



VILNIUS UNIVERSITY
FACULTY OF MATHEMATICS AND INFORMATICS
DATA SCIENCE STUDY PROGRAMME

Master's thesis

**Data Selection Strategies for Multi-Speaker
Text-to-Speech Synthesis in Lithuanian**
Work Title in Lithuanian

Aleksandr Jan Smoliakov

Supervisor : Dr. Gerda Ana Melnik-Leroy

Scientific advisor : Dr. Gražina Korvel

Reviewer : pedagogical/scientific title Name Surname

Vilnius
2025

List of Figures

1	Visual representation of Analog-to-Digital conversion	7
2	Raw waveform, Spectrogram, and Mel-spectrogram	9
3	Tacotron 2 architecture	14
4	Glow-TTS architecture	15
5	Text-to-Speech synthesis pipeline	16
6	General architecture of a Speaker Encoder	17
7	Latin square design for TTS evaluation	19

List of Tables

1	Lithuanian homographs with accentuation ambiguity	12
2	Experimental Data Subsets	23
3	Mel-spectrogram extraction parameters.	24
4	Tacotron 2 DCA training configuration.	26
5	Glow-TTS training configuration.	27
6	Objective evaluation results. Lower is better for both MCD and F0 RMSE. Bold indicates the best performance per architecture; <u>underline</u> indicates the global best. . .	29
7	Mean Opinion Score (MOS) results with 95% confidence intervals. Ratings scale from 1 (Bad) to 5 (Excellent).	30
8	Complete Mel-spectrogram extraction parameters.	35

Contents

List of Figures	2
List of Tables	3
Introduction	6
1 Literature review	7
1.1 Digital representation of audio	7
1.2 Time-Frequency Analysis	8
1.2.1 Fourier Transform	8
1.2.2 Spectrogram and Mel-spectrogram	8
1.3 Text-to-speech synthesis	9
1.3.1 Traditional TTS approaches	9
1.3.2 Concatenative synthesis	9
1.3.3 Parametric synthesis	10
1.4 Linguistic Representation (Text Processing)	10
1.4.1 Text normalization	11
1.4.2 Graphemes vs. Phonemes	11
1.4.3 Specific challenges in Lithuanian	11
1.5 Embeddings and Representation Learning	12
1.5.1 The Concept of Embeddings	12
1.5.2 Text Embeddings	13
1.6 Deep learning for TTS	13
1.6.1 Feedforward neural networks	13
1.6.2 Encoder-Decoder architectures	13
1.6.3 Sequence-to-Sequence models and Tacotron 2	14
1.6.4 Non-autoregressive models and Glow-TTS	15
1.6.5 Other notable TTS models	15
1.6.6 Neural vocoders	15
1.7 Multi-speaker TTS	17
1.7.1 Speaker embeddings	17
1.7.2 Challenges	18
1.8 Evaluation metrics	18
1.8.1 Mean Opinion Score (MOS)	18
1.8.2 Latin square design	19
1.9 Research gap	20
1.10 Summary	21
2 Methodology	22
2.1 Research Design	22
2.1.1 Variables	22
2.2 Data and Preprocessing	22
2.2.1 Liepa 2 Dataset	22
2.2.2 Speaker Selection and Filtering	23
2.2.3 Experimental Data Subsets	23
2.2.4 Text Normalization and Accentuation	23
2.2.5 Audio Preprocessing	24
2.3 Model Architectures	24

2.3.1	Speaker Conditioning	24
2.3.2	Tacotron 2 (Autoregressive)	25
2.3.3	Glow-TTS (Non-autoregressive)	25
2.3.4	Vocoder	25
2.4	Model Training Configurations	25
2.4.1	Environment and Framework	25
2.4.2	Tacotron 2 Configuration	26
2.4.3	Glow-TTS Configuration	26
2.5	Evaluation Protocol	26
2.5.1	Objective Evaluation	27
2.5.2	Subjective Evaluation (MOS)	27
3	Results and Analysis	29
3.1	Objective Evaluation	29
3.1.1	Alignment Convergence	29
3.1.2	Pitch and Spectral Accuracy	29
3.2	Subjective Evaluation (MOS)	30
3.2.1	The Trade-off: Stability vs. Peak Quality	30
3.2.2	Optimal Composition	30
3.3	Discussion	31
4	Conclusion	32
4.1	Summary of findings	32
4.2	Contributions	32
4.3	Limitations of the study	32
4.4	Future work	32
5	References	32
6	Appendix: Audio preprocessing parameters	35

Introduction

The goal of creating machines that can speak like humans has captivated researchers for centuries. One of the earliest known attempts dates back to the 18th century, with Wolfgang von Kempelen’s mechanical speech machine that utilized a bellows-driven lung and physical models of the tongue and lips.

Over the centuries, advancements in technology and understanding of human speech have driven significant progress in this field. Today’s state-of-the-art systems, dominated by end-to-end (E2E) neural models, have achieved highly naturalistic speech with unprecedented acoustic quality. Notably, these end-to-end systems have unified the entire synthesis process into a single neural network, eliminating the need for complex multi-stage pipelines.

Training high-quality TTS models typically requires large amounts of annotated speech data. The common recommendation is to use at least 10 hours of recorded speech from a single speaker to achieve good results.

Liepa 2 [1] is a recently released Lithuanian speech corpus that contains 1000 hours of annotated speech; however, this data is distributed across 2621 speakers, with most speakers contributing under 20 minutes each. The top speaker has around 2.5 hours of recorded speech.

Training a high-quality single-speaker TTS model on such limited data poses a challenge. Multi-speaker TTS models can utilize data from multiple speakers to improve performance. However, training on all available data is a time-consuming and computationally expensive process, especially in the context of a master’s thesis.

Therefore, it makes sense to explore strategies for selecting smaller subsets of the available data for training. The question that arises is, what is the best way to sample multi-speaker data for training TTS models?

This thesis aims to answer the following research questions:

- TODO

Scope: This study is exclusively focused on the Lithuanian language and the Liepa 2 speech corpus. It investigates a fixed total training data size of 22.5 hours. The models are limited to Tacotron 2 with DCA and Glow-TTS architectures within the Coqui TTS framework, using a pre-trained WaveGlow vocoder for waveform generation.

Limitations: The findings may not generalize to other languages, datasets with different characteristics, or other TTS architectures. The 22.5-hour training data size is a practical constraint and may not reflect performance at larger scales.

