# Thesis intermediate Presentation

Moritz Wolter

August 11, 2016

## Outline

Speech in text out.

$$B(f) = 1125 ln(1 + f/700) \tag{1}$$

$$H_m = 0 \qquad\qquad \text{if } k < f[m-1] \tag{2}$$

$$H_m = \frac{k - f[m-1]}{f[m] - f[m-1]} \qquad \text{if } f[m-1] \leq k \leq f[m] \tag{3}$$

$$H_m = \frac{f[m+1] - k}{f[m+1] - f[m]} \qquad \text{if } f[m] \leq k \leq f[m+1] \tag{4}$$

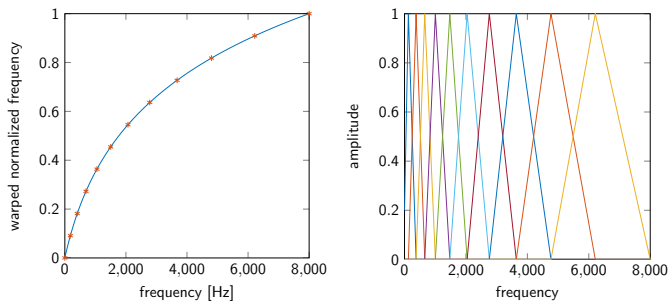$$H_m = 0 \qquad\qquad \text{if } k > f[m+1] \tag{5}$$

Figure 1 : The Mel scale and mel banks.
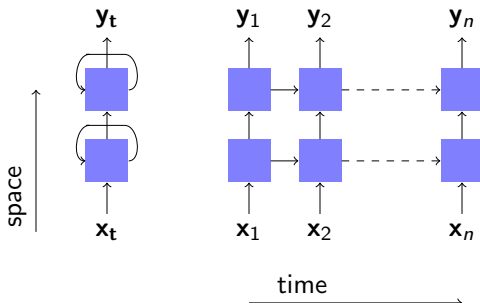
# Recurrent neural nets



Figure 2 : Unrolling a recurrent neural net.
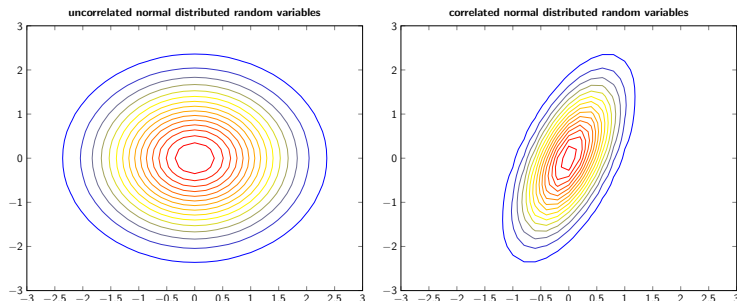
# Gradient updates on correlated data



Figure 3 : Normally distributed random variables with $\Sigma_1 = \mathbf{I}$ and $\Sigma_2 = (0.25\ 0.3;\ 0.3\ 1)$.

# Long Short Term Memory (LSTM)

$$\mathbf{i_t} = \sigma(\mathbf{W}_i[\mathbf{x}_t\ \mathbf{h_{t-1}}\ \mathbf{c_{t-1}}]^T + \mathbf{b}_i) \qquad (6)$$

$$\mathbf{f_t} = \sigma(\mathbf{W}_f[\mathbf{x}_t\ \mathbf{h_{t-1}}\ \mathbf{c_{t-1}}]^T + \mathbf{b}_f) \qquad (7)$$

$$\mathbf{c_t} = \mathbf{f_t}\mathbf{c_{t-1}} + \mathbf{i_t}\tanh(\mathbf{W}_c[\mathbf{x}_t\ \mathbf{h_{t-1}}]^T + \mathbf{b}_c) \qquad (8)$$

$$\mathbf{o_t} = \sigma(\mathbf{W}_o[\mathbf{x}_t\ \mathbf{h_{t-1}}\ \mathbf{c_t}]^T + \mathbf{b}_o) \qquad (9)$$

$$\mathbf{h_t} = \mathbf{o_t}\tanh(\mathbf{c_t}) \qquad (10)$$
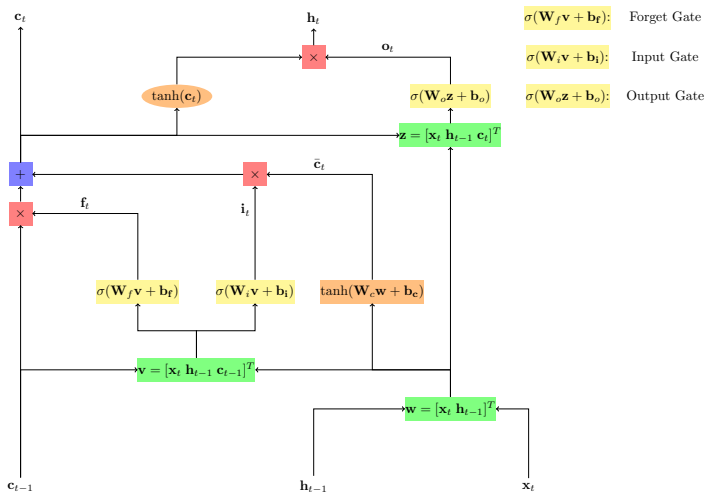
# Long Short Term Memory (LSTM)



