

Building a state of the art speech recogniser

Moritz Wolter

Thesis submitted for the degree of
Master of Science in Mathematical
Engineering

Thesis supervisor:

Prof. dr. ir. Patrick Wambacq

Assessor:

Prof. dr. ir. Johan Suykens

Mentor:

Ir. Vincent Renkens

© Copyright KU Leuven

Without written permission of the thesis supervisor and the author it is forbidden to reproduce or adapt in any form or by any means any part of this publication. Requests for obtaining the right to reproduce or utilize parts of this publication should be addressed to the Departement Computerwetenschappen, Celestijnenlaan 200A bus 2402, B-3001 Heverlee, +32-16-327700 or by email info@cs.kuleuven.be.

A written permission of the thesis supervisor is also required to use the methods, products, schematics and programs described in this work for industrial or commercial use, and for submitting this publication in scientific contests.

Preface

I would like to thank everybody who kept me busy the last year, especially my promotor and my assistants. I would also like to thank the jury for reading the text. My sincere gratitude also goes to my family for supporting me trough my studies.

Moritz Wolter

Contents

Preface	i
Abstract	iii
List of Figures and Tables	iv
List of Abbreviations and Symbols	v
1 Literature Study	1
1.1 Preprocessing and feature extraction	1
1.2 Deep Neural Networks	2
1.3 Recurrent Neural Networks	4
1.4 Tensor-flow	8
1.5 Listen, Attend and Spell	8
Bibliography	11

Abstract

The `abstract` environment contains a more extensive overview of the work. But it should be limited to one page.

List of Figures and Tables

List of Figures

1.1	Frequency Bank input computed from a sentence contained in the <i>TIMIT</i> dataset. Time is shown on x and Frequency on the y-Axis.	1
1.2	The Mel-scale (blue) with Mel-Frequency Cepstrum Coefficients (red) on the left. Filterbank with Mel-spaced filters (right).	2
1.3	Example function network with partial derivatives.	3
1.4	Reverse sweep.	3
1.5	Rolled (left) and unrolled (right) recurrent neural net with two units. .	4
1.6	Visualization of the LSTM architecture following.	6
1.7	A bidirectional Long short term memory layer, according to [8]	7
1.8	The LAS architecture [4, page 3]. BLSTM blocks are shown in red. LSTM blocks in blue and attention nets in green.	10

List of Tables

List of Abbreviations and Symbols

Abbreviations

ConvNet	Convolutional neural network
MSE	Mean Square error
PSNR	Peak Signal-to-Noise ratio

Symbols

42	“The Answer to the Ultimate Question of Life, the Universe, and Everything” according to [?]
c	Speed of light
E	Energy
m	Mass
π	The number pi

Chapter 1

Literature Study

The first contains a general introduction to the work. The goals are defined and the modus operandi is explained. TODO: describe the problem.

1.1 Preprocessing and feature extraction

Filter-Bank features

Filter-banks are collections of filters. These filters are spread out over the whole frequency band [12, page 251]. Filter-bank output is commonly used as input for speech analysis [12][4]. The number of filter-banks depends on the required resolution, 32 is a common choice [13]. The energy within the part of the signal spectrum described by all individual filters is measured. Figure 1.1 shows the resulting energy measurements using 23 filters, for a sentence recording contained in the *TIMIT* data set. The general argument for filter banks is speech recognitions is that the cochlea, in the human ear, resembles a filter bank [12, page 30]. However humans do not perceive frequency linearly. Experimental evidence suggests, that our perceptions is scaled according to [12, page 34]:

$$B(f) = 1125 \ln(1 + f/700) \quad (1.1)$$

A normalized plot of this function is shown in figure 1.2 on the left. Mel-scaling suggests, that human are able to distinguish more lower frequencies than higher frequencies. In the plot the first four thousand Herz occupy roughly eighty percent of the scale. The band from four thousand to eight thousand Herz is left with only twenty percent of the scale, even tough half of the considered frequencies are in this band. Mel spaced filter-banks are an attempt to include the human perception in

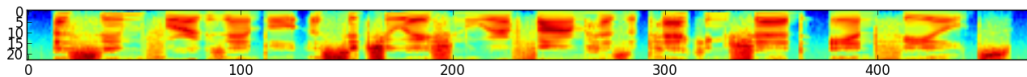


Figure 1.1: Frequency Bank input computed from a sentence contained in the *TIMIT* dataset. Time is shown on x and Frequency on the y-Axis.

