

Forget to remember
Remember to forget

Long Short Term Memories and Gated Recurrent Units

Jonathon Hare

Vision, Learning and Control
University of Southampton

Many of the images and animations used here were made by Adam Prügel-Bennett.

Recap: An RNN is just a recursive function invocation

- $\mathbf{y}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- and the state $\mathbf{c}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- If the output $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-2)$, then prediction will be

$$\mathbf{f}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t-1), \mathbf{g}(\mathbf{x}(t-2) | \mathbf{W}), \mathbf{W}), \mathbf{W}), \mathbf{W})$$

- it should be clear that the gradients of this with respect to the weights can be found with the chain rule

Recap: An RNN is just a recursive function invocation

- $\mathbf{y}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- and the state $\mathbf{c}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- If the output $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-2)$, then prediction will be

$$\mathbf{f}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t-1), \mathbf{g}(\mathbf{x}(t-2) | \mathbf{W}), \mathbf{W}), \mathbf{W}), \mathbf{W})$$

- it should be clear that the gradients of this with respect to the weights can be found with the chain rule
- The back-propagated error will involve applying \mathbf{f} multiple times

Recap: An RNN is just a recursive function invocation

- $\mathbf{y}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- and the state $\mathbf{c}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- If the output $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-2)$, then prediction will be

$$\mathbf{f}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t-1), \mathbf{g}(\mathbf{x}(t-2) | \mathbf{W}), \mathbf{W}), \mathbf{W}), \mathbf{W})$$

- it should be clear that the gradients of this with respect to the weights can be found with the chain rule
- The back-propagated error will involve applying \mathbf{f} multiple times
- Each time the error will get multiplied by some factor a

Recap: An RNN is just a recursive function invocation

- $\mathbf{y}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- and the state $\mathbf{c}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- If the output $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-2)$, then prediction will be

$$\mathbf{f}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t-1), \mathbf{g}(\mathbf{x}(t-2) | \mathbf{W}), \mathbf{W}), \mathbf{W}), \mathbf{W})$$

- it should be clear that the gradients of this with respect to the weights can be found with the chain rule
- The back-propagated error will involve applying \mathbf{f} multiple times
- Each time the error will get multiplied by some factor a
- If $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-\tau)$ then the back-propagated signal will be proportional to $a^{\tau-1}$

Recap: An RNN is just a recursive function invocation

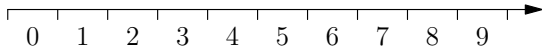
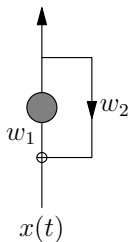
- $\mathbf{y}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- and the state $\mathbf{c}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{c}(t-1) | \mathbf{W})$
- If the output $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-2)$, then prediction will be

$$\mathbf{f}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t), \mathbf{g}(\mathbf{x}(t-1), \mathbf{g}(\mathbf{x}(t-2) | \mathbf{W}), \mathbf{W}), \mathbf{W}), \mathbf{W})$$

- it should be clear that the gradients of this with respect to the weights can be found with the chain rule
- The back-propagated error will involve applying \mathbf{f} multiple times
- Each time the error will get multiplied by some factor a
- If $\mathbf{y}(t)$ depends on the input $\mathbf{x}(t-\tau)$ then the back-propagated signal will be proportional to $a^{\tau-1}$
- This either vanishes or explodes when τ becomes large

Vanishing and Exploding Gradients

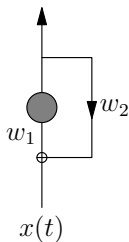
$$y(t) = w_1 (x(t) + w_2 y(t-1))$$



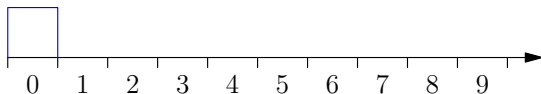
Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

$$w_1 = w_2 = 0.9$$



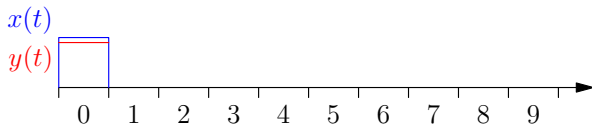
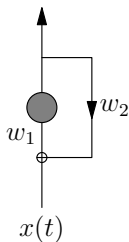
$x(t)$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

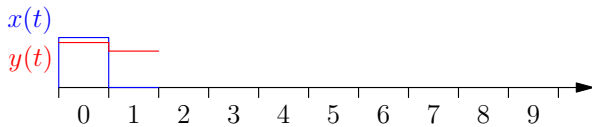
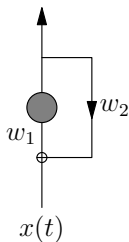
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

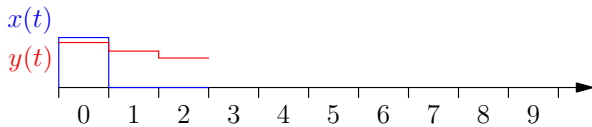
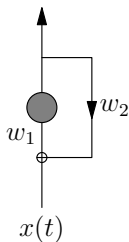
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

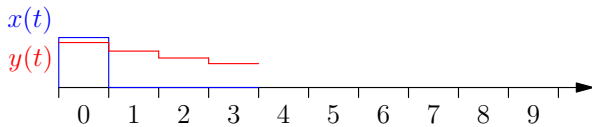
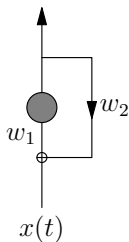
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