1 Literature review

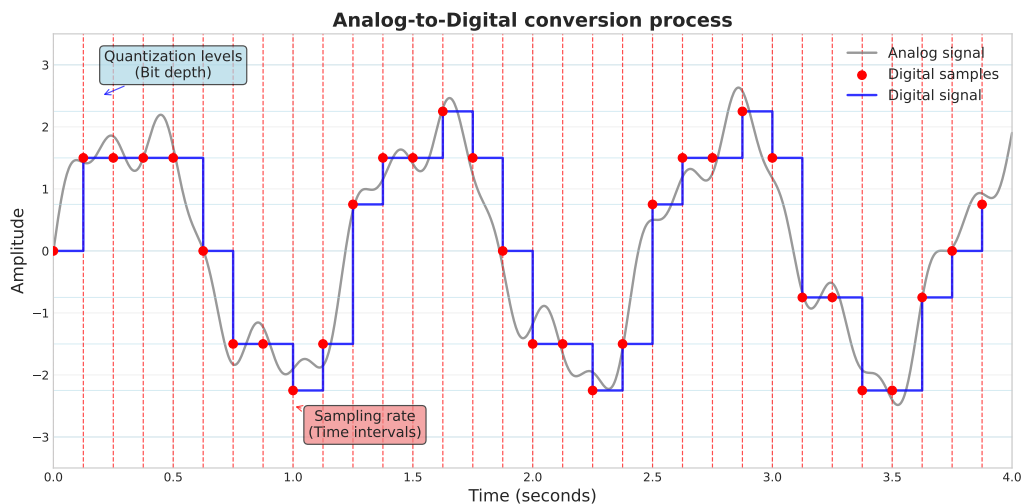
1.1 Digital representation of audio

Speech, or sound in general, is a continuous pressure wave that propagates through a medium, such as air. The key properties of sound waves include frequency (pitch), amplitude (loudness), and phase.

Converting continuous sound waves into a digital format suitable for computer processing involves two main steps: sampling and quantization.

Sampling is the process of measuring the amplitude of the sound wave at regular time intervals. The rate at which these samples are taken is called the sampling rate. According to the Nyquist-Shannon [2] sampling theorem, accurate reconstruction of a continuous signal requires a sampling rate that is strictly greater than twice the highest frequency present in the signal. Frequencies in the range between 300 Hz and 3400 Hz contribute most to human speech intelligibility and recognition. [3]. In text-to-speech applications, common sampling rates for audio are 22.05 kHz and 24 kHz, which can capture frequencies up to approx. 11 kHz and 12 kHz, respectively.

Quantization (also known as bit depth) is the mapping of continuous amplitude values to discrete levels for digital representation, which determines the precision of the representation. Common bit depths for audio are 16-bit and 24-bit formats. A visual representation of both sampling and quantization is provided in Figure 1.



1 Visual representation of Analog-to-Digital conversion. The continuous grey line represents the analog signal. The vertical lines represent the **sampling rate** (time intervals), and the horizontal grid lines represent **quantization levels** (bit depth).

Pre-emphasis is a high-frequency filtering technique applied to audio signals before further processing. Natural speech signals tend to have more energy in the lower frequencies, with a gradual drop-off towards higher frequencies (typically around -6 dB per octave). Pre-emphasis compensates for this spectral tilt by boosting high frequencies using a first-order high-pass filter, which is defined as:

$$y[n] = x[n] - \alpha x[n - 1] \quad (1)$$

where $y[n]$ is the pre-emphasized signal, $x[n]$ is the original signal, α is the pre-emphasis coefficient (typically between 0.9 and 1.0, and often set to 0.97), and n is the sample index.

This transformation balances the frequency spectrum, improving the signal-to-noise ratio for higher frequencies and preventing the model from optimizing only for low-frequency components.

1.2 Time-Frequency Analysis

1.2.1 Fourier Transform

Fourier Transform (FT) is a mathematical technique that transforms a time-domain signal (such as an audio waveform) into its frequency-domain representation. The signal is decomposed into a sum of sine and cosine waves at various frequencies, each with a specific amplitude and phase. This allows us to analyze the frequency content of the signal.

Short-Time Fourier Transform [4] (STFT) extends the FT by applying it to short, overlapping segments (frames) of the signal. This transformation provides a time-frequency representation, showing how the frequency content of the signal changes over time.

In TTS applications, the STFT is computed by dividing the audio signal into short frames (usually, 20–50 ms) with a certain overlap (usually, 50–75%) between frames, windowed by a Hamming or Hann function to reduce the spectral leakage.

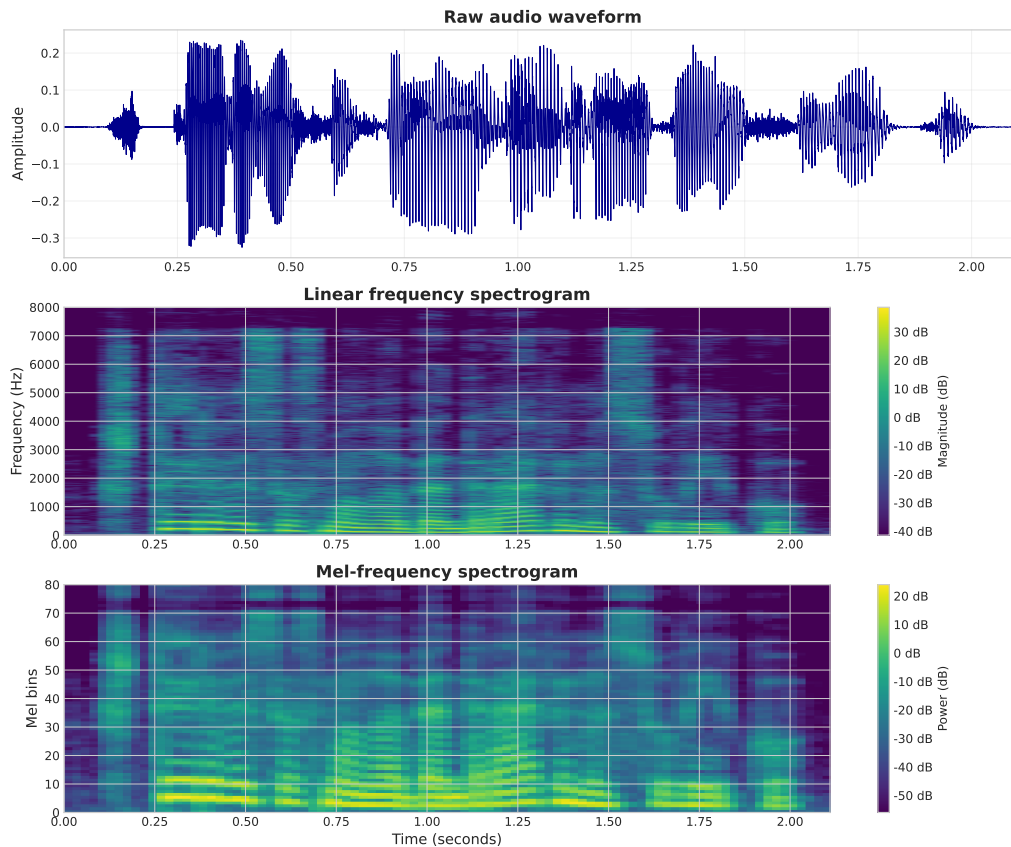
1.2.2 Spectrogram and Mel-spectrogram

The spectrogram is a visual representation of the STFT, displaying frequency on the vertical axis, time on the horizontal axis, and amplitude represented by the color intensity.

