



TAMPEREEN TEKNILLINEN YLIOPISTO
TAMPERE UNIVERSITY OF TECHNOLOGY

LAURI NISKANEN
MULTI-LABEL SPEAKER RECOGNITION USING
RECURRENT NEURAL NETWORKS

Master of Science thesis

Examiners: Professor Tuomas Virtanen and
Professor Hannu-Matti Järvinen

The examiners and topic of the thesis were
approved on 9 August 2017

ABSTRACT

LAURI NISKANEN: Multi-Label Speaker Recognition using Recurrent Neural Networks

Tampere University of Technology

Master of Science thesis, 43 pages

November 2018

Master's Degree Programme in Information Technology

Major: Pervasive Systems

Examiners: Professor Tuomas Virtanen and Professor Hannu-Matti Järvinen

Keywords: speaker recognition, multi-label audio classification, recurrent neural networks

Speech recognition is a popular research topic that analyzes human speech. In addition to understanding the spoken message, it is beneficial to know who is speaking. This thesis studies speaker recognition and presents a machine learning based system for identifying the speakers from audio streams. Our implementation is based on Mel-frequency cepstral coefficients (MFCC) and recurrent neural networks.

The system is developed and evaluated on AMI Meeting Corpus dataset. The dataset contains annotated meeting recordings with typically four participants in each. Our system processes the audio files of the recordings in 20 millisecond slices and produces a list of active speakers at each time step.

We measure the performance of our system using various metrics. The results indicate that our system is capable of identifying the speakers with decent accuracy. The best classifier model that we examined is a 1-layer long short-term memory (LSTM) neural network with layer size 256. Neural networks that are more complex than it do not seem to improve the classification results, but they suffer from increased training times. We also suggest alternative classifications methods for future research.

TIIVISTELMÄ

LAURI NISKANEN: Monen puhujan tunnistaminen takaisinkytkettyillä neuroverkoilla

Tampereen teknillinen yliopisto

Diplomityö, 43 sivua

Marraskuu 2018

Tietotekniikan diplomi-insinöörin tutkinto-ohjelma

Pääaine: Pervasive Systems

Tarkastajat: professori Tuomas Virtanen ja professori Hannu-Matti Järvinen

Avainsanat: puhujantunnistus, äänen multi-label-luokittelu, takaisinkytketyvät neuroverkot

Puheentunnistus on suosittu tutkimusalue, jossa analysoidaan ihmisten puhetta. Puhutun viestin ymmärtämisen lisäksi on hyödyllistä tietää kuka puhuu. Tässä diplomityössä tutkitaan puhujantunnistusta ja esitellään koneoppimiseen perustuva järjestelmä puhujien tunnistamiseksi äänivirroista. Toteutus perustuu MFCC-piirteisiin ja takaisinkytkettyihin neuroverkkoihin.

Järjestelmän kehittämiseen ja testaukseen käytetään AMI Meeting Corpus -aineistoa, jossa on aikaleimallisesti litteroituja kokousäänitteitä. Yhdessä kokouksessa on tyypillisesti neljä osallistujaa. Järjestelmä käsittelee äänitallenteita kahdenkymmenen millisekunnin siivuissa ja tuottaa jokaiselle ajanhetkelle listan aktiivisista puhujista.

Järjestelmän suorituskykyä mitataan erilaisilla metriikoilla. Tulokset osoittavat, että järjestelmä kykenee tunnistamaan puhujia kohtuullisella tarkkuudella. Paras tarkastelluista luokitinmalleista on yksikerroksinen LSTM-neuroverkko, jossa kerroksen koko on 256. Tätä monimutkaisemmat neuroverkot eivät vaikuta parantavan luokitus tuloksia, mutta niiden opettamiseen kuluu enemmän aikaa. Ehdotamme myös vaihtoehtoisia luokitusmenetelmiä jatkotutkimuskohteiksi.

PREFACE

This thesis was written between August 2017 and November 2018 while I was employed by Bitwise Oy in Tampere. I am grateful for the opportunity to work with the interesting topics of this thesis. I enjoyed creating the implementation in practice and finding out how it would actually be able to recognize the speakers.

I would like to thank all my colleagues and friends who helped and supported me. Special thanks to Ville Nukarinen who reviewed the thesis extensively and provided invaluable feedback.

Lauri Niskanen

Tampere, 8 November 2018

CONTENTS

1	INTRODUCTION	1
2	SPEAKER RECOGNITION	3
2.1	Frequency spectrum	4
2.2	Mel-frequency cepstrum	5
3	MACHINE LEARNING	7
3.1	Supervised learning	7
3.2	Artificial neural networks	7
3.3	Backpropagation	10
3.4	Multi-label classification	11
3.5	Overfitting and regularization	12
3.6	Recurrent neural networks	14
4	IMPLEMENTED METHOD	16
5	EVALUATION	19
5.1	Datasets	19
5.2	Evaluation metrics	20
5.3	Hyperparameters	22
5.4	Results	22
6	CONCLUSION	36
	REFERENCES	38

LIST OF ABBREVIATIONS

AMI	Augmented Multi-party Interaction
ANN	artificial neural network
BPTT	backpropagation through time
BRNN	bidirectional recurrent neural network
CNN	convolutional neural network
DCT	discrete cosine transform
DFT	discrete Fourier transform
DTW	dynamic time warping
FN	false negative
FP	false positive
GMM	Gaussian mixture model
GPU	graphics processing unit
GRU	gated recurrent unit
HMM	hidden Markov model
LPC	linear prediction coding
LSTM	long short-term memory
MFCC	Mel-frequency cepstral coefficients
PCM	pulse-code modulation
RNN	recurrent neural network
ROC	receiver operating characteristic
ROC AUC	area under the receiver operating characteristic curve
TN	true negative
TP	true positive
VQ	vector quantization
XML	Extensible Markup Language

1 INTRODUCTION

For a long time computers have been better than any human in doing calculations or simple repetitive tasks. However, tasks that require deeper understanding, creativity, or imagination have been very hard for computers to do. Only recently computers have begun to conquer many of these problems. For example, research topics like computer vision, robotics, natural language processing, and automated medical diagnosis have been greatly advanced with the help of modern machine learning [12, 18, 35, 50].

These hard problems are typically so complex that it becomes impossible for a programmer to manage all possible cases in a systematic way. The solution is to use data-driven statistical methods to reduce the dimensionality of the task. Machine learning is the study of algorithms that make predictions from collected data. In contrast to classical computer algorithms, machine learning algorithms are typically somewhat general and data is in a very important role. The quality, quantity, and representation of the data can have huge impact on the predictions. The recent success of machine learning is a combination of theoretical advances, more and more extensive data collection, and the availability of high performance graphics processing units (GPU). [48]

One important problem area where computer systems have become better with machine learning is processing human speech. Speech recognition is an especially popular research topic in which the textual message of speech is analyzed. However, human speech also contains information about the age, gender, emotion, and identity of the speaker. In this thesis, we focus on speaker recognition and study how the identity of a speaker correlates with audible features. [49, 54]

Speaker recognition has a wide range of possible applications. One of them is annotating who speaks when in meeting recordings or other audio tracks. It could also be used for improving speech controlled home automation systems with multiple users. Voice commands could select their default parameters and preferences based on who gave the command. For example, a specific personal calendar could be selected when a new event is being created. Speaker recognition has also been used as a biometrical component for customer verification in financial services and in criminal investigations [43, 51].

We implement a machine learning system that can identify speakers from given audio samples based on how the voices of individual speakers sound different. The system is applied to a dataset containing recordings of meetings. Each meeting has typically four attendees having natural conversations. Our system analyzes the

audio recording in slices and tries to determine who is speaking at each time point. People mostly speak in turns, but it is not uncommon that two or more people are speaking over each other. Meetings also have brief parts where nobody is speaking. To accommodate these situations, our system does not simply name a single speaker per time point, but instead outputs a list of active speakers. Since this list of speakers is given at each point in time the output for the whole recording is two dimensional as shown in Figure 1.

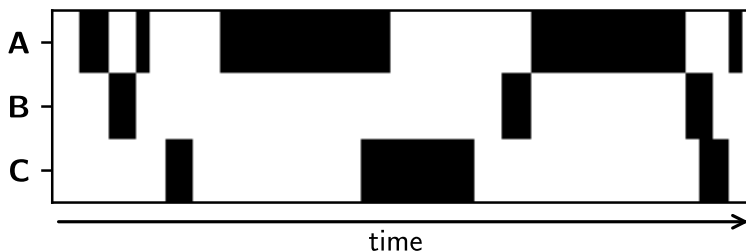


Figure 1: A simplified example of the program output, where A, B, and C are different speakers. The filled segments represent times where each speaker is active during the audio track.

Our system is based on a recurrent neural network classifier. Recurrent neural networks have the ability to remember past information to aid future predictions. This way the system can learn that one speaker is often active for some time before the speaker changes. During difficult points in the audio track the system can support its decisions by the understanding of the previous moments. We compare different neural network models and show how their hyperparameters affect the predictions.

Chapter 2 has an introduction to speaker recognition. Chapter 3 explains how machine learning and neural networks can be used to achieve our goal. Chapter 4 presents our speaker recognition implementation and Chapter 5 analyzes the performance of our methods. Chapter 6 has concluding discussion and proposed topics for future research.

2 SPEAKER RECOGNITION

Speaker recognition aims to detect on which parts of an audio track someone is speaking and to identify who the speakers are on each of those parts. There are many ways to recognize speakers. One approach is to use multiple microphones with known locations in relation to the speakers [61]. It is also possible to do speech recognition on the text of the speech and then recognize one or more specific phrases, for example have the speakers say their name or a password. However, in this thesis we are studying text-independent methods, where we only use acoustic characteristics of speech irrespective of what is being said, and without using the location of the microphones. [28]

This kind of recognition is possible, because people have individual voices. The differences in the voices are mainly caused by anatomical differences in the vocal tract. The shape of the vocal tract produces different resonances in the voice, also called formants. Other affecting factors include the anatomy of the lungs and the trachea. These differences in the voices can be analyzed using the frequency spectrum of the audio. [7]

Traditionally speaker recognition has been split to two phases: speaker diarization and speaker identification. The purpose of speaker diarization is to split an audio track into segments with only one speaker in each and also separate parts where nobody is speaking [65]. The goal of speaker identification is to then detect the identity of the speaker in each segment [54].

However, the traditional approach has some limitations. First, diarization systems often need to process a whole file at a time and speaker identification is typically done only after segmentation. This means that speaker recognition cannot be done live on an audio stream. Second, people often interrupt each other or talk simultaneously when trying to take the floor. This cannot be accurately represented with one-speaker segments.

To resolve these issues, it is possible to do diarization and identification jointly in one pass. Instead of telling who is speaking on each segment, we recognize the speaker at each time frame. The idea of live recognition has been researched by Vinyals and Friedland [67]. As we avoid needing the one-speaker segments we can take a step further and do multi-label speaker recognition by giving a list of simultaneous speakers for each time frame.

Speaker recognition systems contain two main components: feature extraction and feature matching. The purpose of feature extraction is to represent the audio signal

of the examined speech in a compressed form where useless information is filtered out. There are many speech feature extraction methods such as linear prediction coding (LPC) [46] and Mel-frequency cepstral coefficients (MFCC) [13]. The purpose of feature matching is to connect the extracted speech features to the speaker identity or otherwise classify or cluster them. Feature matching techniques that are used with speech include dynamic time warping (DTW) [58, 49], hidden Markov models (HMM) [2, 53], Gaussian mixture models (GMM) [57, 54], vector quantization (VQ) [14, 45], and artificial neural networks (ANN) [16, 38, 39, 47]. [1, 28]

This thesis presents and analyzes a live multi-label text-independent speaker recognition system with Mel-frequency cepstral coefficients (MFCC) as the speech feature extraction method and recurrent neural network (RNN) as the feature matching method. Next, we introduce the theory behind the vocal feature extraction methods that are needed for our implementation. The neural network components are covered in Chapter 3.

2.1 Frequency spectrum

Digital audio streams are typically represented with pulse-code modulation (PCM) [4]. It is a time domain representation where the amplitude of the audio signal is sampled with regular intervals. High quality audio signals are typically stored with 44.1 kHz or 48 kHz sampling rate.

However, in speaker analysis we are interested in the high-level features of the audio signal. Individual amplitude samples of the one-dimensional PCM audio signal do not directly correlate with anything that would be useful for speaker recognition. The solution is to transform the signal to a more useful format. Practically all vocal feature extraction methods used in speaker recognition are based on distinguishing the individual frequency modes, or formants, using the frequency spectrum of the speech sample [15].

The spectrum can be calculated using discrete Fourier transform (DFT) [27], which converts the audio samples from the time domain to the frequency domain. In the time domain we can see how the amplitude changes over time, but in the frequency domain we can analyze how the audio frequencies of the signal behave. Discrete Fourier transform is defined by

$$X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-i2\pi kn/N}, \quad \text{for } k = 0, \dots, N - 1,$$

where x_0, x_1, \dots, x_{N-1} represents uniformly spaced time domain samples and X_0, X_1, \dots, X_{N-1} is a sequence of complex numbers containing information about the amplitude and phase of the frequencies in the signal.

2.2 Mel-frequency cepstrum

One step further to make the signal more compact and better suited for speaker classification is to calculate the Mel-frequency cepstrum [13]. The Mel-frequency cepstral coefficients (MFCC) are widely used as a feature vector in the field of automated speech analysis [40, 55, 70, 72]. The Mel scale [64] is a scale of audio pitches derived from listening experiments with the purpose of mimicking how humans perceive audio signals [40]. A frequency f given in hertz can be converted to mels using the formula [52]

$$f_m = 1127 \ln \left(1 + \frac{f}{700 \text{ Hz}} \right) \text{ mel.}$$

Cepstrum of signal \bar{x} can be defined as

$$\bar{C} = DCT\{\log |DFT[\bar{x}]|^2\},$$

where DFT is the discrete Fourier transform and DCT is the discrete cosine transform [5]. In Mel-frequency cepstrum, the frequency bands are spaced based on the Mel scale [40]. Mel-frequency cepstrum is calculated by applying the Mel-weighting function w_m before the discrete cosine transform [70]:

$$\bar{C}_m = DCT\{w_m(\log |DFT[\bar{x}]|^2)\}.$$

The computation flow for the Mel cepstrum is illustrated in Figure 2. The discrete Fourier transform is used to calculate the power spectrum of the audio signal. Phase information can be discarded because it has been shown to be not as important. The

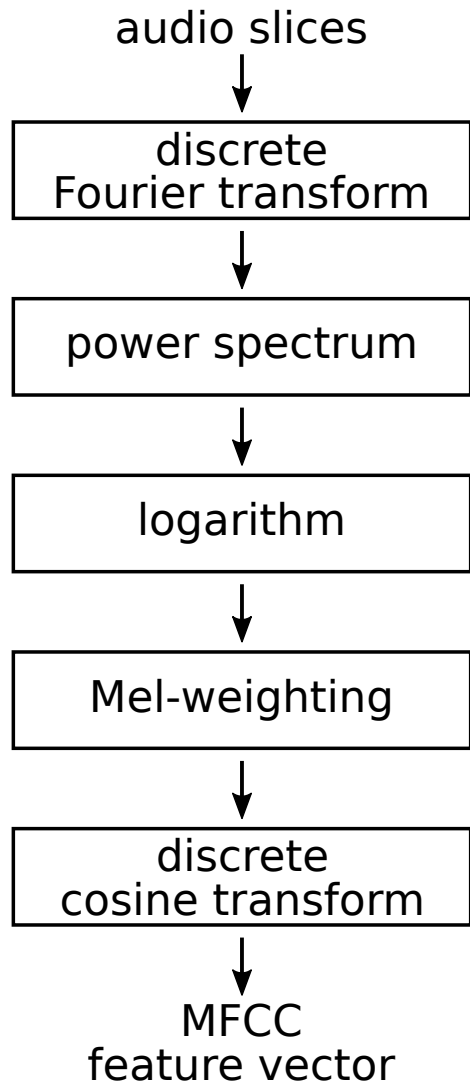


Figure 2: Pipeline for MFCC feature vector calculation.

logarithm of the spectrum is taken, because it approximately matches the perceived loudness of the signal. Finally, after scaling the signal to the Mel scale, discrete cosine transform is taken to reduce the number of parameters. This is useful because the calculated Mel-spectral vectors consist of highly correlated components. Karhunen–Loève transform [33] would be more precise, but with speech signals discrete cosine transform is commonly used to approximate it. [40]

3 MACHINE LEARNING

Machine learning can be applied to a wide range of problems and there are multiple ways to use the collected data to solve the machine learning problems. In this thesis, we are mainly interested in classification, where the goal is to assign a category for each given item. Next, we will explain how classification can be implemented with supervised learning using artificial neural networks.

3.1 Supervised learning

In classification, the goal is to train a prediction model, a classifier, that tells to which category a sample belongs based on its features. The classifier can be trained using example data. Each example has a vector of features and a target label. Features are the input attributes that describe the sample. Labels are the categories to which each sample belongs. [48]

The data must be split to three sets: training samples, validation samples, and test samples. Training samples are used for training the classifier model. Validation samples are used to compare different methods and to adjust the model parameters. Test samples are used to test the performance of the trained classifier. It is important that test samples are not available for the algorithm during the learning stage. The classifier model can be tested by using the model to predict where the test samples should belong based on its features and comparing the result with the sample label that represents the truth. [3]

In supervised learning the algorithm is given access to both the training features and training labels. In contrast, in unsupervised learning the algorithm can only see unlabeled features. There are also other scenarios that differ in how the data is available to the algorithm. [48]

There are many algorithms and models that can be trained to make classification predictions [3]. Next, we will present one of them: artificial neural networks.

