

Building an end-to-end speech recognizer

Moritz Wolter

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

Thesis supervisors:
Prof. dr. ir. Patrick Wambacq
Prof. dr. ir. Hugo van Hamme

Assessor:
Prof. dr. ir. Johan Suykens

Mentor:
Ir. Vincent Renkens

© Copyright KU Leuven

Without written permission of the thesis supervisors 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 supervisors 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 through my studies.

Moritz Wolter

Contents

Preface	i
Abstract	iii
List of Figures and Tables	iv
List of Abbreviations and Symbols	vi
1 Problem statement	1
2 Literature Study	3
2.1 Preprocessing and feature extraction	3
2.2 Neural Networks	5
2.3 Connectionist temporal classification	18
2.4 Listen, Attend and Spell [7]	19
2.5 Tensor-flow	24
3 Methodology	25
3.1 Pylint	25
4 Experiments	27
4.1 TIMIT	27
4.2 BLSTM-CTC	27
4.3 Listen attend and spell experiments	30
Bibliography	37

Abstract

Speech recognition is concerned with transcribing what is said in a recoding of spoken language. In machine learning terms this process is called sequence labeling. A recoding consists of a chain of frames, this chain can be split up into several sequences, these make up words or phonemes, which must be labeled. The sequence of labels forms the transcription.

The meaning of speech depends on context, therefore a good system needs to take it into account. Classical feed-forward networks fail to do that, which is why this system will mainly consist of recurrently connected Long Short Term Memory (LSTM) blocks. Inspired by the recurrent connects of neurons in the human brain, LSTM-RNNs have the ability to store information over long time periods, an important requirement in take context into account.

In order to train machine learning systems, speech and transcription text pairs are used. The text contains the exact information of what is said in the recording, but where in the recoding which word or sound is said is unknown. In other words text to speech alignment is missing. Aligning the data will be an important issue in this thesis.

List of Figures and Tables

List of Figures

2.1	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.	3
2.2	The Mel-scale (blue) with Mel-Frequency Cepstrum Coefficients (red) on the left. Filterbank with Mel-spaced filters (right).	4
2.3	Plot of $y = -\ln(x)$ for $x \in (0, 1)$	7
2.4	Example function network with partial derivatives.	9
2.5	Reverse sweep.	9
2.6	“Dropout Neural Net Model. Left: A standard neural net with 2 hidden layers. Right: An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.” [24], image and caption taken from the same source.	12
2.7	Visualization of a single recurrent cell.	14
2.8	Rolled (left) and unrolled (right) recurrent neural net with two units.	14
2.9	Visualization of the LSTM architecture.	16
2.10	A bidirectional Long short term memory layer, according to [15]	18
2.11	The LAS architecture [7, page 3]. BLSTM blocks are shown in red. LSTM blocks in blue and attention nets in green.	22
4.1	48 to 39 phoneme folding as shown in [21].	28
4.2	Raw time domain signal and feature vectors. Features include 40 mel banks as well as the first and second derivatives. Raw data was taken from TIMIT utterance faem0-si2022 “What outfit does she drive for?”.	28
4.3	Validation and Test set error while training a two layer BLSTM network with CTC output layer on the TIMIT speech corpus with 39 folded phonemes.	29
4.4	The training progress shown for the Listener with added CTC layer. The loss values shown in blue and green have been scaled with $\frac{1}{100}$. On the right a closeup on the last two training processes is shown.	31
4.5	The training progress shown for the full las architecture with greedy decoding. The loss values shown in blue and green have been scaled with $\frac{1}{100}$. On the right a closeup on the last two training processes is shown.	32

LIST OF FIGURES AND TABLES

4.6	Repetitions of the same experiment with network output reuse probabilities 0.2, 0.4, 0.6, 0.8.	33
4.7	Training full las over 40 epochs with a network output reuse probability of 0.5	34
4.8	Two different attend and spell cell configurations.	35

List of Tables

2.1	An example of the edit distance matrix for the dutch word huis and its english translation house. The edit distance can be found in the cell in the lower right of the table.	23
-----	---	----

List of Abbreviations and Symbols

Abbreviations

Con-	Convolutional neural network
vNet	
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

Problem statement

Chapter 2

Literature Study

2.1 Preprocessing and feature extraction

Filter-Bank features

Filter-banks are collections of filters. These filters can be spread out over audible frequencies¹. Filter-bank output is commonly used as input for speech analysis [18][7]. The number of filter-banks depends on the required resolution, 32 is a common choice [19]. The energy within the part of the signal spectrum described by all individual filters is measured. Figure 2.1 shows the resulting energy measurements using 23 filters, for a sentence recording contained in the *TIMIT* data set. The general argument for filter banks in speech recognition is that the cochlea, in the human ear, resembles a filter bank [18, page 30]. Humans do not perceive frequency linearly. Experimental evidence suggests, that our perception of is scaled according to the Mel-Scale [18, page 34]:

$$B(f) = 1125 \ln(1 + f/700) \quad (2.1)$$

A normalized plot of this function is shown in figure 2.2 on the left. According to the Mel-scale, humans 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 about 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

¹ Approximately 16 to 16000 Hz.

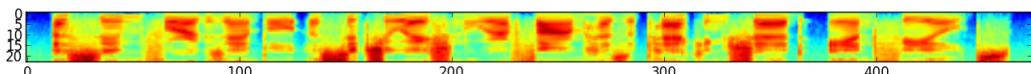


Figure 2.1: Frequency Bank input computed from a sentence contained in the *TIMIT* dataset. Time is shown on x and Frequency on the y-Axis.

2. LITERATURE STUDY

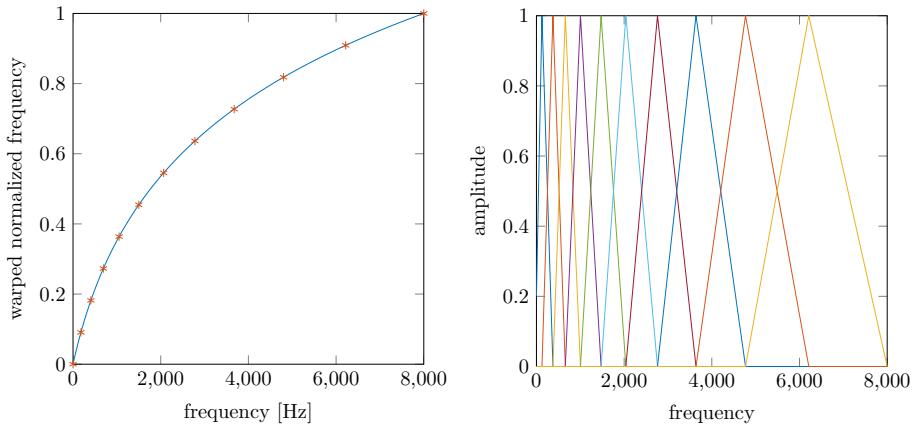


Figure 2.2: The Mel-scale (blue) with Mel-Frequency Cepstrum Coefficients (red) on the left. Filterbank with Mel-spaced filters (right).

in speech recognition. The filter functions are defined by [18, page 317]:

$$H_m = 0 \quad \text{if } k < f[m - 1] \quad (2.2)$$

$$H_m = \frac{k - f[m - 1]}{f[m] - f[m - 1]} \quad \text{if } f[m - 1] \leq k \leq f[m] \quad (2.3)$$

$$H_m = \frac{f[m + 1] - k}{f[m + 1] - f[m]} \quad \text{if } f[m] \leq k \leq f[m + 1] \quad (2.4)$$

$$H_m = 0 \quad \text{if } k > f[m + 1] \quad (2.5)$$

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 left of figure 2.2. The right plot shows the triangular filter banks. These banks are spaced 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.

In addition to perceptually motivated instantaneous features, derivative information can be appended to each feature vector. These so called delta features add information regarding the speed of change taking into account $2 * N$ feature vectors. N denotes the number of frames back and forward in time. The derivative information can be computed using [27]. The mel feature vectors are denoted by \mathbf{c} :

$$\mathbf{d}_t = \frac{\sum_{n=1}^N n(\mathbf{c}_{t+n} - \mathbf{c}_{t-n})}{2 \sum_{n=1}^N n^2} \quad (2.6)$$

The formula above is a central difference, because the value at time t is not taken into account, in contrast to forward or backward differences, where the value at t is part of the difference. A common choice for N is two. In this case the formula above

simplifies to:

$$\mathbf{d}_t = \frac{2\mathbf{c}_{t+2} + \mathbf{c}_{t+1} - \mathbf{c}_{t-1} - 2\mathbf{c}_{t-2}}{10} \quad (2.7)$$

In this case two future frames are used to compute the derivative at time t , which implies, that at computation time a full recording is available. Should that not be the case a backward difference scheme should be considered instead. To further augment the features second derivative or acceleration information can be added as well. These can be computed by simply applying the central difference to the deltas one more time. It is important to note that augmenting the features with derivative information significantly increases their dimension. If 40 mel-filters are used then the first and second derivatives will contain 40 entries each, therefore increasing the input dimension from 40 to 120.

Before using the features it is best practice to standardize the components of the input vectors. It is desirable to use features with an overall mean of zero and a standard deviation of 1 over the entire training set [11, page 30]. The feature normalization process starts by computing the mean,

$$\mathbf{m}_i = \frac{1}{\|S\|} \sum_{\mathbf{x} \in S} \mathbf{x}_i \quad (2.8)$$

and standard deviation

$$\sigma_i = \sqrt{\frac{1}{|S|-1} \sum_{x \in S} (\mathbf{m}_i - \mathbf{x}_i)^2} \quad (2.9)$$

Of a data set S consisting of input \mathbf{x} and target \mathbf{t} pairs. Standardized input vectors can then be computed using:

$$\hat{\mathbf{x}}_i = \frac{\mathbf{x}_i - \mathbf{m}_i}{\sigma_i} \text{ for all } \mathbf{x}_i \in S. \quad (2.10)$$

Input standardization is beneficial to network performance, because it puts the input values in ranges, which standard activation functions such as sigmoids or hyperbolic tangents can handle better [11, page 30].

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 less preprocessing, is that the network might be able to come up with something better on its own.

2.2 Neural Networks

Artificial neural networks were originally designed to model biological neurons. Even though it is now known that artificial neural networks have little in common with their counterparts in biology, the models are still popular due to their useful pattern recognition properties [11, page 13]. All neural networks consist of elementary functions with weighted and biased inputs, the operations necessary to evaluate them

2. LITERATURE STUDY

form a computational graph. In a nutshell the idea is to choose the weights and biases such that the network output corresponds to patterns seen in the input data. The process of finding good weights is called training. Training in turn is done using gradient descent.

