	
2c-4!	STX	STX	S	1	STX	S	STX	S	STX	none	none	none	h1	pa	X.11	
2c-5	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-6	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-7	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-8	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-9	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-10	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-11	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-12	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-13	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-14	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-15	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2c-16	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	none	none	none	h1	pa	X.11	
2d	1XX	1XX	S	1	1	1	1XX	S	1XX	-E,-H,-S	-1,-E,-S	H1	h1	pa	1.X	
2d!	1XX	1XX	S	1	1	1	1XX	S	1XX	-E,-H,-S	-1,-E,-S	H1	h1	pa	1.X	
2d-	1XX	1XX	S	1	1	1	1XX	S	1XX	-E,-H,-S	-1,-E,-S	H1	h1	pa	1.X	
2d-1	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	r	-E,-H	all	h1	pa	X.2	
2d-2	STX	STX	S	1	STX	S	STX	S	STX	r	-E,-H	all	h1	pa	X.2	
2d-3	1XXXX	1XXXX	S	1	1XXXX	S	1XXXX	S	1XXXX	r	-E,-H	all	h1	pa	X.2	
2e	1XX	1XX	S	1	1	1	1XX	S	1XX	r	r	all	h1	pa	1.1	
2e!	1XX	1XX	S	1	1	1	1XX	S	1XX	r	-E,-H	all	h1	pa	1.1	
2e-	1XX	1XX	S	1	1	1	1XX	S	1XX	r	-E,-H	all	h1	pa	1.1	

1	2	3	4	5	6
TestID	subset	interval	test set size	stopping criterion	success
1	t1	[-0.2, 0.2]	2500	training 100 test correctly pred.	see test
2a	t1	[-0.2, 0.2]	no test set	after 5 million exemplars	<b>AFS[0.25]</b>
2B	t2	[-0.2, 0.2]	10000	after 5 million exemplars	<b>AFS[0.25]</b>
2c	t2	[-0.2, 0.2]	10000	after 5 million exemplars	<b>AFS[0.25]</b>
2d	t3	[-0.1, 0.1]	2500	ST1 and ST2 have test	<b>AFS[0.25]</b>
2E	t3	[-0.1, 0.1]	2500	ST1 and ST2 have test	<b>AFS[0.25]</b>
2f	t3	[-0.1, 0.1]	2500	ST1 and ST2 have test	see test
2g	t3	[-0.1, 0.1]	2500	ST3[0.1]	<b>AFS[0.25]</b>
2h	t3	[-0.1, 0.1]	2500	see test	<b>AFS[0.25]</b>
2i	t3	[-0.1, 0.1]	2500	ST3[0.1]	<b>AFS[0.25]</b>
2j	t3	[-0.1, 0.1]	2500	ST3[0.1]	<b>AFS[0.25]</b>

Table 11: **Summary of experimental conditions** from LSTM. **First column:** task name. **Each column:** including exemption selection, where “**W**” stands for “uniformly chosen from underlying set”, “**R**” for “uniformly chosen from  $\mathbb{R}$  classes”, and “**B**” for “uniformly generated numbers”. **Each column:** memory bottleneck. **WU column:** task cell size. **BU column:** skipping criterion for learning, where “**ONE**” stands for “memory learning over token  $\mathcal{E}$  with the PMM until recent sequences meet precision criteria”. **BU column:** success [*current classification*] condition, where “**THREE**” stands for “absolute error of all output units of sequence with  $\mathcal{E}$ ”.

## 3 DISCUSSION

### Limitations of LSTM.

- The particularly efficient truncated backprop version of the LSTM algorithm will not easily solve problems similar to “strongly delayed XOR problems”, where the goal is to compute the XOR of two suddenly segmented inputs that previously occurred somewhere in a noisy sequence. The reason is that storing only one of the inputs will not help to reduce the expected error — the task is non-decomposable in the sense that it is impossible to incrementally reduce the error by first solving an easier subgoal.
- In theory, this limitation can be circumvented by using the full gradient [including with additional conventional hidden units receiving input from the memory cells]. But was the not recommended computing the full gradient for the following reasons: [1] It increases computational complexity. [2] Constant error flow through RNNs can be shown only for truncated LSTM. [3] We actually did conduct a few experiments with non-truncated LSTM. There was no significant difference to truncated LSTM, meaning the error outside the RNNs error flow tends to vanish quickly. For this same reason full BIETT does not outperform truncated BIETT.
- Recall incoming cell block needs two additional units [input and output gate]. In comparison to standard recurrent nets, however, this does not increase the number of weights by more than a factor of 2: each conventional hidden unit is replaced by at most 8 units in the LSTM architecture, increasing the number of weights by a factor of  $2^k$  in the fully connected sense. Note, however, that our experiments use quite comparable weight numbers for the architectures of LSTM and competing approaches.
- Considerably speaking, due to its constant error flow through RNNs within incoming cells, LSTM runs into problems similar to those of feedforward nets saving the entire input string at once. For instance, there are tasks that can be quickly solved by random weight guessing but not by the truncated LSTM algorithm with small weight initializations, such as the 500-state parity problem [see introduction to Section 5]. Here, LSTM’s problems are similar to the ones of a feedforward net with 500 inputs, trying to save 500-bit parity. Indeed LSTM typically behaves much like a feedforward net trained by backprop that sees the entire input. But that’s also precisely why it so clearly outperforms previous approaches on many non-trivial tasks with significant recall errors.
- LSTM does not have any problems with the notion of “memory” that go beyond those of other approaches. All gradient-based approaches, however, suffer from practical inability to precisely count discrete time steps. If it makes a difference whether a certain signal occurred 9 or 100 steps ago, then an additional accounting mechanism seems necessary. Dasier tasks, however, such as ones that only requires to make a difference between, say, 8 and 11 steps, do not pose any problems to LSTM. For instance, by generating an appropriate negative connection between incoming cell output and input, LSTM can give more weight to recent inputs and learn things without necessary