3.2 Artificial neural networks

Neural networks were first researched as a way to represent biological information processing with mathematics by McCulloch in 1943 [47]. The term has since been associated with numerous different models, most of which are only remotely related to biology if at all. Neural networks consist of a fixed number of interconnected para-

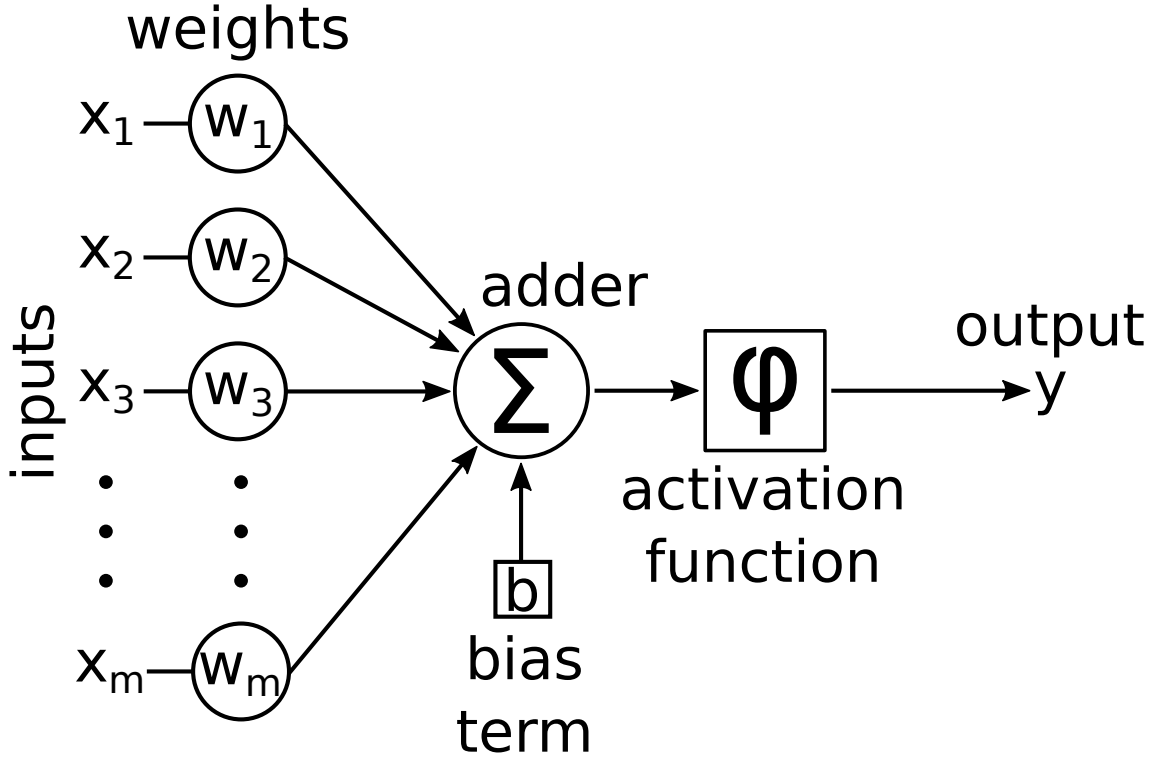


Figure 3: The components of an artificial neuron.

metric units, or artificial neurons, whose parameters can be adjusted in the training phase so that they and the network as a whole produce the desired output. [3]

There are multiple models for artificial neurons, but in the common basic case a neuron can be defined by the parts shown in Figure 3. Each neuron has a set of connecting links to predecessor neurons with weights associated to each link. The adder component takes the output values from the connected predecessor neurons, multiplies them by the link weights, and calculates the sum of these values. A bias term can also be added to the sum. Activation function is a function that takes the calculated sum as an input and produces an output value for the neuron. Some commonly used activation functions are listed in Table 1. [29]

rectified linear unit	$\begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
logistic sigmoid	$\frac{1}{1 + e^{-x}}$
hyperbolic tangent	$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
softplus	$\ln(1 + e^x)$

Table 1: Common activation functions [21].

This basic neuron model can be described mathematically as the equation

$$y = \varphi \left(b + \sum_{j=1}^m w_j x_j \right),$$

where y is the output value of the neuron, φ is the activation function, b is the bias term, w_1, w_2, \dots, w_m are the link weights, and x_1, x_2, \dots, x_m are the connected predecessor neuron values. [29]

A feedforward neural network is composed of individual neurons arranged in layers as shown in Figure 4. The links between the neurons are defined so that the predecessor of a neuron is in the preceding layer. Some of the neurons are selected to be in the output layer of the network. Their value is visible as the output vector of the whole network. Similarly, some of the neurons are used as the input to the network. The input neurons do not have any predecessors neurons, but instead their value is given from outside of the network. Neurons that are neither input nor output units are called hidden units as they are not directly exposed to the outside. [3]

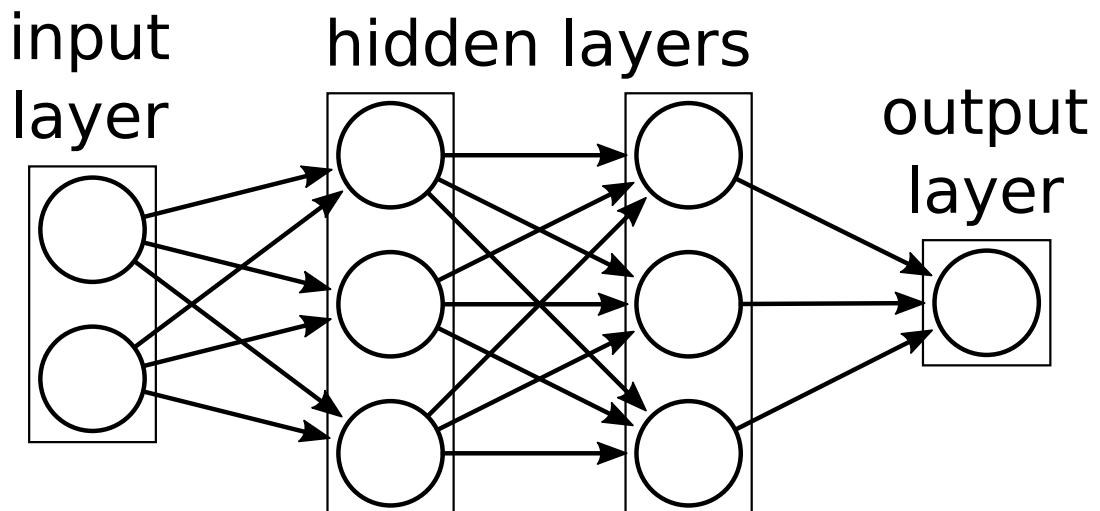


Figure 4: A simplified feedforward network composed of neurons arranged in layers.

Typically the network structure is expressed by defining the number of hidden layers in the network and the number of units in each layer, also called the size of the layer [3]. Different layers may have different sizes. The input and output layer sizes are usually defined by the intended use of the network. The number and sizes of hidden layers can be adjusted based on the desired network complexity. Different layers may also use different activation functions or extended neuron models, some of which we will examine later. These parameters related to the structure of the network are called the hyperparameters of the network. They are typically not changed during the training process. However, next we see how other parameters,

especially the neuron link weights, are not fixed in place but instead being adjusted dynamically during neural network training.

3.3 Backpropagation

Training a neural network is a process where the neuron link weights, and sometimes other parameters, are adjusted so that the network can produce the desired output values for given input values [29]. In other words, the internal parameters of a neural network are modified so that the network can approximate a given function. Backpropagation [68] is a popular and effective algorithm for training neural networks.

In supervised learning, the neural network is trained with a set of labeled training examples. These examples, or training samples, are pairs of input and output values. When the network is trained to do classification, the input values are called features, which describe a given sample in some way. The output vector of the network is used for encoding the class of the sample. For example, in a binary classification case there are positive and negative samples, and the class of the sample can be encoded simply by using one neuron in the output layer and by differentiating the class by the value of the neuron (e.g. 0 or 1). [29]

Before starting the training process it is not known what values the neuron link weights should have. They can be initialized randomly. The backpropagation algorithm works by iterating training samples making a forward pass and a backward pass for each of them. In the forward pass, the values of the input units are set according to the features of the sample. Then, the neuron functions are evaluated propagating the signal towards the output layer. This process produces values for the output neurons of the network. The output values are compared with the expected output vector, which in supervised learning is known for each sample. In the backward pass the neuron link weights are adjusted so that the difference between the expected output vector and the produced output vector decreases. The goal of the training process is to iteratively tweak the parameters in the network so that the response produced in the forward pass matches the desired one more and more closely. [29]

There are multiple ways to determine how the parameters should be adjusted. Mathematically speaking, we need to define a loss function that the training process is trying to minimize. For binary classification the loss function can be the logistic sigmoid function on an output neuron [3]. We will later examine the choice of the loss function and the encoding scheme for output neurons for various classification

problems. The gradient of the loss function is useful in calculating the needed change to the parameters. To decrease the loss, a small step can be taken in the direction of the negative gradient [3]:

$$a_{n+1} = a_n - \eta \nabla E(a_n),$$

where a_n is the value of the examined parameter before the adjustment and a_{n+1} after it, ∇E is the gradient of the loss function, and η is a parameter called the learning rate. Learning rate determines how great the adjustment steps are. This way of adjusting the network is known as stochastic gradient descent if the training samples are iterated randomly evaluating the loss function one sample at a time [56].

3.4 Multi-label classification

Basic single-label binary classification answers the yes–no question "Does the sample match the criteria?" In practice the question can be for example "Is the object in the given image a cat?" In these cases, the neural network output is typically encoded as a single neuron whose value is either one or zero. Values in between may represent varying levels of uncertainty. This encoding scheme can be extended to support multiple mutually exclusive classes, answering for example the question "Is the animal in the image a cat, a dog, or a horse?" With n classes the output is often encoded with n neurons so that one neuron has the value one and the rest have value zero. Uncertainty can be encoded by assigning weighted values to the neurons so that the sum of the values of all output neurons is one. This is called multiclass or multinomial classification with one-hot encoding. [66]

Another way to extend the scheme is to allow multiple labels that are not mutually exclusive. The classification question could be "Which of the listed animals, if any, are present in the image?" Multi-label classification answers to n questions at the same time. Often these are binary questions and the output can be encoded with n neurons whose values are all separately between one and zero. It is also possible to further extend the output encoding to allow more complicated structures. [66]

We have argued that the structure of the output layer depends on the type of classification. In addition, the activation function of the output layer neurons and the classifier loss function must be chosen accordingly. In single-label binary classification, the activation function is typically logistic sigmoid and the loss function is logistic loss. There are also many alternative activation and loss functions for this case. In multiclass classification, the natural choice is to use softmax activation and

categorical cross entropy loss, which is an extension of logistic loss. Softmax is a function that squashes a vector of values in between zero and one so that their sum is one. It is mathematically defined as

$$\varphi_i(\mathbf{x}) = \frac{e^{x_i}}{\sum_{k=1}^n e^{x_k}} \quad \text{for } i = 1, \dots, n,$$

where x_1, x_2, \dots, x_n are the input values [3]. In binary multi-label classification logistic sigmoid can be used as the activation function for all the output neurons separately. The multi-label extension of logistic loss is called binary cross entropy loss. These activation and loss functions for different classification problems are listed in Table 2. [22]

classification type	output activation	loss function
binary single-label	logistic sigmoid	logistic loss
multiclass	softmax	categorical cross entropy
binary multi-label	logistic sigmoid	binary cross entropy

Table 2: Common output layer activation functions and loss functions for different classification problems.

Choosing the activation and loss functions correctly is important. For example, if softmax and categorical cross entropy loss were used in a multi-label classification problem where the labels are not actually mutually exclusive, the network would be unable to make an output with two simultaneously active labels and naïvely interpreting the output values would produce somewhat correlated but as a matter of fact meaningless results. The other way around, if a mutually exclusive multiclass classifier used separate logistic sigmoid activations instead of softmax, the produced probability distribution would be invalid because the total probability would not be exactly one.

3.5 Overfitting and regularization

A supervised model should be able to repeat predictions that it was given in the training phase. If the model is too simple, it might fail to do this because it is not able to capture all the details of the training samples. This phenomenon is called underfitting. However, it is also very important that the model can generalize and make good predictions about new samples. It is not enough to remember all the details of the training samples since the number of training samples is in often limited and the samples may contain noise. By noise we mean random details in the features that do not represent the underlying properties of the data. When the model learns too much of this noise instead of the intended structure, it is said to be overfitting. An example of this behavior is shown in Figure 5. [6]

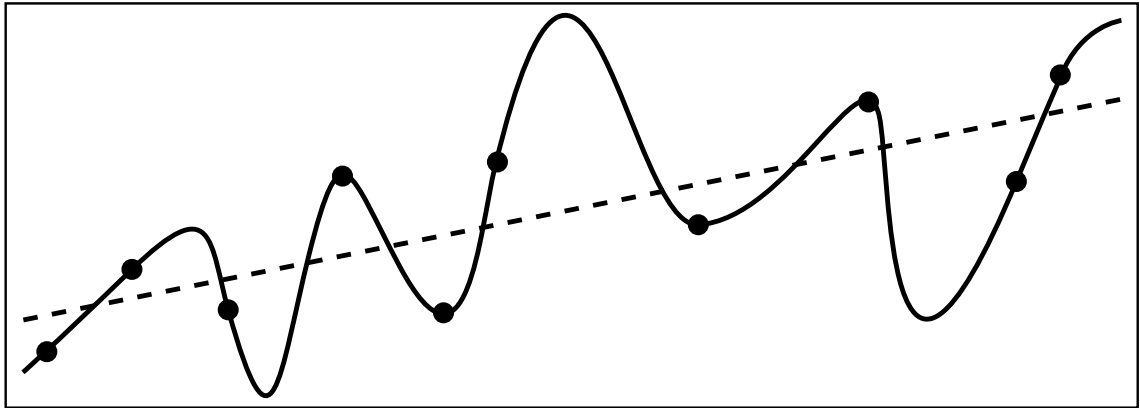


Figure 5: The solid line is a precise fit to the data points, but the dashed line may generalize better when making predictions.

It is beneficial to use a model structure that is so simple that it cannot learn all the noise in the data so it must instead approximate the samples it was given in order to minimize the loss function. This approximate often generalizes better producing good predictions for unknown samples. In neural network models a natural way to adjust the complexity and learning capacity is to choose the number of neurons accordingly. One way to prevent overfitting is to train multiple models and average the predictions. However, there are regularization techniques that are designed to prevent overfitting with less computation. [36]

Overfitting can be visualized using learning curves. These curves show validation loss as a function of training time. The validation loss of an ideal model would monotonically converge to a global minimum. If the model overfits, the validation loss starts to increase after a point. Early stopping is a popular regularization method that leverages this behavior by stopping the training when the loss starts to go up. Figure 6 has learning curves for two example models. The dashed line converges and the solid line represents a model that overfits. [60]

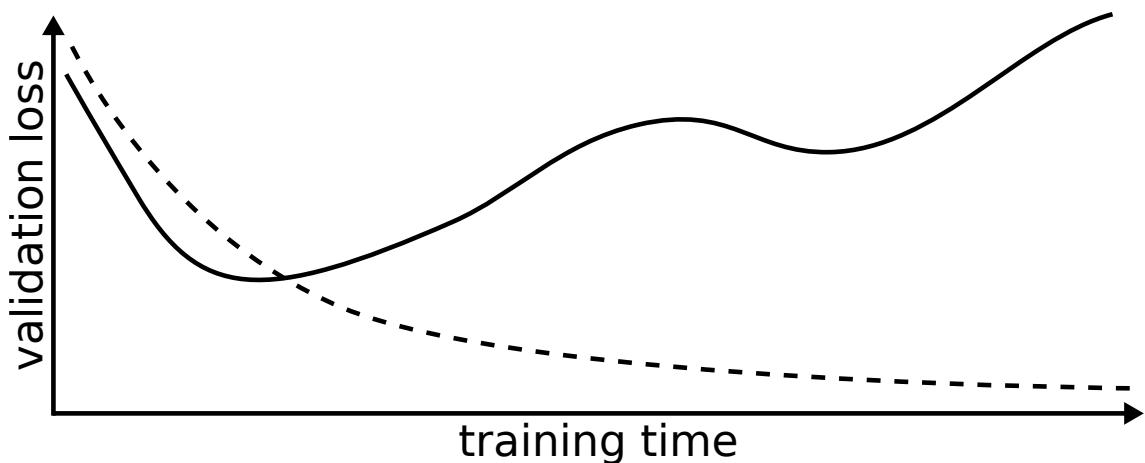


Figure 6: Learning curves for two models showing loss convergence and overfitting.