However, the human ear does not perceive frequencies linearly — it is more sensitive to lower frequencies than higher ones. To mimic this perceptual characteristic, the Mel scale [5] maps linear frequency f (in Hz) to a perceptual scale m (in Mels) using the following formula:

$$m = 2595 \cdot \log_{10} \left(1 + \frac{f}{700} \right) \quad (2)$$

Mel-spectrograms are computed by applying a Mel filterbank of overlapping triangular filters (or kernels) to the magnitude spectrogram obtained from the STFT. This results in a compressed representation of the audio signal that aligns more closely with human auditory perception. Such Mel-spectrograms are commonly used as input features for modern TTS systems. The differences between the raw waveform, the standard spectrogram, and the Mel-spectrogram are illustrated in Figure 2. Note how the Mel-spectrogram has a higher resolution in the lower frequencies, where the majority of the speech energy is concentrated.



2 Raw audio waveform (top), its spectrogram (middle), and Mel-spectrogram (bottom) representations for the utterance “Štai ir visas mano bendravimas su vaiku”.

1.3 Text-to-speech synthesis

Text-to-Speech (TTS) synthesis, also known as speech synthesis, is the process of converting written text into human-like spoken words. Nowadays TTS is a key technology in numerous applications, including virtual assistants, accessibility tools, and language learning platforms.

1.3.1 Traditional TTS approaches

Early attempts at artificial speech synthesis evolved from the first mechanical devices in the 18th century to electronic systems. Wolfgang von Kempelen’s mechanical speech machine demonstrated basic phoneme production using a physical model of the vocal tract. In 1939, Homer Dudley’s invention of the Voder [6] became the first electronic speech synthesizer that could produce intelligible speech through operator-controlled acoustic parameters, establishing the foundation for modern electronic synthesis methods.

In the decades that followed, two main approaches for speech synthesis emerged: concatenative synthesis and parametric synthesis.

1.3.2 Concatenative synthesis

The concatenative synthesis approach [7] synthesizes speech by piecing together pre-recorded segments of human speech. This method involves several steps. First, it requires pre-recording a

large database of speech segments spoken by a human voice actor in pristine, highly controlled studio conditions to ensure consistent audio quality and minimize background noise. Each segment is labeled and indexed based on its phonetic and prosodic properties.

During synthesis, the system breaks down the input text into short linguistic units (such as phonemes or syllables) using a text analysis module. Then, it queries the speech database to find the best-matching segments for each unit using selection cost functions [8]. The retrieved segments are blended and concatenated to form a continuous speech waveform. Finally, the system uses signal processing techniques to smooth the transitions between segments and adjust pitch and duration to match the desired output characteristics.

Concatenative synthesis can produce natural-sounding individual speech units, but the final audio often has noticeable audible continuity distortions at the concatenation points [8]. The segments may not blend smoothly due to differences in pitch, duration, and timbre. The prosody also tends to sound “choppy” and unnatural, since stringing disjointed segments together does not capture the natural rhythm and intonation patterns of connected speech.

Finally, concatenative synthesis requires language-specific expertise to design and maintain the underlying speech database and selection algorithms. This need for extensive data can make it challenging to develop concatenative TTS systems for low-resource languages or dialects.

1.3.3 Parametric synthesis

In contrast, statistical parametric speech synthesis [9] (SPSS) uses statistical models, typically Hidden Markov Models (HMMs) [10], to generate the parameters that control a speech waveform.

This method involves training a statistical model on a large corpus of recorded speech. The model learns the relationship between linguistic features (like phonemes and prosody) and the acoustic features of the speech signal, such as spectral envelope and fundamental frequency. During synthesis, the system takes text as input, converts it to a sequence of linguistic features, and then uses the trained model to generate a corresponding sequence of acoustic parameters.

Compared to concatenative synthesis, the statistical approach allows for more flexibility and control over the speech synthesis process, enabling the generation of a wider variety of voices and speaking styles. However, HMM-based synthesis [10] had a persistent problem: the statistical averaging built into the models tended to over-smooth the acoustic features, creating the characteristic “buzzy” or “muffled” sound that lacked the sharpness and detail of natural human speech.

1.4 Linguistic Representation (Text Processing)

In TTS systems, the input text must be pre-processed and converted into a suitable linguistic representation that the synthesis model can use. The main goal is to map the raw sentences into a sequence of symbols that can be more closely mapped to the acoustic features of speech.

Although theoretically an end-to-end TTS model could learn to map raw text directly to audio, in practice, pre-processing the text makes the model convergence easier and improves the quality of the synthesized speech.

This process typically involves several steps, such as text normalization, grapheme-to-phoneme conversion, and possibly prosody prediction.

1.4.1 Text normalization

Text normalization [11] is the process of converting raw written text with non-standard words (NSWs) into a more standardized “spoken” form. Typical steps include expanding abbreviations (e.g., expanding “Dr.” to “Doctor”), punctuation removal, number normalization (e.g., converting “123” to “one hundred twenty-three”), and lowercasing.

As an example, the input text “Dr. Smith has 2 cats.” could be normalized to “doctor smith has two cats”.

Text normalization helps reduce the variability and complexity in the input text, decreases the number of unique symbols, and removes the ambiguities that could confuse the TTS model. The resulting normalized text is not only easier for the model to process, but can also be further converted into phonemes, which provide an even closer representation of the spoken language.

1.4.2 Graphemes vs. Phonemes

Text-to-speech systems use a discrete input representation derived from text, generally divided into grapheme-based or phoneme-based sequences.

Grapheme-based models ingest raw character sequences (orthography). This approach simplifies the inference pipeline by eliminating the dependency on external grapheme-to-phoneme (G2P) converters. However, it forces the model to implicitly learn pronunciation rules from data, which can be a significant challenge for languages with complex orthographies or inconsistent grapheme-to-phoneme mappings (e.g., “read” vs. “read”).

In contrast, the phoneme-based approach uses a phonetic transcription of the text, typically in the International Phonetic Alphabet (IPA) or ARPABET form. By resolving pronunciation ambiguities prior to training, phonemes provide a more direct mapping to acoustic features, simplifying the model’s task of learning alignment. The downside is that this approach requires an external grapheme-to-phoneme (G2P) conversion step [12]. Additionally, errors in the G2P conversion can propagate to the TTS model, affecting the quality of the synthesized speech.

There is another approach that augments the grapheme-based representation with explicit lexical stress markers or diacritics (e.g., tilde, acute, grave accents). This intermediate method helps the model disambiguate pronunciation of homographs and easier learn prosodic patterns without requiring a full phonetic transcription, particularly in languages where stress placement alters meaning.

1.4.3 Specific challenges in Lithuanian

Lithuanian is a Baltic language with a rich inflectional morphology and complex prosodic structure. It is a pitch-accent language with free stress, meaning the stress can fall on any syllable in a word, and can change the position depending on the grammatical form.

Challenges in Lithuanian TTS synthesis include:

- **High OOV rate:** Due to extensive word inflection, the number of unique word forms is significantly higher than in English. This leads to data sparsity issues where many valid word forms may not appear in the training set.
- **Ambiguity without accentuation:** Typically, stress marks are omitted in written Lithuanian. However, stress position and tone (acute, circumflex, or short) determine the meaning of monographic words. Examples are shown in Table 1. A grapheme-based model with accentuation marks has been shown to improve synthesis quality in Lithuanian. [13].