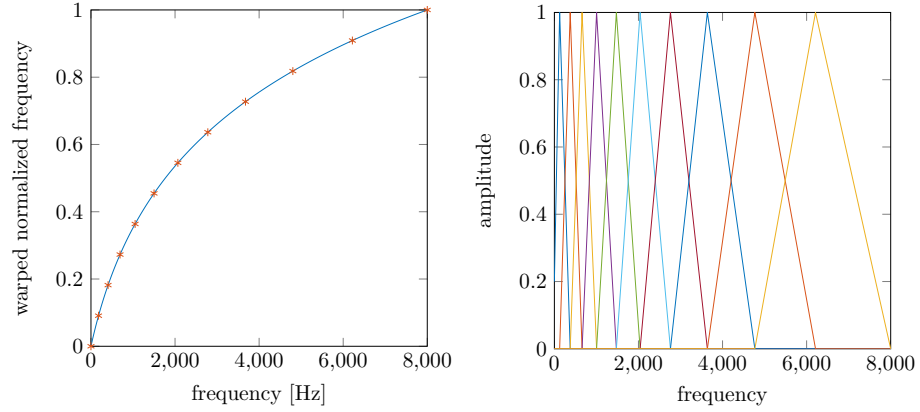


Figure 1.2: The Mel-scale (blue) with Mel-Frequency Cepstrum Coefficients (red) on the left. Filterbank with Mel-spaced filters (right).

speech recognition. The filter functions are defined by [12, page 317]:

$$H_m = 0 \quad \text{if } k < f[m-1] \quad (1.2)$$

$$H_m = \frac{k - f[m-1]}{f[m] - f[m-1]} \quad \text{if } f[m-1] \leq k \leq f[m] \quad (1.3)$$

$$H_m = \frac{f[m+1] - k}{f[m+1] - f[m]} \quad \text{if } f[m] \leq k \leq f[m+1] \quad (1.4)$$

$$H_m = 0 \quad \text{if } k > f[m+1] \quad (1.5)$$

In the equations above H_m denotes the magnitude of filter m with a total of M filters. The frequency is denoted by k , the vector f contains $M+2$ linearly spaced filter border values. These are the red stars on the right of figure 1.2. The left plot shows the triangular filter bank using, where the spacing is done according to the same values. Roughly speaking using mel-filter banks means using a high filter resolution where human hearing is good and a low resolution where it is bad.

Mel-Frequency banks are considered high level feature inputs. When the recognition system is found to work with these, features on a lower level or even raw data could be used as input. The idea behind doing fewer preprocessing, is that the network might be able to come up with something better.

1.1.1 Text and Phoneme output

1.2 Deep Neural Networks

1.2.1 Gradient descent

Stochastic gradient descent

When training networks on very large training sets, working with the full data set to compute the current gradient becomes very inefficient. As a remedy its is good

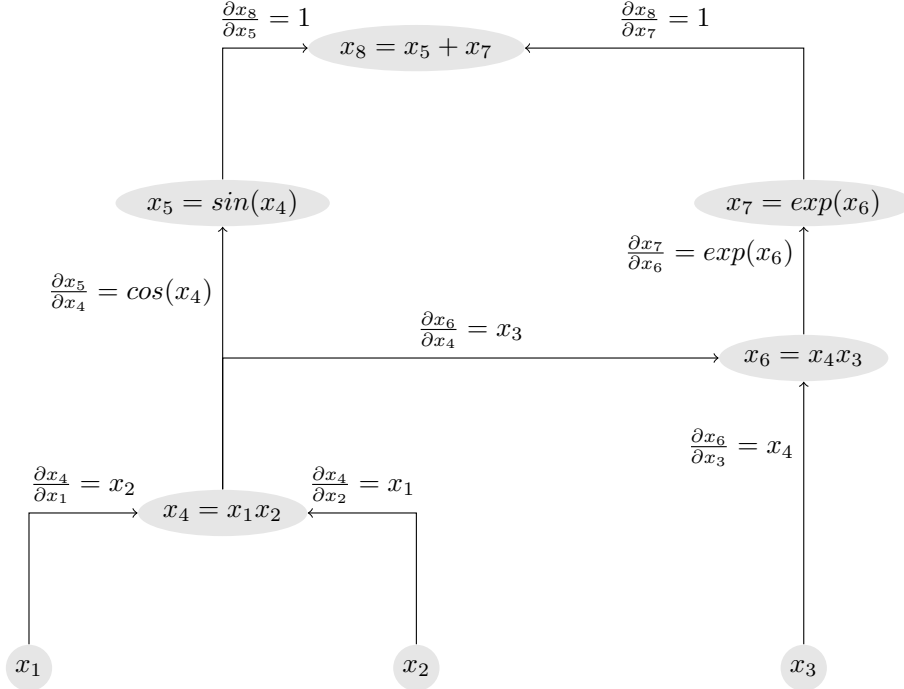


Figure 1.3: Example function network with partial derivatives.