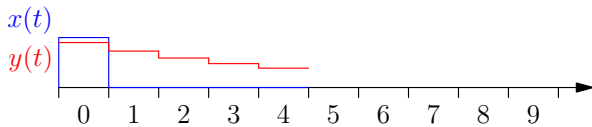
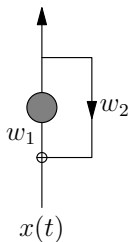
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

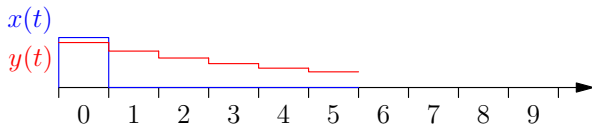
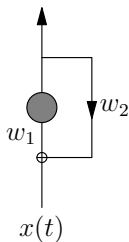
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

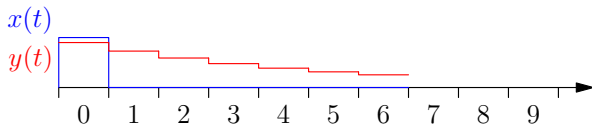
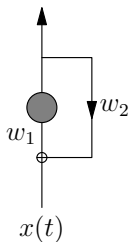
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

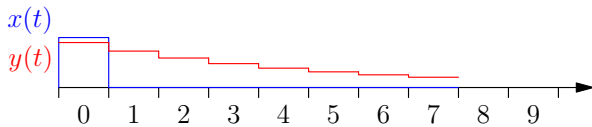
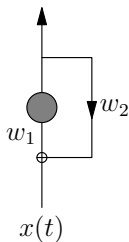
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

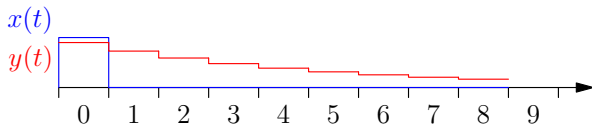
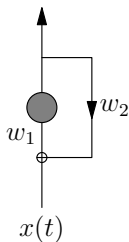
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

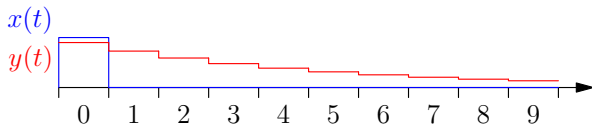
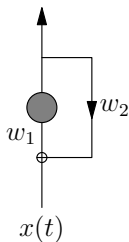
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

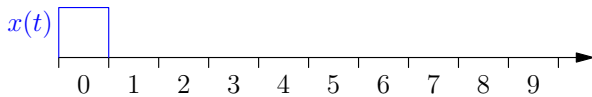
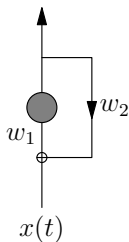
$$w_1 = w_2 = 0.9$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

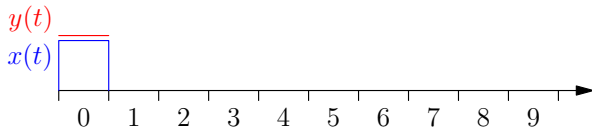
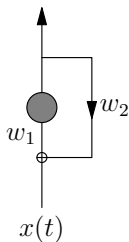
$$w_1 = w_2 = 1.1$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

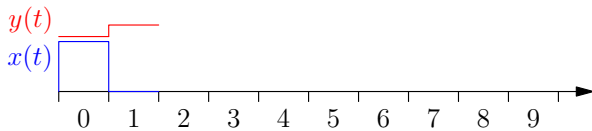
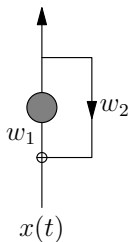
$$w_1 = w_2 = 1.1$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

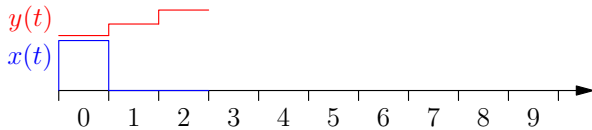
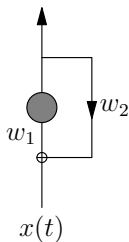
$$w_1 = w_2 = 1.1$$



Vanishing and Exploding Gradients

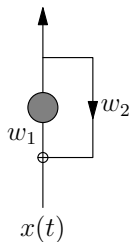
$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

$$w_1 = w_2 = 1.1$$

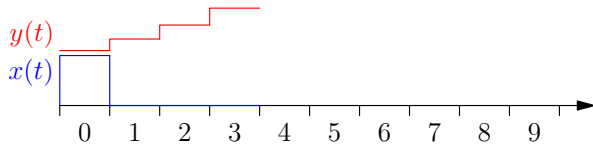


Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$

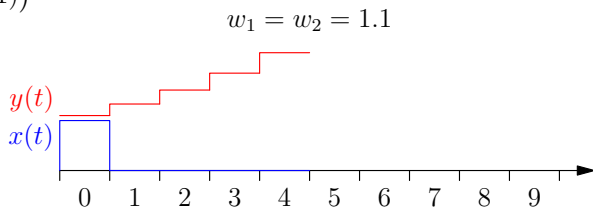
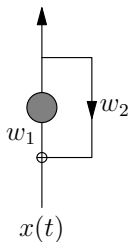


$$w_1 = w_2 = 1.1$$



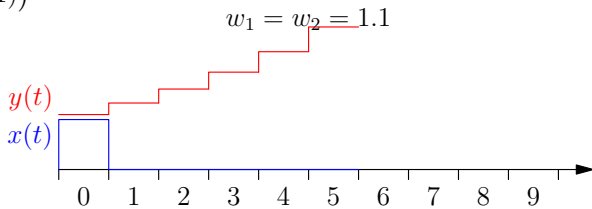
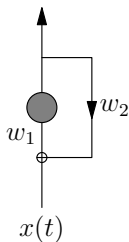
Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$