Word	Accentuation	Meaning
Antis	<i>ántis</i> (Acute)	A duck (noun)
	<i>añtis</i> (Circumflex)	Bosom/Chest (noun)
Kasa	<i>kãsa</i> (Circumflex)	He/she digs (verb)
	<i>kasà</i> (Short)	Braid/Pancreas (noun)

1 Examples of Lithuanian homographs where accentuation determines meaning. A grapheme-only model cannot distinguish these without context or explicit stress marks.

To overcome these challenges, tools like **Kirčiuoklis** [14] (Vytautas Magnus University) are often employed in the text normalization pipeline. Kirčiuoklis automatically assigns stress marks to raw text. One weakness of Kirčiuoklis is that it relies on simple word-dictionary based lookup, which does not take into account the context of the word. Thus, it suggests multiple possible accentuation variants for homographs, leaving it up to the user to select the correct one.

In the absence of a high-quality, context-aware Grapheme-to-Phoneme (G2P) converter for Lithuanian, this thesis will focus on grapheme-based TTS synthesis with accentuation marks provided by Kirčiuoklis. In cases where Kirčiuoklis suggests multiple accentuation variants for a word, no stress marks will be added, leaving the TTS model to infer the correct prosody from context.

1.5 Embeddings and Representation Learning

1.5.1 The Concept of Embeddings

In machine learning, embeddings are dense vector representations of discrete entities (such as words, characters, or speakers) to a high-dimensional continuous vector space. Unlike one-hot encodings, which are sparse and highly dimensional, embeddings provide a dense, lower-dimensional representation that captures semantic relationships between underlying entities. For instance, in word embeddings, similar words tend to have more similar (correlated) vector representations, while dissimilar words map to more distant points in the vector space. [15]

1.5.2 Text Embeddings

The “Encoder” part of a TTS model is responsible for converting a sequence of input symbols (characters or phonemes) into a sequence of feature vectors. Usually, this is done using an embedding layer, which maps each “categorical” input symbol to a learnable fixed-size vector representation. During training, these embeddings are learned jointly with the rest of the TTS model.

1.6 Deep learning for TTS

The limitations of complex, multi-stage pipelines motivated the creation of the end-to-end (E2E) model. E2E systems learn the entire speech synthesis process — from input text directly to acoustic output — using a single neural network. This approach promised to eliminate the need for hand-crafted pipelines that were difficult to design, required extensive expertise, and suffered from errors that accumulated across multiple components. By learning directly from text-audio pairs, E2E models showed they could produce speech with higher naturalness and expressiveness than previous methods, representing a significant leap in TTS technology.

Although deep learning TTS models are more robust to variations in data quality compared to concatenative approaches, they are essentially “data-hungry” systems that require large amounts of training data to achieve optimal performance. Extrapolating from results in language modeling, it is observed that model performance follows general scaling laws [16], improving as the amount of training data increases.

However, in the context of multi-speaker synthesis, there is a trade-off between the breadth of the data (number of distinct speakers) and the depth of the data (duration of audio per speaker). In theory, training on a dataset with a massive number of speakers, even with limited data per speaker, may allow the model to learn a more generalized latent space of voice characteristics. This high variance in the training data could act as a form of regularization, preventing overfitting to noise and idiosyncrasies of individual speakers. In contrast, datasets with fewer speakers but high duration per speaker allow the model to capture fine-grained prosodic details specific to those voices, potentially achieving higher stability but lower generalization capabilities.

1.6.1 Feedforward neural networks

Feedforward Neural Networks (FNNs) are the simplest type of artificial neural networks, consisting of layers of interconnected nodes (neurons) where information flows in one direction — from the input, through hidden layers, to the output. While FNNs can be useful for basic regression or classification tasks, they lack the memory and context-awareness needed for processing sequential data like text and speech. Therefore, FNNs are not suitable for modelling TTS tasks that require understanding of temporal dependencies.

1.6.2 Encoder-Decoder architectures

The Encoder-Decoder architecture is a neural network architecture consisting of two components, namely an encoder and a decoder. The encoder processes the input data and compresses it

into a high-dimensional latent representation. This vector captures the meaningful features of the input. The decoder uses this latent representation as context to generate the final output. This architecture is commonly used in sequence-to-sequence tasks, such as machine translation (text-to-text) and text-to-speech synthesis (text-to-audio frames).

1.6.3 Sequence-to-Sequence models and Tacotron 2

A common approach in modern neural TTS is the sequence-to-sequence (seq2seq) framework [17], which uses an encoder-decoder architecture with an attention mechanism to map input (text) sequences to output (audio) frames.

Tacotron [18] and Tacotron 2 [19] are two notable TTS models based on the sequence-to-sequence architecture. The variant that this thesis primarily focuses on is Tacotron 2 with Dynamic Convolutional Attention (DCA). The complete architecture of Tacotron 2 is depicted in Figure 3.

This architecture has three main components:

1. **Encoder:** The encoder's input is a character or phoneme sequence. A stack of convolutional layers followed by a bidirectional LSTM converts the character sequence into a high-level hidden feature representation.
2. **Attention mechanism:** A location-sensitive attention [20] mechanism learns to align the high-level text representation with the decoder steps. This alignment determines which parts of the input text should be attended to when generating each of the output audio frames.
3. **Decoder and Post-net:** The autoregressive LSTM decoder generates a coarse Mel-spectrogram frame. This output is then passed through a convolutional **Post-net** which predicts a residual to refine the spectral details and improve reconstruction quality.
4. **Stopnet:** A linear layer projects the LSTM decoder's output to a scalar, predicting the probability that the current frame is the "stop token", which halts the synthesis process. This allows the model to dynamically determine the output duration.

The model is optimized by minimizing a combination of losses: the mean squared error (MSE) between the predicted and ground truth Mel-spectrograms (both the Decoder and Post-net outputs), the spectral similarity index (SSIM) loss for improving spectral similarity, the "guided attention" loss to encourage diagonal attention alignments, and the binary cross-entropy loss for the stop token prediction.

While Tacotron 2 can generate high-quality speech, its autoregressive nature makes inference slow, as each audio frame must be generated sequentially. Additionally, the attention mechanism can sometimes fail, leading to issues like skipped or repeated words in the synthesized speech.

3 The Tacotron 2 architecture. Note the recurrent connections in the decoder and the attention mechanism aligning encoder outputs to decoder steps. [19]

1.6.4 Non-autoregressive models and Glow-TTS

To address the slow inference speed and stability issues of autoregressive models, non-autoregressive (parallel) models were developed. Glow-TTS [21] is a notable example of such a model that applies flow-based generative models to text-to-speech tasks.