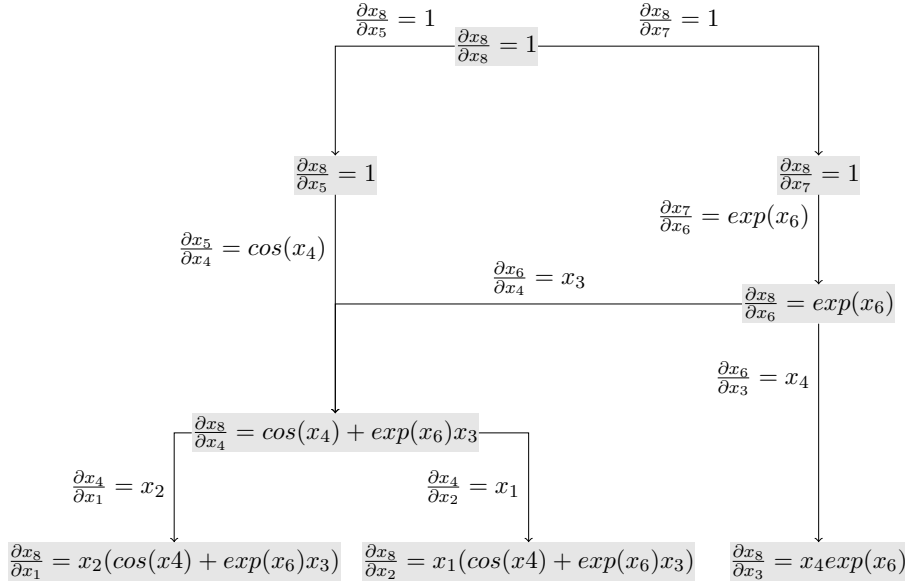


Figure 1.4: Reverse sweep.

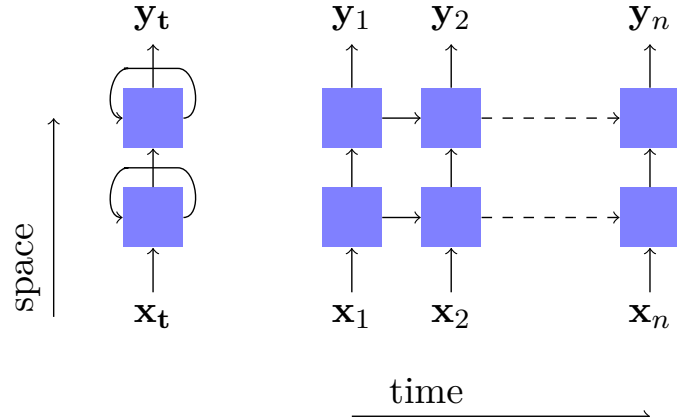


Figure 1.5: Rolled (left) and unrolled (right) recurrent neural net with two units.

practice in machine learning to work with so called mini-batches. A mini-batch includes a random subset of the training data set. This procedure is known as randomized gradient descent. With an added momentum term it can be formalized as [9, page 4]:

$$\Delta w_{ij}(t) = \alpha \Delta w_{ij}(t-1) - \epsilon \frac{\partial C}{\partial w_{ij}(t)} \quad (1.6)$$

C is the cost, which is computed by comparing the current network output to the desired output. $\alpha \in (0, 1)$ is a momentum coefficient. The weight of a connection from unit i in the layer under consideration to unit j in the layer below is given by the expression w_{ij} .

1.2.2 Convnets

When looking for a signal in different parts of a recording it is not always advisable to relearn recognition of the signal in different locations. With a conventional fully connected structure every sample on the time axis gets its own input weight. For shifted version of this signal the network will not be able to reuse weights it used to recognize the same signal at another time point in the recording. Convolutional neural nets aim to solve this problem [6, page 6], while preserving essential ordering information.

1.3 Recurrent Neural Networks

When processing speech it is important to take context into account. When spelling the letters, which make up a word, it is important to know what the previous letter was, in order to make the right decision. Feed-forward neural nets do not possess memory. These networks make decisions, starting from zero every time. In order to fix this a cell state variable can be introduced. This state is fed back into the cell together with new inputs every time step. Such a layout is shown in figure 1.5 on

the left. Another way to depict the same network is to not only consider the spacial dimension, but add the time axis as well. Figures, which show the spacial and time dimension are called unrolled network diagrams, shown in figure 1.5 on the right. When looking at the

1.3.1 The exploding and vanishing gradient problem

Even though past information is available in theory, learning long time dependencies is problematic with classical neural nets. The back-propagated derivative can sometimes become waker and weaker until it ultimately vanishes [10]. Another problem is that sometimes classical recurrent neural nets produce a gradient that blows up [14]. The exploding gradients can be fixed by clipping, but vanishing gradients require more sophisticated treatment [3].