Dropout [63] is a commonly used regularization method for neural networks [19]. Dropout works by randomly dropping neurons from the network during training as illustrated in Figure 7. This prevents the neurons from relying too much on each other and forces the network to create redundant paths. After the training phase all neurons are enabled. The output of single neurons and ultimately the output of the whole network is averaged over multiple paths making the results more robust against noise. Using dropout will increase the training loss of the network, but if working properly it will lower the validation loss. [63]

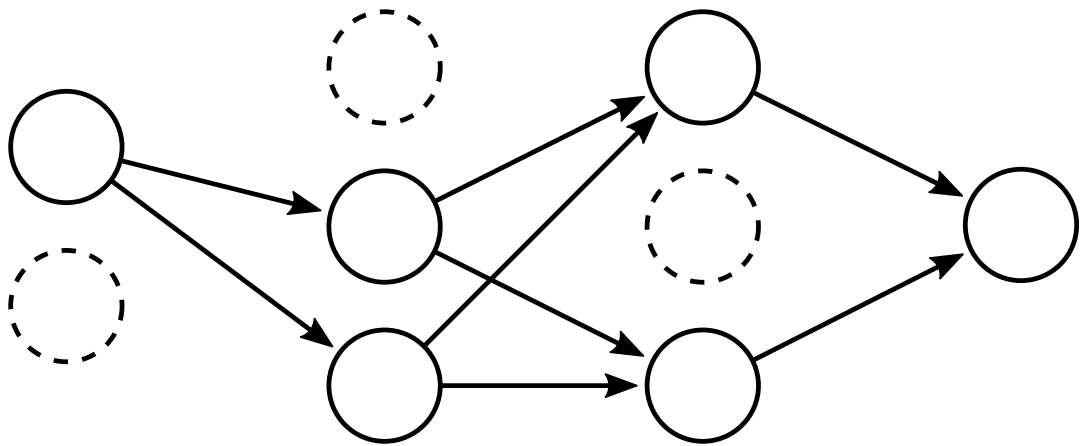


Figure 7: With dropout random neurons are disabled in the training phase.

3.6 Recurrent neural networks

Recurrent neural networks (RNN) [32] were first introduced by Hopfield in 1982. They gained popularity after Hochreiter and Schmidhuber discovered long short-term memory (LSTM) [30] networks in 1997. Long short-term memory networks outperform many other models in multiple fields including natural language text processing and speech recognition [17, 44]. Gated recurrent units (GRU) [9] have also been proved to have similar performance with LSTM networks.

Recurrent neural networks are applied to sequences like a stream of text, audio, or time series data points. In contrast, traditional neural networks must be applied to fixed-length vectors and cannot handle variable-length sequences. Recurrent neural networks are also able to output variable-length sequences. Recurrent neural networks are based on recurrent neurons that extend the basic neuron model by adding connections to earlier neurons along the sequence. This effectively enables recurrent neurons to remember values over time. [24]

An LSTM neuron contains an input gate, a memory cell, a forget gate, and an output gate as shown in Figure 8. The dashed lines represent connections along the

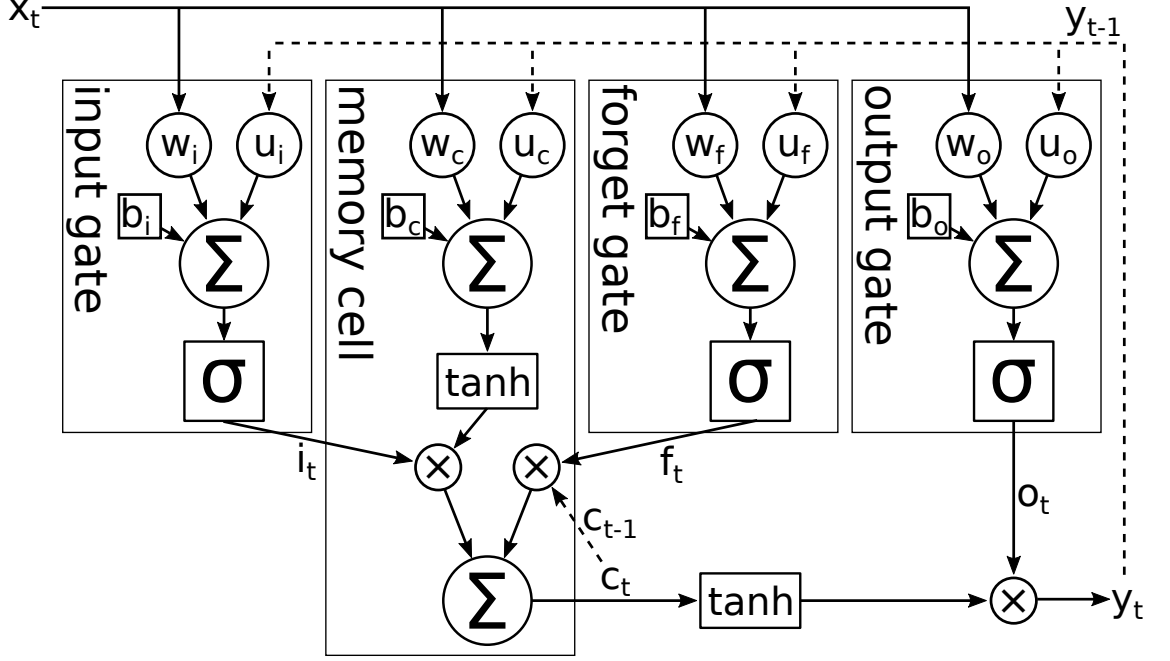


Figure 8: The structure of a long short-term memory neuron.

sequence. The LSTM gates regulate how the information is stored and accessed from the cell allowing the neuron to remember a value for a time it is useful but no longer than necessary. This structure tries to avoid the vanishing gradient problems that often limit the usefulness of recurrent neural network models [31]. An extension of the backpropagation algorithm, called backpropagation through time (BPTT) [69], can be used to efficiently train recurrent neural networks. [20]

The forward pass of an LSTM neuron is defined by equations [20]

$$\begin{aligned}
 i_t &= \sigma(w_i x_t + u_i y_{t-1} + b_i) \\
 f_t &= \sigma(w_f x_t + u_f y_{t-1} + b_f) \\
 o_t &= \sigma(w_o x_t + u_o y_{t-1} + b_o) \\
 c_t &= i_t \tanh(w_c x_t + u_c y_{t-1} + b_c) + f_t c_{t-1} \\
 y_t &= o_t \tanh(c_t),
 \end{aligned}$$

where x_t is the input vector, f_t , i_t , and o_t are the activation vectors of the forget gate, the input gate, and the output gate, respectively, c_t is the cell state vector, y_t is the output vector of the neuron, σ is the logistic sigmoid function, and w , u , and b are the weight matrices and the bias vectors.

4 IMPLEMENTED METHOD

We implemented a multi-label speaker recognition model that can identify known speakers in a given audio track. For each trained speaker, the model outputs a probability of speaking at each time frame in the audio track. The output is visualized in Figure 9. In this model, it is possible that multiple speakers are speaking at the same time. The model is trained with an audio track and a synchronized binary label track for each speaker. The label track used in training has the same structure as the program output, but with only binary values at each time point as shown in Figure 1 in the introduction chapter.

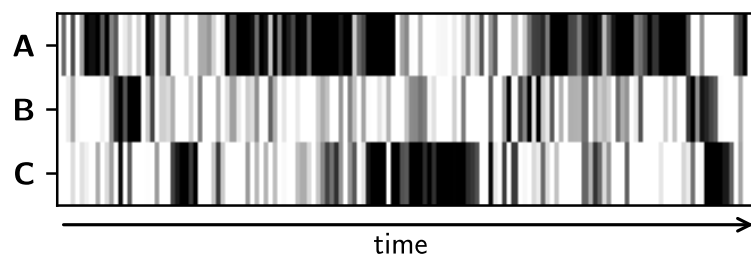


Figure 9: Program output with predicted probabilities shown as grayscale values.

The pipeline for making speaker activity predictions is shown in Figure 11 (see the next page). The input audio track is given in PCM format. It is sliced to time frames of 20 milliseconds and 18 MFCC audio feature values are calculated for each frame. These MFCC features frames are given to the classifier model and it will output a probability vector for each frame. The probabilities represent the likelihood of each speaker being active during the frame. With an ideal classifier the output probabilities would mimic the binary ground truth labels where each speaker is either speaking or not on each time frame. However, in practice the output contains time frames with uncertain predictions which may be incorrect.

In the training phase (Figure 10) the pipeline has the same preprocessing and feature extraction stages, but instead of making predictions in the classification stage the model is given both the MFCC features vectors and the ground truth labels. The

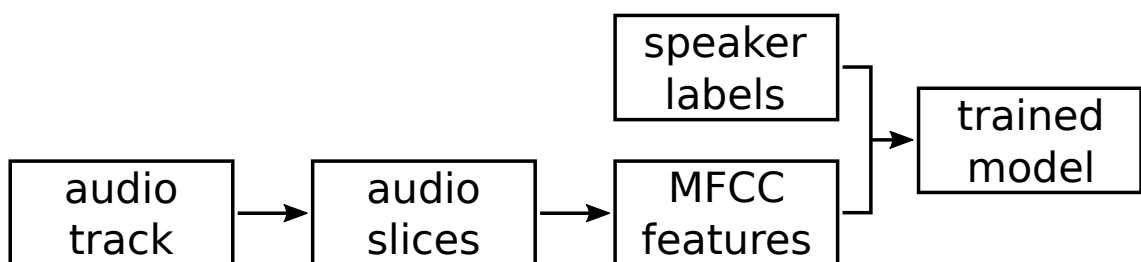


Figure 10: Training the classifier model.

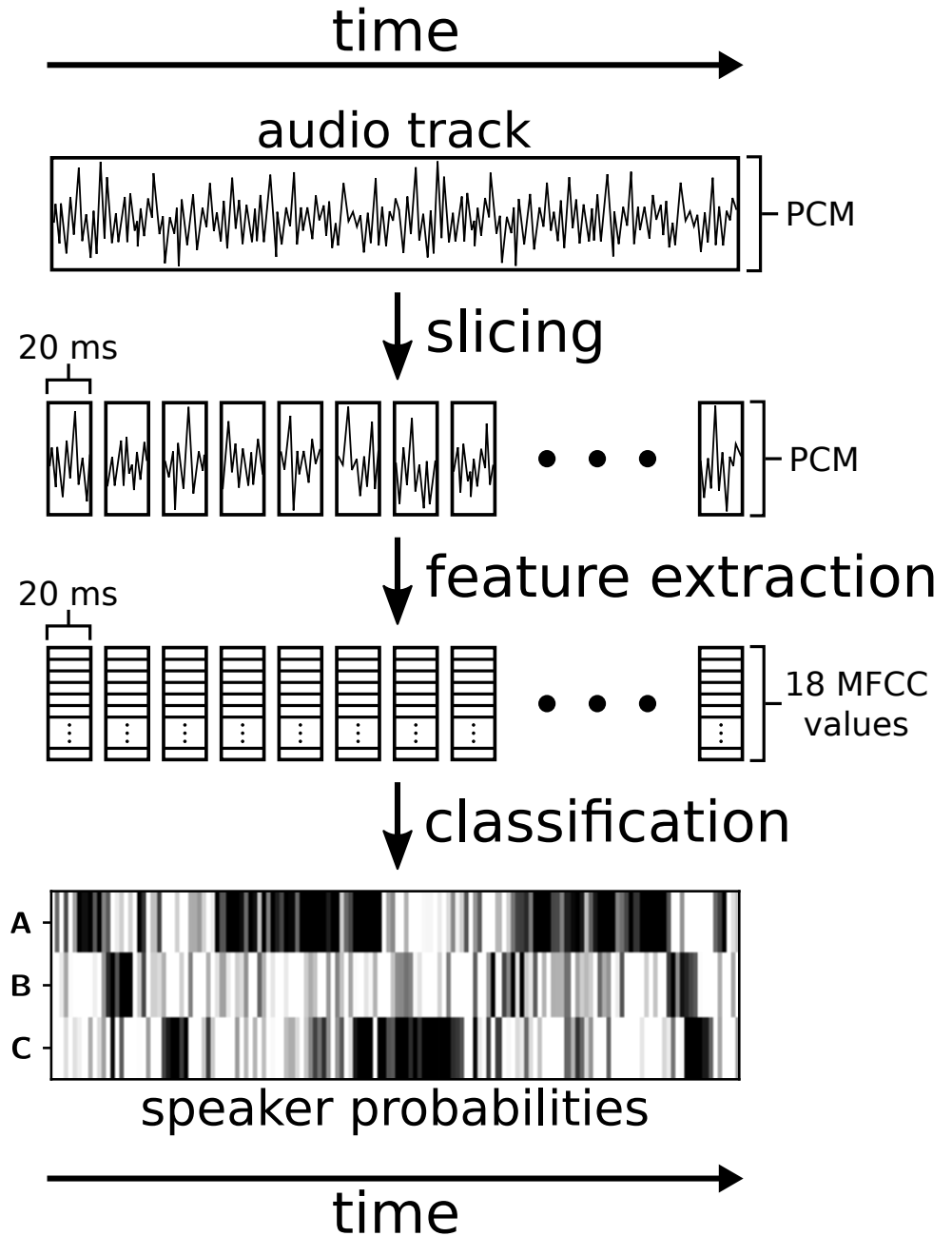


Figure 11: Pipeline for making speaker activity predictions on a given audio track.

model will try to fit its internal parameters so that it will be able to make useful predictions.

The classifier model, as shown in Figure 12, is a neural network with LSTM layers and a densely connected output layer. The MFCC feature frames are given as the input sequence for the first LSTM layer. During training we limit the sequence to 200 time steps. The number and dimensionality of the LSTM layers can be adjusted. The LSTM layers also have an adjustable dropout parameter.

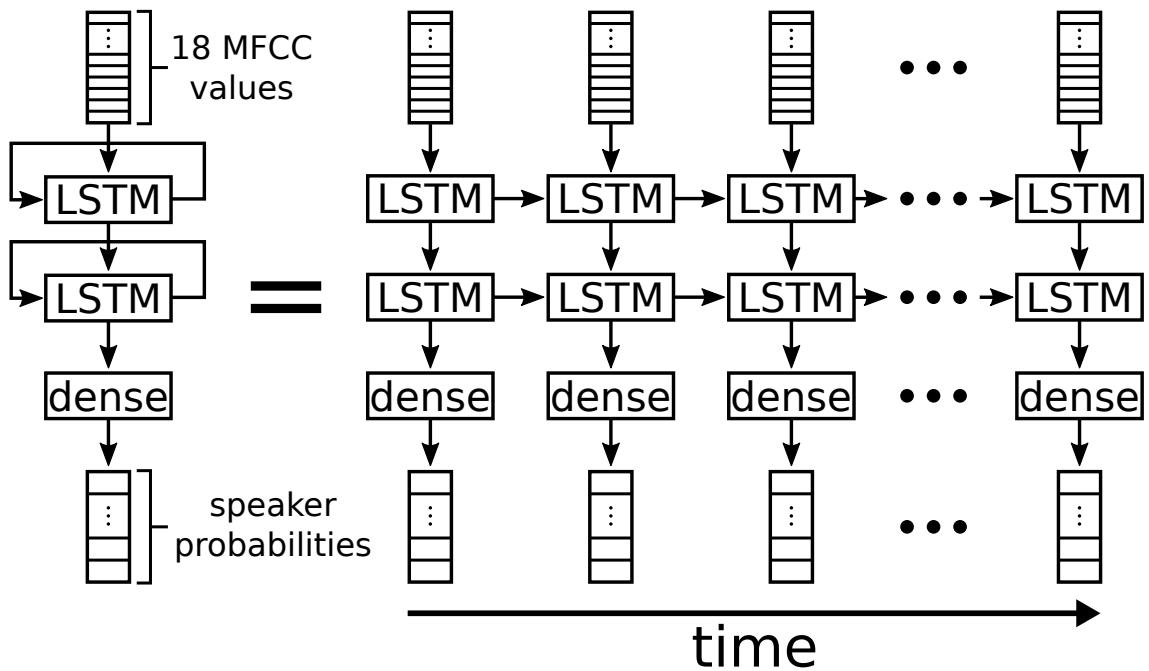


Figure 12: The structure of the classifier model with two LSTM layers.

The LSTM model can be given a bit more information by allowing it to see some time steps into the future. This can be done by delaying the output targets by a fixed number of time steps during training [23]. The output will be delayed by the same amount during evaluation, but the target delay does not need to be very long and so it does not greatly affect the ability of doing live prediction. If live operation is not required or if greater delays are acceptable, a bidirectional recurrent neural network (BRNN) [62] could be used to utilize future information more extensively. However, bidirectional recurrent neural networks are not studied in this thesis.

We have implemented the model in Python using Keras [11], an open-source neural network library. MFCC audio features are generated with `python_speech_features` library [42].

5 EVALUATION

In this chapter, we present the datasets and evaluation metrics we used and then analyze the performance of our model with different hyperparameter combinations. We try to select the best model using our validation dataset and then compare the results with a test dataset that was not used at all before the final tests.

5.1 Datasets

Our main dataset for model development is AMI Meeting Corpus [8]. It is a publicly available, Creative Commons licensed set of recorded meetings with multiple audio tracks along with various additional signals and annotations. The meetings were recorded in English by native and non-native speakers in acoustically different rooms using various kinds of microphones. Our interest is mainly in the high quality synchronized transcript that we can use to train and test speaker recognition models. The transcript annotations in the dataset are in an XML-based format that requires some preprocessing to fit our use case.