Unlike Tacotron 2, Glow-TTS model generates the entire Mel-spectrogram in parallel, significantly speeding up the inference. It utilizes an invertible flow-based decoder. A key innovation of Glow-TTS is its ability to learn alignment internally without requiring external priors (aligner), achieved through the following components:

- **Monotonic Alignment Search (MAS):** Unlike FastPitch [22], which relies on external aligners to train a duration predictor, Glow-TTS treats alignment as a latent variable. It employs MAS to find the most probable monotonic path between the input text and the target Mel-spectrogram, maximizing the log-likelihood of the data.
- **Flow-based decoder:** The model uses a stack of normalizing flows—specifically invertible 1×1 convolutions and affine coupling layers. This decoder transforms a simple prior distribution (conditioned on the text input) into the complex distribution of Mel-spectrograms.

By utilizing the properties of flows, Glow-TTS allows for varied speech synthesis by sampling from the latent space and manipulating the noise temperature. Furthermore, the duration of the speech can be controlled by modifying the predicted duration of the alignment.

The high-level architecture of Glow-TTS is shown in Figure 4.

4 The Glow-TTS architecture. It utilizes a Transformer-based text encoder and a flow-based decoder, connected via Monotonic Alignment Search (MAS) to enable parallel Mel-spectrogram generation. [21]

The primary advantages of Glow-TTS over Tacotron 2 are inference speed (due to non-autoregressive generation), robustness (the monotonic constraint prevents skipping or repeating words), and the elimination of the need for an external aligner during training.

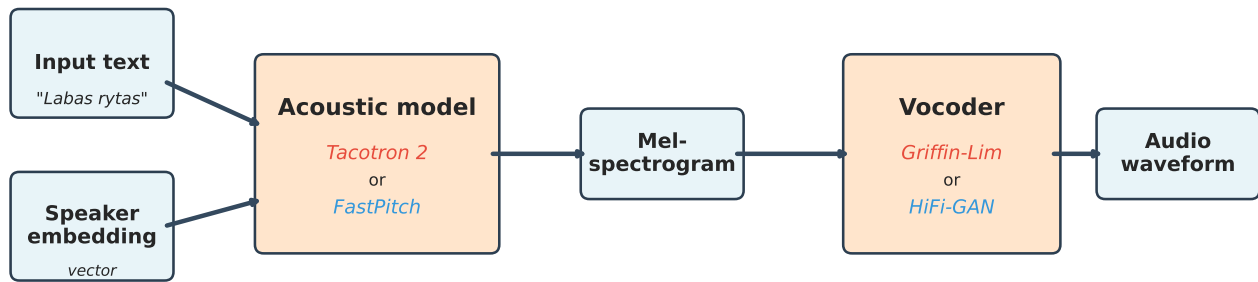
1.6.5 Other notable TTS models

Besides Tacotron 2 and Glow-TTS, other notable TTS architectures include **FastPitch** [22], which uses a feed-forward Transformer architecture with explicit pitch and duration predictors, and **VITS** [23] (Conditional Variational Autoencoder with Adversarial Learning), which combines the acoustic TTS model (Glow-TTS) with a neural vocoder (HiFi-GAN) into a single end-to-end architecture.

1.6.6 Neural vocoders

As illustrated in Figure 5, modern TTS systems typically employ a two-stage pipeline: an acoustic model like Tacotron 2 and Glow-TTS generates intermediate acoustic features (Mel-spectrograms),

End-to-End Text-to-Speech pipeline



5 Text-to-Speech synthesis pipeline. The TTS model generates Mel-spectrograms from input text, which are then converted to raw audio waveforms by a neural vocoder.

but do not directly produce raw audio waveforms. Spectrograms are lossy representations that only capture the magnitude of the sound frequency bands, discarding phase information. Converting a lossy spectrogram into audio is a non-trivial task, as the phase information must be estimated. This challenge is known as the *inversion problem*.

An additional component called a vocoder is required to reconstruct the raw waveform from the Mel-spectrogram.

Traditionally, the Griffin-Lim algorithm [24] has been used to iteratively estimate and reconstruct the phase information from the magnitude spectrogram. However, this method often produces audio with noticeable artifacts and lower quality compared to natural speech.

Modern TTS systems use neural vocoders, which are deep generative models trained to map acoustic features to raw waveforms. **WaveNet** [25] was one of the first autoregressive models to produce high-fidelity audio, but its sequential generation process made it prohibitively slow for real-time applications.

To address the speed limitations, Generative Adversarial Network (GAN) based vocoders were introduced. **HiFi-GAN** [26] is currently one of the state-of-the-art neural vocoders. It consists of a Generator that upsamples the Mel-spectrograms using transposed convolutions and a set of Discriminators (multi-scale and multi-period discriminators) that ensure the generated audio is indistinguishable from real human speech. HiFi-GANs are highly efficient and capable of faster-than-real-time synthesis on consumer hardware while maintaining high perceptual quality.

The framework used in this thesis, Coqui TTS [27], comes with a pre-trained HiFi-GAN v2 vocoder trained on a large multi-speaker dataset (VCTK [28]) with 110 English speakers.

The use of a vocoder trained on English data for Lithuanian synthesis is justified by the language-agnostic nature of the phase reconstruction task. Neural vocoders' primary function is to model the physics of human speech production rather than linguistic features. While language-dependent phonetic nuances exist, studies have shown that vocoders trained on large, diverse datasets can effectively generalize to unseen speakers and languages [29]. Therefore, this thesis will utilize the pre-trained HiFi-GAN v2 model for waveform generation.

This should allow the vocoder model to effectively generalize to unseen speakers and languages, provided that the acoustic feature parameters (including sampling rate, FFT size, Mel-filterbank limits) of the input Mel-spectrograms match those used during the vocoder’s training.

Therefore, this thesis will utilize the pre-trained HiFi-GAN v2 model for waveform generation. The TTS models’ acoustic parameters will be configured to the exact same acoustic parameters used during the vocoder’s training to ensure compatibility.

1.7 Multi-speaker TTS

Multi-speaker TTS models are designed to synthesize speech in the voices of multiple speakers. In order to achieve this, these models are indeed trained on data from many different speakers, allowing them to learn the characteristics of each voice and synthesize speech that sounds like a specific individual, while still being able to generalize the shared linguistic and acoustic patterns across speakers.

1.7.1 Speaker embeddings

To enable multi-speaker synthesis, TTS models require a representation of the speaker’s identity. In multi-speaker models, the network is conditioned on a speaker embedding. The model learns a shared representation of phonetics (how text maps to sound generally) while using an additional input — the speaker embedding — to adjust the timbre and prosodic characteristics specific to a voice.

Early successful implementations of this approach include Deep Voice 2 [30], which demonstrated effective multi-speaker synthesis by learning speaker-specific embeddings.

Nowadays there are several techniques for incorporating speaker embeddings into TTS models:

Lookup Tables (LUT): Early multi-speaker approaches used simple, learnable embeddings where each speaker ID is mapped to a unique vector. The vectors are initialized randomly and learned jointly with the TTS model. While this method is straightforward and efficient, it cannot generalize to speakers not seen during training.

d-vectors and x-vectors: Transfer learning approaches [31] have demonstrated adapting speaker verification models for multispeaker TTS synthesis, enabling better speaker adaptation and higher voice quality. The general architecture of such a speaker encoder is illustrated in Figure 6. A speaker encoder model pre-trained on a massive, noisy dataset with thousands of speakers (e.g., the VoxCeleb dataset [32]) learns the general speaker space. Its pre-trained weights are frozen and used to extract embeddings for the TTS training data, allowing the TTS model to effectively account for multi-speaker variation.