2.2.1 Gradient descent

The optimization process of neural networks works with a training data set $\{\{\mathbf{x}_1, \mathbf{t}_1\}, \dots, \{\mathbf{x}_p, \mathbf{t}_p\}\}$ [23, page 156]. The elements of this set are the input- and output-patters \mathbf{x} and \mathbf{t} respectively. For each input vector the network output \mathbf{o} is computed. Ideally the network output should be the same as the desired target \mathbf{t} one for all data pairs. The difference between target and current outputs could be measured by the cost function [23, page 156]:

$$E = \frac{1}{2} \sum_{i=1}^p \|\mathbf{o}_i - \mathbf{t}_i\|^2. \quad (2.11)$$

However generally the more complex cross-entropy is used as a cost function. For classification tasks, cross entropy is defined as [6, page 245]:

$$E = - \sum_{i=1}^p \mathbf{t}_i \ln(\mathbf{o}_i) + (1 - \mathbf{t}_i) \ln(1 - \mathbf{o}_i) \quad (2.12)$$

In order to understand why the cross entropy is a measure of output error, the two scalar cases $t = 0$ and $t = 1$ are considered. For $t = 0$ the cross entropy simplifies to $-\ln(1 - o)$. Looking at figure 2.3, this expression turns out to be a good cost function, because for values close to the desired output $o = 0$, $-\ln(1 - o)$ will be close to zero. If the output moves away from the desired value zero, the cost grows asymptotically. If $t = 1$ is considered the cost function simplifies to $-\ln(o)$ which will be zero for the desired output and once more displays asymptotic growth for undesired values far away from one. If the target probability is somewhere between these two extreme cases the cost will be the sum of the two cases considered earlier. Cross entropy is statistically motivated, the output vectors \mathbf{o} describe probability distributions over the possible labels, their elements must therefore be $\in (0, 1)$ and sum up to one.

During the training process a local minimum of the error function E is sought. At this minimum the difference between the network output \mathbf{o} and the target values \mathbf{t} is small. This value might not be the smallest possible value, termination can happen at the global or a local minimum. After the training process has completed the network is expected to identify similarities to data seen during the training process and produce a similar output. In order to reach a local minimum, the gradient of the error function is needed. The key idea of gradient descent is to follow the negative gradient until a local minimum is reached. Neural networks can be considered as large composite functions, which are made up of elementary operations. The evaluation of the network can be written as a graph. Computations are done at each node and information travels through the network along directed edges from

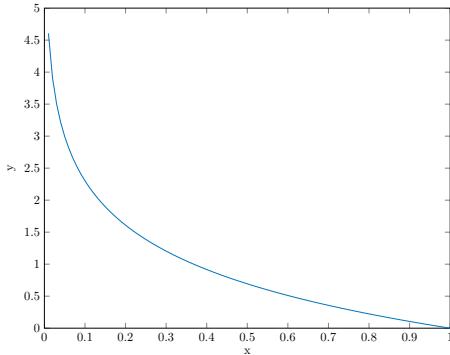


Figure 2.3: Plot of $y = -\ln(x)$ for $x \in (0, 1)$.

node to node. In order to create a computational graph for the training process each of the output units of the network under consideration are connected to a new node which computes $\frac{1}{2}(o_{ij} - t_{ij})^2$ [23, page 157] or the cross-entropy term of one data target pair. These new nodes in turn are connected to one more node, which sums up all error values and produces E_i . The process described above must be repeated for all training data pairs. One final node is added, which sums up all values E_i . Its output gives the value for the error function E which is now in the form of a large graph.

Reverse mode algorithmic differentiation or back-propagation is an algorithm to compute the gradient of a graph consisting of basic elementary operations. As an example its operation is now illustrated using the function [9, page 69]:

$$f(w_1, w_2, w_3) = \sin(w_1 w_2) + \exp(w_1 w_2 w_3) \quad (2.13)$$

This function written in terms of five elementary operations as [9, page 70]:

$$w_4 = w_1 w_2 \quad (2.14)$$

$$w_5 = \sin(w_4) \quad (2.15)$$

$$w_6 = w_4 w_3 \quad (2.16)$$

$$w_7 = \exp(w_6) \quad (2.17)$$

$$w_8 = w_5 + w_7 \quad (2.18)$$

$$o = w_8 \quad (2.19)$$

Computing the gradient means computing the partial derivatives of f with respect to all inputs. In this case this means finding:

$$\frac{\partial f(w_1, w_2, w_3)}{\partial w_1} = w_2(\cos(w_1 w_2) + w_3 \exp(w_1 w_2 w_3)) \quad (2.20)$$

$$\frac{\partial f(w_1, w_2, w_3)}{\partial w_2} = w_1(\cos(w_1 w_2) + w_3 \exp(w_1 w_2 w_3)) \quad (2.21)$$

$$\frac{\partial f(w_1, w_2, w_3)}{\partial w_3} = w_1 w_2 \exp(w_1 w_2 w_3) \quad (2.22)$$

2. LITERATURE STUDY

Above the derivatives have been found by hand using the chain rule. Now these will be computed using back-propagation. Figure 2.4 shows a graphical representation of equation 2.13. The partial derivatives needed for the backwards sweep can be found on the edges.

The gradient is computed using a forward and backward sweep. During the forward sweep the inputs are fed into the network and the functions at each node are evaluated layer by layer, until the output at the last node is known. In figure 2.4 this means computing w_4 to w_8 .

After the forward sweep the gradient is found by going back through the network from the output to each input node. Using a seed value of 1 at the output node the lower unit values are computed by multiplying the associated partial derivative found on each edge. If a node has more than one incoming value, their sum is computed. The process is illustrated in figure 2.5. At the roots of the tree the partial derivatives of the output with respect to each input can be found. Together these root values make up the gradient. To be able to perform the first forward sweep the network weights are initialized at random. The training data pairs are known and can be added as constants to the graph. If the weights of the network are stored in a weight vector \mathbf{w} the value of the gradient after the n th update may be written as [11, page 27]:

$$\Delta \mathbf{w}(n) = -\alpha \frac{\partial f}{\partial \mathbf{w}(n)}. \quad (2.23)$$

Where $\alpha \in [0, 1]$ denotes the learning rate. Unfortunately simple gradient descent tends to get stuck in local optima. In order to increase the chance of gradient descent to escape from such a local minimum a momentum term can be added to the formulation [6, page 267][11, page 27]:

$$\Delta \mathbf{w}(n) = m \Delta \mathbf{w}(n-1) - \alpha \frac{\partial f}{\partial \mathbf{w}(n)}. \quad (2.24)$$

Above $m \in [0, 1]$ denotes the momentum. In order to understand the effect of the momentum term consider a weight space region of very low curvature, where it can be assumed, that the gradient stays constant. Using the momentum formula above yields [6, page 267]:

$$\Delta \mathbf{w}(n) = -\alpha \frac{\partial f}{\partial \mathbf{w}} (1 + m + m^2 + \dots) \quad (2.25)$$

$$= -\frac{\alpha}{1-m} \frac{\partial f}{\partial \mathbf{w}}. \quad (2.26)$$

Using the fact that the sum of the geometric series is known and $|m| < 1$ must hold. The momentum term therefore changes the effective learning rate from α to $\frac{\alpha}{1-m}$ in regions of low curvature. Elsewhere the gradient is oscillatory, which leads to cancellation of successive contribution of the momentum term [6, page 267].

2.2.2 Regularization

The optimization process works on the training set, but later the trained network should be able to recognize patterns seen in the training data set in other inputs

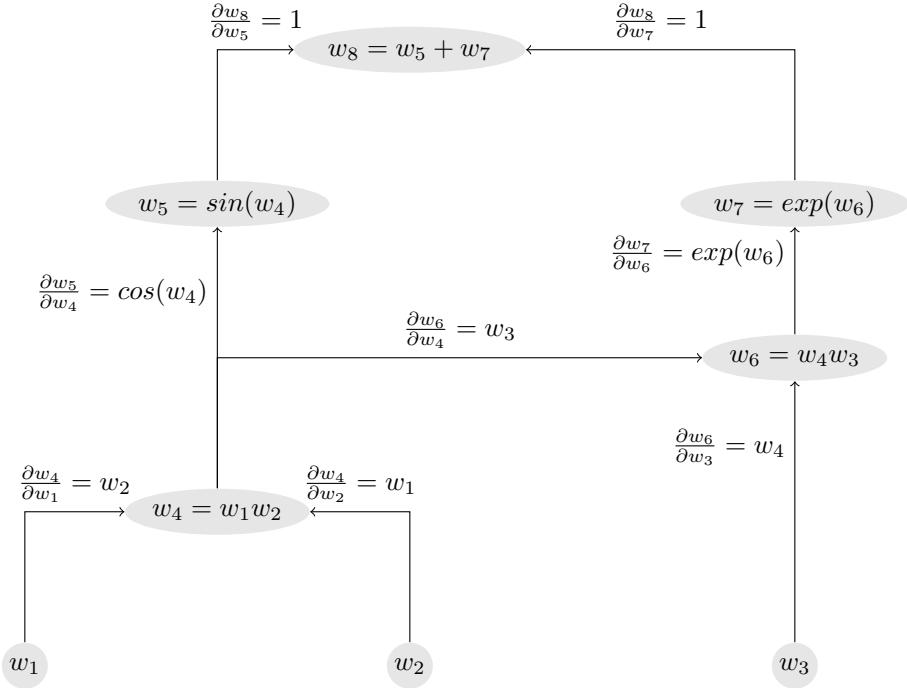


Figure 2.4: Example function network with partial derivatives.

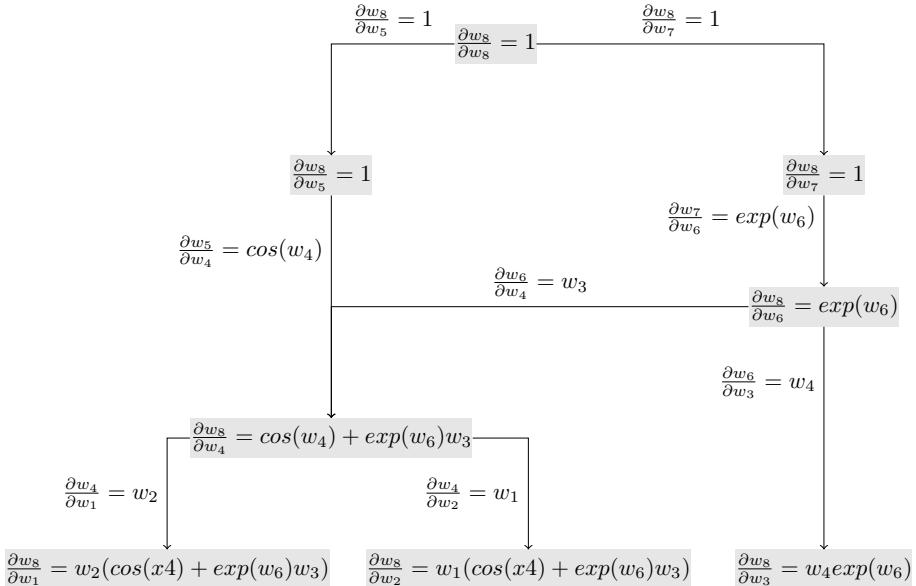


Figure 2.5: Reverse sweep.