AMI Meeting Corpus consists of recording sessions. The sessions have unique identifiers, for example ES2016. Each session has a fixed set of participants and meetings they recorded. Typically a session has four speakers in four meetings. In the dataset each speaker have been assigned a unique identifier, for example MEO062 or FEE064. The first letter in the identifier corresponds to the gender of the speaker (M=male, F=female). These identifiers are shown in some of the result visualizations later. Each session has a new group of participants, but some participants have joined multiple sessions. Most of the meetings are simulated for the dataset, but with the aim to have natural, uncontrolled conversations. [8]

Each meeting in the dataset is recorded with multiple sets of different microphones. The dataset includes individual audio tracks for each speaker that could be used for speaker recognition based on the track volumes alone. However, they are not used in this thesis, because we are interested in a classifier that can distinguish different speakers from a single audio track. It is possible that the classifier would overfit to the properties of a specific microphone type. We want to test if our system can make good predictions without relying on the learned microphones and so we use a different set of microphones in the training phase. Training uses sound mixed from lapel microphones and evaluation uses mixed headset microphone tracks.

We train our speaker recognition model for each group of speakers separately. Out of the four meetings in a session, meetings B, C, and D are used for training and the re-

maintaining meeting A for validation. Figure 13 illustrates the structure of one recording session. Our goal is to find a model and hyperparameter values that perform well on all recording sessions with different groups of speakers. We have reserved sessions with even numbered identifiers for model selection and hyperparameter adjusting and the remaining ones are used only for testing.

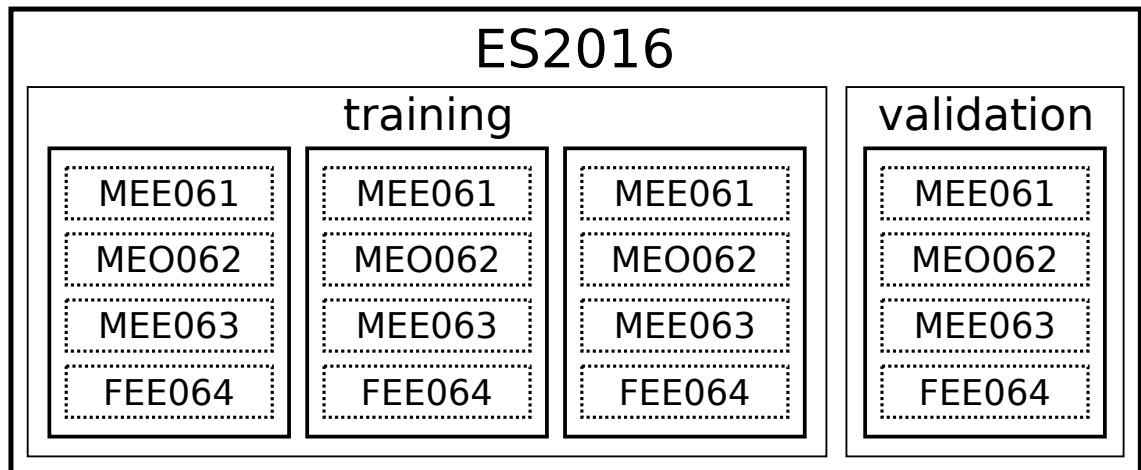


Figure 13: The structure of one recording session, which consists of four meetings. Each meeting has the same four participants.

5.2 Evaluation metrics

We want to analyze the performance of our models. This is needed for two purposes. First, we want to compare different models and hyperparameters so that we can select the best ones. This is done on the validation set. After selecting the best model, we want to analyze its performance on an independent test set. We need one or more evaluation metrics to produce numerical values that we can then use to assess the quality of the model and for comparing the models with each other.

One metric we can use is the loss metric that the neural network optimizer uses to train the network. In our multi-label case this metric is binary cross entropy loss as discussed in Chapter 3.4. When training the network, this loss is calculated on the training set. However, we can also calculate this metric on the validation set or the test set. This is a good basic metric that is especially useful for analyzing if the model starts to overfit.

In order to define the other metrics we use we first need to define some supporting terms. Table 3 has a confusion matrix where binary prediction outcomes are divided into four categories. The column is chosen based on the true label of the sample and the row based on the predicted label. True positives and true negatives represent correct predictions. [59]

	positive sample	negative sample
predicted positive	true positive (TP)	false positive (FP)
predicted negative	false negative (FN)	true negative (TN)

Table 3: Confusion matrix with four prediction outcome categories.

Based on these outcome categories we can now define three useful measures

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}},$$

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}},$$

$$\text{fall-out} = \frac{\text{FP}}{\text{FP} + \text{TN}},$$

where TP, FP, FN, and TN are the number of samples in each category. [59]

The area under the receiver operating characteristic curve (ROC AUC) is a metric that is commonly used for model comparison [26]. The ROC curve is calculated by varying the binary class discrimination threshold and plotting recall as a function of fall-out. An example of this curve is shown in Figure 14. The area under the curve summarizes the quality of a classifier as a single numerical value. The ROC AUC score has been criticized as being noisy and having some other issues when used in model comparison, but nevertheless still remains popular in machine learning. [25]

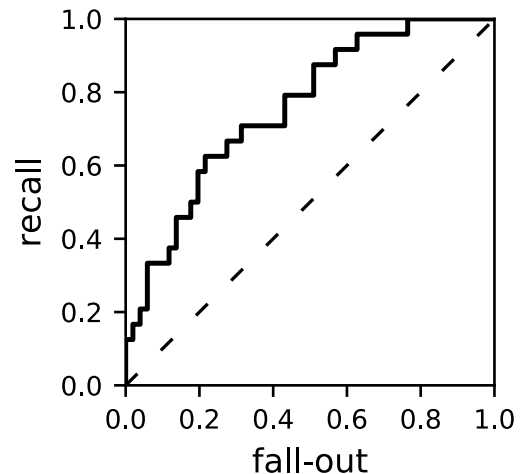


Figure 14: An example receiver operating characteristic curve. The dashed line represents random guessing.

We use both binary cross entropy loss and ROC AUC score to compare the models on the validation set. When evaluating the best model on the test set, we add one more metric. Precision and recall are good measures that we could use as a pair. Our final metric is F_1 score, which combines precision and recall to one value by taking their harmonic mean [59]:

$$F_1 = \left(\frac{\text{precision}^{-1} + \text{recall}^{-1}}{2} \right)^{-1}$$

To extend ROC AUC score and F_1 score for multi-label classification we need to decide how the results of each label are combined. There are two common possibilities for this. We could add the total number of true positives, false positives,

false negatives, and true negatives together from all of the labels and then use these values to calculate the metrics. Doing it this way is called taking the micro average. The alternative is to calculate the macro average, where the metrics are first calculated for each label separately and then the mean of those metrics becomes the joined metric. Micro and macro averages calculate slightly different things and we cannot say that one is necessarily better than the other. In this thesis we use macro averages for the metrics. [71]

5.3 Hyperparameters

Our classifier model is not remarkably deep or otherwise complex. Yet there are still many tweakable hyperparameters that have significant impact on the classification result. The main hyperparameters are the number of stacked LSTM layers and the size of those layers. As discussed in Chapter 3, larger networks are prone to overfitting. To combat this, our third major hyperparameter is dropout. In this part of the thesis, we are evaluating the different values for these hyperparameters to see which combinations works the best. The hyperparameters to be tested and the values we are going to choose from are listed in Table 4.

hyperparameter	examined values
number of LSTM layers	1, 2
size of LSTM layers	16, 64, 256, 1024
dropout	0 % (disabled), 50 %

Table 4: Examined hyperparameters and their values.

The number of training steps can also be seen as a hyperparameter for the model. It is expected that good classification results require a certain amount of training steps. In an ideal case the classification accuracy would improve over time and converge to some level. However, mainly due to overfitting, the accuracy may start to decrease at some point with more training. We are trying to analyze these behaviors with learning curves where classification accuracy is plotted as a function of training epochs.

5.4 Results

First we examine a model with one LSTM layer and no dropout. The layer size is varied to see how it affects the results. In Figure 15, there are two loss curves for each tested LSTM layer size. The dashed lines represent training losses and the solid lines validation losses. Initially, the losses decrease rapidly, but after around epoch five the validation loss for all models except the one with layer size 16 start to increase while training losses continue to decrease. This most likely indicates

overfitting. Figure 16 has a ROC AUC score curve for each tested layer size. We can see that scores increase on the very first epochs. On layer size 16 the score continues to increase slowly, but the scores of the models with greater layer sizes start to decrease.

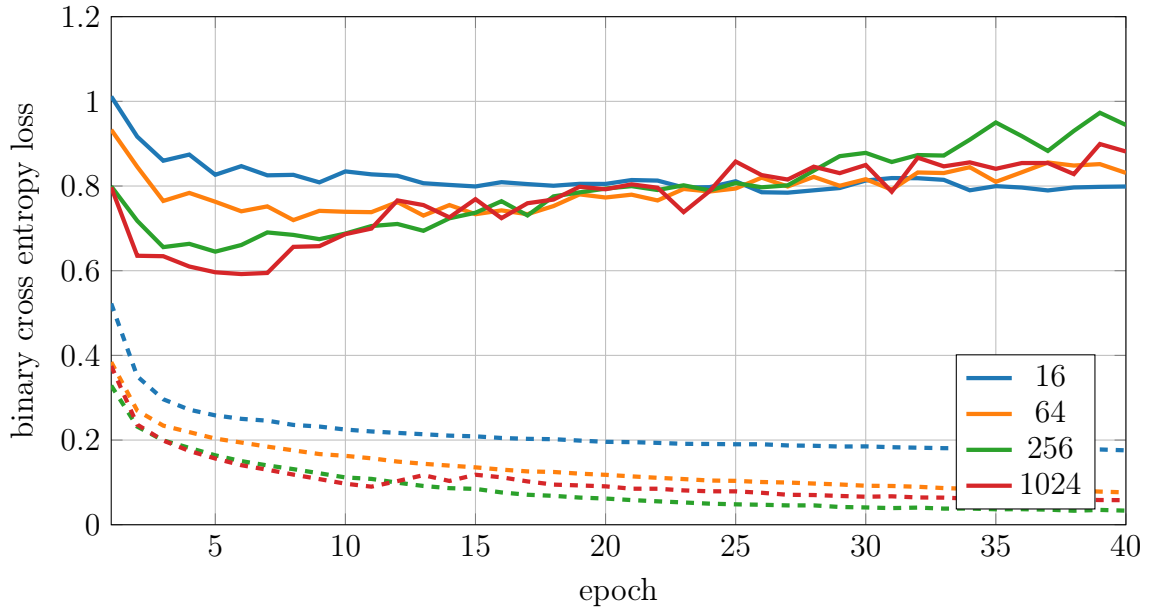


Figure 15: The effect of layer size on training loss (dashed lines) and validation loss (solid lines).

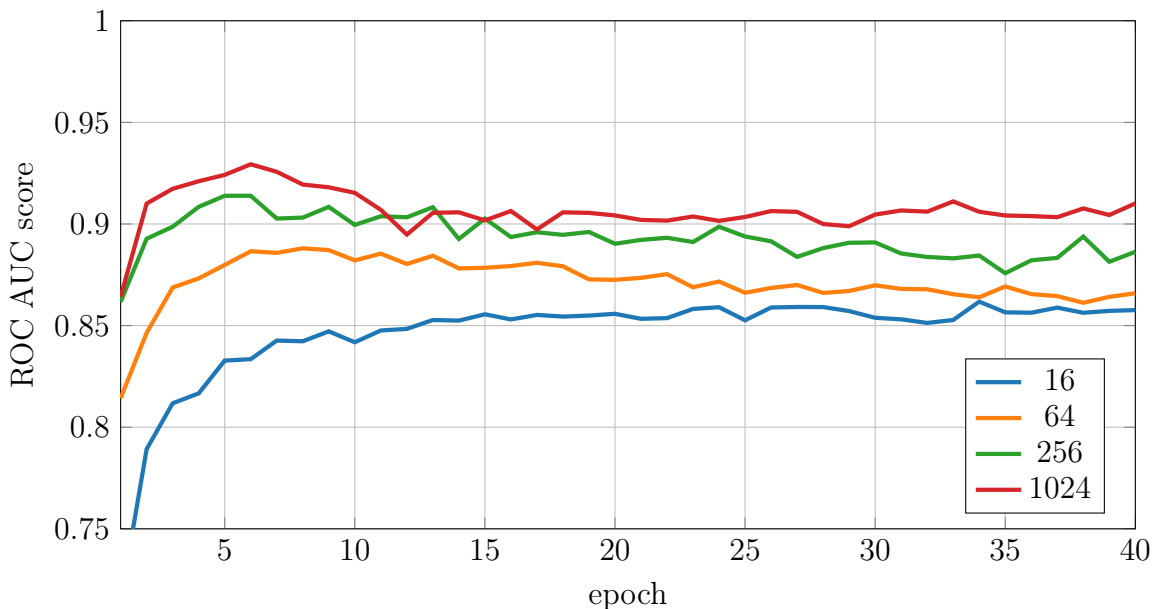


Figure 16: The effect of layer size on ROC AUC score.

Next, we test how adding dropout changes the results. Dropout may reduce overfitting allowing us to use larger network sizes. Figure 17 has learning curves with training and validation losses for models with layer size 256 and dropout both enabled and disabled. Shaded regions are one standard deviation error bands for the

variance between meetings. Enabling dropout increased training loss. As discussed earlier in Chapter 3, this is only an artifact caused by dropout itself. Dropout makes fitting to training data more difficult by design. The model with dropout has lower validation loss on all epochs compared to the model without dropout. The curve with dropout is converging and has less variance. In contrast, the validation loss for the model without dropout has significantly more variance and starts to increase after around five epochs of training showing overfitting. At this layer size, these observations would support selecting the model with dropout over the one without. We can further see the difference in Figure 18 where ROC AUC score is compared.

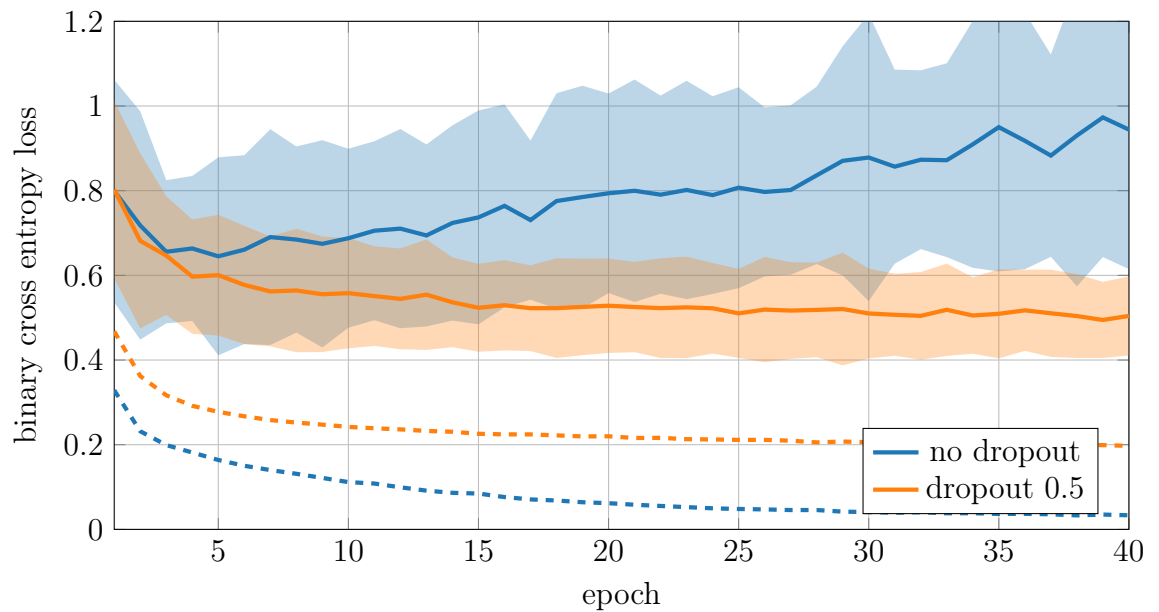


Figure 17: The effect of dropout on training and validation losses (LSTM, layer size 256).

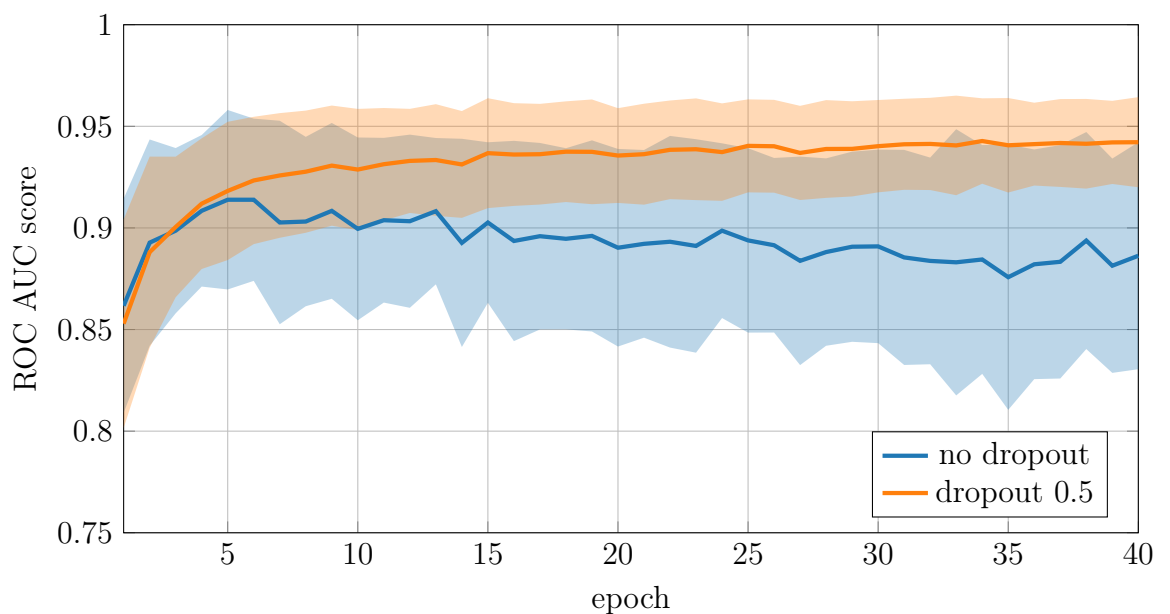


Figure 18: The effect of dropout on ROC AUC score (LSTM, layer size 256).