6 General architecture of a Speaker Encoder. A reference audio of arbitrary length is processed (typically by LSTM or TDNN layers) and pooled to produce a fixed-length embedding vector (e.g., d-vector) representing the speaker identity.

d-vectors: d-vectors [33] are fixed-length speaker embeddings derived from a separate speaker

verification model. A reference encoder network takes a reference audio recording of arbitrary length and compresses it into a fixed-length vector known as a d-vector, that summarizes the speaker’s timbral and prosodic characteristics. These d-vectors are then provided as additional input to the TTS model, and are kept fixed during TTS training.

x-vectors: An evolution of d-vectors, x-vectors [34] use a Time Delay Neural Network (TDNN) architecture to capture the temporal context more effectively. These embeddings have shown an improved ability in zero-shot TTS scenarios.

One limitation of d-vectors and x-vectors is that if the reference audio is of poor quality or contains background noise, the resulting speaker embedding may not accurately represent the speaker’s identity, leading to degraded synthesis quality.

1.7.2 Challenges

One key challenge in multi-speaker TTS is ensuring that the model can generalize across many speakers while still maintaining high quality for each. There is a trade-off between the *breadth* of the dataset (number of speakers) and the *depth* (minutes of audio per speaker).

Standard TTS systems historically required 10 to 20 hours of recorded speech for a single professional speaker. However, deep learning models capable of *transfer learning* can produce intelligible speech for a new speaker with significantly less data, potentially as little as a few minutes — if the base model has been pre-trained on a sufficiently diverse multi-speaker dataset.

1.8 Evaluation metrics

Evaluating Text-to-Speech systems is notoriously difficult because “quality” is a subjective metric defined by human perception. There is no single mathematical objective function that perfectly correlates with human judgement of naturalness and intelligibility. Therefore, TTS systems are typically evaluated using subjective listening tests.

1.8.1 Mean Opinion Score (MOS)

The most standard metric for evaluating speech synthesis quality is the Mean Opinion Score (MOS), originally derived from telecommunications quality standards (ITU-T P.800). [35].

In a MOS test, human listeners (raters) are presented with a set of synthesized speech audio samples and asked to rate them on a 5-point Likert scale. The standard scale for “Naturalness” is:

- **5:** Excellent (Imperceptible difference from real speech)
- **4:** Good (Perceptible but not annoying)
- **3:** Fair (Slightly annoying)
- **2:** Poor (Annoying)
- **1:** Bad (Very annoying / Unintelligible)

The final score is the arithmetic mean of all ratings collected for a specific TTS system. Although MOS is subjective, with a sufficient number of raters (typically, at least 15–20), the scores tend to converge and provide a reliable ranking between different models.

1.8.2 Latin square design

A major challenge in subjective listening tests is controlling for biases. If a rater hears the same sentence produced by different TTS systems in a row, their ratings may be influenced by the repetition (repetition effect) or by the relative order of presentation (order effect). For instance, a “Slightly annoying” sample may be rated more harshly if it follows an “Excellent” sample (contrast effect).

In order to mitigate these biases, a Latin square design [36] is often employed for MOS tests. In this experimental design:

1. A set of test sentences (utterances) is selected.
2. The listeners are divided into groups.
3. The presentation is balanced such that each listener hears every test sentence exactly once, and every TTS system (model) exactly once per block of trials, but never the same sentence-system combination twice.

Latin square design for TTS evaluation

Listener group	Group 1	Model A S1	Model B S2	Model C S3	Model D S4
	Group 2	Model B S3	Model C S4	Model D S1	Model A S2
	Group 3	Model C S2	Model D S1	Model A S4	Model B S3
	Group 4	Model D S4	Model A S3	Model B S1	Model C S2
		Order 1	Order 2	Order 3	Order 4

Presentation order

Model A
 Model B
 Model C
 Model D

7 Latin square design for TTS evaluation. Each listener group hears each sentence exactly once, and each TTS system exactly once per block, ensuring balanced exposure and mitigating order/repetition biases.

An example Latin square design for 4 TTS systems and 4 test sentences is illustrated in Figure 7.

For a multi-speaker TTS evaluation (as is the case in this thesis), the Latin square design ensures that the ratings reflect the quality of the model rather than the linguistic content of the sentence or listener fatigue. By rotating the systems and sentences across listener groups, the influence of specific difficult sentences is averaged out across all models.

1.9 Research gap

While the literature demonstrates the capabilities of modern deep learning TTS architectures like Tacotron 2 and Glow-TTS to produce highly natural-sounding speech, several questions remain unanswered regarding their application to low-resource, morphologically complex languages like Lithuanian.

Firstly, although neural TTS models may follow general neural model scaling laws [16], implying that performance improves with more data, there is limited understanding of the optimal composition of training data under a fixed budget. In low-resource settings, scaling up the dataset size is not always feasible, and this may be further constrained by the computational resources required for training large models. A critical question is whether it is more beneficial to train on a smaller number of speakers with more data per speaker (high depth) or a larger number of speakers with less data per speaker (high breadth).

Current research primarily focuses on high-resource languages like English, where the availability of large, balanced multi-speaker datasets masks the nuances of this trade-off. For a pitch-accent language like Lithuanian, the requirements may be different. It is hypothesized that high-diversity datasets may help the model learn a richer representation of prosodic patterns, while high-depth datasets may improve the model’s naturalness for the target speakers.

Secondly, most multi-speaker TTS research assumes access to large-scale datasets with thousands of utterances per speaker. There is a lack of research exploring how different TTS architectures (autoregressive vs. non-autoregressive) perform when the data per speaker is scarce (e.g., under 10 minutes).

This thesis aims to fill the research gap by systematically evaluating the efficiency of Tacotron 2 and Glow-TTS models trained on Lithuanian speech data. By controlling the total dataset size and varying the distribution of speakers and data per speaker, this study will provide insights into the optimal data composition for training multi-speaker TTS models in low-resource settings.

To summarize, the key research questions this thesis seeks to answer are:

- How does the trade-off between data breadth (number of speakers) and data depth (minutes per speaker) affect the performance of multi-speaker TTS models for Lithuanian?
- How do different TTS architectures (Tacotron 2 vs. Glow-TTS) perform under varying data selection strategies in low-resource settings?

The experiments will involve training models on three distinct data selection strategies:

- **High-resource per speaker:** Fewer speakers (30), but high fidelity (45 min each), total: 22.5 hours.
- **Balanced:** Moderate diversity (60 speakers), moderate data (22.5 min each), total: 22.5 hours.
- **High-diversity:** Many speakers (180), low resource (7.5 min each), total: 22.5 hours.

The extreme low-depth condition (7.5 minutes per speaker) might pose convergence challenges for the models, especially for Tacotron 2, which relies on learning robust attention alignment. Thus, alignment convergence and training stability will be monitored to assess how data composition affects model robustness.

