# Realtime control of sequence generation with character based Long Short Term Memory Recurrent Neural Networks

**Memo Akten, Goldsmiths University of London**

Mick Grierson, Goldsmiths University of London

Huosheng Hu, University of Essex

## Abstract

Recurrent Neural Networks (RNNs) — particularly Long Short Term Memory (LSTM) RNNs — are a popular and very successful model for generating sequences. However, most LSTM based sequence generation techniques are currently not interactive and do not allow continuous control of the sequence generation, let alone in a gestural or expressive manner. This research investigates methods of realtime continuous control and steering of RNN sequence generation, as well as ways of expressively controlling the output.

## Introduction

Recurrent Neural Networks (RNN) are Artificial Neural Networks with recurrent connections, allowing them to model sequences. Long Short Term Memory (LSTM) networks (Hochreiter and Schmidhuber 1997) are a kind of RNN with differentiable gates, allowing them to overcome the so-called *exploding and vanishing gradients* problem and preserve information from many time-steps in the past.

Now with increased compute power and large training sets, LSTMs and variants are having lots of success not only in sequence *classification* (Greff et al. 2015), but also in sequence *generation* (Eck and Schmidhuber 2002; Sutskever, Martens, and Hinton 2011; Sutskever 2013; Graves 2013; Nayebi and Vitelli 2015; Gregor et al. 2015; Wu and King 2016; Crnkovic-friis and Crnkovic-friis 2016).

However, sequence generation with RNNs is currently not an interactive process. Some recent implementations have used a turn-based approach, such as the online text editor Word Synth (Goodwin 2016), which allows a user to enter a seed for the RNN before it generates the next phrase. Though this is still not *realtime continuous* control of the output, let alone an *expressive* experience.

## Method

We're using character based text models as it has been demonstrated that LSTMs are successful in generating character sequences (Graves 2013; Karpathy 2015). First we train a number of models on varied corpus of text. Then the system generates text character by character, and while it's outputting, we can gesturally steer the output towards the style of different models.

**1. Training** We train numerous LSTM models, each on a different corpus of text. These include the works of Shakespeare, Baudelaire, Nietzsche, Jane Austen, Donald Trump, the Dalai Lama, the King James Bible, assorted love song lyrics, Linux kernel C code, LaTeX source and more. Each corpus is picked due to its nature of being easily recognizable in style and content.

**2. Prediction, visualization and interaction** Once trained, the system generates text by: running each of the models, mixing their output probability distributions $p_i$ via model mixture weights $w_m$, sampling a character from the mixed probability distribution $p_{mix}$, and feeding the sampled character back in as an input.

The system generates and outputs characters at roughly 10-20 chars/second. While it's outputting, we can steer the output towards different models by interacting with the system and dynamically shifting the mixture weights $w_m$ for each model. Interaction is through moving the mouse cursor, moving our hand (tracked by a LeapMotion), or playing with midi sliders.

We are thus able to guide the system to morph between the different models' output with relatively smooth transitions.

## Results and discussion

In this study we trained multiple models on different corpora and mixed their predicted $p_i$ via $w_m$, interactively controlled at every step by a user's gestures. This allows a user to continuously control the output of an LSTM as it generates a sequence, and seamlessly morph it between styles while it's generating. Figure 1 shows an example output.

The system works well and demonstrates interesting behaviour. When multiple models are active, containing words or phrases that are common to all models, the probabilities for the common next characters accumulate whilst probabilities specific to individual models are suppressed, i.e. when multiple models are active the system tends towards common words and phrases.