## Advantages of LSTM.

- The constant error flow propagation within memory cells results in LSTM's ability to bridge long time lags in case of problems similar to those discussed above.
- For long time lag problems such as those discussed in this paper, LSTM can handle noise, distributed representations, and continuous values. In contrast to finite state automata or hidden Markov models LSTM does not require an *a priori* choice of a finite number of states. In principle it can deal with unlimited state numbers.
- The problems discussed in this paper LSTM generalizes well — even if the positions of widely separated, relevant inputs in the input sequence do not matter. Unlike previous approaches, ours quickly learns to distinguish between two or more widely separated occurrences of a particular element in an input sequence, without depending on appropriate short time lag training examples.
- There appears to be no need for parameter fine tuning. LSTM works well over a broad range of parameters such as learning rates, input gate bias and output gate bias. For instance, to some readers the learning rates used in our experiments may seem large. However, a large learning rate pushes the output gates towards zero, thus automatically counteracting its own negative effects.
- The LSTM algorithm's update complexity per weight and time step is essentially that of RNN, namely  $O(1)$ . This is excellent in comparison to other approaches such as BRNN. Unlike full BRNN, however, LSTM is *local by field space with time*.

## 7 CONCLUSION

Each memory cell's internal architecture guarantees constant error flow within its constant error surrounded CDR, provided that truncated backprop cuts off error flow trying to leak out of memory cells. This represents the basis for bridging long time lags. Two gate units learn to open and close access to error flow within each memory cell's CDR. The multiplicities input gate affords protection of the CDR from perturbation by irrelevant inputs. Likewise, the multiplicities output gate protects writer units from perturbation by currently irrelevant memory contents.

**Future work.** To find out about LSTM's practical limitations we intend to apply it to real world data. Application areas will include [1] time series prediction, [2] music composition, and [3] speech processing. It will also be interesting to augment sequence chunkers [Sallmi/Huller 1992], [1993] by LSTM to examine the advantages of that.

## 8 ACKNOWLEDGMENTS

Thanks to Miles Moon, Wilfried Breuer, Nico Sallmehelph, and several anonymous referees for valuable comments and suggestions that helped to improve a previous version of this paper [Heidtmann and Sallmehelph 1995]. This work was supported by DFG grant SCHA 142/3-1 from "Bundesforschungsministerium".

## APPENDIX

### A.1 ALGORITHM DETAILS

In what follows, the index  $k$  ranges over output units,  $l$  ranges over hidden units,  $\eta_k$  stands for the  $k$ -th memory cell block,  $\eta_k^i$  denotes the  $i$ -th unit of memory cell block  $\eta_k$ ,  $m, l, m$  stand for multiplying units,  $t$  ranges over all time steps of a given input sequence.

The gate unit logistic sigmoid (with range [0, 1]) used in the experiments is

$$\sigma(x) = \frac{1}{1 + \exp(-x)}. \quad [2]$$

The function  $M$  (with range [-1, 1]) used in the experiments is

$$M(x) = \frac{x}{1 + \exp(-|x|)} - 1. \quad [3]$$

The function  $g$  (with range [-2, 2]) used in the experiments is

$$g(x) = \frac{x}{1 + \exp(-|x|)} - 2. \quad [4]$$

#### Forward pass.

The net input and the activation of hidden unit  $i$  are

$$\begin{aligned} m_{hi}[t] &= \sum_m m_{hi,m} g^m[t-1] \\ g^i[t] &= \sigma(m_{hi}[t]). \end{aligned} \quad [5]$$

The net input and the activation of  $h_{\text{out}}$  are

$$\begin{aligned} m_{h_{\text{out}},i}[t] &= \sum_m m_{h_{\text{out}},im} g^m[t-1] \\ g^{h_{\text{out}}}[t] &= \sigma_{h_{\text{out}}}(m_{h_{\text{out}},i}[t]). \end{aligned} \quad [6]$$

The net input and the activation of  $m_{ij}$  are

$$\begin{aligned} m_{m_{ij}}[t] &= \sum_m m_{m_{ij},m} g^m[t-1] \\ g^{m_{ij}}[t] &= \sigma_{m_{ij}}(m_{m_{ij}}[t]). \end{aligned} \quad [7]$$

The net input  $m_{a_j}$ , the internal state  $s_{a_j}$ , and the output activation  $y^{a_j}$  of the  $a$ -th memory cell of memory cell block  $a_j$  are:

$$\begin{aligned} m_{a_j}[t] &= \sum_m m_{a_j,m} g^m[t-1] \\ s_{a_j}[t] &= s_{a_j}[t-1] + g^{h_{\text{out}}}[t] g\left(m_{a_j}[t]\right) \\ y^{a_j}[t] &= g^{m_{a_j}}[t] g(s_{a_j}[t]). \end{aligned} \quad [8]$$

The net input and the activation of output unit  $k$  are

$$\begin{aligned} m_{k_k}[t] &= \sum_{m: m \text{ is not a gate}} m_{k_k,m} g^m[t-1] \\ g^k[t] &= \sigma_k(m_{k_k}[t]). \end{aligned}$$

The backward pass to the described later is based on the following truncated backprop formulae. Approximate derivatives for truncated backprop. The truncated version (see Section II) only approximates the partial derivatives, which is reflected by the “≈” signs in the notation below. It truncates error flow since it leaves memory cells as gate units. Truncation ensures that there are no loops across which an error that left some memory cell through its input or input gate can re-enter this cell through its output or output gate. This in turn ensures constant error flow through the memory cell’s block.