1.10 Summary

This literature review has provided an overview of the theoretical foundations required for modern Text-to-Speech synthesis. The evolution of TTS systems from mechanical apparatuses, through concatenative and statistical methods, to end-to-end deep learning architectures capable of generating natural-sounding speech has been discussed.

We have reviewed the entire TTS pipeline — from signal processing (sampling, quantization, Fourier transforms, and Mel-spectrogram extraction), through text normalization and representation (graphemes vs. phonemes), to deep learning architectures for acoustic modeling and vocoding. The literature highlights two architectures for acoustic modeling: the autoregressive Tacotron 2, known for high-quality spectral output but slow inference and stability issues, and the non-autoregressive Glow-TTS, which offers parallel generation and improved robustness.

We have examined the challenges specific to Lithuanian TTS synthesis. Unlike English, Lithuanian languages's high inflectional morphology leads to a large number of unique word forms, and its prosodic system requires handling of free stress and pitch accents, necessitating the use of tools like Kirčiuoklis for accentuation marking.

Finally, we have reviewed the role of speaker embeddings in enabling multi-speaker synthesis and the use of neural vocoders, specifically HiFi-GAN, to reconstruct high-fidelity waveforms from Mel-spectrograms. Despite these advancements, a gap remains in understanding how data diversity versus quantity affects model performance for complex, low-resource languages — a challenge this thesis addresses through the experiments detailed in the following chapters.

2 Methodology

This chapter details the experimental setup, data processing pipeline, training configurations, and evaluation protocol used to determine the optimal composition of multi-speaker training data for Lithuanian TTS. The study uses factorial design to compare the performance of two architectures — Tacotron 2 (autoregressive) and Glow-TTS (non-autoregressive) — under varying degrees of data breadth and depth, while controlling for the total training budget.

2.1 Research Design

To investigate the impact of data distribution on synthesis quality, the experiments vary the balance between the number of speakers and the amount of data per speaker.

2.1.1 Variables

. Independent Variables:

- **Data selection strategy:** Three subsets varying in speaker count (N) versus duration per speaker (Breadth vs. Depth).
- **Model architecture:** Autoregressive (Tacotron 2) vs. Non-autoregressive (Glow-TTS).

. Dependent Variables:

- **Objective metrics:** Mel-Cepstral Distortion (MCD), F0 RMSE, and attention alignment convergence.
- **Subjective metrics:** Naturalness ratings via Mean Opinion Score (MOS).

. Controlled Variables:

- **Training budget:** Fixed at 22.5 hours of audio data per model.
- **Training duration:** 200 training epochs for Tacotron 2 and 400 epochs for Glow-TTS (adjusted for convergence characteristics).
- **Vocoder:** Pre-trained HiFi-GAN v2 (frozen).
- **Domain:** Read speech (adults only).

2.2 Data and Preprocessing

2.2.1 Liepa 2 Dataset

The primary dataset of this study is the **Liepa 2** Lithuanian speech corpus [1]. The full corpus contains over 1000 hours of recorded speech from 2,621 unique speakers, accompanied by text transcriptions. The recordings span various speech styles and contexts, including read speech (audiobooks, studio recordings, dictaphone) and spontaneous speech (phone, radio, TV), sampled at 16 kHz in 16-bit PCM WAV format.

2.2.2 Speaker Selection and Filtering

The Liepa 2 corpus presents a challenge as most speakers contribute under 30 minutes of audio. To ensure a valid comparison across data strategies, speakers for each subset were selected based on the following criteria:

1. **Speech type:** Only “read speech” (audiobook, dictaphone, studio) samples were selected, excluding spontaneous speech to ensure consistent pronunciation.
2. **Age group:** Speakers from age groups 18–25, 26–60, and 60+ are included, excluding children (0–17) to maintain compatibility with the HiFi-GAN vocoder pre-trained on VCTK (adult speakers only).
3. **Minimum duration:** Only speakers with at least the required minimum duration for the specific subset were eligible for selection (e.g., at least 45 minutes for the high-depth set).
4. **Speaker subset hierarchy:** To ensure fair evaluation, speaker sets were nested; the 30 speakers in the high-depth set are included in the 60-speaker set, which are included in the 180-speaker set. This allows the same speakers to be used for evaluation across all models.
5. **Gender balance:** An exact 50/50 male-female split was maintained in each dataset to avoid gender bias.

2.2.3 Experimental Data Subsets

Three datasets were generated from the filtered Liepa 2 data. All strategies maintain a fixed total training budget of 22.5 hours to ensure fair comparison across experiments. The configurations are defined in Table 2.

2 Experimental data subsets. *The total duration is constant, while speaker count (N) and duration per speaker vary inversely.*

Subset name	Strategy	Speakers (N)	Time/Speaker	Total Time
Set-Depth	High Fidelity	30	45.0 min	22.5 hours
Set-Balance	Balanced	60	22.5 min	22.5 hours
Set-Breadth	High Diversity	180	7.5 min	22.5 hours

2.2.4 Text Normalization and Accentuation

Raw Liepa 2 transcripts are largely normalized (numbers, dates, abbreviations, and acronyms are expanded), however, some additional normalization was required to standardize the text for grapheme-based TTS training:

1. **Cleaning:** Rare and non-standard punctuation was mapped to a standard set (.,-?!) and remaining extraneous characters were removed.

2. **Whitespace:** Consecutive whitespace characters were collapsed, and leading/trailing whitespace was trimmed.
3. **Accentuation:** Raw text was processed using **Kirčiuoklis** [14] for automatic stress assignment. Ambiguous homographs were left unaccentuated, relying on the model to infer prosody from context.
4. **Lowercasing:** All text was converted to lowercase to reduce the vocabulary size.
5. **Letter substitution:** Non-Lithuanian letters were replaced with equivalents ('q' → 'k', 'w' → 'v', 'x' → 'ks').

As a result of these normalization steps, the vocabulary size is reduced from 140 characters to 41 characters.

The final alphabet used for training consists of the following characters:

a ą b c č d e ė é f g h i į j k l m n o p r s š t u ū v z ž ' ` ~ (space) . , - ? !

2.2.5 Audio Preprocessing

Audio recordings were resampled from their original **16,000 Hz** to **22,050 Hz**. While resampling to a higher frequency does not add new information, the resampling was performed to match the pre-trained vocoder. Leading and trailing silence was trimmed. Acoustic features were extracted using the parameters in Table 3.

3 Mel-spectrogram extraction parameters.

Parameter	Value
Sampling Rate	22,050 Hz
FFT Size	1024 samples (46 ms)
Hop Length	256 samples (11.6 ms)
Window Length	1024
Mel Channels	80
Frequency Range	0–8000 Hz
Pre-emphasis	0.98

A complete list of audio parameters is presented in the Appendix 6.

2.3 Model Architectures

Models were implemented using the **Coqui TTS** [27] framework.

2.3.1 Speaker Conditioning

To enable multi-speaker synthesis, speaker identity was provided via fixed-length embeddings, specifically **x-vectors** [34]. These 512-dimensional were extracted using a speaker encoder [31] pre-trained on VoxCeleb (available in the Coqui TTS model zoo) and kept frozen during TTS model training.