2. LITERATURE STUDY

and produce similar outputs. In order to achieve this goal it is important not to memorize training data samples. Another problem with complex networks is that they are able to track noise in the training data. Regularization techniques aim to counteract sample memorization and noise tracking in order to ensure that progress on the training data carries over to new data sets or in other words regularization aims to improve generalization.

Early stopping

The idea of early stopping is to observe the training and validation loss during training. As long as both are falling simultaneously, we assume that the network will also generalize well to other data. This is a reasonable assumption, because the validation data set plays no part in computing the weight updates. When the training set loss continues fall, but the validation loss rises, the network is overfitting [12, page 31]. The training process should be stopped early in these cases. A test set should be used to guard against parameter choices that work well only on the training and validation set. Ideally the performance on the test set should only be evaluated, once the training process has finished.

Care must be taken, when dividing the entire data set into training, validation and test subsets. All three sets must be good representations of the overall data set. Differences for example regarding the noise level should be avoided. Furthermore the sizes of all sets must be chosen carefully. Using too few samples for validation and testing, makes it very hard to guard against over-fitting. If the training set becomes too small, overall performance might suffer. In practice a good trade off must be found experimentally.

When data collection is expensive cross validation should be used. Cross validation divides the merged training and validation data sets into k subsets. During each iteration one of the k subsets is used for validation and the remaining $k - 1$ sets are used for training. Cross-validation has the advantage that each subset is used for training and validation, thus making more efficient use of the total data set. The test set is very important when working with cross validation, because the validation data is part of the training process. It should be sufficiently large to detect over-fitting.

Input Noise

Training with input noise means that fixed variance, zero mean Gaussian noise is added to the inputs during training [12, page 32]. The targets remain unchanged. The idea is to reuse the training data set more efficiently. When the artificial noise is similar to the noise already present in the data set, reusing the noisy samples is better than reusing the original samples. Input noise decreases the reliability of the inputs during training, it covers small variations caused by noise in the data, which the network could remember without it during training, if the noise was not present. Furthermore the artificial noise will help the network to deal with similar noise in the future, making the trained system more robust. A problem with input noise is

that it is difficult to determine how large the variance should be, one option is to manually tune the parameter using the validation performance as a metric.

Weight noise

When using weight noise zero mean fixed variance Gaussian noise is added to the weights during training. The idea is to reduce the precision, with which the network weights can be described [12, page 32]. The reduced precision makes it impossible for the network to model small variations in the training samples, which are likely to be noise. Instead it will only be able to model larger changes and generalize better. When using input noise the following algorithm should be used [12, page 33]:

```
while stopping condition is False:  
    Shuffle training data pairs  
    for each training sample pair:  
        Add Gaussian noise to weights  
        compute the gradient  
        restore original weights  
        update the original weights
```

However weight noise should be used with care, because it can lead to very slow convergence.

Regularized loss minimization

Regularized loss minimization includes a model complexity penalty term in the cost function. Instead of minimizing the loss E an augmented loss [6, p. 338][3, p. 171]:

$$\hat{E} = E + v\Omega \quad (2.27)$$

is used. Ω stands for a penalty term, which should grow, when the network is likely to over-fit the training data. $v \in (0, 1)$ denotes the weight parameter, which determines how much emphasis the optimization algorithm should place on regularization. A straightforward choice for Ω is the l_2 norm of the model weights:

$$\Omega_{l2} = \|\mathbf{w}\|^2 = (\sqrt{\sum w^2})^2 \quad (2.28)$$

The weight norm serves as a measure for the model complexity. The idea of choosing the simplest possible model which can explain the data, often referred to as Occam's razor is expressed in the modified cost function \hat{E} , the result of the optimization process will be a compromise between fit to the training data and model complexity.

Dropout

Like input noise and weight noise dropout is part of the group of regularization techniques, that artificially corrupt training data to improve prediction performance [26]. When applying dropout between two layers, some units and their connections are zeroed out at random during training. The process is illustrated in figure 2.6. One

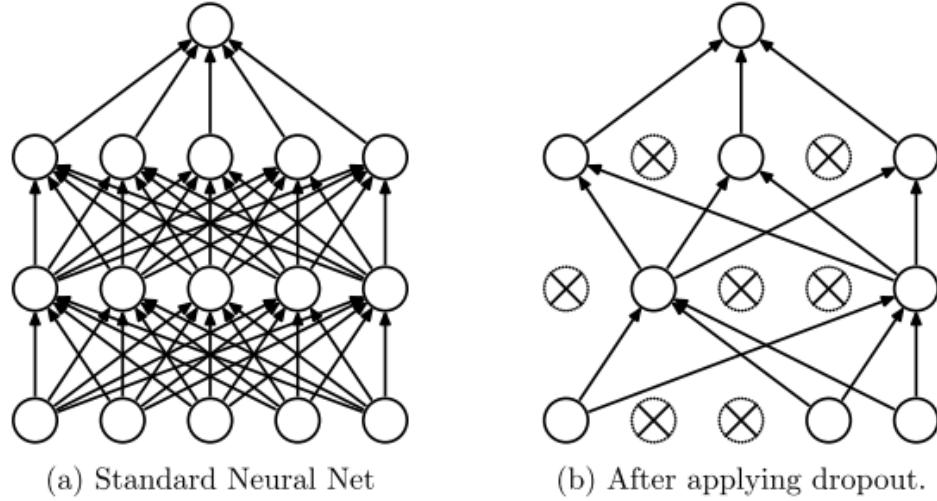


Figure 2.6: “Dropout Neural Net Model. Left: A standard neural net with 2 hidden layers. Right: An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.” [24], image and caption taken from the same source.

dropout interpretation relates dropout networks to committee networks. Network committees consist of several independently trained networks, each of which makes predictions on its own. The committee prediction is then produced by using the mean output of the committee members or by feeding the committee outputs into another neural net, which decides on the final output. When training large deep networks however committee networks become infeasible due to their large memory requirements. Dropout networks are related to committee networks, because when some units are dropped the remaining units will have to learn to work with the remaining units. Just like committee networks, dropout networks contain multiple mechanism to recognize the same pattern, and thus gains some of the prediction robustness that committee networks are known to have. Another motivation for dropout comes from the theory of sexual selection [24]. The theory assumes that the superiority of sexual reproduction is caused by the fact that removes strong interconnections between genes. After all a specific gene might not be there a generation from now. To function over many generations, genes must therefore develop redundant ways that work with different genes. Following this interpretation dropout forces network units to do the same thing.

2.2.3 Deep Neural Networks

Deep learning deals with the design of computational nets, consisting of multiple layers, which learn representations of data. Deep networks are often used when abstract patterns must be found. Simple networks are often unable to cope with

the problem complexity, and fail to recognize the underlying structure. Human speech features such complex patterns, therefore deep networks must be considered. Increasing the network depth leads to many additional weights, for which a suitable value must be found. In addition to the added memory load caused by the extra weights, the optimization process must look through large data sets until a suitable optimum is reached.

Stochastic gradient descent

When training deep networks on very large data sets, working with the full data set to compute the current gradient becomes very inefficient. It is common practice to use stochastic gradient descent instead of the classic algorithm. The idea is to select small subsets of the full data set. These so called mini-batches are then used to compute the outputs and errors, as well as the gradients for these smaller examples. Instead of working with the gradient of the union of these mini-batches the average gradient is computed and used to update the network weights. This approach dramatically reduces the memory requirements of the training process, which makes it possible to work with more sophisticated network designs.

2.2.4 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. A simple recurrent cell where the current state is set to the previous output is shown in figure 2.7. Another way to depict the same architecture is shown in Figure 2.8, here the cell state is not labeled explicitly, the single cell on the left is just a simplified version of the figure above. If the network is unrolled in time the flow of the state becomes apparent, which is done in figure 2.8 on the right. The unrolled form shows a direct dependency of the output at time t on the previous output at $t - 1$, which in turn depends on the previous outputs. This causes y_t and y_{t-1} to change together. In other words: The introduction of recurrent connections leads to correlation of the two outputs.

The exploding and vanishing gradient problem

Even though past information is available in theory, RNNs have only limited access to contextual information in practice [11, page 1]. Due to problems with gradient descent on correlated data, the back-propagated derivative can sometimes become weaker and weaker until it ultimately vanishes [16]. Another problem is that sometimes classical recurrent neural nets produce a gradient that blows up [22]. The exploding gradients can be fixed by clipping, but vanishing gradients require more sophisticated treatment [5].

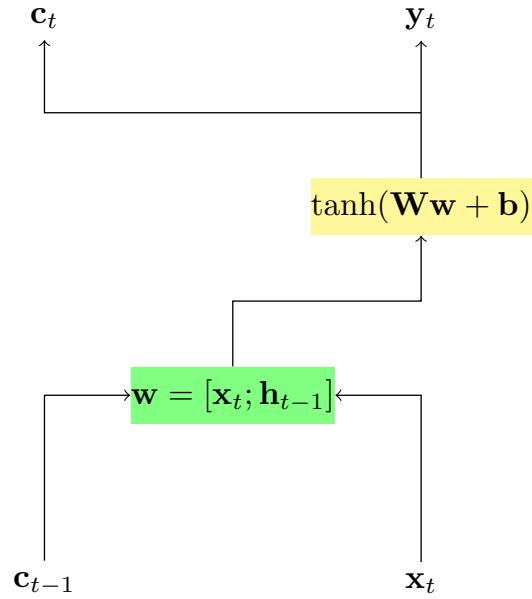


Figure 2.7: Visualization of a single recurrent cell.

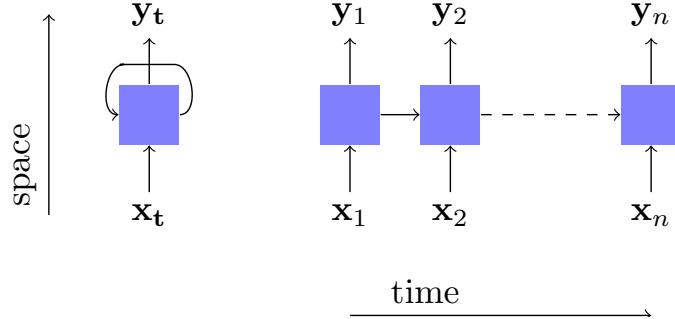


Figure 2.8: Rolled (left) and unrolled (right) recurrent neural net with two units.