1.3.2 Long short-term memory

Research seemed to be focused on solving the problem by making changes to the back-propagation algorithm. However a good solution to the problem turned out to be changing the network instead. Long short-term memory (LSTM) cells as proposed in [11] are more complex network units. These cells use the equation system [7, page 5]:

$$\mathbf{i}_t = \sigma(\mathbf{W}_{ix}\mathbf{x}_t + \mathbf{W}_{ih}\mathbf{h}_{t-1} + \mathbf{W}_{ic}\mathbf{c}_{t-1} + \mathbf{b}_i) \quad (1.7)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{fx}\mathbf{x}_t + \mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fc}\mathbf{c}_{t-1} + \mathbf{b}_f) \quad (1.8)$$

$$\mathbf{c}_t = \mathbf{f}_t\mathbf{c}_{t-1} + \mathbf{i}_t \tanh(\mathbf{W}_{cx}\mathbf{x}_t + \mathbf{W}_{ch}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (1.9)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{ox}\mathbf{x}_t + \mathbf{W}_{oh}\mathbf{h}_{t-1} + \mathbf{W}_{oc}\mathbf{c}_t + \mathbf{b}_o) \quad (1.10)$$

$$\mathbf{h}_t = \mathbf{o}_t \tanh(\mathbf{c}_t) \quad (1.11)$$

$$(1.12)$$

From the definition of the matrix product follows that

$$\mathbf{A}\mathbf{x}_1 + \mathbf{B}\mathbf{x}_2 = \begin{bmatrix} \mathbf{A} & \mathbf{B} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}. \quad (1.13)$$

Which this relation in mind the equations above can be rewritten, by creating column wise concatenated weight matrices for every neuron gate W_i , W_f , W_o , as well as for the state W_c . These matrices can then be multiplied by a row wise concatenated vector $[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}]^T$, which leads to the slightly simplified system of equations below:

$$\mathbf{i}_t = \sigma(\mathbf{W}_i[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}_{t-1}]^T + \mathbf{b}_i) \quad (1.14)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_f[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}_{t-1}]^T + \mathbf{b}_f) \quad (1.15)$$

$$\mathbf{c}_t = \mathbf{f}_t\mathbf{c}_{t-1} + \mathbf{i}_t \tanh(\mathbf{W}_c[\mathbf{x}_t \ \mathbf{h}_{t-1}]^T + \mathbf{b}_c) \quad (1.16)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}_t]^T + \mathbf{b}_o) \quad (1.17)$$