In this truncated backprop version, the following shortcomings are replaced by zeros:

$$\frac{\text{Grad}_{\text{out}_i}[t]}{\text{Eg}^n[t-1]} \approx_{\text{in}} 0 \quad \forall i,$$

$$\frac{\text{Grad}_{\text{out}_{i,j}}[t]}{\text{Eg}^n[t-1]} \approx_{\text{in}} 0 \quad \forall i, j,$$

$$\frac{\text{Grad}_{\text{out}_j}[t]}{\text{Eg}^n[t-1]} \approx_{\text{in}} 0 \quad \forall j.$$

Therefore we get

$$\frac{\text{Eg}^{\text{out}_i}[t]}{\text{Eg}^n[t-1]} = \mathbb{E}'_{\text{out}_i}[\text{Grad}_{\text{out}_i}[t]] \frac{\text{Grad}_{\text{out}_i}[t]}{\text{Eg}^n[t-1]} \approx_{\text{in}} 0 \quad \forall i,$$

$$\frac{\text{Eg}^{\text{out}_{i,j}}[t]}{\text{Eg}^n[t-1]} = \mathbb{E}'_{\text{out}_{i,j}}[\text{Grad}_{\text{out}_{i,j}}[t]] \frac{\text{Grad}_{\text{out}_{i,j}}[t]}{\text{Eg}^n[t-1]} \approx_{\text{in}} 0 \quad \forall i, j,$$

and

$$\frac{\text{Eg}^{\text{out}}[t]}{\text{Eg}^n[t-1]} = \frac{\text{Eg}^{\text{out}_i}[t]}{\text{Grad}_{\text{out}_i}[t]} \frac{\text{Grad}_{\text{out}_i}[t]}{\text{Eg}^n[t-1]} + \frac{\text{Eg}^{\text{out}_j}[t]}{\text{Grad}_{\text{out}_j}[t]} \frac{\text{Grad}_{\text{out}_j}[t]}{\text{Eg}^n[t-1]} + \frac{\text{Eg}^{\text{out}_i,j}[t]}{\text{Grad}_{\text{out}_i,j}[t]} \frac{\text{Grad}_{\text{out}_i,j}[t]}{\text{Eg}^n[t-1]} \approx_{\text{in}} 0 \quad \forall t.$$

This implies for all  $m_{lm}$  not on connections to  $\text{out}_i^2$ ,  $\text{out}_j$ , and  $\text{out}_{i,j}$  (that is,  $b \notin \{\text{out}_i^2, \text{out}_j, \text{out}_{i,j}\}$ ):

$$\frac{\text{Eg}^{\text{out}}[t]}{\text{Eg}_{lm}[t]} = \sum_m \frac{\text{Eg}^{\text{out}_i}[t]}{\text{Grad}_{\text{out}_i}[t]} \frac{\text{Grad}_{\text{out}_i}[t-1]}{\text{Eg}_{lm}[t]} \approx_{\text{in}} 0.$$

The truncated derivatives of output unit  $k$  are:

$$\begin{aligned} \frac{\text{Eg}^k[t]}{\text{Eg}_{lm}} &= \mathbb{E}'_{\text{out}_k}[\text{out}_k[t]] \left( \sum_{m_l: m_l \text{ not in gate}} m_{km_l} \frac{\text{Eg}^m[t-1]}{\text{Eg}_{lm}} + \mathbb{E}_{\text{out}_k}[\text{out}_k[t-1]] \right) \approx_{\text{in}} 0 \quad [10] \\ &= \mathbb{E}'_{\text{out}_k}[\text{out}_k[t]] \left( \sum_i \sum_{j=1}^{S_j} \mathbb{E}_{\text{out}_j} m_{kj_j} \frac{\text{Eg}^{\text{out}_j}[t-1]}{\text{Eg}_{lm}} + \sum_i [\mathbb{E}_{\text{out}_i} + \mathbb{E}_{\text{out}_{i,j}}] \sum_{j=1}^{S_j} m_{kj_j} \frac{\text{Eg}^{\text{out}_j}[t-1]}{\text{Eg}_{lm}} + \right. \\ &\quad \left. \sum_{k: k \text{ hidden unit}} m_{kk} \frac{\text{Eg}^k[t-1]}{\text{Eg}_{lm}} + \mathbb{E}_{\text{out}_k}[\text{out}_k[t-1]] \right) = \\ &= \mathbb{E}'_{\text{out}_k}[\text{out}_k[t]] \begin{cases} \mathbb{E}_{\text{out}_k}[\text{out}_k[t-1]] & b = k \\ m_{kj_j} \frac{\text{Eg}^{\text{out}_j}[t-1]}{\text{Eg}_{lm}} & b = \text{out}_j \\ \sum_{j=1}^{S_j} m_{kj_j} \frac{\text{Eg}^{\text{out}_j}[t-1]}{\text{Eg}_{lm}} & b = \text{out}_{i,j} \text{ OR } b = \text{out}_{i,j} \\ \sum_{k: k \text{ hidden unit}} m_{kk} \frac{\text{Eg}^k[t-1]}{\text{Eg}_{lm}} & b \text{ otherwise} \end{cases}, \end{aligned}$$

where  $\mathbb{E}$  is the Kronnecker Delta ( $\mathbb{E}_{ab} = 1$  if  $a = b$  and 0 otherwise), and  $S_k$  is the size of memory cell block  $m_k$ . The truncated derivatives of a hidden unit  $b$  that is not part of a memory cell are:

$$\frac{\text{Eg}^b[t]}{\text{Eg}_{lm}} = \mathbb{E}'_{\text{out}_k}[\text{out}_k[t]] \frac{\text{Grad}_{\text{out}_k}[t]}{\text{Eg}_{lm}} \approx_{\text{in}} 0, \quad [11]$$

[Note: Here it would be possible to use the full gradient without affecting constant error flow through internal states of memory cells.]