Sometimes, while a sequence is being generated, a particular model might have a spiking $p_i$ (very high confidence for a particular character). If at that point other models have wider $p_i$ (lower confidence across multiple characters), then the first model will overpower and dominate the se-
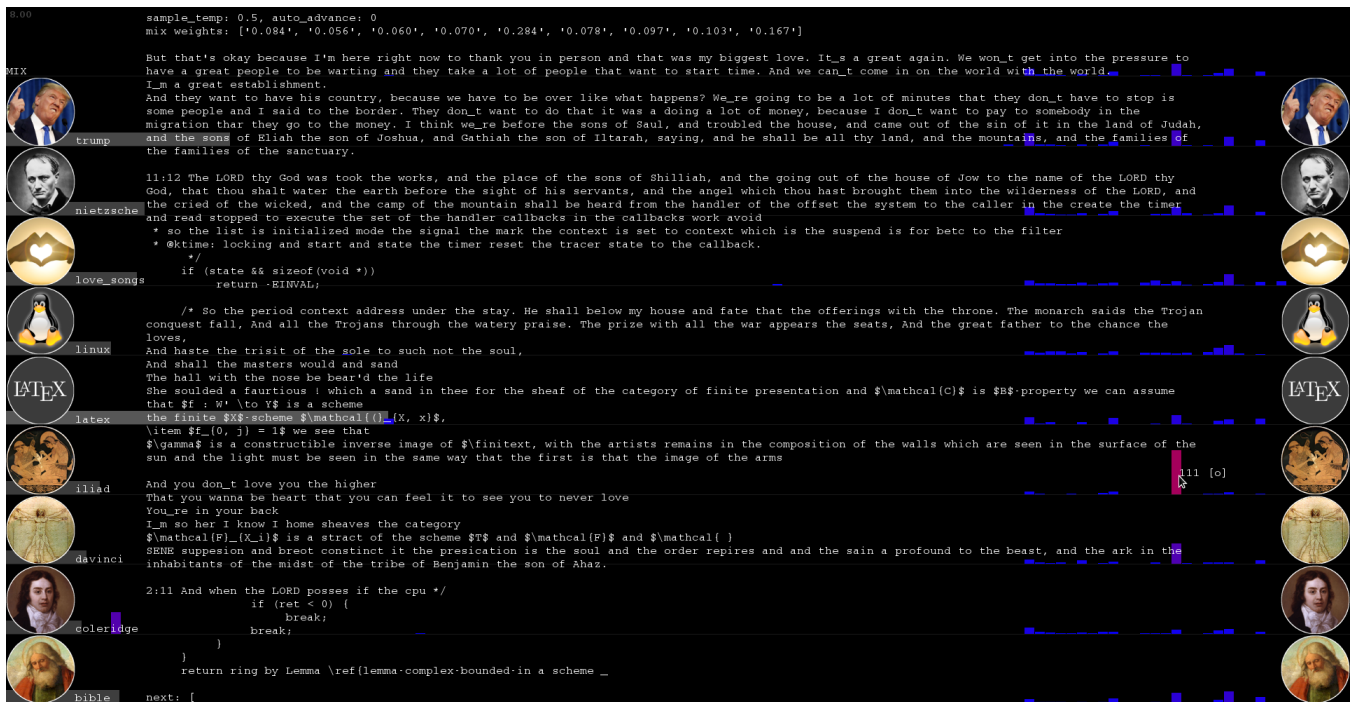
Figure 1: An example of the screen and output. Each model's output probability distributions $\mathbf{p_i}$ can be seen, as well as the model mixture weights $\mathbf{w_m}$ (horizontal grey bars) and final mixed probability distribution $\mathbf{p_{mix}}$ (top row).

quence generation. This is quite likely to start a positive feedback loop and that model will stay in control of the sequence generation until it reaches a point where its probability distribution widens, and another model spikes. So it's very possible to see hints of love songs, philosophy or poetry within C comments and variable names or LaTeX equations. It seems there are 'hand-over' words or sequences which are common to many models, but have stronger connotations in some models over others. This is of course further guided by the user's actions, who can choose to push further towards the emerging theme, or pull towards another style and seamlessly go from one style to another over these hand-over words.

As opposed to training many models independently on different corpora, another approach we are looking at is using a single model trained on the entire corpora. Then controlling output via manipulating the internal state of the LSTM. This has advantages and disadvantages, particularly when it comes to adding a new corpus (style) to the system.

## References

Crnkovic-friis, L., and Crnkovic-friis, L. 2016. Generative Choreography using Deep Learning. *arXiv preprint arXiv:1605.06921*.

Eck, D., and Schmidhuber, J. 2002. A First Look at Music Composition using LSTM Recurrent Neural Networks. *Istituto Dalle Molle Di Studi Sull Intelligenza Artificiale* 103.

Goodwin, R. 2016. Word Synth.

Graves, A. 2013. Generating sequences with Recurrent Neural Networks. *arXiv preprint arXiv:1308.0850*.

Greff, K.; Srivastava, R. K.; Koutník, J.; Steunebrink, B. R.; and Schmidhuber, J. 2015. LSTM: A Search Space Odyssey. *arXiv preprint arXiv:1503.04069*.

Gregor, K.; Danihelka, I.; Graves, A.; and Wierstra, D. 2015. DRAW: A Recurrent Neural Network For Image Generation. *arXiv preprint arXiv:1502.04623*.

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

Karpathy, A. 2015. The Unreasonable Effectiveness of Recurrent Neural Networks.

Nayebi, A., and Vitelli, M. 2015. GRUV : Algorithmic Music Generation using Recurrent Neural Networks.

Sutskever, I.; Martens, J.; and Hinton, G. 2011. Generating Text with Recurrent Neural Networks. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, volume 131. 1017–1024.

Sutskever, I.; Vinyals, O.; and Le, Q. V. 2014. Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems*, 3104–3112.

Sutskever, I. 2013. *Training Recurrent Neural Networks*. Ph.D. Dissertation, University of Toronto.

Wu, Z., and King, S. 2016. Investigating Gated Recurrent Neural Networks for Speech Synthesis. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 5140–5144. IEEE.