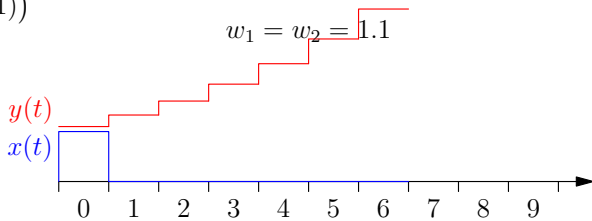
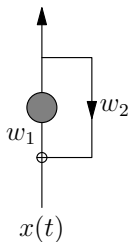
Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$



Vanishing and Exploding Gradients

$$y(t) = w_1 (x(t) + w_2 y(t-1))$$



- LSTMs (long-short term memory) was designed to solve this problem

- LSTMs (long-short term memory) was designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\mathbf{c}(t) = \mathbf{c}(t - 1)$$

- LSTMs (long-short term memory) was designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\mathbf{c}(t) = \mathbf{c}(t - 1)$$

- Sometimes we have to forget and sometimes we have to change a memory

- LSTMs (long-short term memory) was designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\mathbf{c}(t) = \mathbf{c}(t - 1)$$

- Sometimes we have to forget and sometimes we have to change a memory
- To do this we should use 'gates' that saturate at 0 and 1

- LSTMs (long-short term memory) was designed to solve this problem
- Key ideas: to retain a 'long-term memory' requires

$$\mathbf{c}(t) = \mathbf{c}(t - 1)$$

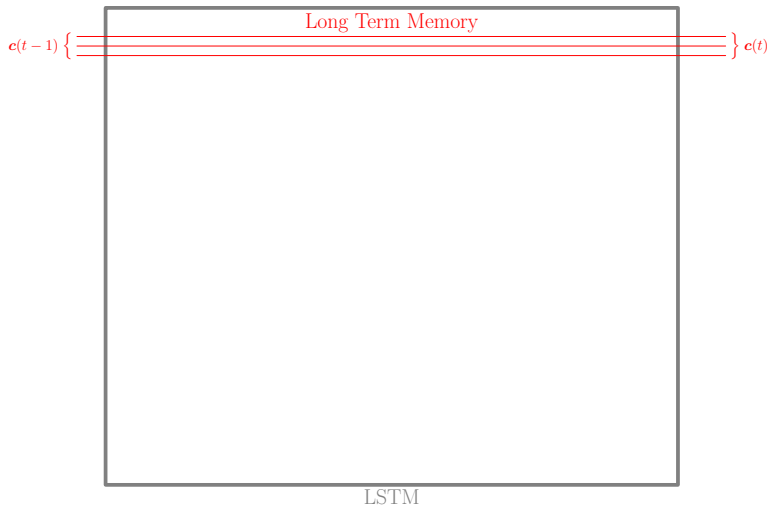
- Sometimes we have to forget and sometimes we have to change a memory
- To do this we should use 'gates' that saturate at 0 and 1
- Sigmoid functions naturally saturate at 0 and 1