Cell Block  $\alpha_j$ 's truncated derivatives are:

$$\frac{\mathbf{E}g^{l_{m_j}}[t]}{\mathbf{E}m_{lm}} = \mathbf{E}'_{m_j}[\mathbf{mult}_{m_j}[t]] \frac{\mathbf{E}m_{lm}}{\mathbf{E}m_{lm}} \approx_{lm} \mathbf{E}_{m_j}[\mathbf{E}'_{m_j}[\mathbf{mult}_{m_j}[t]]] y^m[t-1]. \quad [12]$$

$$\frac{\mathbf{E}g^{x^{out_j}}[t]}{\mathbf{E}m_{lm}} = \mathbf{E}'_{m_j}[\mathbf{mult}_{x^{out_j}}[t]] \frac{\mathbf{E}m_{lm}}{\mathbf{E}m_{lm}} \approx_{lm} \mathbf{E}_{m_j}[\mathbf{E}'_{m_j}[\mathbf{mult}_{x^{out_j}}[t]]] y^m[t-1]. \quad [13]$$

$$\frac{\mathbf{E}s_{xy}[t]}{\mathbf{E}m_{lm}} = \frac{\mathbf{E}s_{xy}[t-1]}{\mathbf{E}m_{lm}} + \frac{\mathbf{E}g^{l_{m_j}}[t]}{\mathbf{E}m_{lm}} g'[\mathbf{mult}_{s_j}[t]] + y^{l_{m_j}}[t] g'[\mathbf{mult}_{s_j}[t]] \frac{\mathbf{E}m_{lm}}{\mathbf{E}m_{lm}} \approx_{lm} \quad [14]$$

$$(\mathbf{E}_{m_j}[\mathbf{I}] + \mathbf{E}_{s_j}[\mathbf{I}]) \frac{\mathbf{E}s_{xy}[t-1]}{\mathbf{E}m_{lm}} + \mathbf{E}_{m_j}[\mathbf{I}] \frac{\mathbf{E}g^{l_{m_j}}[t]}{\mathbf{E}m_{lm}} g'[\mathbf{mult}_{s_j}[t]] +$$

$$\mathbf{E}_{s_j}[\mathbf{I}] y^{l_{m_j}}[t] g'[\mathbf{mult}_{s_j}[t]] \frac{\mathbf{E}m_{lm}}{\mathbf{E}m_{lm}} = \\ (\mathbf{E}_{m_j}[\mathbf{I}] + \mathbf{E}_{s_j}[\mathbf{I}]) \frac{\mathbf{E}s_{xy}[t-1]}{\mathbf{E}m_{lm}} + \mathbf{E}_{m_j}[\mathbf{I}] \mathbf{E}'_{m_j}[\mathbf{mult}_{m_j}[t]] g'[\mathbf{mult}_{s_j}[t]] y^m[t-1] + \\ \mathbf{E}_{s_j}[\mathbf{I}] y^{l_{m_j}}[t] g'[\mathbf{mult}_{s_j}[t]] y^m[t-1].$$

$$\frac{\mathbf{E}g^{s_j^2}[t]}{\mathbf{E}m_{lm}} = \frac{\mathbf{E}g^{x^{out_j}}[t] W[\mathbf{s}_{s_j}[t]] + W[\mathbf{s}_{s_j}[t]] \frac{\mathbf{E}s_{xy}[t]}{\mathbf{E}m_{lm}} y^{x^{out_j}}[t]}{\mathbf{E}m_{lm}} \approx_{lm} \quad [15]$$

$$\mathbf{E}_{m_j}[\mathbf{I}] \frac{\mathbf{E}g^{x^{out_j}}[t] W[\mathbf{s}_{s_j}[t]]}{\mathbf{E}m_{lm}} + (\mathbf{E}_{m_j}[\mathbf{I}] + \mathbf{E}_{s_j}[\mathbf{I}]) W[\mathbf{s}_{s_j}[t]] \frac{\mathbf{E}s_{xy}[t]}{\mathbf{E}m_{lm}} y^{x^{out_j}}[t].$$

To efficiently update the system at time  $t$ , the only (truncated) derivatives that need to be stored at time  $t-1$  are  $\frac{\mathbf{E}s_{xy}[t-1]}{\mathbf{E}m_{lm}}$ , where  $b = \frac{\mathbf{E}s_{xy}[t-1]}{\mathbf{E}m_{lm}}$ , unless  $b = \mathbf{E}_s^2$  or  $b = \mathbf{E}_{s_j}$ .

**Backward pass.** We will describe the backward pass only for the particularly efficient “truncated gradient version” of the LSTM algorithm. For simplicity we will use equal signs even without approximations and make according to the truncated backprop equations above.

The squared error at time  $t$  is given by

$$E[t] = \sum_{k: k \text{ output unit}} [y^k[t] - y^k[t]]^2, \quad [16]$$

where  $y^k[t]$  is output unit  $k$ 's target at time  $t$ .

Time  $t$ 's contribution to  $m_{lm}$ 's gradient-based update with learning rate  $\alpha$  is

$$\Delta m_{lm}[t] = -\alpha \frac{\partial E[t]}{\partial m_{lm}}. \quad [17]$$

We define some unit  $k$ 's error at time step  $t$  by

$$e_k[t] := -\frac{\partial E[t]}{\partial m_{k,t}[t]}, \quad [18]$$

Using (truncated) standard backprop, we first compute updates for weights to output units  $[b = \mathbf{I}]$ , weights to hidden units  $[b = \mathbf{I}]$  and weights to output gates  $[b = \mathbf{mult}_b]$ . We obtain (compute formulas [10], [11], [18]):

$$b = \mathbf{I} [\text{output}] : \quad w_b[t] = \mathbf{E}'[\mathbf{mult}_b[t]] [y^k[t] - y^k[t]], \quad [19]$$

$$b = \mathbf{I} [\text{hidden}] : \quad w_b[t] = \mathbf{E}'[\mathbf{mult}_b[t]] \sum_{k: k \text{ output unit}} m_{k,t} e_k[t], \quad [20]$$

$b = \text{out}_k$  [output gates] : [21]

$$\Delta m_{\text{out}_j}[t] = b'_{\text{out}_j} [\text{out}_{\text{out}_j}[t]] \left( \sum_{v=1}^{S_j} w[\mathbf{x}_{v,j}[t]] - \sum_{k: K \text{ output unit}} m_{k,j} w_k[t] \right).$$