Long short-term memory

Initially researchers tried to solve the vanishing gradient problem by making changes to the training algorithm using simulated annealing, time delays or compression [11, page 32]. However a good solution to the problem turned out to be changing the RNN cell architecture. Long short-term memory (LSTM) cells as proposed in [17] are more complex network units. LSTMs are differentiable versions of memory chips in a digital computer. Memory chips generally have read, write, and erase ports which can be set in order to allow the chips state to be read, modified or emptied. In the LSTM case the input, output, and forget gates serve the same function [11, page 33]. Throughout the literature these gates are generally denoted as i , f and o . The content of the memory or state is written as c , the output as h , while a t subscript denotes the time step. These memory cells use the differentiable equation system

[13, page 5]²:

$$\mathbf{i}_t = \sigma(\mathbf{W}_{ix}\mathbf{x}_t + \mathbf{W}_{ih}\mathbf{h}_{t-1} + \mathbf{W}_{ic}\mathbf{c}_{t-1} + \mathbf{b}_i) \quad (2.29)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{fx}\mathbf{x}_t + \mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fc}\mathbf{c}_{t-1} + \mathbf{b}_f) \quad (2.30)$$

$$\mathbf{c}_t = \mathbf{f}_t \mathbf{c}_{t-1} + \mathbf{i}_t \tanh(\mathbf{W}_{cx}\mathbf{x}_t + \mathbf{W}_{ch}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (2.31)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{ox}\mathbf{x}_t + \mathbf{W}_{oh}\mathbf{h}_{t-1} + \mathbf{W}_{oc}\mathbf{c}_t + \mathbf{b}_o) \quad (2.32)$$

$$\mathbf{h}_t = \mathbf{o}_t \tanh(\mathbf{c}_t) \quad (2.33)$$

$$(2.34)$$

From the definition of the matrix product follows that

$$\mathbf{Ax}_1 + \mathbf{Bx}_2 = [\mathbf{A} \quad \mathbf{B}] \cdot \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}. \quad (2.35)$$

With 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:

$$\mathbf{i}_t = \sigma(\mathbf{W}_i[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}_{t-1}]^T + \mathbf{b}_i) \quad (2.36)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_f[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}_{t-1}]^T + \mathbf{b}_f) \quad (2.37)$$

$$\mathbf{c}_t = \mathbf{f}_t \mathbf{c}_{t-1} + \mathbf{i}_t \tanh(\mathbf{W}_c[\mathbf{x}_t \ \mathbf{h}_{t-1}]^T + \mathbf{b}_c) \quad (2.38)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o[\mathbf{x}_t \ \mathbf{h}_{t-1} \ \mathbf{c}_t]^T + \mathbf{b}_o) \quad (2.39)$$

$$\mathbf{h}_t = \mathbf{o}_t \tanh(\mathbf{c}_t) \quad (2.40)$$

This system of equations is visualized in figure 2.9. Just like diagram 2.7 this depiction is read from bottom to top. The most important part is the line from \mathbf{c}_{t-1} to \mathbf{c}_t [8]. 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 $\sigma(\cdot)$ 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 $\in (0, 1)$ well 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 [8]. 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 $\tanh(\cdot)$ function makes sure all vector elements are between -1 and 1 . The neuron computing the candidate state values $\bar{\mathbf{c}}_t$ looks at input data and the past outputs. Both are labeled \mathbf{w} in figure 2.9, \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

² Various versions of LSTM cells exist. This one is commonly referred to as the “peephole” variant.

2. LITERATURE STUDY

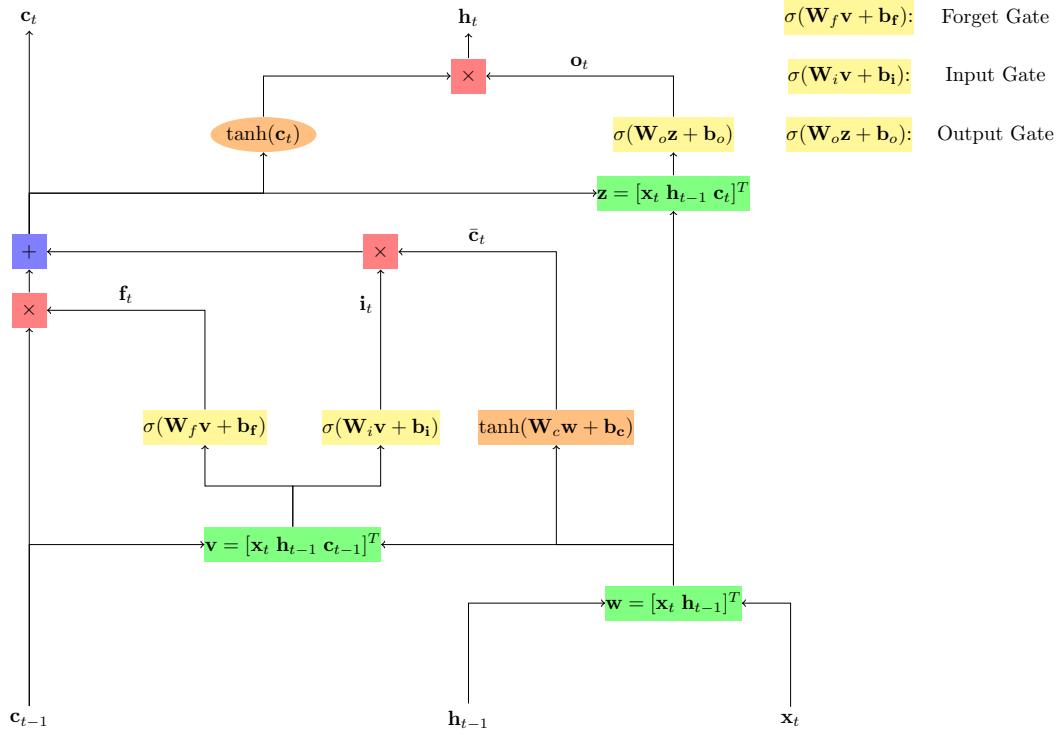


Figure 2.9: Visualization of the LSTM architecture.

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 the output gate. Its output 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 o_t then yields the new output h_t .

At this point it is interesting to note that the graphs considered in this thesis will almost exclusively consist out the exponential, logarithmic, sigmoid and hyperbolic tangent functions. These functions are intimately related. It is well known that the exponential and logarithmic functions are connected by:

$$\exp(\ln(x)) = x \quad \text{if } x > 0. \quad (2.41)$$

The sigmoid and hyperbolic tangent functions in turn consist out of exponentials, and are related trough [11, page 15]:

$$\tanh(x) = 2\sigma(2x) - 1. \quad (2.42)$$

The equation above can be proven by using the definitions of the sigmoidal $\sigma =$

$\frac{1}{1+\exp(-x)}$, and hyperbolic tangent $\tanh = \frac{\exp(2x)-1}{\exp(2x)+1}$:

$$\tanh(x) = \frac{\exp(2x) - 1}{\exp(2x) + 1} \quad (2.43)$$

$$= \frac{2\exp(2x) - (\exp(2x) + 1)}{\exp(2x) + 1} \quad (2.44)$$

$$= \frac{2\exp(2x)}{\exp(2x) + 1} - \frac{\exp(2x) + 1}{\exp(2x) + 1} \quad (2.45)$$

$$= \frac{2\exp(2x)}{\exp(2x) + 1} - 1 \quad (2.46)$$

$$= \frac{2}{\frac{\exp(2x)+1}{\exp(2x)}} - 1 \quad (2.47)$$

$$= \frac{2}{1 + \frac{1}{\exp(2x)}} - 1 \quad (2.48)$$

$$= 2\frac{1}{1 + \exp(-2x)} - 1 = 2\sigma(2x) - 1 \quad (2.49)$$

These two non-linear functions are thus very similar, in the LSTM-cell their different function value ranges decide how they are used. Sigmoids determine how much of something should be used, stored, or deleted, because its output values range between zero and one, while computations done on the cell state use the hyperbolic tangent, due to its wider range of function values.

Bidirectional Long Short Term Memory

With the advent of LSTMs deep recurrent networks became feasible in speech recognition [15]. 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 2.10. It is important to note, that several LSTM layers are often followed by a single linear layer, as is the case in [15]:

$$\vec{\mathbf{h}}_t = \text{LSTM}(\mathbf{W}_{\vec{\mathbf{h}}_t}^n [\mathbf{x}_t \ \mathbf{h}_{t-1}]^T + \mathbf{b}_{\vec{\mathbf{h}}_t}^n) \quad (2.50)$$

$$\overleftarrow{\mathbf{h}}_t = \text{LSTM}(\mathbf{W}_{\overleftarrow{\mathbf{h}}_t}^n [\mathbf{x}_t \ \mathbf{h}_{t+1}]^T + \mathbf{b}_{\overleftarrow{\mathbf{h}}_t}^n) \quad (2.51)$$

$$\mathbf{y}_t = \mathbf{W}_y^n [\vec{\mathbf{h}}_t \ \overleftarrow{\mathbf{h}}_t]^T + \mathbf{b}_y \quad (2.52)$$

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:

$$\vec{\mathbf{h}}_t^n = \text{LSTM}(\mathbf{W}_{\vec{\mathbf{h}}_t}^n [\mathbf{h}_t^{n-1} \ \mathbf{h}_{t-1}]^T + \mathbf{b}_{\vec{\mathbf{h}}_t}^n) \quad (2.53)$$

$$\overleftarrow{\mathbf{h}}_t^n = \text{LSTM}(\mathbf{W}_{\overleftarrow{\mathbf{h}}_t}^n [\mathbf{h}_t^{n-1} \ \mathbf{h}_{t+1}]^T + \mathbf{b}_{\overleftarrow{\mathbf{h}}_t}^n) \quad (2.54)$$

$$\mathbf{h}_t^N = \mathbf{W}_y^N [\vec{\mathbf{h}}_t \ \overleftarrow{\mathbf{h}}_t]^T + \mathbf{b}_y^N \quad (2.55)$$

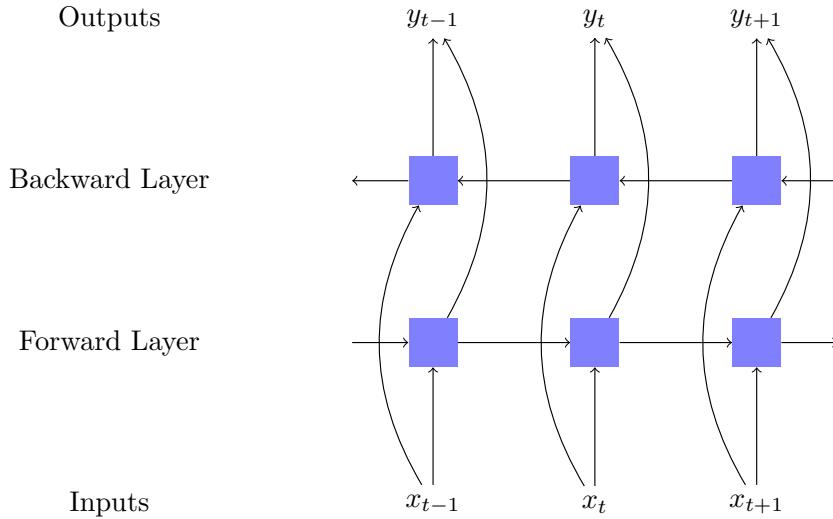


Figure 2.10: A bidirectional Long short term memory layer, according to [15]

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

2.3 Connectionist temporal classification

A bidirectional LSTM layer establishes access to past and future input information, but the alignment of the speech inputs and it's transcription remains unknown. Connectionist temporal classification is an output layer, which in tandem with the BLSTM layer enables the network to work around this problem, while at the same time allowing the gradient to propagate through it. The resulting network can be trained end-to-end. This entire section is based on [11, chapter 7].