$$\mathbf{h}_t = \mathbf{o}_t \tanh(\mathbf{c}_t) \quad (1.18)$$

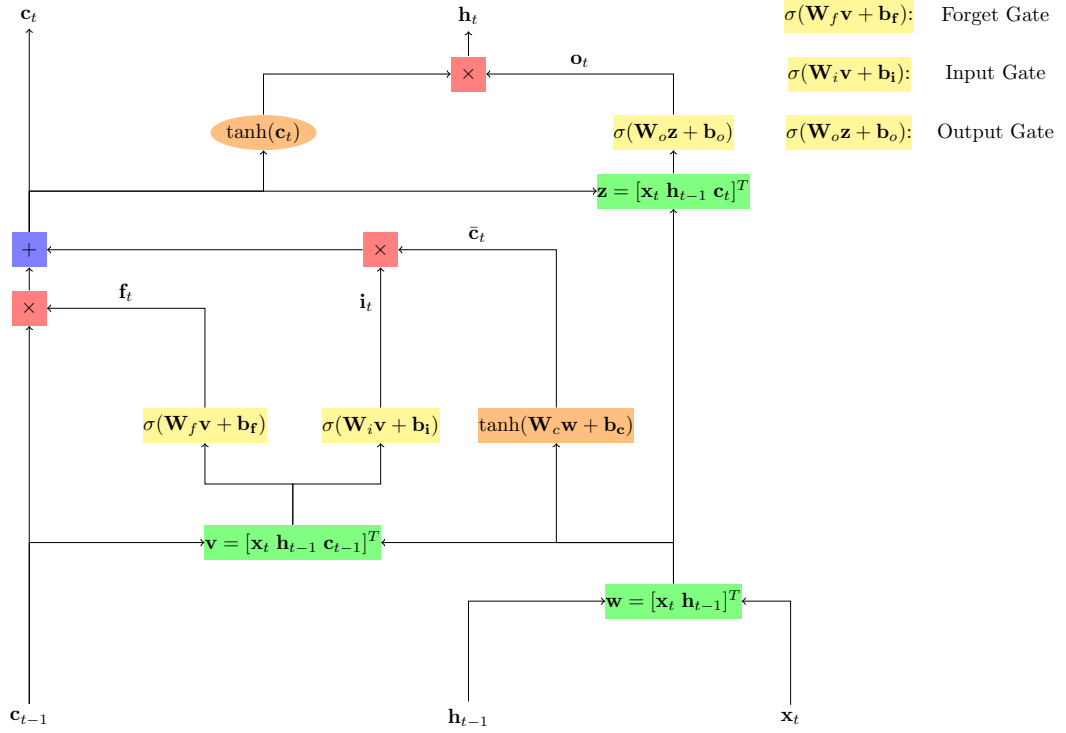


Figure 1.6: Visualization of the LSTM architecture following.

This system of equations is visualized in figure 1.6. The diagram is read from bottom to top. The most important part is the line from \mathbf{c}_{t-1} to \mathbf{c}_t [5]. It records operations on the cell state \mathbf{c}_t . The cell state contains information from the past which helps the block make decisions regarding the current output \mathbf{h}_t . The sigmoid functions $\sigma()$ are applied element wise on the input vectors and produce outputs between zero and one. In the case of the forget gate output \mathbf{f}_t these values $\in (0, 1)$ will serve as a measure of how much of the past state the cell would like to remember. One means keep this variable and zero throw it away [5]. The following task is to determine what should be added to the memory. This information can be found in the input gate result \mathbf{i}_t . \mathbf{i}_t is multiplied element wise with the candidate values $\bar{\mathbf{c}}_t$. These are computed by a hyperbolic tangent neuron. The $\tanh()$ function makes sure all vector elements are between -1 and 1 . The neuron looks at input data and the past outputs. Both are labeled \mathbf{w} in figure 1.6, \mathbf{w} contains all information that could possibly be included in the new state. Finally the weighted candidate values are added to what was previously stored. This operation leads to the updated memory state \mathbf{c}_t . Last but not least the new output value has to be computed, which will be a filtered version of the cell state. The decision of which and how much of each state variable will be send outside is made by output gate. It's output \mathbf{o}_t is multiplied with a rescaled version of the cell state. The rescaling is done using another hyperbolic tangent, which again sets all values between minus one and one. The product of this rescaled state and the weights found in \mathbf{o}_t then yields the new output \mathbf{h}_t .

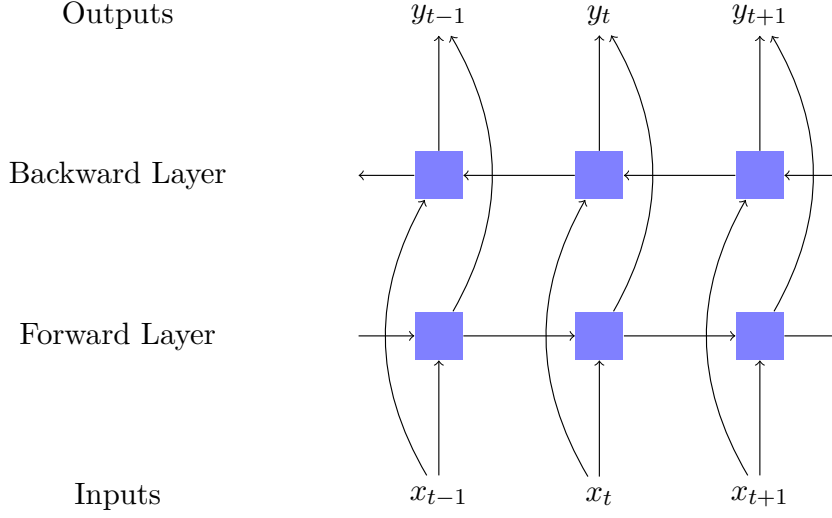


Figure 1.7: A bidirectional Long short term memory layer, according to [8]