LSTM Architecture

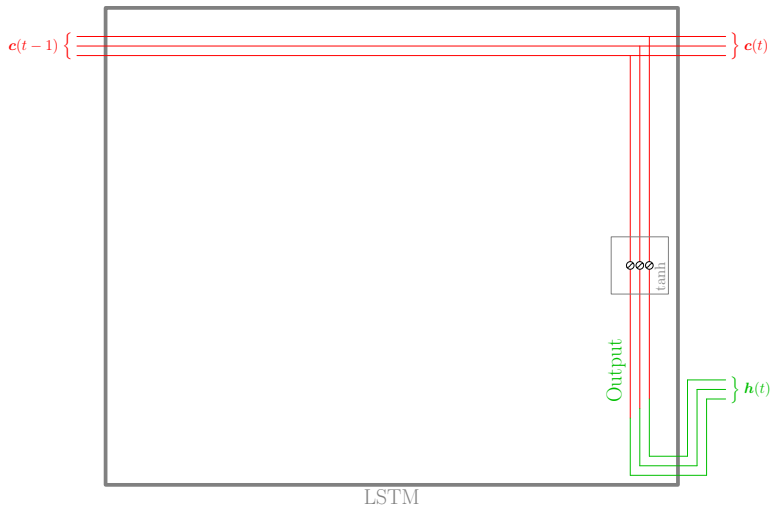


LSTM

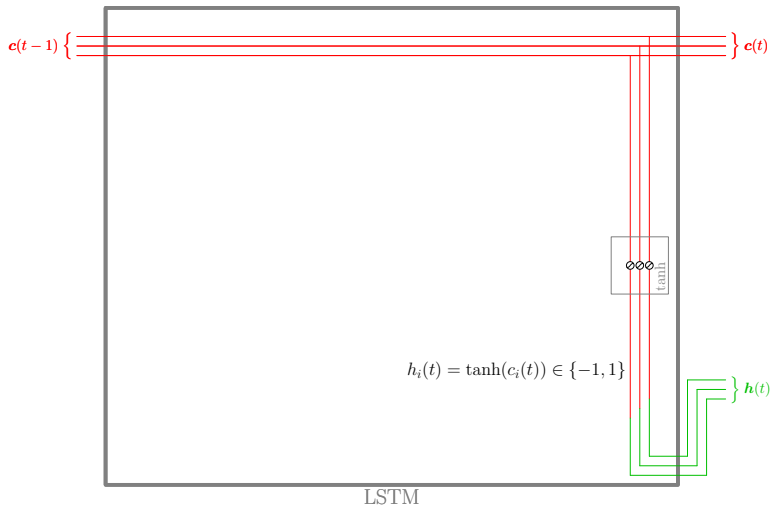
LSTM Architecture



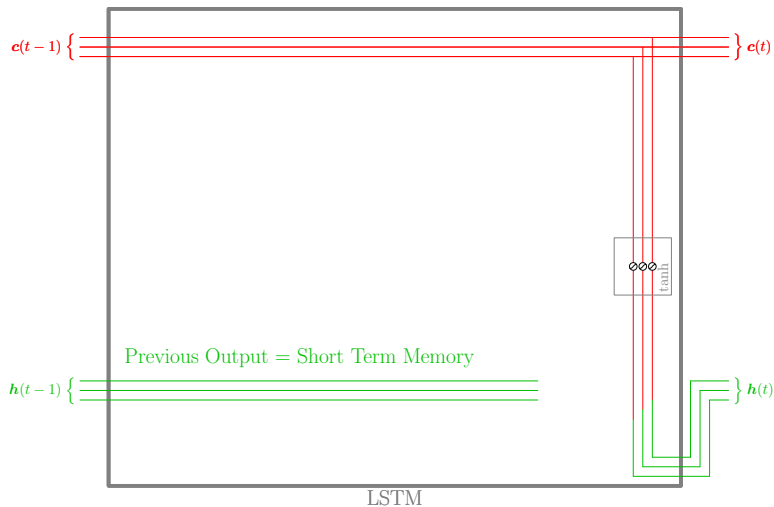
LSTM Architecture



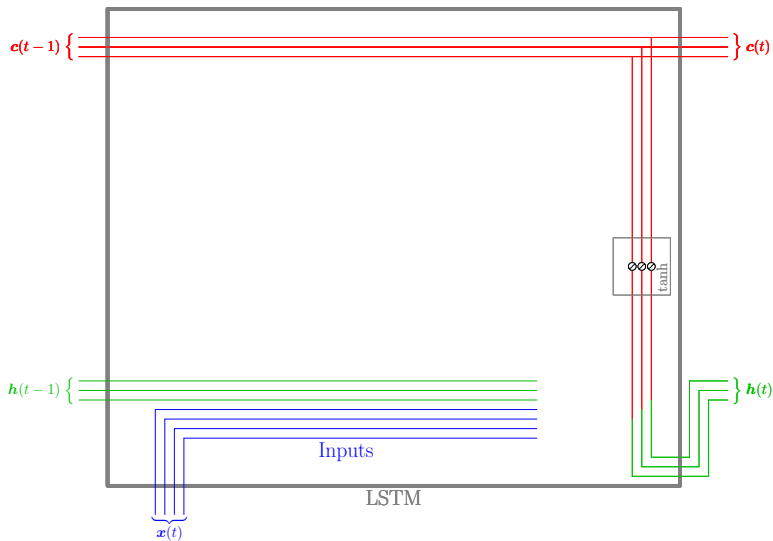
LSTM Architecture



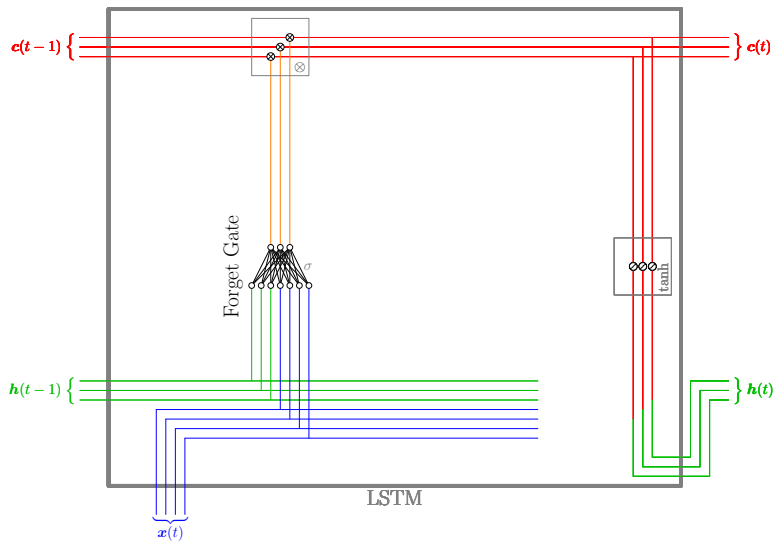
LSTM Architecture



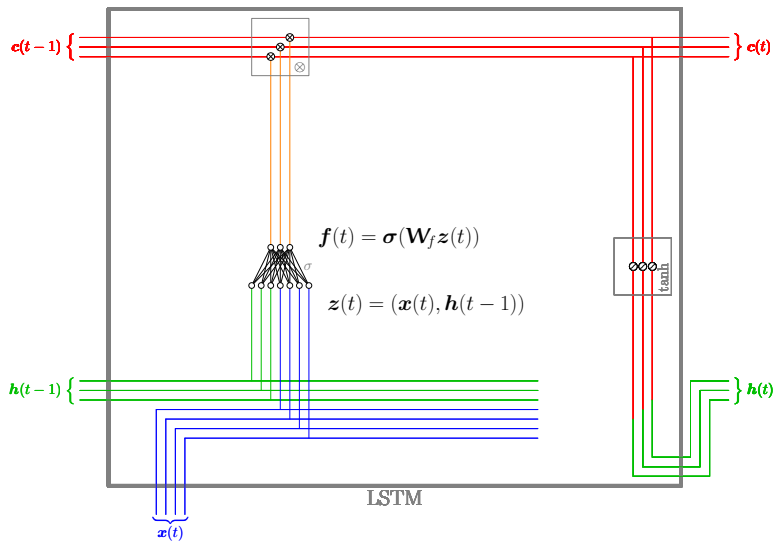
LSTM Architecture



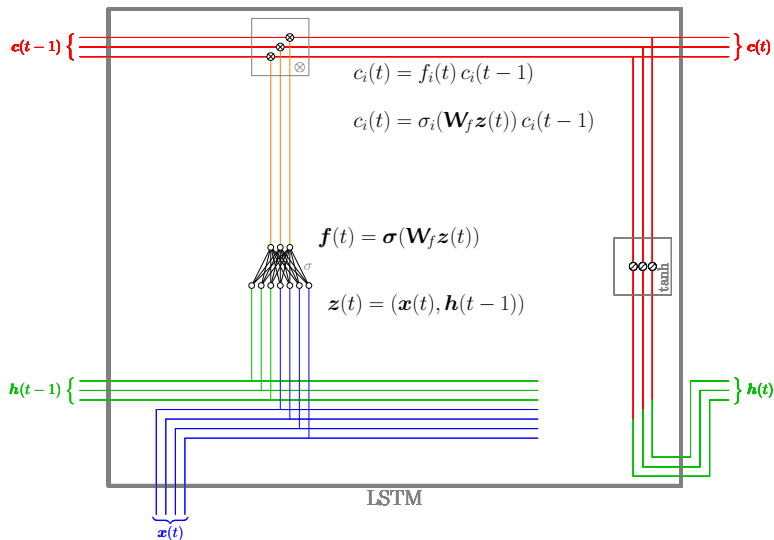
LSTM Architecture



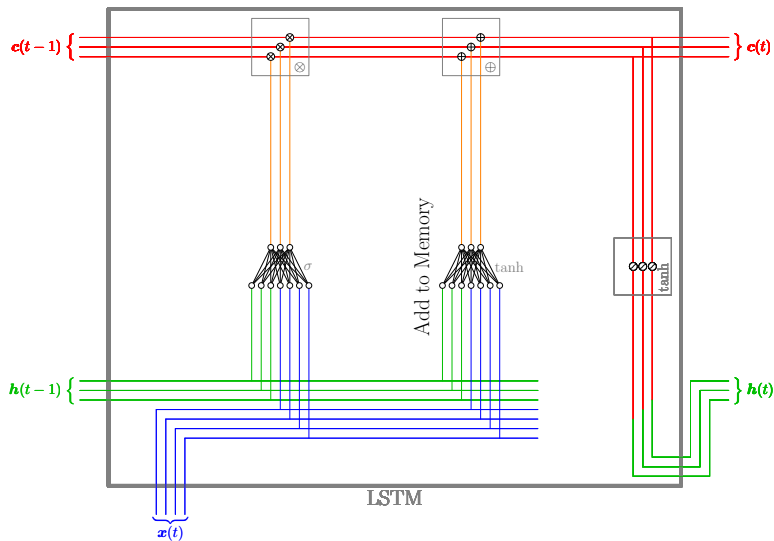