To find the best hyperparameter combination, we now test more layer sizes with dropout enabled. Figure 19 shows how varying the layer size affects training and validation losses when dropout is enabled. With dropout, the loss curves do not start to increase anymore like in Figure 15. ROC AUC scores with dropout enabled in Figure 20 are converging and reaching better values than without dropout in Figure 16 where most of the scores started to decrease. It can be seen that dropout is useful for reducing overfitting and making the models more stable.

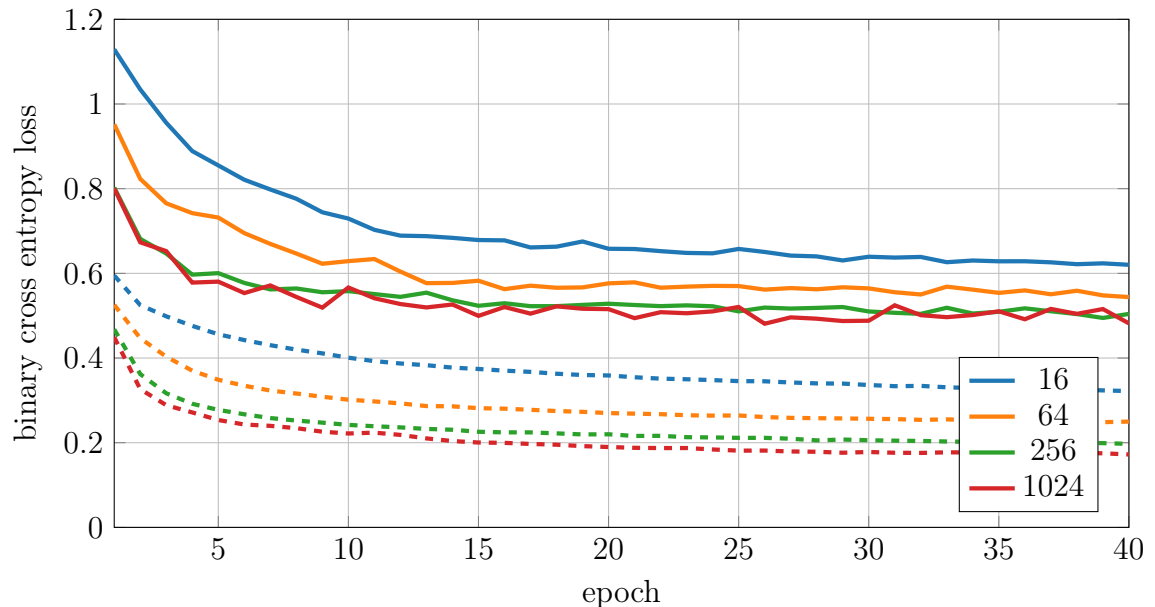


Figure 19: The effect of layer size on training and validation loss with dropout enabled.

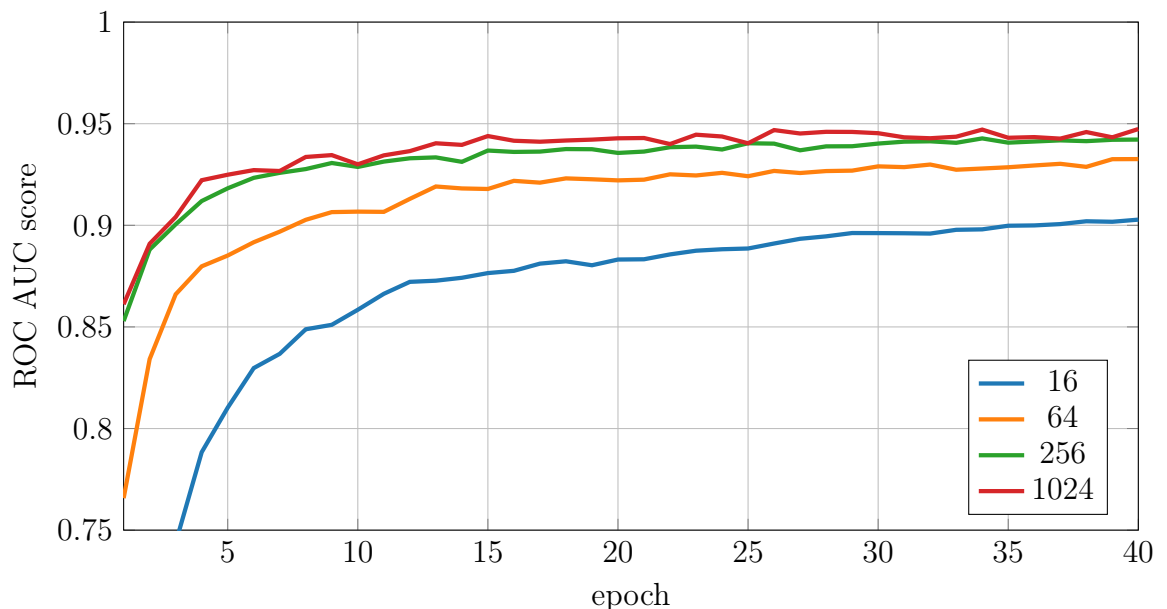


Figure 20: The effect of layer size on ROC AUC score with dropout enabled.

Among the tested layer sizes it would appear that 256 and 1024 are the best models on both validation loss and ROC AUC score. Layer size 1024 may have marginally better results, but the difference is very small. Layer size 64 is not far behind, but

size 16 does not seem to be large enough to fully capture the underlying phenomena. However, even with the layer size 16 the loss and score values are better with dropout than without. This can be seen in Figures 21 and 22 where the model without dropout is compared to models with dropout enabled. A curve for the model with layer size 256 is also included in these figures for reference. At layer size 16, the model learns initially a bit faster without dropout but the model with dropout enabled will reach and exceed the performance with further training.

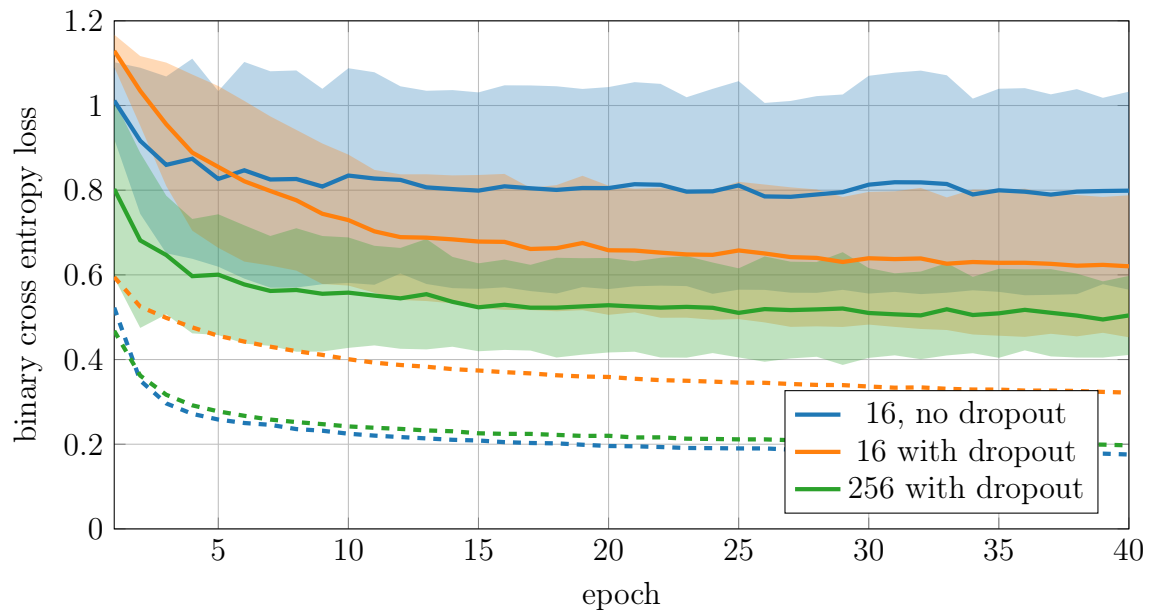


Figure 21: Training and validation losses comparison of some models with dropout enabled and disabled.

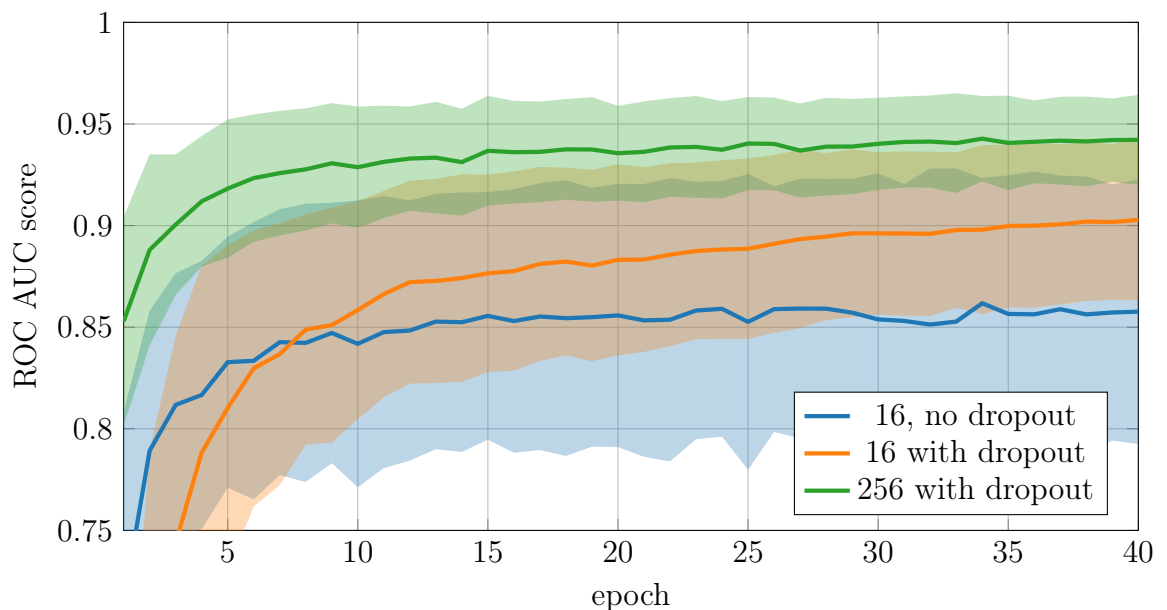


Figure 22: ROC AUC score comparison of some models with dropout enabled and disabled.

Plotting learning curves with the loss metrics and ROC AUC scores is a good way to visualize how the classifier models learn. However, we can also visualize the predictions made by our system. Figures 23, 24, and 25 have two rows for each speaker in the meeting. The first row has ground truth labels showing when the speaker was actually active. The second row has the predictions that the classifier produced. Differences between the rows are mistakes in the prediction. The x-axis is time in the audio track starting from the beginning of the meeting. The meetings are longer, but for visualization reasons these figures are limited to the first ten minutes.

The predictions in Figure 23 are made after only one epoch of training. At that stage the classifier is still underfitted so the predictions are mostly uncertain and include lots of random noise. Figure 24 has the same classifier and configuration, but after 40 epochs of training. The prediction signal is now significantly stronger. There are still some incorrect predictions, but most of the speaking segments are correctly identified. For comparison, Figure 25 has the predictions for the same meeting made by a classifier model with layer size 16 and dropout disabled. This figure looks similar, but the predictions, especially for MEO062, have more mistakes compared to Figure 24.

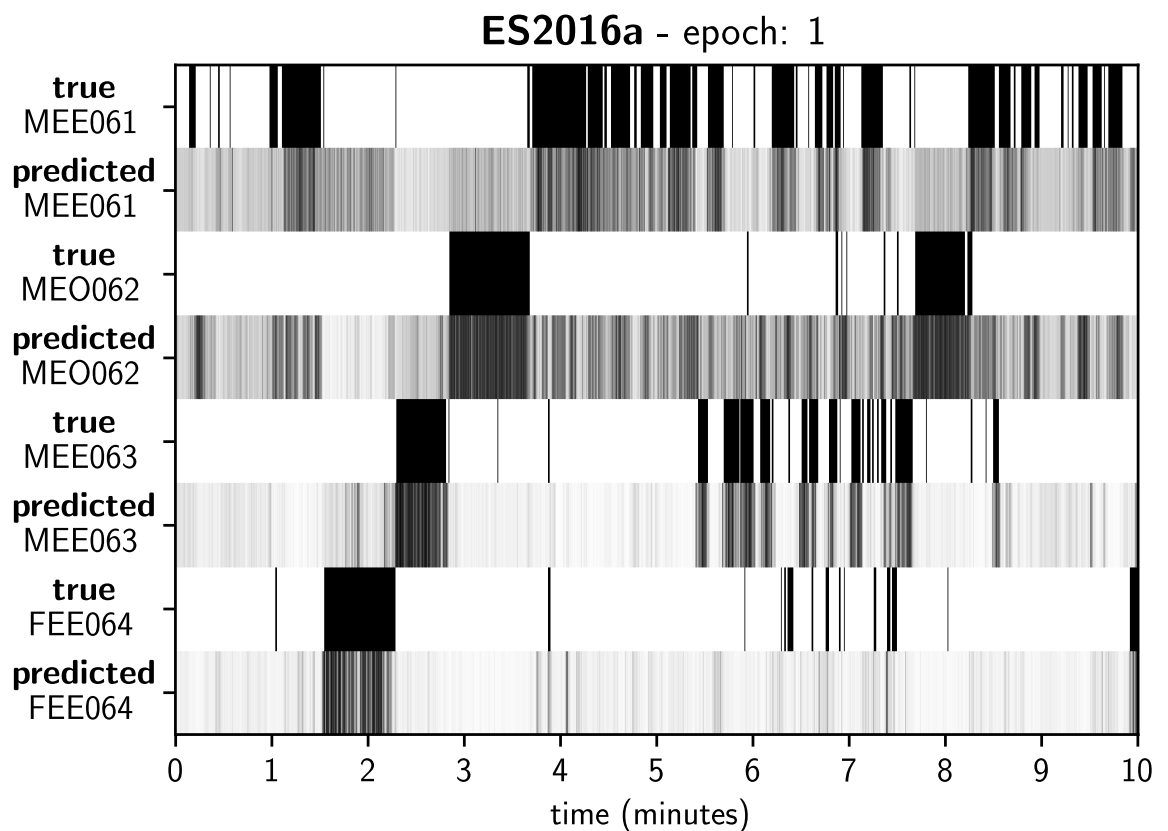


Figure 23: Predictions compared to ground truth labels after only one epoch of training (layer size 256, dropout).

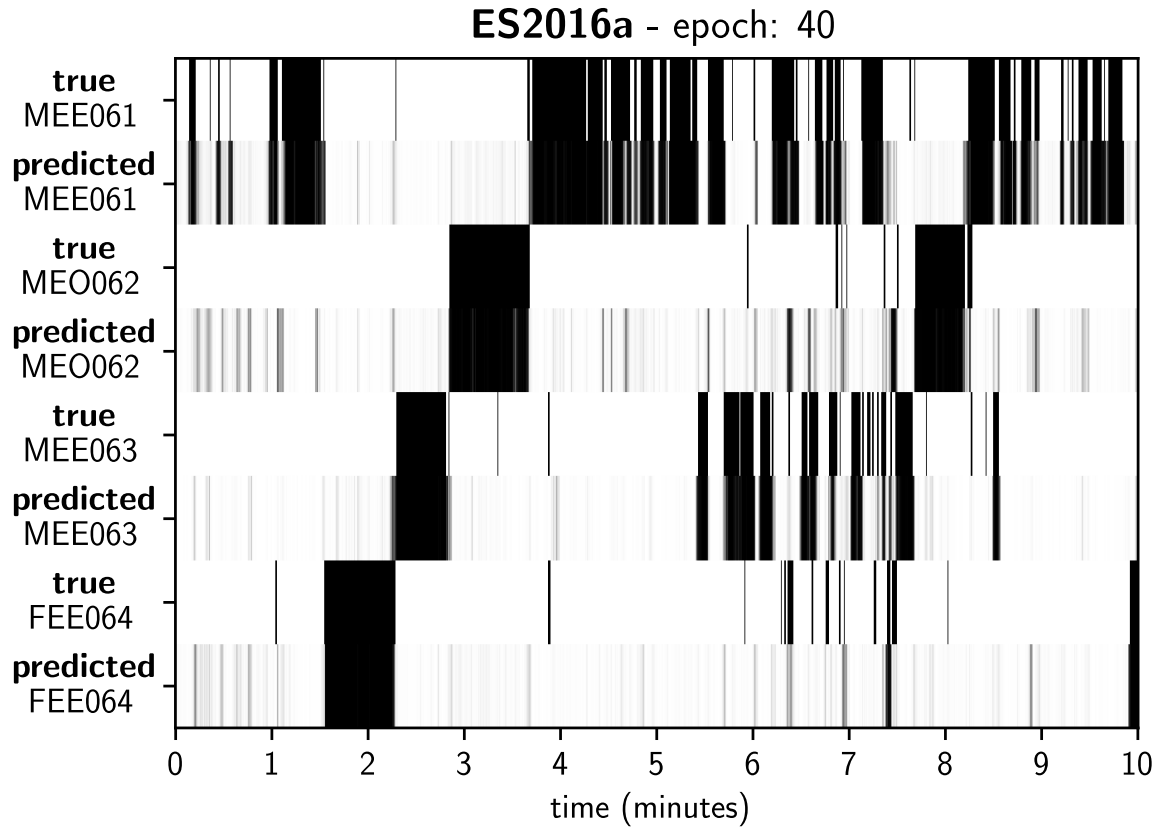


Figure 24: Predictions after 40 epochs of training (layer size 256, dropout).