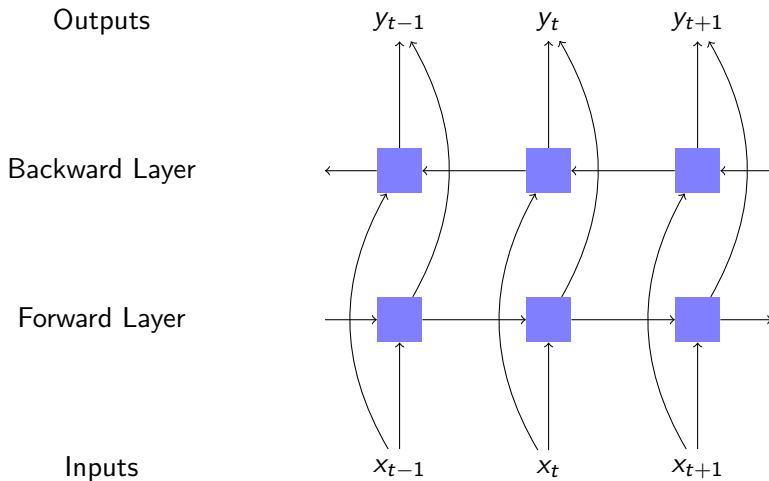Figure 4 : Visualization of the LSTM architecture

# Bidirectional BLSTM



Figure 5 : Bidirectional LSTM architecture

## Listener

- Pyramidal Bidirectional long short term memory (*pBLSTM*).
- Pyramid structure compresses the features.
- Three *pBLSTM*s on top of a *BLSTM* layer $\Rightarrow$ compression factor $2^3 = 8$.
- Pyramidal inputs concatenate the out put from previous layers:

$$\mathbf{h}_i^j = \text{pBLSTM}(\mathbf{h}_{i-1}^j, [\mathbf{h}_{2i}^{j-1}, \mathbf{h}_{2i+1}^{j-1}]) \tag{11}$$

- $i$ denotes the time step (from 0) and $j$ the layer.

## Attend and Spell

- attention based *LSTM* transducer.
- Find the most likely character given the features and previously found letters.

$$\mathbf{c}_i = \text{AttentionContext}(\mathbf{s}_i, \mathbf{H}) \tag{12}$$

$$\mathbf{s}_i = \text{RNN}(\mathbf{s}_i, \mathbf{H}) \tag{13}$$

$$P(\mathbf{y}_i | \mathbf{x}, y_{<i}) = \text{CharacterDistribution}(\mathbf{s}_i, \mathbf{c}_i) \tag{14}$$

## Attention Context

- Produce a context vector $\mathbf{c}_i$, with alignment information.

$$e_{i,u} = \phi(\mathbf{s}_i)^T \psi(\mathbf{h}_u) \tag{15}$$

$$\alpha_{i,u} = \frac{\exp(e_{i,u})}{\sum_u exp(e_{i,u})} \tag{16}$$

$$\mathbf{c}_i = \sum_u \alpha_{i,u} \mathbf{h}_u \tag{17}$$

- $\phi$ and $\psi$ are feed-forward MLP networks.
- $\mathbf{s}_i$ is the decoder state.
- The $\alpha$s work like a sliding window.
- $U$ denotes the total number of feature vectors.
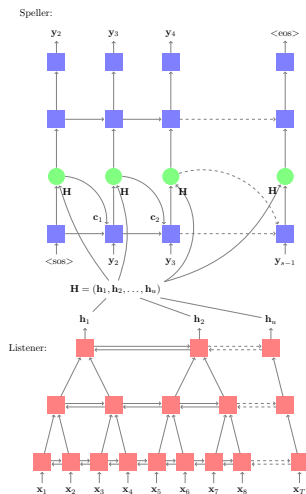
# The LAS-Architecture



Figure 6 : The LAS architecture

## Decoding and Rescoring

- Humans do not read character distributions.
- Left to right beam search turns distributions into text.
- Generate a tree using the $n$ most likely characters.
- Select from the tree according to:

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{LM}(\mathbf{y}) \qquad (18)$$

- The first summand is the total probability found from the tree.
- The second summand is a weighted language model contribution.

## What is tensorflow
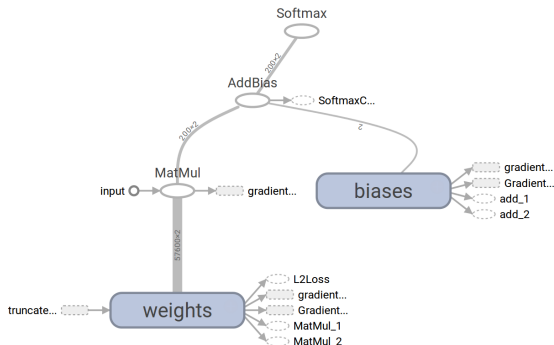
A machine learning toolbox.

# Tensorflow



Figure 7 : A simple linear node in tensorboard

# Summary and questions