LSTM Architecture



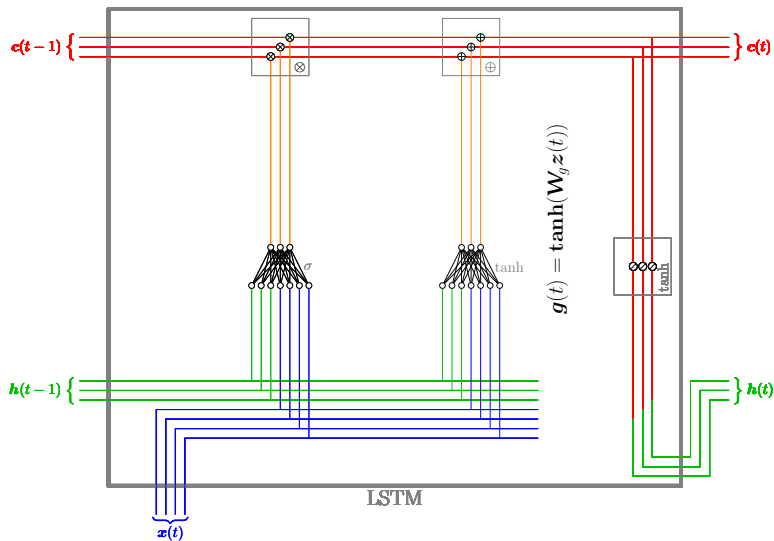
LSTM Architecture



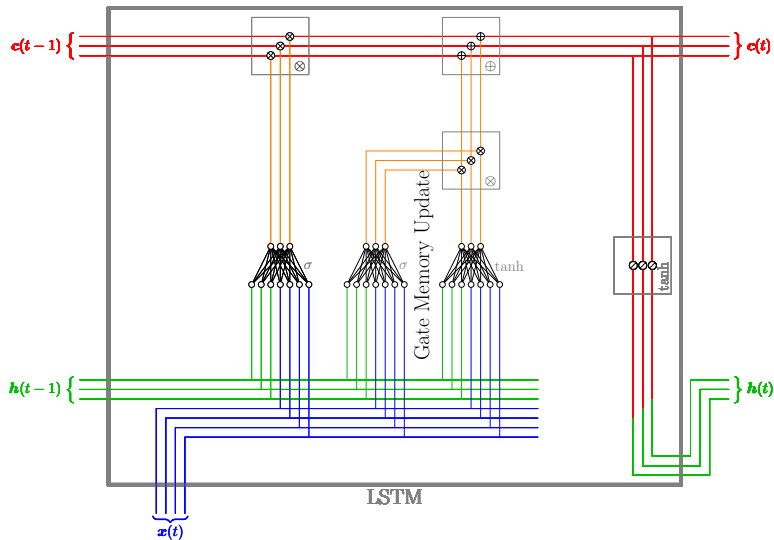
LSTM Architecture



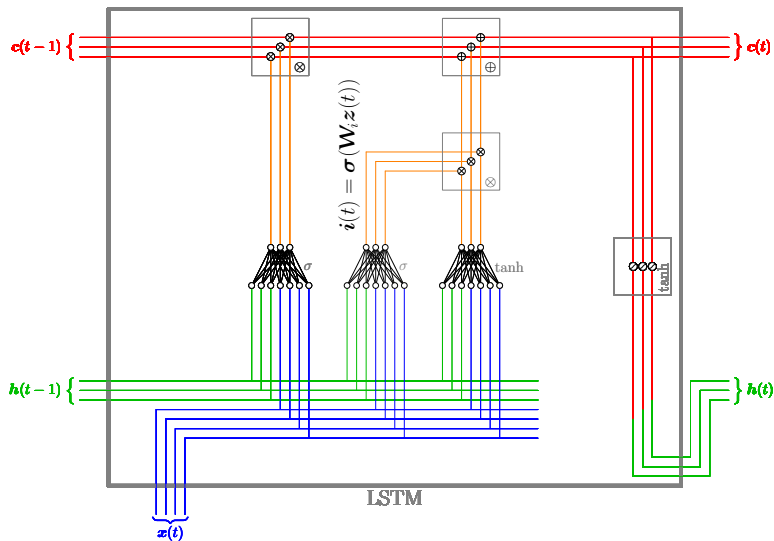
LSTM Architecture



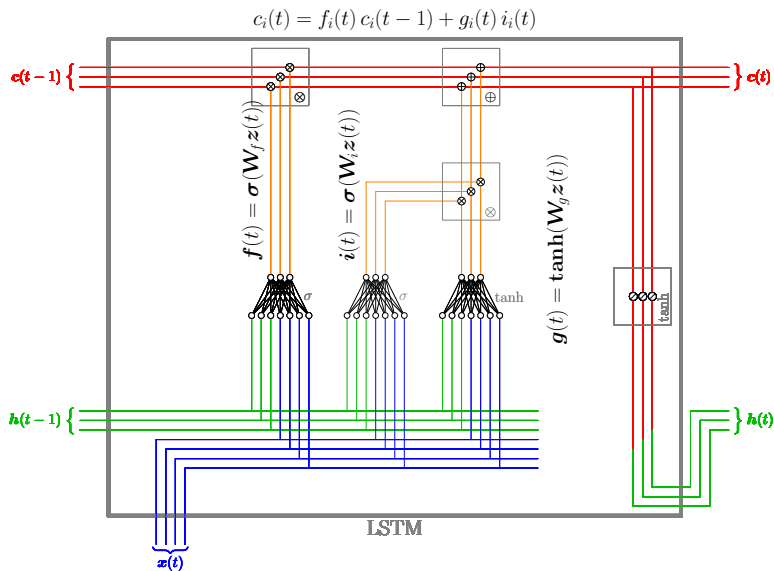
LSTM Architecture



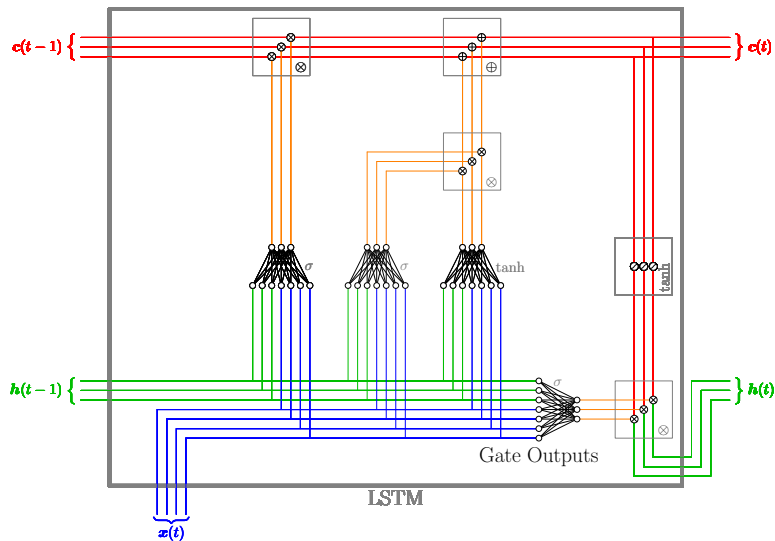
LSTM Architecture



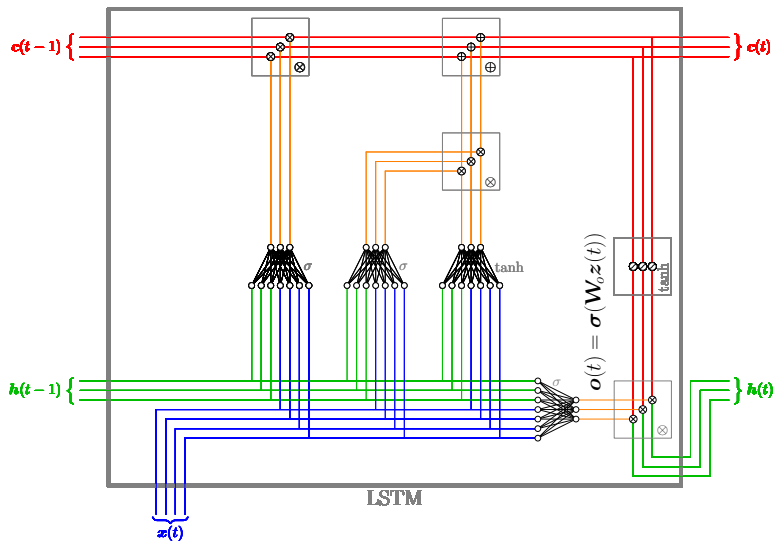
LSTM Architecture



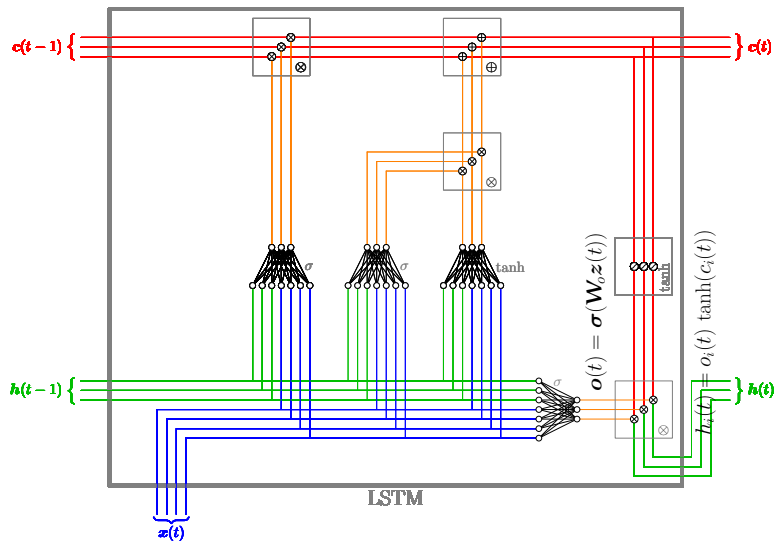
LSTM Architecture



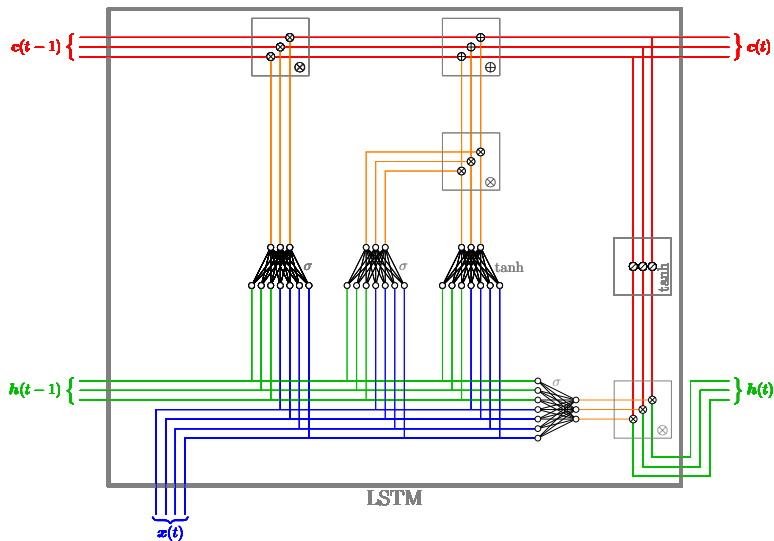