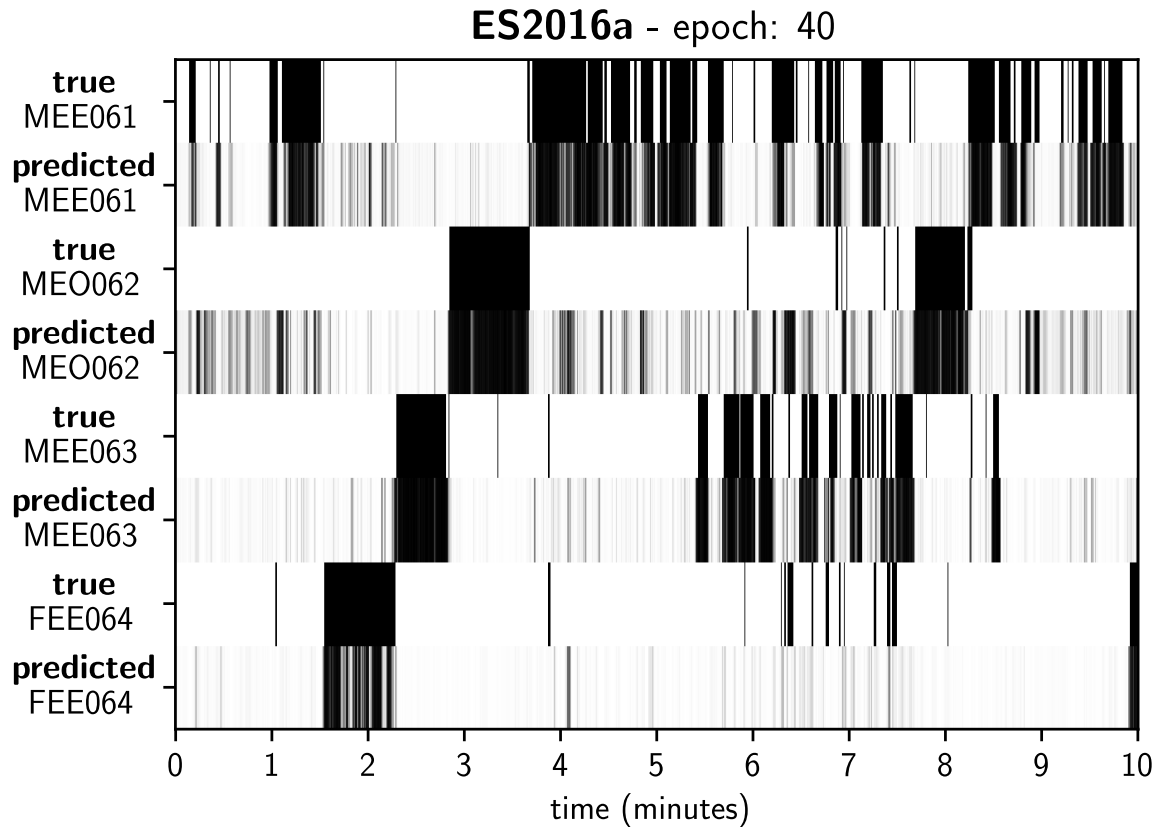


Figure 25: Predictions made by a model with layer size 16 and no dropout.

Stacking another LSTM layer to the network model is a way to increase the learning capacity. Figures 26 and 27 compare 1-layer and 2-layer models with layer size 256 and dropout enabled. The 2-layer model performs almost exactly the same as the 1-layer model with the same parameters. The simpler 1-layer model learns a bit faster, but with more training the 2-layer model reaches the same results. The 2-layer model may be marginally better, but the difference between the results is not clear. It seems that the 1-layer model at this layer size is already complex enough to accurately fit the problem.

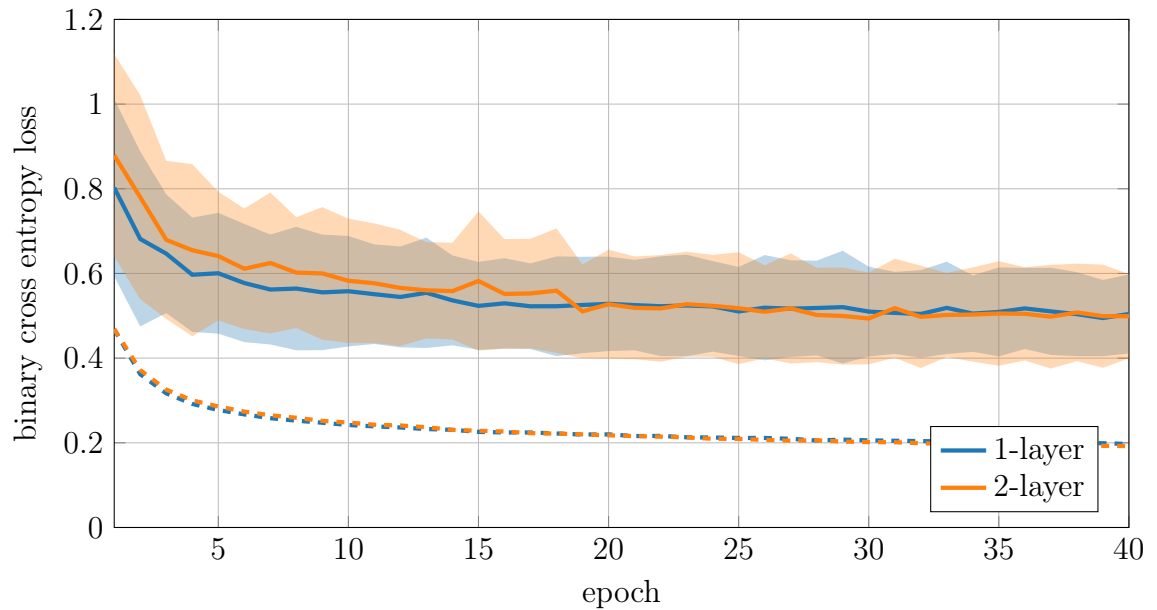


Figure 26: Loss comparison of 1-layer and 2-layer models with layer size 256 and dropout.

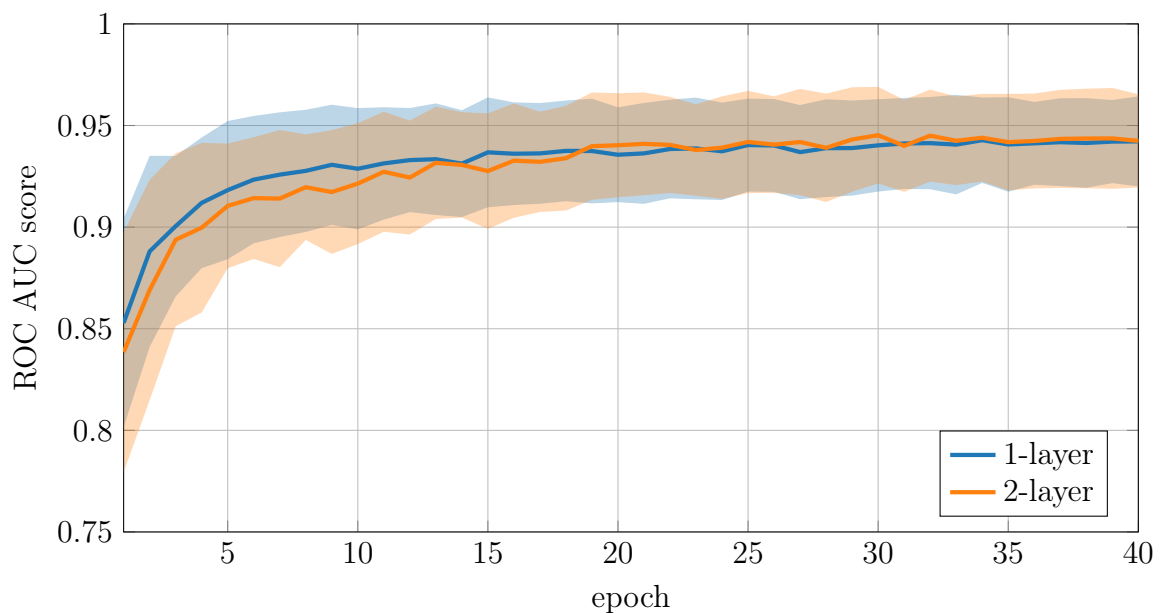


Figure 27: ROC AUC comparison of 1-layer and 2-layer models with layer size 256 and dropout.

In order to determine if adding a second hidden layer increases the learning capacity we compare 1-layer and 2-layer models with layer size 16 and dropout disabled. The loss and ROC AUC score learning curves are shown in Figures 28 and 29. It can be seen that the 2-layer model performs better with these parameters. Again, the 1-layer model learns initially faster, but the 2-layer model achieves better results later.

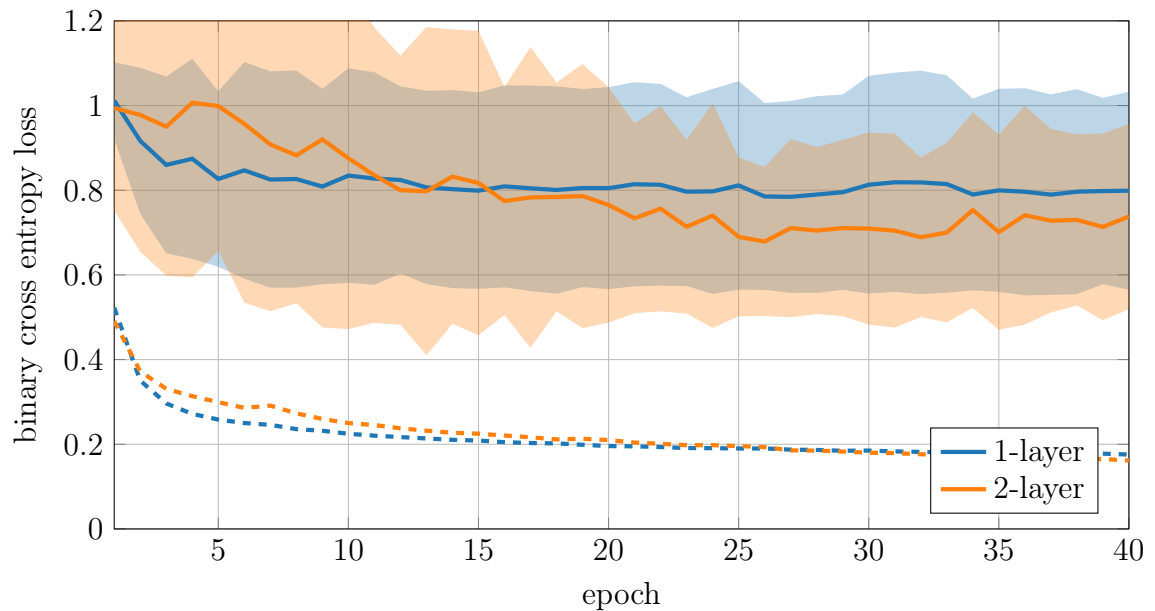


Figure 28: Loss comparison of 1-layer and 2-layer models with layer size 16.

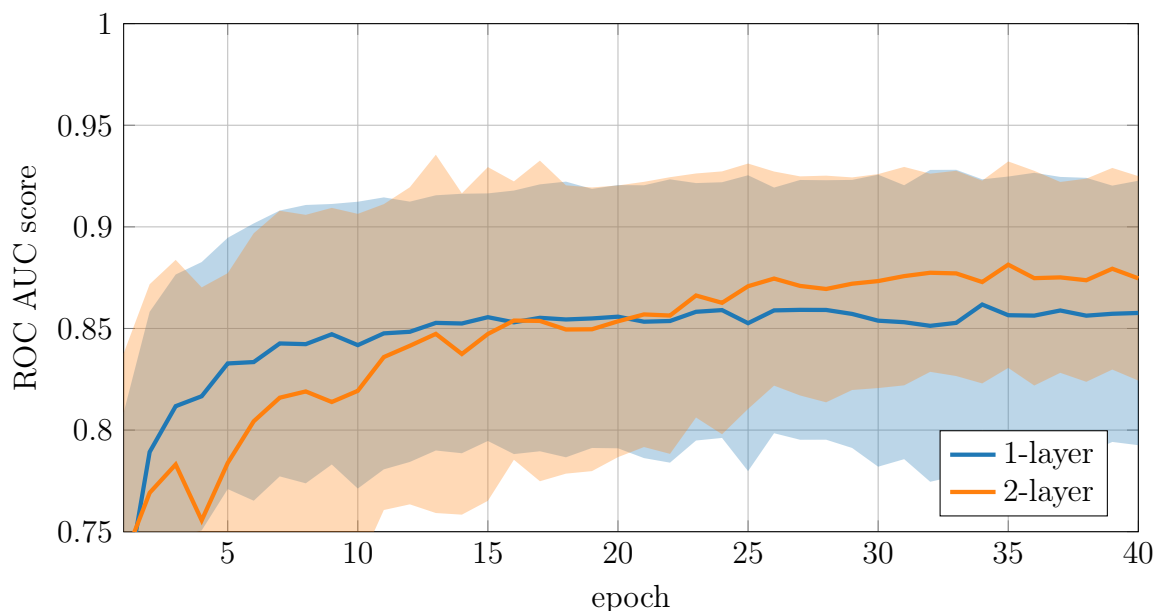


Figure 29: ROC AUC comparison of 1-layer and 2-layer models with layer size 16.

To further compare the parameter combinations, we test 2-layer models with each layer size. The learning curves for each 2-layer model with different layer sizes are shown in Figure 30 and 31. These curves look very similar to the 1-layer model curves in Figures 19 and 20. Layer sizes 256 and 1024 are still the best ones. It

seems that the 2-layer models take more training time to reach the same results and there is no visible benefit in having 2 LSTM layers.

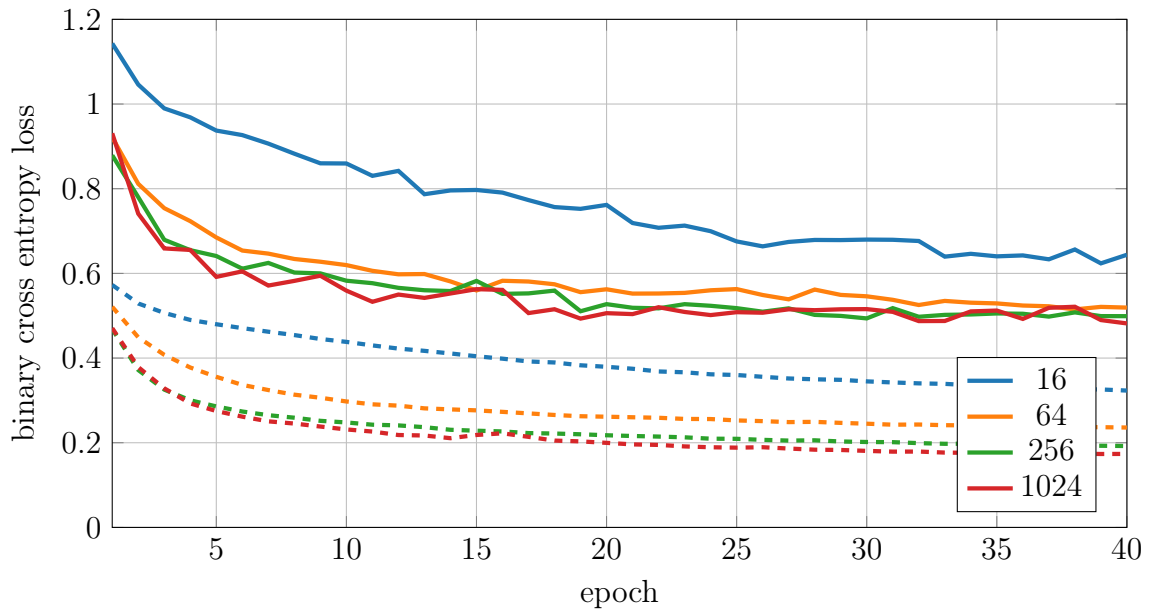


Figure 30: The effect of layer size on training and validation loss with 2-layer models and dropout enabled.