For all possible  $b$  times  $t$ 's contribution to  $m_{\text{out}}$ 's update is

$$\Delta m_{\text{out}}[t] = \alpha w[t] y^m[t-1]. \quad [22]$$

The remaining updates for weights to input gates [ $b = \text{in}_k$ ] and to cell units [ $b = c_k^c$ ] are less conventional. We define some internal states  $\mathbf{x}_{v,j}$ 's error:

$$\begin{aligned} \mathbf{x}_{c,j}^v &:= -\frac{\partial w[t]}{\partial \mathbf{x}_{v,j}[t]} = \\ &= b'_{\text{out}_j} [\text{out}_{\text{out}_j}[t]] \cdot W[\mathbf{x}_{v,j}[t]] - \sum_{k: K \text{ output unit}} m_{k,j} w_k[t]. \end{aligned} \quad [23]$$

We obtain for  $b = \text{in}_k$  or  $b = c_k^c$ ,  $v = 1, \dots, S_k$

$$-\frac{\partial w[t]}{\partial m_{\text{in}_j,m}} = \sum_{v=1}^{S_k} \mathbf{x}_{c,j}^v[t] \frac{\partial \mathbf{x}_{v,j}[t]}{\partial m_{\text{in}_j,m}}. \quad [24]$$

The derivatives of the internal states with respect to weights and the corresponding weight updates are as follows [compute expression [10]]:

$b = \text{in}_k$  [input gates] : [25]

$$\frac{\partial \mathbf{x}_{v,j}[t]}{\partial m_{\text{in}_j,m}} = \frac{\mathbf{x}_{v,j}[t-1]}{\mathbf{x}_{m_{\text{in}_j,m}}} + \alpha' [\text{out}_{v,j}[t]] \cdot b'_{\text{out}_j} [\text{out}_{\text{out}_j}[t]] \cdot y^m[t-1];$$

Therefore time  $t$ 's contribution to  $m_{\text{in}_j,m}$ 's update is [compute expression [10]]:

$$\Delta m_{\text{in}_j,m}[t] = \alpha \sum_{v=1}^{S_k} \mathbf{x}_{c,j}^v[t] \frac{\partial \mathbf{x}_{v,j}[t]}{\partial m_{\text{in}_j,m}}. \quad [26]$$

Similarly we get [compute expression [10]]:

$b = c_k^c$  [memory cells] : [27]

$$\frac{\partial \mathbf{x}_{v,j}[t]}{\partial m_{c,j,m}} = \frac{\mathbf{x}_{v,j}[t-1]}{\mathbf{x}_{m_{c,j,m}}} + \alpha' [\text{out}_{v,j}[t]] \cdot b'_{\text{out}_j} [\text{out}_{\text{out}_j}[t]] \cdot y^m[t-1];$$

Therefore time  $t$ 's contribution to  $m_{c,j,m}$ 's update is [compute expression [10]]:

$$\Delta m_{c,j,m}[t] = \alpha \mathbf{x}_{c,j}^v[t] \frac{\partial \mathbf{x}_{v,j}[t]}{\partial m_{c,j,m}}. \quad [28]$$

All we need to implement for this forward pass are equations [20], [21], [22], [23], [24], [25], [26], [27], [28]. Total weight's total update is the sum of the contributions of all time steps.

**Computational complexity.** LSTM's update complexity per time step is

$$\mathcal{O}(Km + KCN + CM + CS) = \mathcal{O}(m), \quad [29]$$

where  $K$  is the number of output units,  $C$  is the number of memory cell blocks,  $S > 1$  is the size of the memory cell blocks,  $M$  is the number of hidden units,  $N$  is the [internal] number of units forward-connected to memory cells, gates units and hidden units, and

$$CM = KM + KCN + CSN + CM + MM = \mathcal{O}(Km + KCN + CS) + MM$$

is the number of weights. Expression [23] is obtained by considering all computations of the backword pass: equation [19] needs  $M$  steps; [20] needs  $MH$  steps; [21] needs  $MSC$  steps; [22] needs  $M(M+G)$  steps for output units,  $MH$  steps for hidden units,  $MG$  steps for output gates; [23] needs  $MSC$  steps; [24] needs  $GSC$  steps; [25] needs  $GSH$  steps; [26] needs  $GSM$  steps; [27] needs  $GSH$  steps. The total is  $M + 2M H + MG + 2MSC + MH + G + 2GSC$  steps, or  $O(MH + MSC + MH + GSC)$  steps. We conclude: LSTM's update complexity per time step is just like RNN's for a fully recurrent net.

At a given time step, only the  $GSC$  most recent  $\frac{\partial \mathbf{x}_{t,j}}{\partial \mathbf{x}_{t,k}}$  values from equations [25] and [27] need to be stored. Hence LSTM's storage complexity also is  $O(MH)$  — it does not depend on the input sequence length.

## A.2 ERROR FLOW

We compute how much an error signal is scaled while flowing back through a memory cell for  $g$  time steps. As a  $\mathbf{y}_t$ -gradient, this analysis confirms that the error flow within a memory cell's SGD is indeed constant, provided that truncated backprop cuts off error flow trying to leave memory cells (see also Section 3.3). This analysis also highlights a potential for undesirable long-term shifts of  $\mathbf{x}_{t,j}$  (see [28] below), as well as the beneficial, countermeasuring influence of negatively biased input gates (see [28] below).