CTC consists of a softmax layer with one output per alphabet element plus one for blank or no label. The idea is to allow the network to output a label at any input time step, and post process these outputs to get the overall sequence of labels correct, while ignoring the timing of each label. More formally L denotes the used alphabet and $L' = L \cup \emptyset$ the alphabet with the empty label. The network output vectors are denoted by \mathbf{y}^t and will have length $|L'|$. Given the training set S and the input sequence \mathbf{x} the probability for observing a sequence π consisting of elements from L' with length T is given by [11, page 56]:

$$p(\pi | \mathbf{x}, S) = \prod_{t=1}^T y_{\pi_t}^t. \quad (2.56)$$

In a first step all a mapping \mathcal{B} is introduced, which simplifies the allowed paths through the label probabilities, by removing duplicates. For example $\mathcal{B}(x \circ xy \circ) = \mathcal{B}(\emptyset xx \circ xyy) = xxy$. In other words a new label is only produced if the output label changed in the path. The probability of a labeling l can be computed using

[11, page 57]:

$$p(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{l})} p(\pi|\mathbf{x}). \quad (2.57)$$

Generally speaking this collapsing process is what makes it possible to use CTC with unaligned data.

CTC takes into account only the features and (depending on how the beam search is done) ignores the label it previously produced. Additionally the network output is coupled to the input frames. This shortcoming lead to the development of so called transducers. These models decouple the decoding time from the input time and use past labels to improve performance.

2.4 Listen, Attend and Spell [7]

The Listen Attend and Spell architecture is a deep neural network, designed to jointly learn to align and transcribe speech data. In contrast to CTC it belongs to the transducer family of speech recognition algorithms. All transducers have in common that they consist of an encoder and a decoder, with decoupled time steps. The idea also appeared in the automatic translation literature [2]. Working on the source language text, an annotator or encoder produces high level annotations. These annotations are then given to the translator or decoder, which using its own output time step works with the annotations to compute a context, which it then uses to produce a translation. In the speech application the encoder is referred to as the listener. The decoder as the speller. Combined these two form the las-network. The listener is a pyramidal recurrent neural net. It accepts filter bank spectra \mathbf{x}_n as inputs and produces compressed high level output features \mathbf{h}_m . Compression reduces the computational load during feature processing later. The speller in turn accepts the features as input and outputs distributions over character sequences \mathbf{y}_p . Due to the uncoupling of the input and decoding time step, it requires a state, as well as an attention mechanism. The state provides a memory of what happened in the decoder in the past. The attention function determines, which listener-features are relevant at a given decoding time step. Combining attention and state makes it possible to label the input data. An overview of the las-achrcitcure is given in figure 2.11.

2.4.1 The Listener

The listener, consists of Long Short Term Memory blocks. These blocks are arranged in layers. The inputs are first fed into a recurrent bidirectional layer (BLSTM). This choice gives the system access to future data, therefore only fully recorded data can be analyzed. This system is restricted to applications, which do not require transcriptions in real time or where its acceptable to wait with decoding until a sentence has been recorded completely. When going up in figure 2.11, pyramidal layers (pBLSTM) follow the initial BLSTM layer. The pyramidal structure concatenates the hidden values computed previously, such that their time dimension is halved:

$$\mathbf{h}_t^n = \text{BLSTM}(\mathbf{h}_{t-1}^n, [\mathbf{h}_{2t}^{n-1}, \mathbf{h}_{2t+1}^{n-1}]) \quad (2.58)$$

2. LITERATURE STUDY

Technically instead of two, three or more previously computed feature vectors could be concatenated, which increases the compression factor per pyramidal 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, because a longer time span has to be considered to decode a single character. Additionally the compression reduces the problem complexity, which speeds up the training process significantly.

2.4.2 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-1} . The index i denotes decoding time, $i - 1$ is used to refer to results from the last time step. The last decoding step I , at which the system terminates is a learned quantity. While the last input step U , depends on the input features, lower case u denotes the input step. During operation the Attend and spell functions keep track of previous output labels and previously important features, which are contained in the context. This information is stored in the state \mathbf{s}_i . To function the network must determine, which part of the computed features \mathbf{H} are relevant at any given decoding time step i . The context vector \mathbf{c}_i contains a linear combination of relevant features, weighted according to their importance. The AttentionContext function determines these weights based on the state. Finally the speller function finds a probability over possible labels using the on the relevant features in the context and the state. The computing steps are therefore [7, page 4]:

$$\mathbf{s}_i = \text{RNN}(\mathbf{s}_{i-1}, \mathbf{y}_{i-1}, \mathbf{c}_{i-1}) \quad (2.59)$$

$$\mathbf{c}_i = \text{AttentionContext}(\mathbf{s}_i, \mathbf{H}) \quad (2.60)$$

$$P(\mathbf{y}_i | \mathbf{x}, \mathbf{y}_{<i}) = \text{CharacterDistribution}(s_i, \mathbf{c}_i) \quad (2.61)$$

The state follows from a recurrent neural a multilayer LSTM. In contrast to the listener the speller is causal, meaning that it makes decisions only based on information computed during previous decoding steps. LSTMs are necessary here, because past states must be remembered. 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 \mathbf{h}_u . The computation starts with two feedforward neural networks or multilayer perceptrons (MLP), ϕ and ψ [7, page 5]:

$$e_{i,u} = \phi(\mathbf{s}_i)^T \psi(\mathbf{h}_u) \quad (2.62)$$

$$\alpha_{i,u} = \frac{\exp(e_{i,u})}{\sum_u \exp(e_{i,u})} \quad (2.63)$$

$$\mathbf{c}_i = \sum_u \alpha_{i,u} \mathbf{h}_u \quad (2.64)$$

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[7, page 5]. Thus it is justified to think of the alphas as a sliding window. This window contains only the currently relevant parts of the condensed input data set.

2.4.3 Training

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

$$\max_{\theta} \sum_i \log P(y_i | \mathbf{x}, y_{<i}; \theta). \quad (2.65)$$

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. In practice the objective is minimized by working with a cross entropy loss function. 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 [7, page 5]:

$$\hat{y}_i = \text{CharacterDistribution}(s_i, \mathbf{c}_i) \quad (2.66)$$

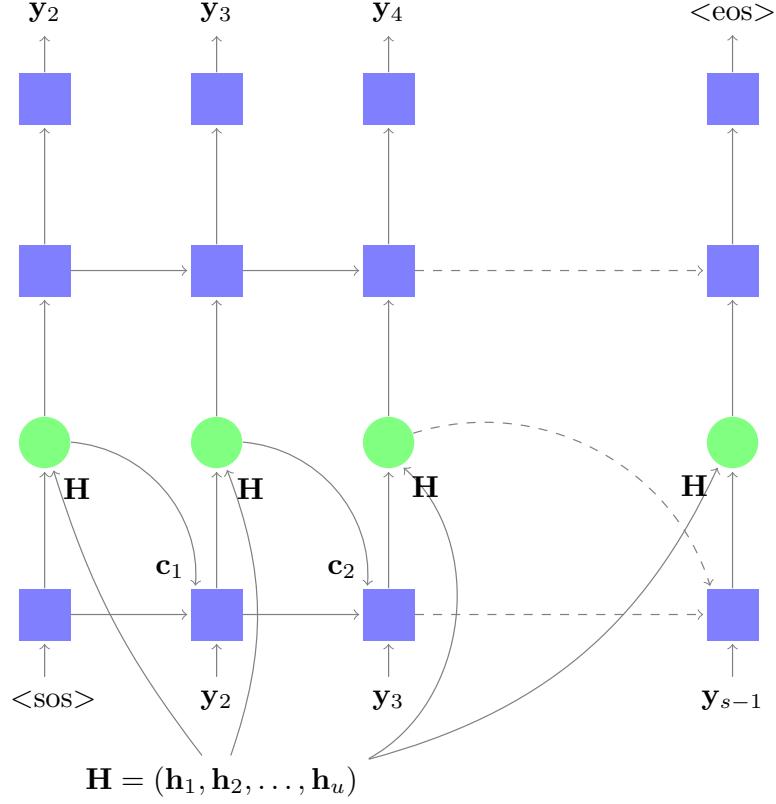
$$\max_{\theta} \sum_i \log P(y_i | \mathbf{x}, \hat{y}_{<i}; \theta) \quad (2.67)$$

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 [4].

2.4.4 Decoding with beam search

In order to generate a readable text, it is necessary to choose characters from the generated character distributions. One way to do this is to simply pick the most likely letter from each distribution. This approach is called greedy decoding. It ignores the possibility of generating better results by also considering less likely options. It is reasonable to expect better results, when considering more than just the most likely label, because the attend and spell decoder takes past labels into account. Consequently a broader search through the most likely options is considered. Unfortunately memory limitations make it impossible to search through all possible combinations. Therefore only the n most likely options are explored and the rest is disregarded. This approach is referred to as beam search. Taking into account the most likely options for each label produces a tree of possible transcriptions. The different routes along this tree are called hypotheses. A score for each hypothesis can be computed, by multiplication of the probability values the las-network assigned

Speller:



Listener:

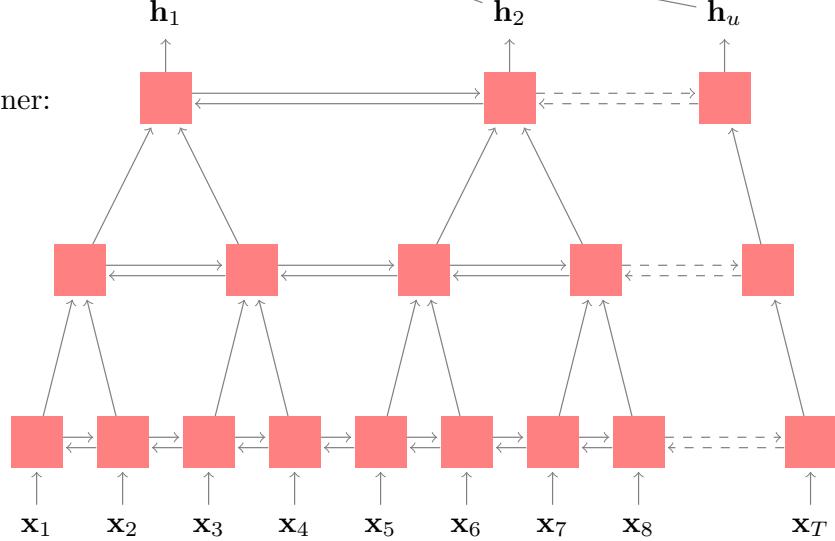


Figure 2.11: The LAS architecture [7, page 3]. BLSTM blocks are shown in red. LSTM blocks in blue and attention nets in green.