LSTM Architecture



LSTM Architecture



LSTM Architecture



Update Equations

- Inputs $\mathbf{z}(t) = (\mathbf{x}(t), \mathbf{h}(t-1))$

Update Equations

- Inputs $\mathbf{z}(t) = (\mathbf{x}(t), \mathbf{h}(t-1))$
- Network updates (\mathbf{W}_* are the free parameters)

$$\mathbf{f}(t) = \sigma(\mathbf{W}_f \mathbf{z}(t))$$

$$\mathbf{i}(t) = \sigma(\mathbf{W}_i \mathbf{z}(t))$$

$$\mathbf{g}(t) = \tanh(\mathbf{W}_g \mathbf{z}(t))$$

$$\mathbf{o}(t) = \sigma(\mathbf{W}_o \mathbf{z}(t))$$

Update Equations

- Inputs $\mathbf{z}(t) = (\mathbf{x}(t), \mathbf{h}(t-1))$
- Network updates (\mathbf{W}_* are the free parameters)

$$\mathbf{f}(t) = \sigma(\mathbf{W}_f \mathbf{z}(t))$$

$$\mathbf{i}(t) = \sigma(\mathbf{W}_i \mathbf{z}(t))$$

$$\mathbf{g}(t) = \tanh(\mathbf{W}_g \mathbf{z}(t))$$

$$\mathbf{o}(t) = \sigma(\mathbf{W}_o \mathbf{z}(t))$$

- Long-term memory update

$$\mathbf{c}(t) = \mathbf{f}(t) \otimes \mathbf{c}(t-1) + \mathbf{g}(t) \otimes \mathbf{i}(t)$$