Using the truncated backprop learning rule, we obtain

$$\begin{aligned} \frac{\partial \mathbf{x}_{t,j}[t-M]}{\partial \mathbf{x}_{t,j}[t-M-1]} &= \quad [20] \\ 1 + \frac{\partial \mathbf{y}^{out}[t-M]}{\partial \mathbf{x}_{t,j}[t-M-1]} g'[\text{mult}_{t,j}[t-M]] + \mathbf{y}^{out}[t-M] g'[\text{mult}_{t,j}[t-M]] \frac{\partial \text{mult}_{t,j}[t-M]}{\partial \mathbf{x}_{t,j}[t-M-1]} &= \\ 1 + \sum_m \left[ \frac{\partial \mathbf{y}^{out}[t-M]}{\partial \mathbf{y}^u[t-M-1]} \frac{\partial \mathbf{y}^u[t-M-1]}{\partial \mathbf{x}_{t,j}[t-M-1]} \right] g'[\text{mult}_{t,j}[t-M]] &+ \\ \mathbf{y}^{out}[t-M] g'[\text{mult}_{t,j}[t-M]] \sum_m \left[ \frac{\partial \text{mult}_{t,j}[t-M]}{\partial \mathbf{y}^u[t-M-1]} \frac{\partial \mathbf{y}^u[t-M-1]}{\partial \mathbf{x}_{t,j}[t-M-1]} \right] &\approx_{t_M} 1. \end{aligned}$$

The  $\approx_{t_M}$  sign indicates equality plus the fact that truncated backprop replaces the zero till the following derivatives:  $\frac{\partial \mathbf{y}^{out}[t-M]}{\partial \mathbf{y}^u[t-M-1]} \approx_{t_M}$  and  $\frac{\partial \text{mult}_{t,j}[t-M]}{\partial \mathbf{y}^u[t-M-1]} \approx_{t_M}$ .

In what follows, an error  $\mathbf{E}_t[t]$  starts flowing back at  $\mathbf{y}_t$ 's output. We redefine

$$\mathbf{E}_t[t] := \sum_k m_{t,k} \mathbf{E}_t[t+1]. \quad [21]$$

Following the definitions/connections of Section 3.1, we compute error flow for the truncated backprop learning rule. The error occurring at the output gate is

$$\mathbf{E}_{\text{out},j}[t] \approx_{t_M} \frac{\partial \mathbf{y}^{out}[t]}{\partial \text{mult}_{\text{out},j}[t]} \frac{\partial \mathbf{y}^{out}[t]}{\partial \mathbf{y}^{out}[t]} \mathbf{E}_t[t]. \quad [22]$$

The error occurring at the internal state is

$$\mathbf{E}_{\text{c}_j}[t] = \frac{\mathbf{E}_{\mathbf{x}_{t,j}}[t+1]}{\mathbf{E}_{\mathbf{x}_{t,j}}[t]} \mathbf{E}_{\text{c}_j}[t+1] + \frac{\partial \mathbf{y}^{in}[t]}{\partial \mathbf{x}_{t,j}[t]} \mathbf{E}_t[t]. \quad [23]$$

Since we use truncated backprop we have  $\mathbf{E}_t[t] = \sum_k m_{t,k} \mathbf{E}_t[t+1]$ ; therefore we get

$$\frac{\mathbf{E}_{\text{c}_j}[t]}{\mathbf{E}_{\text{c}_j}[t+1]} = \sum_k m_{t,k} \frac{\mathbf{E}_{\text{c}_j}[t+1]}{\mathbf{E}_{\text{c}_j}[t+1]} \approx_{t_M} 1. \quad [24]$$

The previous equations [88] and [89] imply constant error flow through internal states of memory cells:

$$\frac{\partial E_{c_{ij}}[t]}{\partial E_{c_{ij}}[t+1]} = \frac{E_{x_{ij}}[t+1]}{E_{x_{ij}}[t]} \approx_{ta} 1. \quad [90]$$

The error occurring at the memory cell input is

$$E_{inj}[t] = \frac{E_{y^{outj}}[t]}{E_{outj}[t]} \frac{E_{x_{ij}}[t]}{E_{y^{outj}}[t]} E_{c_{ij}}[t]. \quad [91]$$

The error occurring at the input gates is

$$E_{inj}[t] \approx_{ta} \frac{E_{y^{outj}}[t]}{E_{outj}[t]} \frac{E_{x_{ij}}[t]}{E_{y^{outj}}[t]} E_{c_{ij}}[t]. \quad [92]$$

**No external error flow.** Errors are propagated back from units  $b$  to unit  $a$  along outgoing connections with weights  $m_{ab}$ . This “external error” [note that for conventional units there is nothing but external error] at time  $t$  is

$$E_a[t] = \frac{E_{y^a}[t]}{E_{outa}[t]} \sum_b \frac{E_{outb}[t+1]}{E_{y^a}[t]} m_{ba}[t+1]. \quad [93]$$

We obtain

$$\begin{aligned} \frac{\partial E_a[t-1]}{\partial E_a[t]} &= \quad [94] \\ \frac{E_{y^a}[t-1]}{E_{outa}[t-1]} \left( \frac{\partial E_{outj}[t]}{\partial E_y[t]} \frac{\partial E_{outj}[t]}{\partial E_{y^a}[t-1]} + \frac{\partial E_{outj}[t]}{\partial E_x[t]} \frac{\partial E_{outj}[t]}{\partial E_{y^a}[t-1]} + \frac{\partial E_{outj}[t]}{\partial E_c[t]} \frac{\partial E_{outj}[t]}{\partial E_{y^a}[t-1]} \right) &\approx_{ta} 1. \end{aligned}$$

We observe: the error  $E_y$  arriving at the memory cell output is not backpropagated to units  $a$  via external connections to  $b_{outj}, b_{outk}, b_g$ .