	ϵ	h	u	i	s
ϵ	0	1	2	3	4
h	1	0	1	2	3
o	2	1	1	2	3
u	3	2	1	2	3
s	4	3	2	2	2
e	5	4	3	3	3

Table 2.1: An example of the edit distance matrix for the dutch word huis and its english translation house. The edit distance can be found in the cell in the lower right of the table.

to each branch along its path. In order to account for different hypotheses lengths the total probability must be divided by the hypothesis length. In beam search only the m most likely hypotheses are kept. Using only las-probabilities is be equivalent taking only the acoustic data and their labels into account. In general text data is far more abundantly available then speech data. To take advantage of these large text corpora a language model trained on these can be used to make a more informed decision when choosing a hypothesis in the beam. A selection can then be made according to [7, page 6]:

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{LM}(\mathbf{y}) \quad (2.68)$$

Here P_{LM} denotes the weight the language model assigns to each hypothesis. And λ is a weight factor, which determines the language model importance. The formula above describes beam selection using a language model to re-score the attend and spell probabilities with a language model.

2.4.5 Levenshtein distance

A metric is required to measure the quality is label sequences. The Levenshtein or edit distance is such a metric. It can be thought of as the minimum number insertions, deletions or substitutions required to transform one sequence into another. An example computation for the two words huis and house is shown in table 2.1. The table cells count replacements, deletions or insertions. At position i, j in the table the first sequence is taken into account until element i , likewise the second sequence is considered up to position j . In order to compute the entries in the table the first row and column of the distance matrix D are initialized to:

$$D[0, j] = j \text{ and } D[i, 0] = 0. \quad (2.69)$$

These entries are known beforehand, because if one index is kept at zero, no elements of the corresponding sequence are considered and i or j insertions will always be

2. LITERATURE STUDY

necessary. If $i \neq 0$ and $j \neq 0$ the following update rule is used [20, slide 10]:

$$D[i, j] = \min \begin{cases} D[i - 1, j] + 1, & \text{deletion} \\ D[i, j - 1] + 1, & \text{insertion} \\ D[i - 1, j - 1] + \delta(x[i - 1], y[j - 1]) & \text{match or replacement} \end{cases} \quad (2.70)$$

$$\text{with } \delta(a, b) = 0 \text{ if } a = b \text{ else } 1 \quad (2.71)$$

Based in these rules the table 2.1, has been computed. Intuitively the distance three makes sense, because to go from huis to house, an o has to be inserted, the i replaced with an s and finally and an extra e has to be added, which amounts to three edits.

2.5 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”.

In tensorflow a directed graph is constructed based on the input operations which define a computational model. The graph represents how data flows trough the model. Each node of the graph has zero or more inputs and zero or more outputs. The nodes represent operations. Tensors flow along the edges of the graph. Graphs are executed in sessions. When running a graph the desired output variables or operations must be specified. Tensorflow then checks if all inputs required to compute the requested output values or perform the desired operation are present. Afterwards its evaluates only the operations including the requested node to find the desired output value or complete the requested operation. Generally tensors do not outlive a single execution of the graph. Variables however are persistent. The model parameters must be stored in variables. The gradient can then be applied to these variables to improve the model parameters over multiple graph executions.

Chapter 3

Methodology

3.1 Pylint

When writing code it is important to ensure that not only the author, but also future readers and coworkers are able to understand it. An important part of ensuring future usability is to document functionality. However readable code is not just well documented, but also respects well established coding conventions. An important collection of python code conventions is the python foundations PEP 8 style guide [25]. However following these conventions is tedious and easily forgotten. Additionally software development is team work and the best way to ruin the esprit de corps is an environment where team members constantly complain to each other about code convention violations. A good way to enforce decent coding style is to get the computer to help. In the python world a useful tool for this purpose is `pylint`¹. `Pylint` is a code analysis tool that looks for patterns in the code that might indicate bad style, for example incorrect variable names, missing documentation, wrong indentation and so forth. Being an interpreted language python code does not need to be compiled. This advantage comes with an important drawback, because during compiling many errors in the code are found and reported to the programmer. Python can do the same thing during runtime, but it crashes when an error is encountered and already computed data is often lost during those crashes. `Pylint` therefore taps into the python interpreter to scan the code for errors like a compiler. By alerting the programmer to potential issues in the code before execution, `pylint` can decrease runtime crashes and increase development efficiency.

¹<https://pylint.readthedocs.io/en/latest/intro.html>

Chapter 4

Experiments

4.1 TIMIT

The *timit* speech corpus [10], contains recordings of ten phonetically rich sentences, for example:

“She had your dark suit in greasy wash water all year. ”

For each sentence a transcription of the spoken phonemes is also available. Phonemes are sets of sounds, which are considered equivalent in a given language. In alphabetic writing systems such as the latin one the phoneme to letter mappings should ideally be one. Due to the fact that the Latin script was devised for classical Latin as well as the fact that when pronunciation changes the spelling often remains the same, the phoneme to letter mappings are often far from one. Therefore the **timit** data set comes with phonetic transcriptions for all sentences. For the sentence considered above the spoken phonemes are:

“h# sh ix hv eh dcl jh ih dcl d ah kcl k s ux q en gcl g r ix s ix w ao
sh epi w ao dx axr ao l y ih axr h# ”

The transcriptions contain a total of 64 possible phonetic labels, in the literature the full set and foldings with 48 and 39 labels are considered [21]. In all following experiments the 39 labels shown in 4.1 will be considered.

For all experiments the timit data set is split into a training, validation and test set. Containing 3696, 400 and 192 sentences [12, page 80].

4.2 BLSTM-CTC

This section explores bidirectional long short term memory layers with CTC output on the timit corpus [12, 14]. Training began after transforming of the speech data into Mel frequency cepstral coefficients (*MFCCCs*) augmented with first and second derivative information. The raw time domain signal as well as the feature vectors are shown in figure 4.2. The phonemes where folded as described above. Two *BLSTM* layers where stacked on top of each other, both using LSTM cells with a

4. EXPERIMENTS

TABLE I
LIST OF THE PHONES USED IN OUR PHONE RECOGNITION TASK

Phone	Example	Folded	Phone	Example	Folded
iy	<u>beat</u>		en	<u>button</u>	
ih	<u>bit</u>		ng	<u>sing</u>	eng
eh	<u>bet</u>		ch	<u>church</u>	
ae	<u>bat</u>		jh	<u>judge</u>	
ix	<u>roses</u>		dh	<u>they</u>	
ax	<u>the</u>		b	<u>bob</u>	
ah	<u>butt</u>		d	<u>dad</u>	
uw	<u>boot</u>	ux	dx	(butter)	
uh	<u>book</u>		g	<u>gag</u>	
ao	<u>about</u>		p	<u>pop</u>	
aa	<u>cot</u>		t	<u>tot</u>	
ey	<u>bait</u>		k	<u>kick</u>	
ay	<u>bite</u>		z	<u>zoo</u>	
oy	<u>boy</u>		zh	<u>measure</u>	
aw	<u>bough</u>		v	<u>very</u>	
ow	<u>boat</u>		f	<u>fief</u>	
l	<u>led</u>		th	<u>thief</u>	
el	<u>bottle</u>		s	<u>sis</u>	
r	<u>red</u>		sh	<u>shoe</u>	
y	<u>yet</u>		hh	<u>hay</u>	hv
w	<u>wet</u>		cl (sil)	(unvoiced closure)	pcl,tcl,kcl,qcl
er	<u>bird</u>	axr	vcl (sil)	(voiced closure)	bcl,dcl,gcl
m	<u>mom</u>	em	epi (sil)	(epenthetic closure)	
n	<u>non</u>	nx	sil	(silence)	h#,#h,pau

Figure 4.1: 48 to 39 phoneme folding as shown in [21].

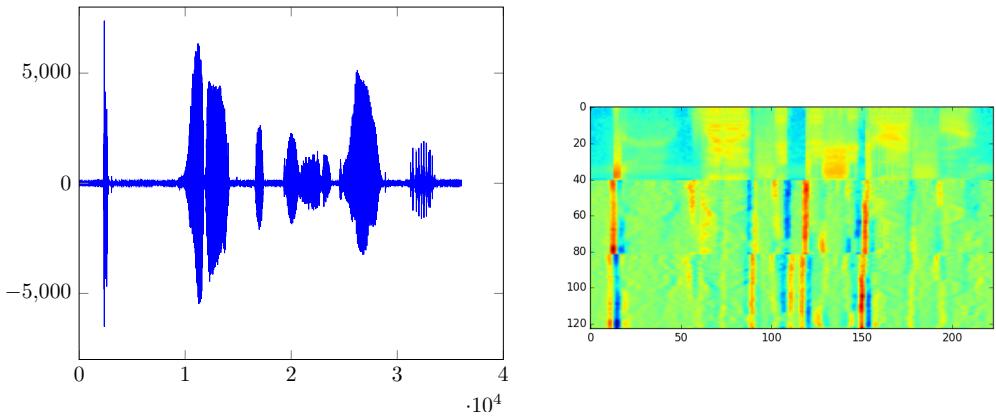


Figure 4.2: Raw time domain signal and feature vectors. Features include 40 mel banks as well as the first and second derivatives. Raw data was taken from TIMIT utterance faem0-si2022 “What outfit does she drive for?”.

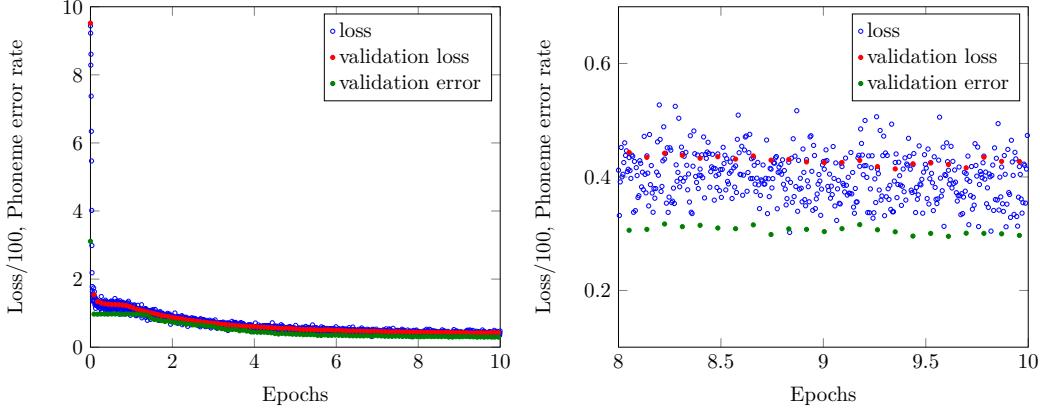


Figure 4.3: Validation and Test set error while training a two layer BLSTM network with CTC output layer on the TIMIT speech corpus with 39 folded phonemes.

dimension of 64 in each direction. The logits computed by the two layers were then fed into a CTC output layer, finally beam search decoding with a beam width of 100 was used to find the phoneme predictions. The whole network was optimized using gradient descent and a learning rate of 10^{-3} . To ensure generalization normal input noise $\mathcal{N}(\mu = 0, \sigma = 0.65)$ was added to the inputs. All feed-forward weights were initialized using another Gaussian $\mathcal{N}(0, 0.1)$, lstm weights were initialized using a random uniform distribution $\mathcal{U}(-0.1, 0.1)$. Gradients were clipped such that all gradient elements are within $(-1, 1)$. Results of the training process are shown in figure 4.3. In the experiment the training process was stopped after 10 epochs to be able to compare the result to the pyramidal listener experiment in the next chapter, at this point the phoneme error rate of this layout was at 29%. However if the training process is continued with decreasing learning rates the error rate values end up close to 26% eventually. Performance can be increased further by assigning the LSTM cells a larger dimension.