1.3.3 Bidirectional Long Short Term Memory

With the advent of LSTMs deep recurrent networks became feasible in speech recognition [8]. RNNs are always deep in time, because their hidden state depends on past inputs. To enable abstraction their structure must also be deep in space. A bidirectional LSTM layer is shown in figure 1.7. It is important to note, that linear neurons are used to compute the LSTM input as well as the outputs according to the equations [8]:

$$\vec{\mathbf{h}}_t = \text{LSTM}(\mathbf{W}_{\vec{\mathbf{h}}_t} [\mathbf{x}_t \mathbf{h}_{t-1}]^T + \mathbf{b}_{\vec{\mathbf{h}}_t}) \quad (1.19)$$

$$\overleftarrow{\mathbf{h}}_t = \text{LSTM}(\mathbf{W}_{\overleftarrow{\mathbf{h}}_t} [\mathbf{x}_t \mathbf{h}_{t+1}]^T + \mathbf{b}_{\overleftarrow{\mathbf{h}}_t}) \quad (1.20)$$

$$\mathbf{y}_t = \mathbf{W}_y [\vec{\mathbf{h}}_t \overleftarrow{\mathbf{h}}_t]^T + \mathbf{b}_y \quad (1.21)$$

If stacked on top of each other, these bidirectional LSTM layers form a deep recurrent network. Defining $\mathbf{h}^0 = \mathbf{x}$, $\mathbf{h}^N = \mathbf{y}$ looking at time from $t = 1$ to T and taking N layers leads to:

$$\vec{\mathbf{h}}_t^n = \text{LSTM}(\mathbf{W}_{\vec{\mathbf{h}}_t}^n [\mathbf{h}_t^{n-1} \mathbf{h}_{t-1}^n]^T + \mathbf{b}_{\vec{\mathbf{h}}_t}^n) \quad (1.22)$$

$$\overleftarrow{\mathbf{h}}_t^n = \text{LSTM}(\mathbf{W}_{\overleftarrow{\mathbf{h}}_t}^n [\mathbf{h}_t^{n-1} \mathbf{h}_{t+1}^n]^T + \mathbf{b}_{\overleftarrow{\mathbf{h}}_t}^n) \quad (1.23)$$

$$\mathbf{h}_t^n = \mathbf{W}_y^n [\vec{\mathbf{h}}_t^n \overleftarrow{\mathbf{h}}_t^n]^T + \mathbf{b}_y^n \quad (1.24)$$

In this setting each LSTM cell has access to information from before and after it. For this to work the speech sequence, which is analyzed has to be recoded completely. In this case future information is available and should be used for recognition purposes.

1.4 Tensor-flow

In this section is devoted to the toolbox, which will be used to implement the Listen Attend and spell, architecture. According to the Tensor-flow authors [1]: “TensorFlow is an interface for expressing machine learning algorithms, and an implementation for executing such algorithms”. It was released by Google in 2015 and after installation can be used from within Python or C++.

1.5 Listen, Attend and Spell

The Listen Attend and Spell architecture (LAS) is the main Idea around which this thesis revolves. This entire section is based on [4]. The las-network consists of two mayor parts, the listener and the speller. The listener is a pyramidal recurrent neural net. It accepts filter bank spectra \mathbf{x}_n as inputs and produces high level output features \mathbf{h}_m . The speller in turn accepts the features as input and outputs distributions over Latin character sequences \mathbf{y}_p . An overview of the las-achrcitecture is given in figure 1.8.

1.5.1 The listener

The listener shown in figure 1.8 on the bottom, consists of Bidirectional Long Short Term Memory RNN (BLSTM) blocks. This choice implies that only fully recorded data can be analyzed. These blocks are arranged in a pyramidal structure, such that the time resolution is cut in half in every layer. This operation reduces the length U of the high level features \mathbf{H} . Without this compression the following attend and spell operation has a hard time extracting the relevant information. Additionally the compression reduces the problem complexity, which speeds up the training process significantly [4, page 4].

1.5.2 Attend and spell

The speller takes the features and produces a distribution over Latin character sequences as output. The computation of this output involves the context vector \mathbf{c}_i , the decoder state \mathbf{s}_i , the features \mathbf{H} and the previous output \mathbf{y}_i . The index i denotes time, $i - 1$ is used to refer to results from the last time step.