**Error flow within memory cells.** We now focus on the error flow within a memory cell’s SIMD. This is actually the only type of error flow that can bridge several time steps. Suppose error  $E_y[t]$  arrives at  $y_j$ ’s output at time  $t$  and is propagated back for  $g$  steps until it reaches  $b_{outj}$  or the memory cell input  $y_j^{outj}[t]$ . It is scaled by a factor of  $\frac{E_{y^j}[t-x]}{E_{y^j}[t]}$ , where  $x = b_{outj}, y_j$ . We first compute

$$\frac{\partial E_{c_{ij}}[t-x]}{\partial E_y[t]} \approx_{ta} \begin{cases} \frac{E_{y^{outj}}[t]}{E_{c_{ij}}[t]} & x = 0 \\ \frac{E_{c_{ij}}[t-x+1]}{E_{c_{ij}}[t-x]} \frac{E_{c_{ij}}[t-x+1]}{E_{y^{outj}}[t-x]} & x > 0 \end{cases}. \quad [95]$$

Decomposing equation [95], we obtain

$$\begin{aligned} \frac{\partial E_a[t-x]}{\partial E_y[t]} &\approx_{ta} \frac{\partial E_a[t-x]}{\partial E_{c_{ij}}[t-x]} \frac{\partial E_{c_{ij}}[t-x]}{\partial E_y[t]} \approx_{ta} \\ &\frac{\partial E_a[t-x]}{\partial E_{c_{ij}}[t-x]} \left( \prod_{m=1}^1 \frac{E_{x_{ij}}[t-m+1]}{E_{x_{ij}}[t-m]} \right) \frac{E_{y^a}[t]}{E_{x_{ij}}[t]} \approx_{ta} \\ &E_{y^{outj}}[t] \frac{E_{x_{ij}}[t]}{E_{y^{outj}}[t]} \left\{ \begin{array}{ll} \frac{y_j^{outj}[t-x]}{y_j^{outj}[t]} & x = y_j \\ \frac{y_j^{outj}[t-x]}{b_{outj}[t-x]} & x = b_{outj} \end{array} \right. . \end{aligned} \quad [96]$$

Consider the factors in the previous equation’s last expression. Obviously, error flow is scaled only at times  $t$  [when it enters the cell] and  $t-x$  [when it leaves the cell], but not in between [constant error flow through the SIMD]. We observe:

[1] This output gate's effect is:  $y^{out}[t]$  scales down those errors that can be reduced early during training without using the memory cell. Likewise, it scales down those errors resulting from using [activating/deactivating] the memory cell at later training stages — without this output gate, the memory cell might for instance suddenly start causing avoidable errors in situations that already seemed under control [because it was easy to reduce the corresponding errors without memory cells]. See “output weight conflict” and “efflux problem” in Sections 8/9.

[2] If there are large positive or negative  $s_{h,i}[t]$  values [because  $s_{h,i}$  has shifted since time step  $t - \tau$ ], then  $W[s_{h,i}[t]]$  may be small [assuming that  $W$  is a logistic sigmoid]. See Section 8. Drifts of the memory cell's internal state  $s_{h,i}$  can be counteracted by negatively biasing the input gate  $b_{hi}$  [see Section 9 and next point]. Recall from Section 8 that this precise bias value does not matter much.

[3]  $y^{out}[t - \tau]$  and  $b_{hi}^T[s_{h,i}[t - \tau]]$  are small if the input gate is negatively biased [assume  $b_{hi}$  is a logistic sigmoid]. However, the potential significance of this is negligibly compared to the potential significance of drifts of the internal state  $s_{h,i}$ .

Some of the factors above may weakly bias LSTM's overall error flow, but not in a manner that depends on the length of the time lag. The flow will still be much more effective than an exponentially [cf. earlier  $\gamma$ ] decaying flow without memory cells.

## References

- Almeida, L. R. [1987]. A learning rule for asynchronous perceptrons with feedback in a combinatorial environment. In *1987 International Conference on Neural Networks, San Diego*, volume 2, pages 303–318.
- Baldi, P. and Hinneburg, H. [1991]. Contrastive learning and neural oscillator. *Neural Computation*, 3:523–545.
- Bengio, Y. and Hinton, G. [1995]. Credit assignment through time: Alternatives to backpropagation. In Cowan, J. D., Tesauro, G., and Alspector, J., editors, *Advances in Neural Information Processing Systems 7*, pages 62–68. San Mateo, CA: Morgan Kaufmann.
- Bengio, Y., Simard, P., and Frasconi, P. [1995]. Learning long-term dependencies with gradient descent is difficult. *1995 International Conference on Neural Networks*, 5(2):127–133.
- Chernoutsos, A., Sommerville, B., and McCallum, A. L. [1992]. Finite-state automata and simple recurrent networks. *Neural Computation*, 1:875–881.
- de Araya, E. and Principe, J. C. [1991]. A theory for neural networks with time delays. In Lippmann, R. P., Moody, J. D., and Touretzky, D. S., editors, *Advances in Neural Information Processing Systems 3*, pages 135–148. San Mateo, CA: Morgan Kaufmann.
- Deng, N. [1992]. Bifurcations in the learning of recurrent neural networks. In *Proceedings of 1992 IEEE International Symposium on Circuits and Systems*, pages 2677–2680.
- Deng, N. and Stellebahn, S. [1993]. Adaptive neural oscillator using continuous-time back-propagation learning. *Neural Networks*, 6:875–888.
- Dilman, A. L. [1988]. Binding structures in time. Technical Report CRL Technical Report 88/1, Center for Research in Language, University of California, San Diego.
- Fallside, S. D. [1991]. The recurrent cascade-correlation learning algorithm. In Lippmann, R. P., Moody, J. D., and Touretzky, D. S., editors, *Advances in Neural Information Processing Systems 3*, pages 187–194. San Mateo, CA: Morgan Kaufmann.
- Hochreiter, S. [1991]. Untersuchungen zu dynamischen neuronalen Netzen. Diploma thesis, Institut für Informatik, Lehrstuhl Prof. Brauer, Technische Universität München. See <http://informatik.tu-muenchen.de/~ Hochrait>.