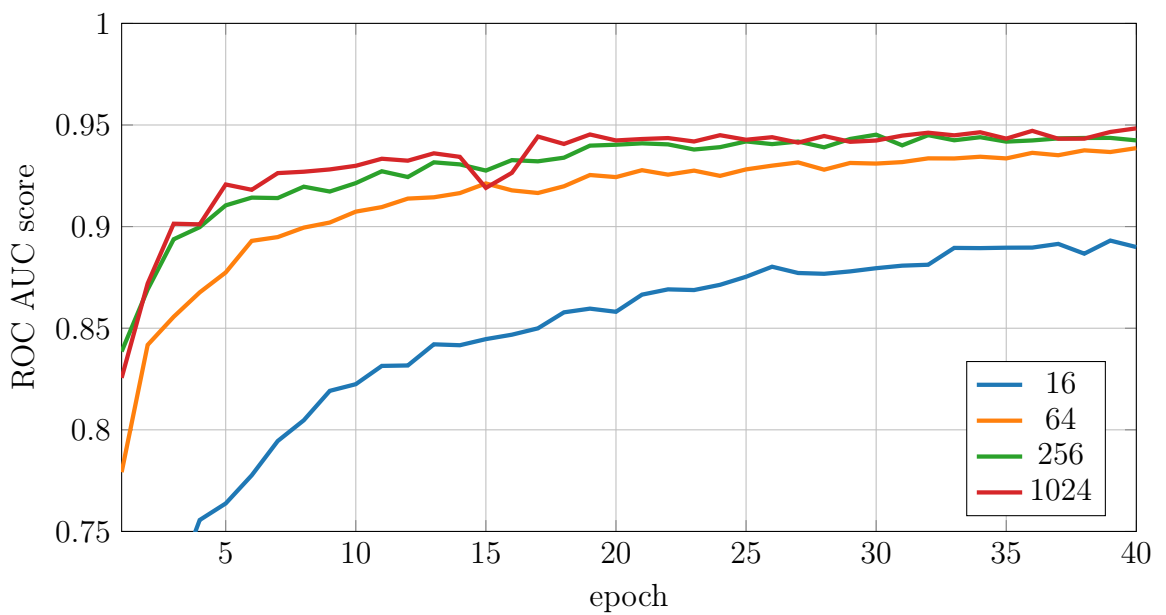


Figure 31: The effect of layer size on ROC AUC score with 2-layer models and dropout enabled.

Different models have different computational costs. Table 5 shows how training speed differs for our models. The training speed depends on the underlying hardware setup. In our case, the measurements were done on a system with Intel Core i5-4670K processor and GeForce GTX 1060 (6GB) graphics processing unit. We can see that models with layer size 1024 and all 2-layer models take significantly more

time for training than the simpler models. The 2-layer model with layer size 1024 takes about ten times the training time compared to the smaller 1-layer models. This gives us reason to select the 1-layer model with layer size 256 since the more complex models did not achieve better accuracy results.

training time per epoch (s)		
	1 layer	2 layers
layer size 16	35.1 ± 0.78	73.8 ± 3.50
layer size 64	36.7 ± 1.94	72.1 ± 2.68
layer size 256	35.0 ± 1.01	74.0 ± 1.84
layer size 1024	121.9 ± 5.22	340.1 ± 13.89

Table 5: Training time for different models (dropout enabled).

Up to this point we have used a validation dataset that combines multiple recording sessions. Using more data is good for achieving generalization as it is not desired to optimize the system for any single case. The meetings in the recording sessions may have differences that affect the classification performance. If we examine predictions made for some specific session, the results may be deviated from the averages shown before. The variance in the results in different sessions can be analyzed by calculating the metrics for each of them separately. This is done in Figure 32 where ROC AUC scores are shown separately for five sessions. The results on some sessions are significantly better compared to others. We can say that, for our classification models, the voices or the speaking patterns in IS1008 seem to be easier to classify than in IS1000. The differences can also be seen by comparing the labels and predictions for meetings ES2016a (Figure 24), IS1008a (Figure 33), and IS1000a (Figure 34). No further analysis on the actual differences in the data is done in this thesis.

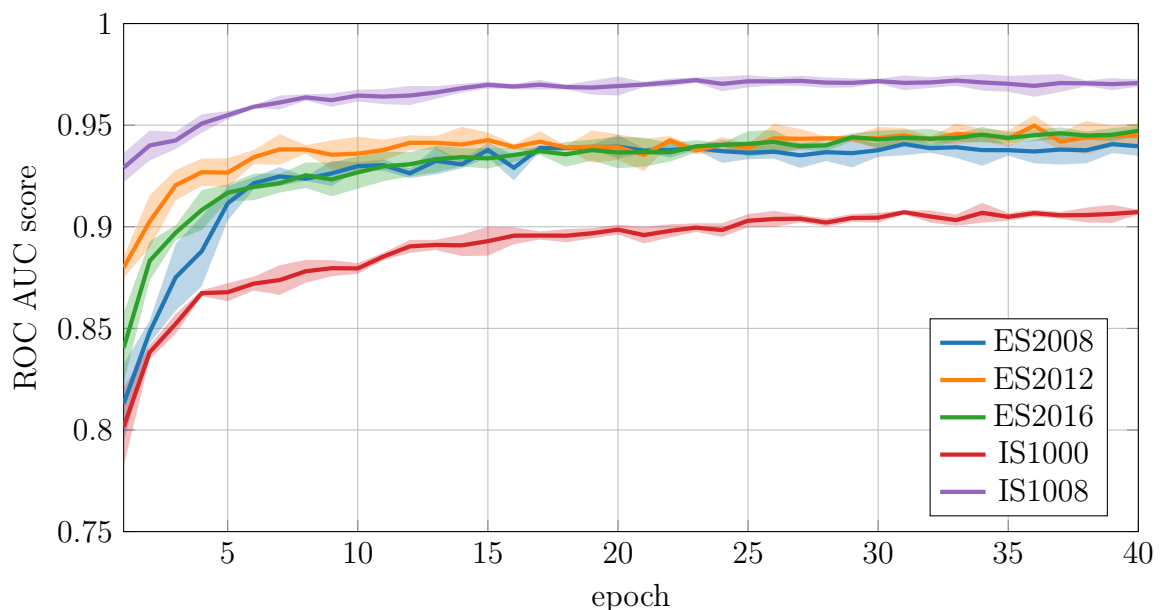


Figure 32: ROC AUC score for individual validation set sessions (LSTM, layer size 256, dropout).

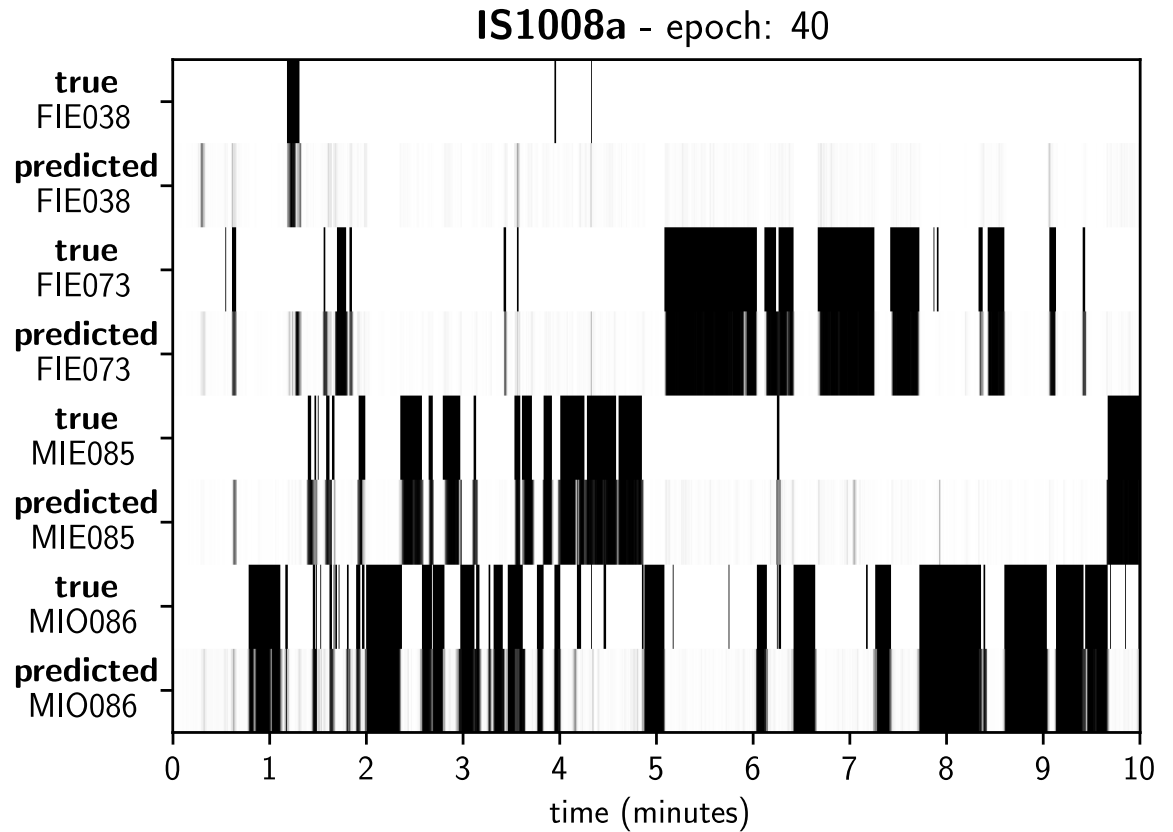


Figure 33: Predictions for meeting IS1008a (layer size 256, dropout).

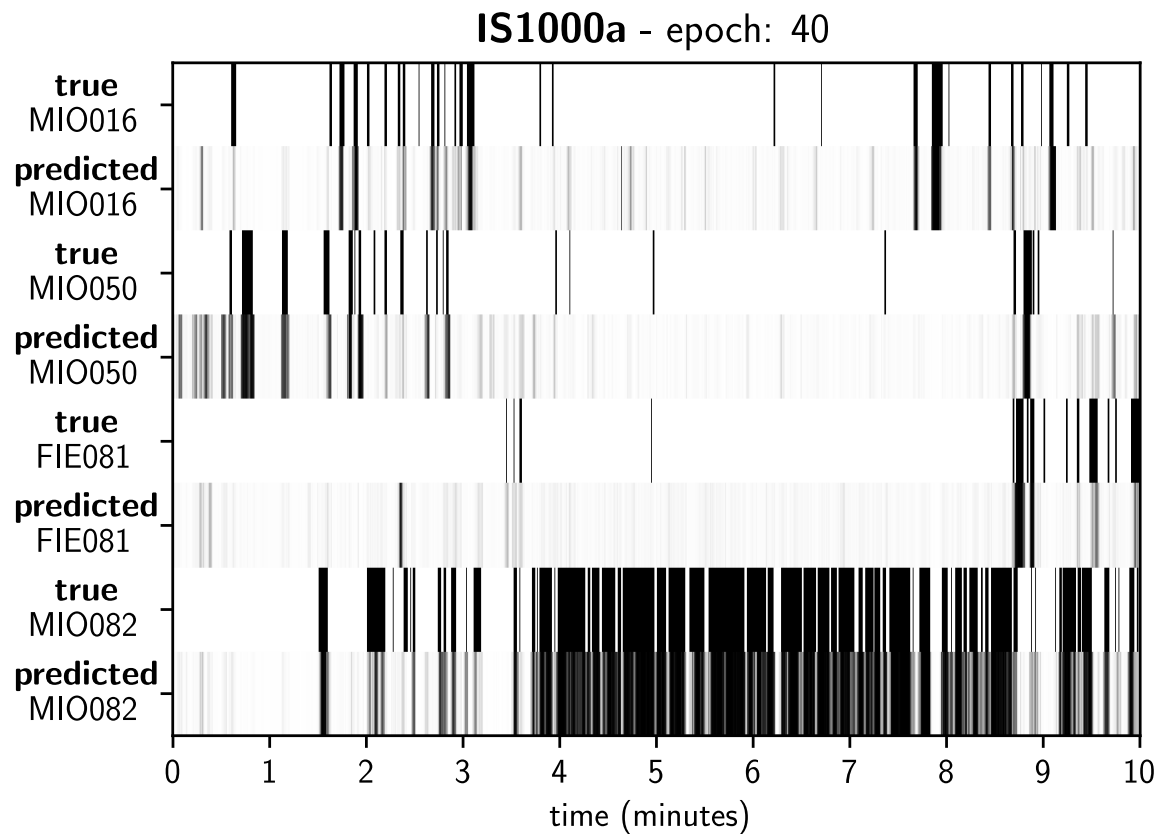


Figure 34: Predictions for meeting IS1000a (layer size 256, dropout).

As we can see, there is variation between the different recording sessions, but when training the network again there is variation even if the exactly same meetings are used. This is because the neural network optimization is a random process. This can be seen in the shaded standard deviation error bands in Figure 32. In previous figures with these shaded error bands, the variance comes mainly from the difference between the recording sessions. Unlike them, the only source of variation in the figure with the individual sessions is the randomness of the training process. However, there is only a small amount of variation between the training instances.

So far we have used our validation dataset for model selection. While the validation data was not available to the network optimizer itself, it is possible that by choosing the best model structure and hyperparameters we have overfitted to the validation data. We can compare the results with a new, unused test set. The best model chosen for testing is a 1-layer LSTM network with layer size 256 and dropout enabled. Table 6 has evaluation results for both the validation set and test set. Test set results can be also seen in Figure 35, which has learning curves for the individual sessions in our test set. This figure can be compared with Figure 32, which has similar learning curves for the validation set sessions.

	cross entropy loss	ROC AUC score	F ₁ score
validation set	0.510 ± 0.2053	0.942 ± 0.0426	0.762 ± 0.1588
test set	0.843 ± 0.2986	0.895 ± 0.0168	0.629 ± 0.0444

Table 6: Comparison of validation and test results (layer size 256, dropout, epoch 40).

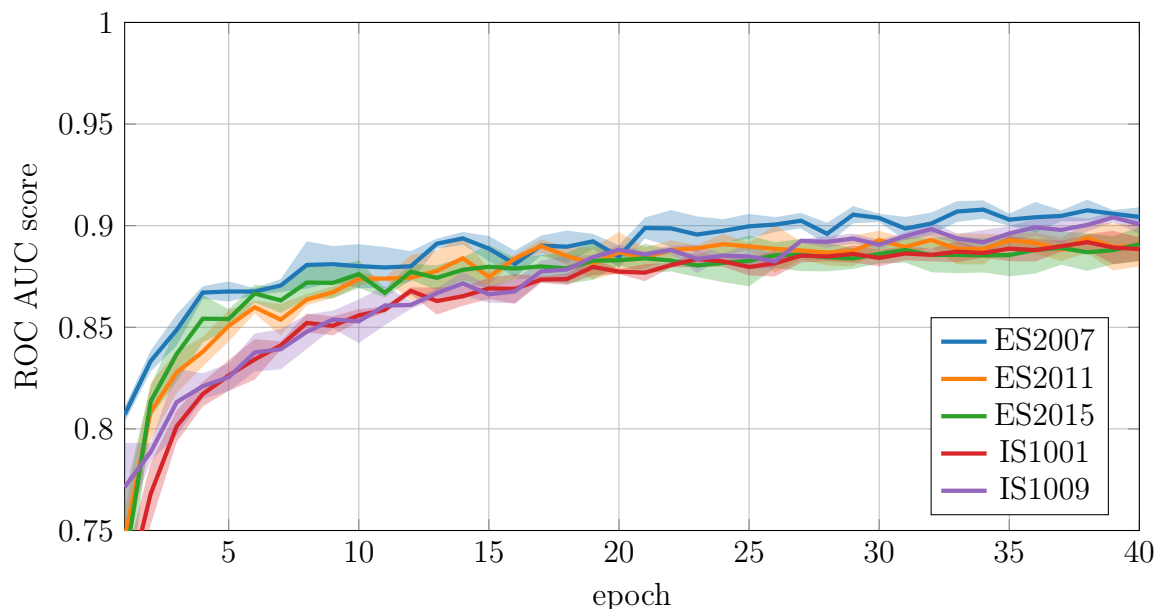


Figure 35: ROC AUC score for individual test set sessions (LSTM, layer size 256, dropout).

We can see that there are noticeable differences in the results. The test set results do not reach the same level of accuracy as the validation results. The difference in the results would suggest that in our model selection process we chose the model and hyperparameters that have the best results on the validation data, but that do not necessarily generalize to new datasets. However, the meetings in the validation set and the test set are completely different and it is possible that the differences in the results are caused by the fact that the validation set happens to have meetings that are easier to classify. In any case, we argue that even the test set results are good showing that our speaker recognition system can make valid predictions.

We can also compare how the predictions made for the test set look. Figure 36 shows predictions made for meeting IS1009a in our test set. The meeting is different so we cannot directly compare it with the meetings showed earlier in Figures 24, 33, and 34, but it can be seen that the predictions are mostly accurate.