4.2.1 Efficient sequence to sequence implementations

When working with sequence to sequence methods such as CTC, data processing problems arise naturally. For efficiency reasons it is beneficial to work with fixed size data tensors. However the number of frames varies significantly between utterances. The same holds true for the target lengths. In order to still be able to work with fixed size tensors, the inputs and targets are padded with zeros or unknown tokens to fill up the length differences. In order to keep the speed benefit of working with tensors it is important to keep track of the sequence lengths of every individual utterance. The sequence length information must be used to avoid actually processing any of the padded data, that the tensors were filled up with in order to allocate memory faster. The listener produces a three dimensional logit tensor of size $[B, T, F]$, where B denotes the batch size, T the largest occurring frame number and F the feature dimension size. When implementing the output layer the naïve approach would be

4. EXPERIMENTS

to loop over the time dimension, and evaluate the linear layer T times with a B times F matrix. Which is simple to implement but very inefficient. Instead the sequence information should be used to remove the padding and concatenate the batch dimension into a matrix \mathbf{X} of size $[S, F]$, where S denotes the sum of all sequence lengths. The linear layer can then be evaluated as:

$$\mathbf{WX} + \mathbf{b} = \mathbf{R} \quad (4.1)$$

The result matrix R can then be converted back into a padded tensor using the sequence length information one more time. An alternative approach which avoids reshaping is to implement a linear recurrent network cell with a state size of zero and use dynamic unrollings to evaluate the linear layer. The resulting memory requirements and speed are comparable.

4.3 Listen attend and spell experiments

In this section experiments use the full listen attend and spell architecture. After verification of the listener using a CTC output layer, the listen attend and spell network is first tested using greedy decoding. In contrast to beam search greedy decoding, does not maintain several hypotheses, instead it works with the most likely label each time step.

4.3.1 Non experimental verification of the code

Testing code implementing artificial neural networks is difficult, because experiments only start to work once all issues in the code have been fixed. A single bug can prevent experiments from working, in such a case the only feedback the programmer gets is a negative experimental result. In order to prevent incorrect implementations in recurrent neural networks it is important to strictly enforce shape invariants of state tensors in loops. This means that the shape of the tensors going into an RNN cell must have the same size when they leave. This way one can be certain, that for example the network does not change the number of label probabilities it predicts over time.

4.3.2 Testing the Listener

A crucial part of the listen attend and spell architecture is formed by the listener. Before working with a fully-fledged `las` a CTC-layer will be attached to the listener. The idea is to verify the implementation. If CTC can extract relevant information from the listener, the attend and spell code should be able to do the same in later experiments. In order to keep memory requirements manageable in later experiments with the same listener all LSTM cells were set up with 64 hidden units. As the listener layers are bidirectional this means 64 in each direction so the hidden dimension is 128 in total. This sum is important, because the feature dimension of the lstm outputs is concatenated for each time step. If not further action is

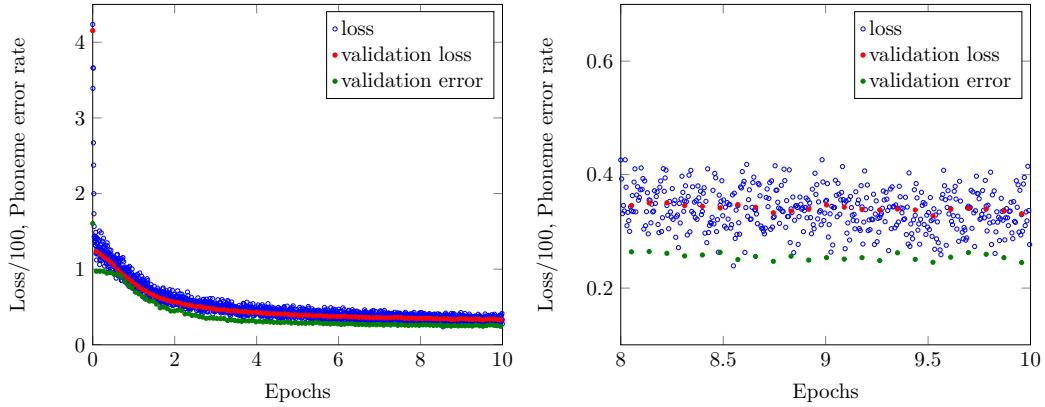


Figure 4.4: The training progress shown for the Listener with added CTC layer. The loss values shown in blue and green have been scaled with $\frac{1}{100}$. On the right a closeup on the last two training process is shown.

taken the listener produces features with a dimension of two times the number of elements per lstm. CTC runs the logits it is given through a softmax layer to compute label probabilities. To function it must therefore be given a logit tensor, where the feature dimension is equal to the number of labels, the system is supposed to output. To meet this requirement an extra linear output layer has been added to the listener which maps the feature dimension to 40, as required. Figure 4.4, shows the optimization algorithms progress, as measured by training loss, validation loss and validation set decoding phoneme error rate. The training was stopped, when the decoding results were no longer improving. During testing a phoneme error rate of 0.268 was observed, which is a pretty good result compared to the twenty four percent error rate of a comparable full BLSTM architecture given that the pyramidal layer compressed the time dimension into half of its original size. During decoding the CTC beam width was once more set to 100.

4.3.3 Design of the decoding loop logic

4.3.4 Greedy Decoding Experiments

Using the tested listener with 64 hidden lstm units per direction and one pyramidal layer, the CTC layer is replaced with attend and spell functions. For computational efficiency these functions have been implemented within a customized RNN cell. Using an RNN framework makes it possible to use optimized code to unroll the attend and spell computations. Among other things dynamically unrolling the second part of the network this way ensures efficient memory utilization and sequence length management. Within the new cell the decoder state size was chosen to be 128, considering the fact that the listener outputs features of size 128, which in turn determines the context vectors to have this same dimension. Hoping to provide sufficient memory to remember past context vectors the decoder state RNN was set to the same dimension. The state and feature networks, ϕ and ψ were given one

4. EXPERIMENTS

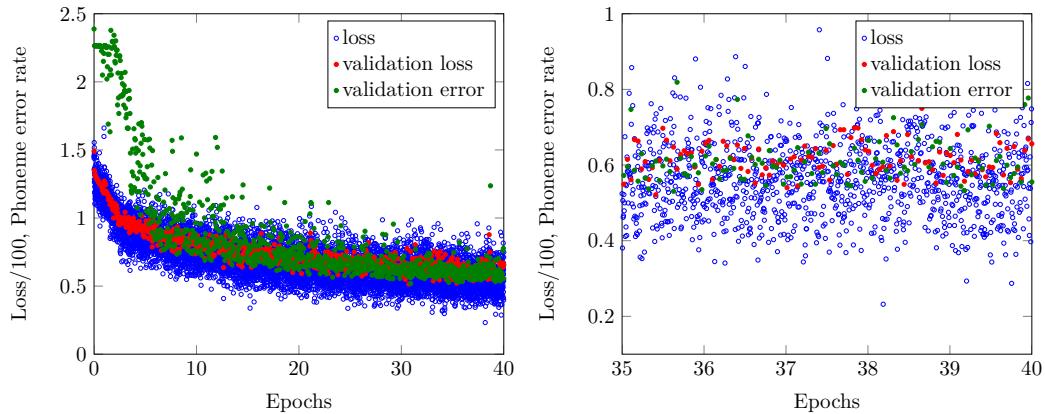


Figure 4.5: The training progress shown for the full las architecture with greedy decoding. The loss values show in blue and green have been scaled with $\frac{1}{100}$. On the right a closeup on the last two training process is shown.

hidden layer each, with a hidden dimension of 64. This choice was made mainly to conserve memory. During training the network output was used instead of the true target with a probability of 0.7. Figure 4.5 shows an overview of the training process. When considering the last five epochs, the decoding error ranges between 0.5 and 0.9. This means that in the best case half of the labels produced by the system must be modified in order to get to the target sequence. Considering timit utterance `fml0_sx295` the folded transcription with additional start and end of sentence tokens is given by:

Listing 4.1: Targets

```
<sos>  sil  ih   f   sil  k   eh   r   l   sil  k   ah   m   z
          sil  t   ah   m   aa   r   ah   hh   ae   v   er   r   ey
          n   jh   f   er   m   iy   dx   iy   ng   ih   sil   t   uw   sil
<eos>
```

From the input features the las network decodes:

Listing 4.2: Network output

```
<sos>  sil  hh  ih   f   sil  k   er   r   ow   ow   sil   sil
          t   ah   m   aa   hh   hh   ae   v   er   r   r   n   n   sil
          f   er   er   m   iy   iy   iy   iy   iy   iy   iy   iy   sil
          sil   t   uw   sil
<eos>
```

The decoding and target sequences clearly bear some resemblance. However significant errors do exist in the network output in particular in the last third. The phoneme sequence `dx iy ng ih` is incorrectly transcribed as `iy iy iy iy iy iy iy`, which has a large impact on the error. The levenshtein distance between the two labellings is 21. Given that the target sequence contains 42 labels including the start and end of sentence tokens, the error rate of this example is 0.5, it is therefore slightly better