- Hochreiter, S. and Schmidhuber, J. [1997]. Long short-term memory. Technical Report 1-97-05, Fakultät für Informatik, Technische Universität München.
- Hochreiter, S. and Schmidhuber, J. [1998]. Bridging long times lags by weight guessing and “Long Short-Term Memory”. In Silver, D. L., Brants, H. G., and Almeida, L. R., editors, *Symposium on models for language and intelligent systems*, pages 35–52. IOS Press, Amsterdam, Netherlands. Series: Frontiers in Artificial Intelligence and Applications, Volume 57.
- Hochreiter, S. and Schmidhuber, J. [1999]. LSTM can solve hard long time lag problems. In *Advances in Neural Information Processing Systems 12*. MIT Press, Cambridge MA. Presented at NIPS 98.
- Lam, W., Waibel, A., and Hinton, G. D. [1990]. A time-delay neural network architecture for isolated word recognition. *Neural Networks*, 3:89–102.
- Miller, G. E. and Giles, C. L. [1998]. Experimental comparison of the effect of order in recurrent neural networks. *International Journal of Human-Computer Studies*, 40(3):513–534.
- Möller, M. C. [1989]. A focused back-propagation algorithm for temporal sequence recognition. *Complex Systems*, 3:379–391.
- Möller, M. C. [1991]. Induction of multiscale temporal structure. In Lippmann, R. P., Moody, J. D., and Touretzky, D. S., editors, *Advances in Neural Information Processing Systems 4*, pages 273–282. San Mateo, CA: Morgan Kaufmann.
- Rumelhart, D. E. [1986]. Learning static space trajectories in recurrent neural networks. *Neural Computation*, 1(2):239–262.
- Rumelhart, D. E. [1995]. Gradient calculations for dynamic recurrent neural networks: A survey. *IEEE Transactions on Neural Networks*, 6(3):1818–1828.
- Sinatra, D. D. [1987]. Generalization of back-propagation to recurrent neural networks. *Highlevel Network Models*, 13(52):2329–2332.
- Sinatra, D. D. [1988]. Dynamics and architecture for neural computation. *Journal of Complexity*, 4:213–232.
- Sohn, T. A. [1993]. Holographic recurrent networks. In S. J. Hanson, D. E. C. and Giles, C. L., editors, *Advances in Neural Information Processing Systems 5*, pages 39–51. San Mateo, CA: Morgan Kaufmann.
- Solla, S. A. [1991]. Language induction by global transition in dynamical recognizers. In Lippmann, R. P., Moody, J. D., and Touretzky, D. S., editors, *Advances in Neural Information Processing Systems 4*, pages 313–320. San Mateo, CA: Morgan Kaufmann.
- Guskevius, G. V. and Polikarpov, L. A. [1995]. Neurocontrol of nonlinear dynamical systems with Kalman filter trained recurrent networks. *IEEE Transactions on Neural Networks*, 5(3):359–367.
- Ring, M. D. [1995]. Learning sequential tasks by incrementally adding higher orders. In S. J. Hanson, D. E. C. and Giles, C. L., editors, *Advances in Neural Information Processing Systems 7*, pages 115–122. Morgan Kaufmann.
- Robinson, A. J. and Falls, D. [1987]. The utility driven dynamic error propagation network. Technical Report COIN/12-INENG/TR.1, Cambridge University Engineering Department.
- Schmidhuber, J. [1992]. The Neural Bucket Brigade: A local learning algorithm for dynamic feedforward and recurrent networks. *Connection Science*, 4(3):303–312.

- Schmidhuber, J. [1992a]. A fixed size storage  $O(n^2)$  time complexity learning algorithm for fully recurrent continually running networks. *Neural Computation*, 4(3):213–235.
- Schmidhuber, J. [1992b]. Learning complete, unlabelled sequences using the principle of history compression. *Neural Computation*, 4(3):237–255.
- Schmidhuber, J. [1992c]. Learning unambiguous reduced sequence descriptions. In Meedya, H. D., Hansen, S. H., and Lippmann, R. P., editors, *Advances in Neural Information Processing Systems 4*, pages 221–235. San Mateo, CA: Morgan Kaufmann.
- Schmidhuber, J. [1992d]. Netzwerkarchitekturen, Zielfunktionen und Netzwerkgod. Habilitationschrift, Institut für Informatik, Technische Universität München.
- Schmidhuber, J. and Hochreiter, S. [1992]. Guessing can outperform many long time lag algorithms. Technical Report DECIM-12-92, DECIM.
- Silva, G. N., Gomez, D. E., Langbein, T., and Almeida, L. R. [1992]. Faster training of recurrent networks. In Silva, G. N., Brinza, D. N., and Almeida, L. R., editors, *Symbolic and parallel systems*, pages 168–175. IOS Press, Amsterdam, Netherlands. Series: Frontiers in Artificial Intelligence and Applications, Volume 55.
- Smith, A. W. and Zippser, D. [1992]. Learning sequential structures with the real-time recurrent learning algorithm. *Unknown and Unknown of Neural Systems*, 1(2):125–131.
- Sun, C., Elton, H., and Law, M. [1992]. Time warping invariant neural networks. In S. H. Hansen, H. D. Meedya, and Giles, C. L., editors, *Advances in Neural Information Processing Systems 4*, pages 180–185. San Mateo, CA: Morgan Kaufmann.
- Whitney, R. L. and Nolin, G. M. [1992]. Induction of finite-state languages using second-order recurrent networks. *Neural Computation*, 4:103–111.
- Whelber, S. J. [1992]. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 5:1.
- Williams, R. J. [1992]. Complexity of exact gradient computation algorithms for recurrent neural networks. Technical Report Technical Report NEU-TRCS-92-37, Boston: Northeastern University, College of Computer Science.
- Williams, R. J. and Peng, Y. [1992]. An efficient gradient-based algorithm for on-line training of recurrent network trajectories. *Neural Computation*, 4:191–201.
- Williams, R. J. and Zipser, D. [1992]. Gradient-based learning algorithms for recurrent networks and their computational complexity. In *Back-propagation: Theory, Architectures and Applications*. Hillsdale, NJ: Erlbaum.