Update Equations

- Inputs $\mathbf{z}(t) = (\mathbf{x}(t), \mathbf{h}(t-1))$
- Network updates (\mathbf{W}_* are the free parameters)

$$\begin{aligned}\mathbf{f}(t) &= \sigma(\mathbf{W}_f \mathbf{z}(t)) & \mathbf{i}(t) &= \sigma(\mathbf{W}_i \mathbf{z}(t)) \\ \mathbf{g}(t) &= \tanh(\mathbf{W}_g \mathbf{z}(t)) & \mathbf{o}(t) &= \sigma(\mathbf{W}_o \mathbf{z}(t))\end{aligned}$$

- Long-term memory update

$$\mathbf{c}(t) = \mathbf{f}(t) \otimes \mathbf{c}(t-1) + \mathbf{g}(t) \otimes \mathbf{i}(t)$$

- Output $\mathbf{h}(t) = \mathbf{o}(t) \otimes \tanh(\mathbf{c}(t))$

- We can train an LSTM by unwrapping it in time

Training LSTM

- We can train an LSTM by unwrapping it in time
- Note that it involves four dense layers with sigmoidal (or tanh) outputs

Training LSTM

- We can train an LSTM by unwrapping it in time
- Note that it involves four dense layers with sigmoidal (or tanh) outputs
- This means that typically it is very slow to train

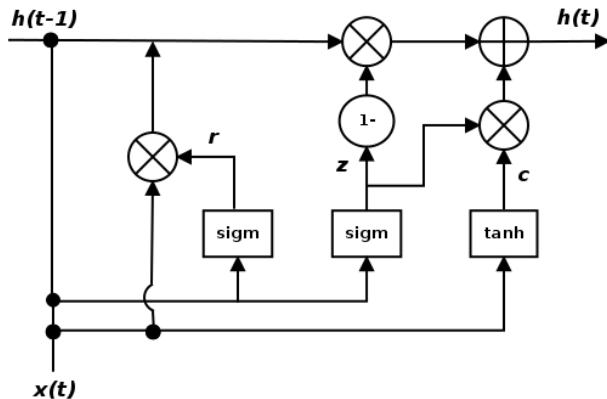
- We can train an LSTM by unwrapping it in time
- Note that it involves four dense layers with sigmoidal (or tanh) outputs
- This means that typically it is very slow to train
- There are a few variants of LSTMs, but all are very similar. The most popular is probably Gated Recurrent Unit (GRU)

LSTM Success Stories

- LSTMs have been used to win many competitions in speech and handwriting recognition
- Major technology companies including Google, Apple, and Microsoft are using LSTMs as fundamental components in new products.
- Google used LSTM for speech recognition on the smartphone, for Google Translate.
- Apple uses LSTM for the "Quicktype" function on the iPhone and for Siri.
- Amazon uses LSTM for Amazon Alexa.
- In 2017, Facebook performed some 4.5 billion automatic translations every day using long short-term memory networks ¹

¹https://en.wikipedia.org/wiki/Long_short-term_memory

Gated Recurrent Unit (GRU)



Gated Recurrent Unit (GRU)

- x_t : input vector
- h_t : output vector
- z_t : update gate vector
- r_t : reset gate vector
- W , U , and b : parameter matrices and vector
- sigm or σ_g is the sigmoid function
- tanh or σ_h is the hyperbolic tangent

Gated Recurrent Unit (GRU)

Initially, for $t = 0$, $h_0 = 0$

$$z_t = \sigma_g(W_z x_t + U_z h_{t-1} + b_z)$$

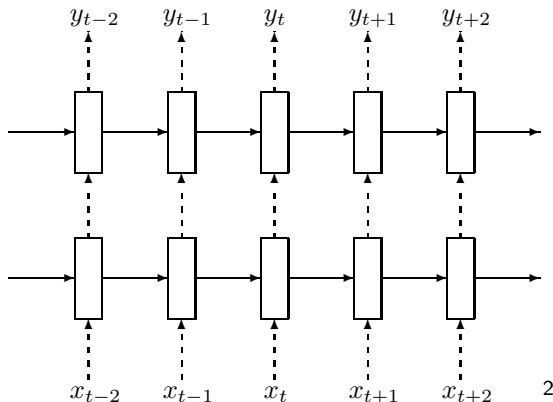
$$r_t = \sigma_g(W_r x_t + U_r h_{t-1} + b_r)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \sigma_h(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h)$$

GRU vs. LSTM?

- GRUs have two gates (reset and update) whereas LSTM has three gates (input/output/forget)
- GRU performance on par with LSTM but computationally more efficient (less complex).
- In general, if you have a very large dataset then LSTMs will likely perform better.
- GRUs are a good choice for smaller datasets.

Regularization in RNNs with LSTM units



2

²<https://arxiv.org/pdf/1409.2329.pdf>

Regularization in RNNs with LSTM units

$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow h_t^l, c_t^l$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} T_{2n,4n} \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} T_{2n,4n} \begin{pmatrix} \mathbf{D}(h_t^{l-1}) \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

3

³<https://arxiv.org/pdf/1409.2329.pdf>