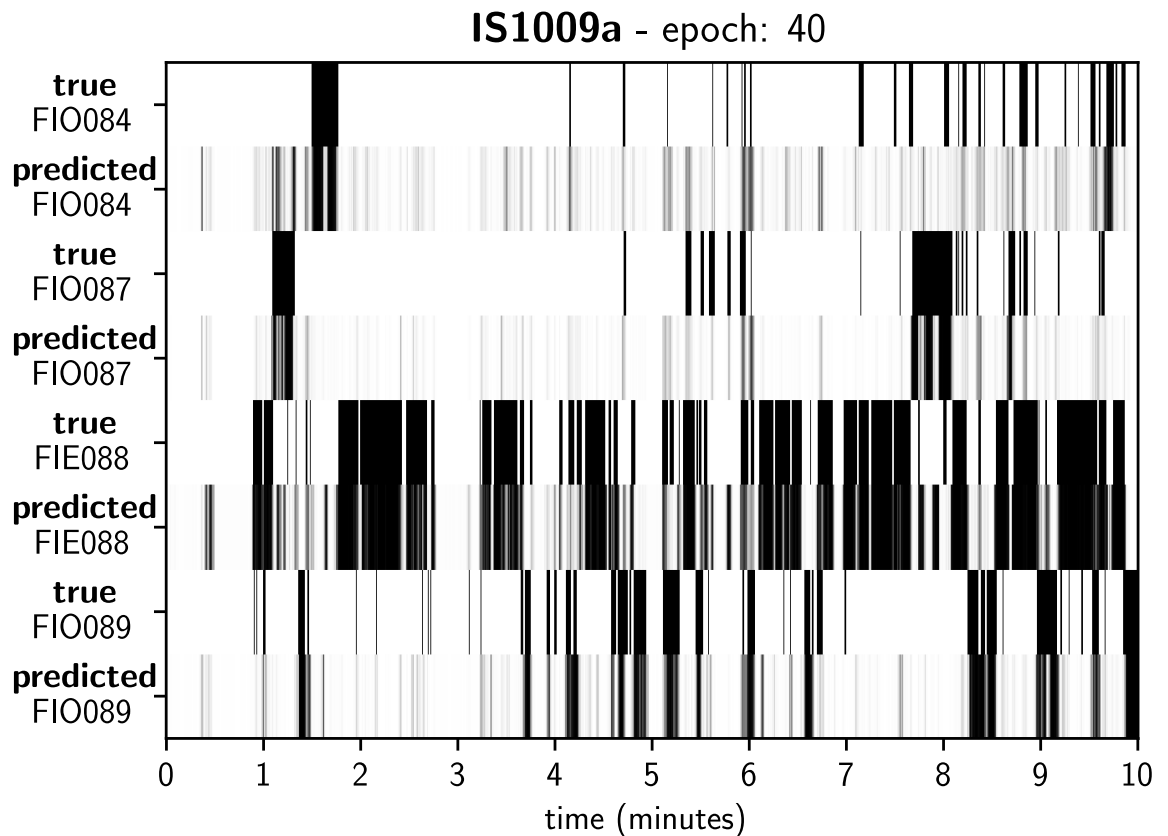


Figure 36: Predictions for a test set meeting IS1009a (layer size 256, dropout).

6 CONCLUSION

This thesis introduced the reader to the theory that is needed to understand our experimental speaker recognition system. We developed a recurrent neural network classifier that successfully performs multi-label speaker recognition on AMI Meeting Corpus dataset. When given a meeting recording the classifier determines when each of the participants in the meeting is speaking. The architecture that we chose can be used for live prediction and it can identify multiple speakers who are active simultaneously. We can see that a rather simple long short-term memory (LSTM) network using MFCC features is sufficient for this task. We compared different hyperparameters and noticed that by choosing good values for them can improve the classification results and reduce computational requirements. The best model was a 1-layer LSTM network with layer size 256. More complex models significantly increase the computation cost, but do not improve the results.

One source of uncertainty in the measured results is the inaccuracies in the speaker label annotations. AMI Meeting Corpus was chosen because the annotations seem to be very good, but we cannot expect them to be absolutely correct. We have also limited the speaker label time resolution to 20 milliseconds.

The main limitation of our implementation is that it only recognizes known speakers that were present in the training material. This is caused by our supervised classification architecture. Unknown speakers could be supported with unsupervised learning that only clusters similar utterances together without actually identifying the speakers. However, the supervised approach may be better if the goal is to only recognize a few designated speakers for whom we have speech samples. Further research could be done on determining the required amount of training samples per speaker and on minimizing that amount.

We used only the AMI Meeting Corpus dataset in our experiments. Further testing could be done with other datasets that do not need to be limited to meeting recordings. In addition, existing data could be augmented to increase the amount of training samples. Possible augmentations include adding noise, distortions, or additional background sounds to the audio track. Artificial meeting scenarios could be created by mixing and matching excerpts from the available material. This could include combining various speakers from different meetings and possibly creating new meetings that have more participants than any meeting in the source material.

Our classifier works on audio streams that would allow real-time speaker recognition on a live audio source, but we did not cover this aspect. Our hypothesis is that real-time evaluation would work without extensive modification to the system. If

evaluation performance turns out to be a problem, gated recurrent units (GRU) may be able to replace our long short-term memory (LSTM) layers using much less computation time as for example Khandelwal et al. found out in their speech recognition experiments [34].

Alternative classification layers, like the gated recurrent units, could also outperform our system in terms of accuracy. Other methods including bidirectional recurrent neural networks (BRNN) and convolutional neural networks (CNN) [37] could also be tested. In addition, attention-based recurrent neural networks have been a popular topic in recent years [10].

Convolutional neural networks might also replace our MFCC based feature extraction layer. Lukic et al. studied this arguing that the commonly used MFCC feature extraction methods do not utilize all the available speaker information in the audio [41]. They also researched speaker recognition as a clustering problem on the latent space of their CNN, which makes it possible to segment speech from unknown speakers.

References

- [1] X. Anguera, S. Bozonnet, N. Evans, C. Fredouille, G. Friedland and O. Vinyals. Speaker Diarization: A Review of Recent Research. *IEEE Transactions on Audio, Speech, and Language Processing* vol. 20, no. 2. 2012. pp. 356–370.
- [2] L.E. Baum and T. Petrie. Statistical Inference for Probabilistic Functions of Finite State Markov Chains. *The Annals of Mathematical Statistics* vol. 37, no. 6. 1966. pp. 1554–1563.
- [3] C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer Science. 2006.
- [4] H.S. Black, and J.O. Edson. Pulse Code Modulation. *Transactions of the American Institute of Electrical Engineers* vol. 66, no. 1. 1947. pp. 895–899.
- [5] L.E. Boucheron and P.L. De Leon. On the Inversion of Mel-Frequency Cepstral Coefficients for Speech Enhancement Applications. *IEEE Signals and Electronic Systems. ICSES'08*. 2008. pp. 485–488.
- [6] K. Burnham and D.R. Anderson. *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. Springer Science & Business Media. 2003.
- [7] J.P. Campbell. Speaker Recognition: A Tutorial. *Proceedings of the IEEE* vol. 85, no. 9. 1997. pp. 1437–1462.
- [8] J. Carletta. Announcing the AMI Meeting Corpus. *The ELRA Newsletter* vol. 11, no. 1, January-March. 2006. pp. 3–5.
- [9] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio. Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation. *arXiv:1406.1078*. 2014.
- [10] J.K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio. Attention-Based Models for Speech Recognition. In *Advances in Neural Information Processing Systems*. 2015. pp. 577–585.
- [11] F. Chollet. Keras. 2015. <https://keras.io/>
- [12] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural Language Processing (Almost) From Scratch. *Journal of Machine Learning Research* vol. 12, Aug. 2011. pp. 2493–2537.

- [13] S. B. Davis, P. Mermelstein. Comparison of Parametric Representations for Monosyllabic Word Recognition in Continuously Spoken Sentences. *IEEE Transactions on Acoustics, Speech and Signal Processing* vol. 28, no. 4. 1980. pp. 357–366.
- [14] H. Dudley. Phonetic Pattern Recognition Vocoder for Narrow-Band Speech Transmission. *The Journal of the Acoustical Society of America* vol. 30, no. 8. 1958. pp. 733–739.
- [15] G. Fant. *Acoustic Theory of Speech Production*. Mouton. 1960.
- [16] K.R. Farrell, R.J. Mammone, and K.T. Assaleh. Speaker Recognition using Neural Networks and Conventional Classifiers. *IEEE Transactions on speech and audio processing* vol. 2, no. 1. 1994. pp. 194–205.
- [17] S. Fernández, A. Graves, and J. Schmidhuber. An Application of Recurrent Neural Networks to Discriminative Keyword Spotting. In *International Conference on Artificial Neural Networks*. Springer, Berlin, Heidelberg. 2007. pp. 220–229.
- [18] K.R. Foster, R. Koprowski, and J.D. Skufca. Machine Learning, Medical Diagnosis, and Biomedical Engineering Research-Commentary. *Biomedical Engineering Online* vol. 13, no. 1. 2014. p. 94.
- [19] Y. Gal and Z. Ghahramani. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. *Advances In Neural Information Processing Systems*. 2016. pp. 1019–1027.
- [20] F.A. Gers, J. Schmidhuber, and J. Cummins. Learning to Forget: Continual Prediction with LSTM. *Neural Computation*, vol. 12, no. 10. 1999. pp. 2451–2471.
- [21] X. Glorot, A. Border, and Y. Bengio. Deep Sparse Rectifier Neural Networks. *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*. 2011. pp. 315–323.
- [22] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. Cambridge, MIT press. 2016.
- [23] A. Graves, and J. Schmidhuber. Framewise Phoneme Classification with Bidirectional LSTM Networks. *Proceedings of International Joint Conference on Neural Networks*, vol. 4. 2005. pp. 2047–2052.
- [24] A. Graves, A.R. Mohamed, and G. Hinton. Speech Recognition with Deep Recurrent Neural Networks. *Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2013. pp. 6645–6649.

- [25] D.J. Hand. Measuring Classifier Performance: A Coherent Alternative to the Area Under the ROC Curve. *Machine learning* vol. 77, no. 1. 2009. pp. 103–123.
- [26] J.A. Hanley, and B.J. McNeil. A Method of Comparing the Areas Under Receiver Operating Characteristic Curves Derived from the Same Cases. *Radiology* vol. 148, no. 3. 1983. pp. 839–843.
- [27] F. J. Harris. On the Use of Windows for Harmonic Analysis with the Discrete Fourier Transform. *Proceedings of the IEEE* 66.1. 1978. pp. 51–83.
- [28] M. R. Hasan, M. Jamil, M. G. Rabbani, & M. S. Rahman. Speaker Identification using Mel Frequency Cepstral Coefficients. *International Conference on Electrical & Computer Engineering*, vol 1, no. 4. 2004. pp. 565–568.
- [29] S. Haykin. *Neural Networks: A Comprehensive Foundation*. Prentice Hall PTR. 1994.
- [30] S. Hochreiter, J. Schmidhuber. Long Short-Term Memory. *Neural Computation*, vol. 9, no. 12. 1997, pp. 1735–1780.
- [31] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber. Gradient Flow in Recurrent Nets: The Difficulty of Learning Long-Term Dependencies. *A Field Guide to Dynamical Recurrent Neural Networks*. IEEE Press. 2001.
- [32] J.J. Hopfield. Neural Networks and Physical Systems with Emergent Collective Computational Abilities. *Proceedings of the National Academy of Sciences* vol. 79, no. 8. 1982. pp. 2554–2558.
- [33] K. Karhunen. Über Lineare Methoden in der Wahrscheinlichkeitsrechnung. *Sana*, vol 37. 1947.
- [34] S. Khandelwal, B. Lecouteux, and L. Besacier. Comparing GRU and LSTM for Automatic Speech Recognition. Ph.D. dissertation, LIG. 2016.
- [35] A. Krizhevsky, I. Sutskever, and G.E. Hinton. Imagenet Classification with Deep Convolutional Neural Networks. *Advances in neural information Processing Systems*. 2012. pp. 1097–1105.
- [36] S. Lawrence, C.L. Giles and A.C. Tsoi. Lessons in Neural Network Training: Overfitting May Be Harder Than Expected. In *AAAI/IAAI*. 1997. pp. 540–545.
- [37] Y. LeCun, B. Boser, J.S. Denker, D. Henderson, R.E. Howard, W. Hubbard, and L.D. Jackel. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural computation* vol. 1, no. 4. 1989. pp 541–551.

- [38] Y. Lei, N. Scheffer, L. Ferrer, and M. McLaren. A Novel Scheme for Speaker Recognition using a Phonetically-Aware Deep Neural Network. *Acoustics, Speech and Signal Processing (ICASSP)*. 2014. pp. 1695–1699.
- [39] R.P. Lippmann. Review of Neural Networks for Speech Recognition. *Neural Computation* vol. 1, no. 1. 1989. pp. 1–38.
- [40] B. Logan. Mel Frequency Cepstral Coefficients for Music Modeling. *ISMIR*, vol. 270. 2000. pp. 1–11.
- [41] Y. Lukic, C. Vogt, O. Dürr, and T. Stadelmann. Speaker Identification and Clustering using Convolutional Neural Networks. *IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP)*, Vietri sul Mare, Italy, 13–16 September. 2016.
- [42] J. Lyons. `python_speech_features`.
https://github.com/jameslyons/python_speech_features
- [43] E. MacAskill. "Did 'Jihadi John' kill Steven Sotloff?". *The Guardian*. 2 September 2014. <https://www.theguardian.com/media/2014/sep/02/steven-sotloff-video-jihadi-john>
- [44] M. Mahoney. Large Text Compression Benchmark.
<http://www.matmahoney.net/dc/text.html>
- [45] J. Makhoul, S. Roucos, and H. Gish. Vector Quantization in Speech Coding. *Proceedings of the IEEE* vol. 73, no. 11. 1985. pp. 1551–1588.
- [46] J. Markel and A. Gray Jr. *Linear Prediction of Speech*. Springer–Verlag. 1976.
- [47] W.S. McCulloch and W. Pitts. A Logical Calculus of the Ideas Immanent in Nervous Activity. *The bulletin of mathematical biophysics* vol. 5, no. 4. 1943. pp. 115–133.
- [48] M. Mohri, A. Rostamizadeh, & A. Talwalkar. *Foundations of Machine Learning*. MIT press. 2012.
- [49] L. Muda, M. Begam, & I. Elamvazuthi. Voice Recognition Algorithms using Mel Frequency Cepstral Coefficient (MFCC) and Dynamic Time Warping (DTW) Techniques. *Journal of Computing*, vol. 2, Issue 3, March. 2010.
- [50] D. Nguyen-Tuong, and J. Peters. Model Learning for Robot Control: A Survey. *Cognitive processing* vol. 12, no. 4. 2011. pp. 319–340.
- [51] Nuance. *Multimodal Voice & Behavioral Biometric Authentication Technology*. Retrieved 18 October 2018.

<https://www.nuance.com/omni-channel-customer-engagement/security/identification-and-verification.html>

- [52] D. O’Shaughnessy. *Speech Communication: Human and Machine*. Universities Press. 1987. p. 150.
- [53] L.R. Rabiner. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE* vol. 77, no. 2. 1989. pp. 257–286.
- [54] D. Reynolds. Speaker Identification and Verification using Gaussian Mixture Speaker Models. *Speech Communication*, vol. 17, no. 1. 1995. pp. 91–108.
- [55] D. Reynolds. *Automatic Speaker Recognition using Gaussian Mixture Speaker Models*. The Lincoln Laboratory Journal. 1995.
- [56] H. Robbins and S. Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics* vol. 22, no 3. 1951. pp. 400–407.
- [57] R.C. Rose and D.A. Reynolds. Text Independent Speaker Identification using Automatic Acoustic Segmentation. *Proceedings of ICASSP*, vol. 90, 1990. pp. 293–296.
- [58] H. Sakoe and S. Chiba. Dynamic Programming Algorithm Optimization for Spoken Word Recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing* vol. 26, no. 1. 1978. pp. 43–49.
- [59] C. Sammut, and G.I. Webb. *Encyclopedia of Machine Learning*. Springer Science & Business Media. 2011.
- [60] W.S. Sarle. Stopped Training and Other Remedies for Overfitting. *Computing Science and Statistics*. 1996. pp. 352–360.
- [61] J. Schmalenstroeer and R. Haeb-Umbach. Online Speaker Change Detection by Combining Bic With Microphone Array Beamforming. *Ninth International Conference on Spoken Language Processing*. 2006.
- [62] M. Schuster, and K.K. Paliwal. Bidirectional Recurrent Neural Networks. *IEEE Transactions on Signal Processing* vol. 45, no. 11. 1997. pp. 2673–2681.
- [63] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The Journal of Machine Learning Research* vol. 15, no. 1. 2014. pp. 1929–1958.
- [64] S.S. Stevens, J. Volkman, & E.B. Newman. A Scale for the Measurement of the Psychological Magnitude Pitch. *The Journal of the Acoustical Society of America* vol. 8 no. 3. 1937. pp. 185–190.

- [65] S. Tranter and D. Reynolds. An Overview of Automatic Speaker Diarization Systems. *IEEE Transactions on Audio, Speech, and Language Processing* vol. 14, no. 5. 2006. pp. 1557–1565.
- [66] G. Tsoumakas, and I. Katakis. Multi-Label Classification: An Overview. *International Journal of Data Warehousing and Mining* vol. 3, no. 3. 2007. pp. 1–13.
- [67] O. Vinyals, G. Friedland. Towards Semantic Analysis of Conversations: A System for the Live Identification of Speakers in Meetings. *Semantic Computing, IEEE International Conference*. 2008. pp. 426–431.
- [68] P. Werbos. *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*. Ph. D. dissertation, Harvard University. 1974.
- [69] P. Werbos. *Backpropagation Through Time: What It Does and How to Do It*. *Proceedings of the IEEE* vol. 78, no. 10. 1990. pp. 1550–1560.
- [70] J. Wu, H. Qin, Y. Hua, & L. Fan. Pitch Estimation and Voicing Classification Using Reconstructed Spectrum from MFCC. *IEICE Transactions on Information and Systems* vol. 101, no. 2. 2018. pp. 556–559.
- [71] Y. Yang, and X. Liu. A Re-Examination of Text Categorization Methods. *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1999. pp. 42–49.
- [72] S.J. Young and S. Young. *The HTK Hidden markov Model Toolkit: Design and Philosophy*. University of Cambridge, Department of Engineering. 1993.