These values are computed using [4, page 4]:

$$\mathbf{s}_i = \text{RNN}(\mathbf{s}_{i-1}, \mathbf{y}_{i-1}, \mathbf{c}_{i-1}) \quad (1.25)$$

$$\mathbf{c}_i = \text{AttentionContext}(\mathbf{s}_i, \mathbf{H}) \quad (1.26)$$

$$P(\mathbf{y}_i | \mathbf{x}, \mathbf{y}_{<i}) = \text{CharacterDistribution}(\mathbf{s}_i, \mathbf{c}_i) \quad (1.27)$$

The state follows from a recurrent neural net (RNN) made of a two layer LSTM. The attention mechanism, called AttentionContext above, computes a new context vector once every time step. This computation starts with the determination of the scalar energy $e_{i,u}$, which will be used as weight for its corresponding feature vector h_u . The

computation starts with two feedforward neural networks or multilayer perceptrons (MLP), ϕ and ψ [4, page 5]:

$$e_{i,u} = \phi(\mathbf{s}_i)^T \psi(\mathbf{h}_u) \quad (1.28)$$

$$\alpha_{i,u} = \frac{\exp(e_{i,u})}{\sum_u \exp(e_{i,u})} \quad (1.29)$$

$$\mathbf{c}_i = \sum_u \alpha_{i,u} \mathbf{h}_u \quad (1.30)$$

α is produced by running \mathbf{e} through a softmax function, which scales \mathbf{e} such that all elements are within $(0, 1)$ and add up to one. These scaled weights, can then be used to form the context vector \mathbf{c}_i . When the training process converges the α_i s typically follow a distribution with sharp edges[4, page 5]. Thus it is justified to think of the alphas as a sliding window. This window contains only those parts of the condensed input data set, which are currently relevant.

1.5.3 Training

For end-to-end speech recognition the all networks must be trained jointly. The objective is to maximize the logarithmic probability:

$$\max_{\theta} \sum_i \log P(y_i | \mathbf{x}, y_{<i}; \theta). \quad (1.31)$$

Here y_i denotes the current output distribution, x the input, θ the various network parameters and finally $y_{<i}$ the ground truth, which is the known true desired output. Using the known output during training creates a situation, where the past outputs are always right. In practice however the situation will be different, as the network is going to make mistakes. As it is desired to create a robust model it is necessary to sometimes include the character distribution generated by the networks being trained. Which leads to the objective [4, page 5]:

$$\hat{y}_i = \text{CharacterDistribution}(s_i, \mathbf{c}_i) \quad (1.32)$$

$$\max_{\theta} \sum_i \log R(y_i | \mathbf{x}, \hat{y}_{<i}; \theta) \quad (1.33)$$

The novelty in comparison to the previous expression is that $\hat{y}_{<i}$ is sometimes taken from the past network outputs instead of the ground truth. An idea which Chan et al. found in [2].

Speller:

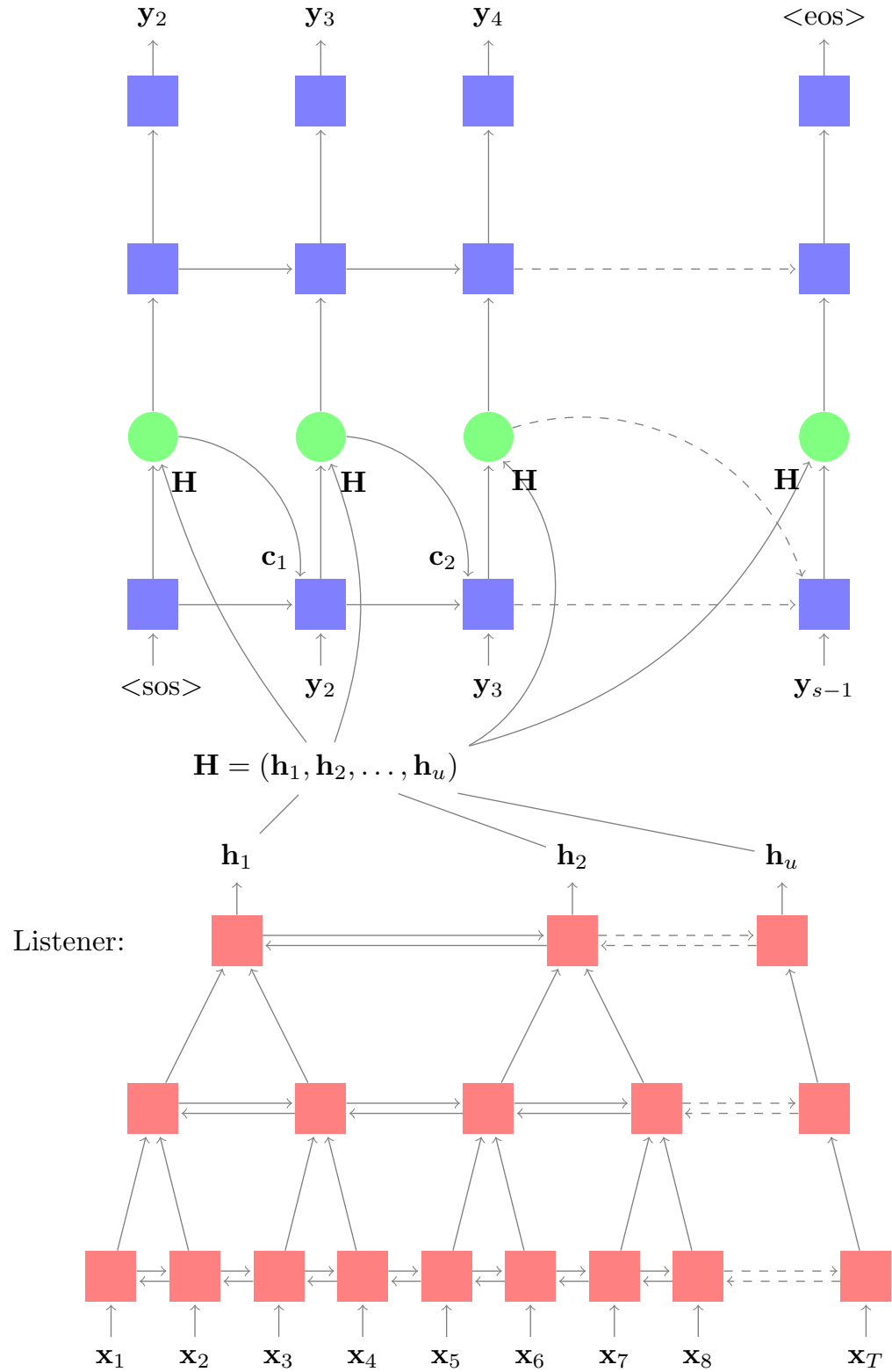


Figure 1.8: The LAS architecture [4, page 3]. BLSTM blocks are shown in red. LSTM blocks in blue and attention nets in green.

Bibliography

- [1] A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. 2015.
- [2] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks. *arXiv*, pages 1–9, 2015.
- [3] Y. Bengio, P. Frasconi, and P. Simard. The problem of learning long-term dependencies in recurrent networks. *IEEE International Conference on Neural Networks - Conference Proceedings*, 1993-Janua:1183–1188, 1993.
- [4] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals. Listen, attend and spell. *arXiv preprint*, pages 1–16, 2015.
- [5] Christopher Olah. Understanding LSTM Networks, 2015.
- [6] V. Dumoulin and F. Visin. A guide to convolution arithmetic for deep learning. pages 1–28, 2016.
- [7] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, pages 1–43, 2013.
- [8] A. Graves, A.-R. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. *2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, (6):6645–6649, 2013.
- [9] G. Hinton, L. Deng, D. Yu, G. Dahl, A.-r. Mohamed, N. Jaitly, V. Vanhoucke, P. Nguyen, T. Sainath, and B. Kingsbury. Deep Neural Networks for Acoustic Modeling in Speech Recognition. *IEEE Signal Processing Magazine*, 29(6):82–97, 2012.
- [10] S. Hochreiter. The Vanishing Gradient Problem During Learning Recurrent Neural Nets and Problem Solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 06(02):107–116, 1998.

- [11] S. Hochreiter and J. Schmidhuber. LONG SHORT TERM MEMORY. *Technical Report FKI-207-95*, pages 1–8, 1995.
- [12] X. Huang, A. Acero, and H.-W. Hon. *Spoken Language Processing: A Guide to Theory, Algorithm and System Development*. 2001.
- [13] B. H. Juang, L. R. Rabiner, and J. G. Wilpon. On the use of bandpass filtering in speech recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, ASSP-35(7), 1987.
- [14] R. Pascanu, T. Mikolov, and Y. Bengio. Understanding the exploding gradient problem. *Proceedings of The 30th International Conference on Machine Learning*, (2):1310–1318, 2012.

Master thesis filing card

Student: Moritz Wolter

Title: Building a state of the art speech recogniser

UDC: 621.3

Abstract:

In the past machine learning relied heavily on algorithms designed by experts to solve a specific task. Which lead to highly sophisticated algorithms, which could be grasped only by small groups of people. The human brain however does not work this way, although specialized areas exist, these areas consist of similar building blocks. Artificial neural networks attempt to mimic this layout. Similar algorithmic structures are used for a wide variety of tasks. This thesis deals with the application of neural networks in speech recognition. Replacing the various subsystems by one integrated network based approach.

Thesis submitted for the degree of Master of Science in Mathematical Engineering

Thesis supervisor: Prof. dr. ir. Patrick Wambacq

Assessor: Prof. dr. ir. Johan Suykens

Mentor: Ir. Vincent Renkens