4.3. Listen attend and spell experiments

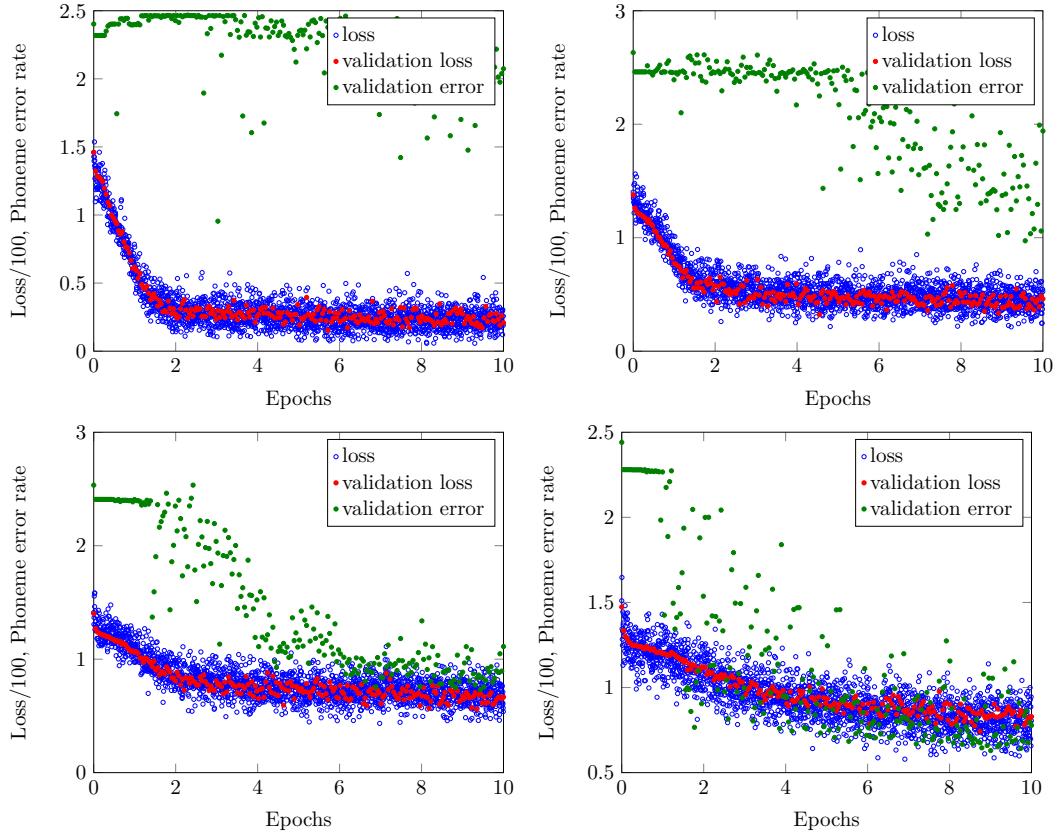


Figure 4.6: Repetitions of the same experiment with network output reuse probabilities 0.2, 0.4, 0.6, 0.8.

than the average of 0.553362, which is the average decoding error rate over the entire validation set.

4.3.5 The effect of the groundtruth selection probability during training

In the previous experiment the groundtruth was only used instead of the network output in 0.3 percent of the cases. As the network output is often incorrect during training, less emphasis will be put on past outputs later during decoding. This section investigates other values. The same las network is retrained for 10 epochs using output probabilities of 0.2, 0.4, 0.6, 0.8. Results are shown in figure 4.6. It can be observed, that decoding results improve when reduced emphasis is placed on past labels during training. This observation is confirmed by the phoneme error rates of 2.08, 1.97, 1.1168, 0.87, which were observed on the validation set. To explore the effect of a low reuse probability over longer training periods the experiment has been repeated one more time with a 0.5 groundtruth probability over 40 epochs. Figure 4.7 depicts the training process. In comparison to the experiment shown in

4. EXPERIMENTS

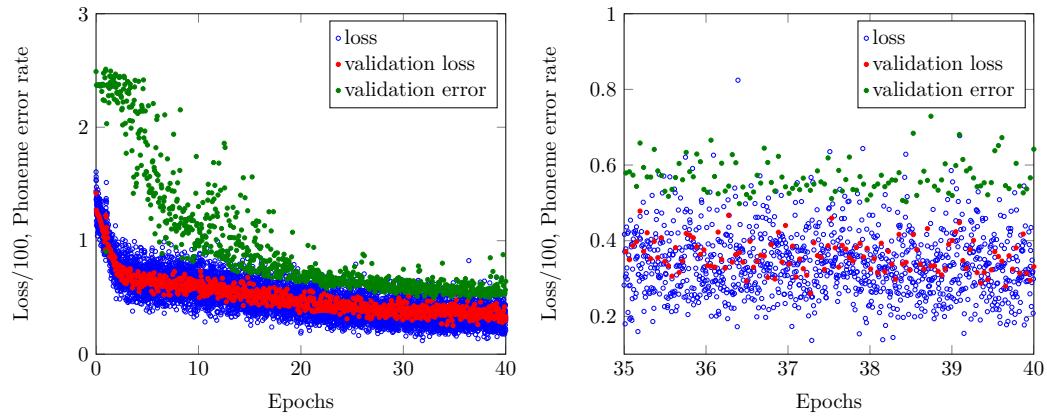


Figure 4.7: Training full las over 40 epochs with a network output reuse probability of 0.5

figure 4.5 the network does significantly better during training, but this improvement does not translate into a better decoding performance.

Ideally one would like to observe the opposite. If the output labels where correct more often during decoding, the networks trained to rely on them should outperform those which where not. Based on the observations above we conclude that beam search should be implemented in order to obtain labels of higher quality during decoding, which in turn should be beneficial to networks, which learned to rely on past outputs.

4.3.6 Comparing two attend and spell variants

4.3.7 Beam-search Decoding

4.3. Listen attend and spell experiments

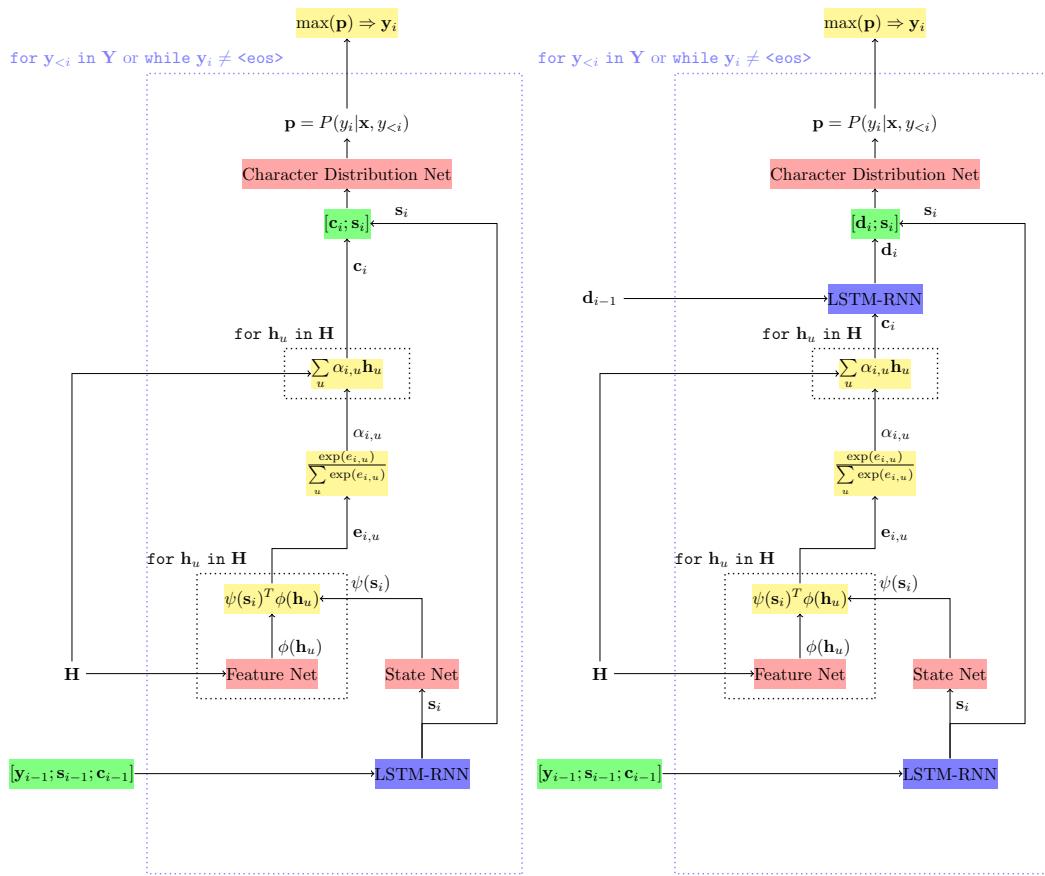


Figure 4.8: Two different attend and spell cell configurations.

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. *CoRR*, abs/1603.0, 2016.
- [2] D. Bahdanau, K. Cho, and Y. Bengio. NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE. *CoRR*, abs/1409.0:1–15, 2014.
- [3] S. Ben-david. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge Cambridge university press, Jerusalem, 2014.
- [4] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks. *arXiv*, pages 1–9, 2015.
- [5] 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.
- [6] C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [7] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals. Listen, attend and spell. *arXiv preprint*, pages 1–16, 2015.
- [8] Christopher Olah. Understanding LSTM Networks, 2015.
- [9] M. Diehl. Script for Numerical Optimization Course B-KUL-H03E3A. Technical report, Optimization in Engineering Center (OPTEC) and Electrical Engineering Department (ESAT-SCD), KU Leuven, Leuven, 2013.
- [10] J. Garofolo, L. Lamel, W. Fisher, J. Fiscus, D. Pallett, N. Dahlgren, and V. Zue. TIMIT Acoustic-Phonetic Continuous Speech Corpus, 1993.
- [11] A. Graves. *Supervised Sequence Labelling with Recurrent Neural Networks*. 2008.

BIBLIOGRAPHY

- [12] A. Graves. *Supervised Sequence Labelling with Recurrent Neural Networks*. 2012.
- [13] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, pages 1–43, 2013.
- [14] A. Graves, S. Fernandez, F. Gomez, and J. Schmidhuber. Connectionist Temporal Classification : Labelling Unsegmented Sequence Data with Recurrent Neural Networks. *Proceedings of the 23rd international conference on Machine Learning*, pages 369–376, 2006.
- [15] 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.
- [16] 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.
- [17] S. Hochreiter and J. Schmidhuber. LONG SHORT TERM MEMORY. *Technical Report FKI-207-95*, pages 1–8, 1995.
- [18] X. Huang, A. Acero, and H.-W. Hon. *Spoken Language Processing: A Guide to Theory, Algorithm and System Development*. Prentice Hall PTR, Redmond, 2001.
- [19] 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.
- [20] B. Langmead. Dynamic programming and edit distance, 2016.
- [21] K. F. Lee and H. W. Hon. Speaker-independent phone recognition using hidden Markov models. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 37(11):1641–1648, 1989.
- [22] 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.
- [23] R. Rojas. *Neural Networks - A Systematic Introduction - Backpropagation*. Springer Berlin Heidelberg, 1996.
- [24] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [25] G. van Rossum, B. Warsaw, and N. Coghlan. PEP 8 – Style Guide for Python Code, 2001.

BIBLIOGRAPHY

- [26] S. Wager, S. Wang, P. Liang, and C. Science. Dropout Training as Adaptive Regularization. *Advances in Neural Information Processing Systems*, 26:351—359, 2013.
- [27] P. Woodland, G. Evermann, and M. Gales. HTKBook, 2006.

Master thesis filing card

Student: Moritz Wolter

Title: Building an end-to-end speech recognizer

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 supervisors: Prof. dr. ir. Patrick Wambacq

Prof. dr. ir. Hugo van Hamme

Assessor: Prof. dr. ir. Johan Suykens

Mentor: Ir. Vincent Renkens