# Learning Nonregular Languages: A Comparison of Simple Recurrent Networks and LSTM

**J. Schmidhuber**
*juergen@idsia.ch*
**F. Gers**
*felix@idsia.ch*
**D. Eck**
*doug@idsia.ch*
*Istituto Dalle Molle di Studi sull'Intelligenza Artificiale, 6928 Manno-Lugano, Switzerland*

**In response to Rodriguez's recent article (2001), we compare the performance of simple recurrent nets and long short-term memory recurrent nets on context-free and context-sensitive languages.**

Rodriguez (2001) examined the learning ability of simple recurrent nets (SRNs) (Elman, 1990) on simple context-sensitive and context-free languages (CFLs and CSLs, respectively). He trained his SRN on short training sequences of length not much greater than 10 and found that the SRN does not generalize well on significantly larger test sets and frequently is unable to store the training set. Similar results were recently reported by Bodén and Wiles (2000), who studied the simple context-sensitive language (CSL) $a^n b^n c^n$. They trained SRNs on sequences defined by $n = 1, 2, \ldots, 10$ and found that SRNs fail to generalize reliably to $n > 13$. Sequential cascaded networks (SCNs) did only moderately better.

We applied long short-term memory (LSTM) recurrent nets (Hochreiter & Schmidhuber, 1997) to similar problems (Gers & Schmidhuber, 2001a, 2001b). Many of our training set sizes were comparable to those of Rodriguez (2001) and Bodén and Wiles (2000). We found that LSTM almost always stores the training set and generalizes well on much larger test sets. For instance, we trained LSTM on strings from the CSL $a^n b^n c^n$ with $n <= 40$. From this training data, LSTM readily generalized up to string size 1500; that is, it learned to accept legal test strings such as $a^{500} b^{500} c^{500}$ while rejecting slightly different illegal strings such as $a^{500} b^{499} c^{500}$. In other experiments, LSTM was able to generalize well having seen only very limited training data. For example, from the training strings with $n = 50, 51$, LSTM generalized to $43 <= n <= 57$. Similarly, LSTM trained on the CFL $a^n b^n$ for $1 <= n <= 30$ generalized in the best case up to $n = 1000$ and on average

up to $n = 408$. We refer readers to Gers and Schmidhuber (2001b) for an analysis of this and similar tasks, as well as for LSTM equations.

True, even LSTM did not learn these languages in the strict sense; we do not have proof that it will generalize to arbitrary $n$. But it is obvious that LSTM exhibits excellent generalization performance, while SRNs do not. Because LSTM was designed to solve hard problems and not specifically to address human cognition, it remains to be seen how well it will model behavioral data.

LSTM also successfully learned CSLs that require embedded counters, for example, the deterministic palindrome language $a^m b^n B^n A^m$ that Rodriguez (2001) discussed; traditional recurrent neural networks (RNNs) did not even learn the training set. We refer readers to Gers and Schmidhuber (2001b) for details. Testing more complex CFLs is also a topic of future research.

How does LSTM solve problems like $a^n b^n c^n$? Counters of potentially unlimited size are automatically and naturally implemented by linear units, the constant error carousels (CECs) of standard LSTM, originally designed to overcome error decay problems plaguing previous RNNs (Hochreiter, 1991; Hochreiter, Bengio, Frasconi, & Schmidhuber, 2001). Each linear CEC is surrounded by a cloud of nonlinear units responsible for controlling the flow of information in and out of the CEC. Whereas standard RNNs have a hard time implementing precise linear counters with their squashing functions—they need finely tuned weights to countermand the nonlinearities—LSTM can simply use the CECs for counting and so is free to focus its robust and precise weight adaptation process on the nonlinear aspects of sequence processing.

In the case of the CFL $a^n b^n$, LSTM uses one CEC to count up for $a$'s and down for $b$'s while using nonlinear gates to protect the CEC from irrelevant signals and generate the appropriate accept or reject response. In the case of the CSL $a^n b^n c^n$, LSTM uses one CEC to count up for $a$'s and down for $b$'s and a second CEC to count up for $b$'s and down on $c$'s. SRN and SCN counting mechanisms are analyzed in Rodriguez (2001) and Bodén and Wiles (2000).

The CEC-based counters not only deal much better with the well-known long time lag problem (Hochreiter, 1991; Hochreiter et al., 2001) but also the problem of installing easily accessible linear counters whose increments and decrements implement natural push and pop operations. This explains in a nutshell why LSTM outperforms other recurrent nets not only on regular (Hochreiter & Schmidhuber, 1997) but also on context-free and context-sensitive languages (Gers & Schmidhuber, 2001b).

**References**

Bodén, M., & Wiles, J. (2000). Context-free and context-sensitive dynamics in recurrent neural networks. *Connection Science, 12*(3/4), 197–210.

Elman, J. L. (1990). Finding structure in time. *Cognitive Science, 14*, 179–211.

Gers, F. A., & Schmidhuber, J. (2001a). Long Short-Term Memory learns context free and context sensitive languages. In V. Kurkova et al. (Eds.), *Proceedings of the ICANNGA 2001 Conference* (Vol. 1, pp. 134–137). New York: Springer-Verlag.

Gers, F. A., & Schmidhuber, J. (2001b). LSTM recurrent networks learn simple context free and context sensitive languages. *IEEE Transactions on Neural Networks, 12*(6), 1333–1340.

Hochreiter, S. (1991). *Untersuchungen zu dynamischen neuronalen Netzen*. Diploma thesis, Institut für Informatik, Technische Universität München. Available on-line: ni.cs.tu-berlin.de/∼hochreit.

Hochreiter, S., Bengio, Y., Frasconi, P., & Schmidhuber, J. (2001). Gradient flow in recurrent nets: The difficulty of learning long-term dependencies. In S. C. Kremer & J. F. Kolen (Eds.), *A field guide to dynamical recurrent neural networks*. Parsippany, NJ: IEEE Press.

Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation, 9*(8), 1735–1780.

Rodriguez, P. (2001). Simple recurrent networks learn context-free and context-sensitive languages by counting. *Neural Computation, 13*(9), 2093–2118.