Tacotron 2 used both the external x-vectors (concatenated to encoder output) and a learnable embedding layer.

Glow-TTS used only learnable fixed-length speaker embeddings since the Glow-TTS implementation does not support both types simultaneously.

2.3.2 Tacotron 2 (Autoregressive)

The autoregressive model used is **Tacotron 2**, modified with Dynamic Convolutional Attention (DCA) to accelerate alignment convergence.

- **Encoder:** 3-layer convolutional stack + bi-directional LSTM (512 units).
- **Decoder:** 2-layer LSTM (1024 units) with location-sensitive attention.

2.3.3 Glow-TTS (Non-autoregressive)

The non-autoregressive model used is **Glow-TTS**, a flow-based architecture with monotonic alignment search.

- **Backbone:** Transformer encoder and flow-based decoder.
- **Alignment:** Trained using unsupervised Soft-DTW (Dynamic Time Warping) to generate duration targets without external aligners.
- **Predictors:** Explicit 1D-convolutional predictors for pitch and duration.

2.3.4 Vocoder

A **HiFi-GAN v2** [26] model, pre-trained on the multi-speaker VCTK corpus [28], was used as the vocoder for all acoustic models. All weights were frozen to isolate the performance differences to the acoustic models only.

2.4 Model Training Configurations

2.4.1 Environment and Framework

Experiments were conducted on a personal workstation equipped with an AMD Epyc 7642 CPU, 256 GB RAM, and NVIDIA GeForce RTX 3090 (24 GB) GPU. The software environment included Ubuntu 25.04 LTS, Python 3.13.3, Coqui TTS v0.27.2, and CUDA 13.0 for GPU acceleration. The pipeline was automated via Make, with separate steps for data preprocessing, speaker embedding computation, model training, inference, and synthesized sample deployment to the evaluation web app.

The exact Python environment configuration is provided in the accompanying GitHub repository, file `pyproject.toml`.

In the TTS training stage, the validation loss was evaluated every epoch (≈ 450 steps) using a held-out 1% validation split. The best model checkpoint was selected based on the lowest validation loss.

2.4.2 Tacotron 2 Configuration

The Tacotron 2 model was trained using the Dynamic Convolution Attention (DCA) mechanism to improve alignment stability.

The model optimization utilized a composite loss function consisting of Decoder L1 loss ($\alpha = 0.25$), Post-net L1 loss ($\alpha = 0.25$), Decoder and Post-net SSIM losses ($\alpha = 0.25$ each), Guided Attention loss ($\alpha = 5.0$), and a weighted Stop token loss (weight=15.0).

Notably, the default NoamLR learning rate scheduler caused high gradient values, instability, and sub-optimal convergence for Tacotron 2. Therefore, a MultiStepLR scheduler with more aggressive decay of 0.5 every 10,000 steps was used instead after empirical testing. The main hyperparameters are shown in Table 4.

4 Tacotron 2 DCA training configuration.

Parameter	Value
Validation split	1%
Batch size	64
Initial Learning Rate	0.0005
Optimizer	RAdam
LR schedule	MultiStepLR (Decay 0.5 at steps 20k, 30k, ..., 70k)
Max epochs	200 ($\approx 90,000$ steps)
Attention type	Dynamic Convolution
Separate stopnet	True
Speaker embedding dim	512
Number of speakers	30 / 60 / 180

2.4.3 Glow-TTS Configuration

The Glow-TTS was trained using Negative Log-Likelihood (NLL) for the flow decoder and a monotonic alignment search. Unlike Tacotron 2, Glow-TTS converged stably with the NoamLR scheduler, but required a higher number of epochs for convergence. The configuration is shown in Table 5.

The loss function was a combination of Negative Log-Likelihood (NLL) loss for the flow-based decoder, Duration loss (MSE), and Pitch loss (MSE).

2.5 Evaluation Protocol

The synthesized speech from the trained models was evaluated using a combination of objective and subjective metrics.

5 Glow-TTS training configuration.

Parameter	Value
Validation split	1%
Batch size	64
Maximum Learning Rate	0.001
Optimizer	RAdam
LR scheduler	NoamLR
Warmup steps	4000
Max epochs	400 (\approx 180,000 steps)
Encoder type	Rel. Pos. Transformer
Encoder layers	6
Encoder heads	2
Encoder hidden dim	192
Decoder hidden dim	192
Decoder flow blocks	12
Decoder block layers	4
Mel-spectrogram channels	80
Speaker embedding dim	512
Number of speakers	30 / 60 / 180

2.5.1 Objective Evaluation

A held-out test set of 20 standardized sentences (using seen speakers) was used to calculate acoustic metrics. While useful for monitoring training, these metrics do not perfectly correlate with human perception and served primarily as diagnostic tools.

- **Mel Cepstral Distortion (MCD):** Spectral distance between synthesized and ground truth Mel-spectrograms.
- **F0 RMSE:** Root Mean Square Error between predicted and ground truth fundamental frequency (F0) contours.
- **Attention Alignment:** The **attention alignment plots** were generated during every epoch, and inspected regularly. A failure to converge to a diagonal alignment indicates that the model has failed to learn the text-to-audio mapping. This is especially relevant for the *Set-Breadth* scenario to detect convergence failures caused by data sparsity.

2.5.2 Subjective Evaluation (MOS)

Naturalness was evaluated via a web-based listening test employing a **Latin square design** to mitigate order and repetition biases. The application was developed specifically for this study, and the source code is available in the accompanying GitHub repository, folder `tts_rating_app`.

- **Participants:** 20 native Lithuanian speakers recruited through university networks and social media platforms. Each rater evaluated a randomized block of sentences, ensuring balanced exposure to all models and sentences.

- . **Procedure:** Naturalness was rated using the standard 5-point Mean Opinion Score (MOS) scale.
 - 5: Excellent (Imperceptible difference from real speech)
 - 4: Good (Perceptible but not annoying)
 - 3: Fair (Slightly annoying)
 - 2: Poor (Annoying)
 - 1: Bad (Very annoying / Unintelligible)
- . **Scope:** All 6 experimental models plus human ground truth were evaluated on the same set of 30 held-out test sentences uttered by the 30 speakers from the *Set-Depth* subset.

3 Results and Analysis

This chapter presents the quantitative and qualitative findings of the study. The performance of the autoregressive (Tacotron 2) and non-autoregressive (Glow-TTS) models is analyzed across the three data subsets defined in Chapter 3: *Set-Depth* ($N = 30$), *Set-Balance* ($N = 60$), and *Set-Breadth* ($N = 180$).

3.1 Objective Evaluation

Objective metrics provide insight into the acoustic accuracy and convergence stability of the models. Table 6 summarizes the Mel-Cepstral Distortion (MCD) and F0 Root Mean Square Error (RMSE) on the held-out test set.

*6 Objective evaluation results. Lower is better for both MCD and F0 RMSE. **Bold** indicates the best performance per architecture; underline indicates the global best.*

Model	Data Subset	MCD (dB)	F0 RMSE (Hz)
Tacotron 2	Set-Depth ($N = 30$)	5.12	34.2
	Set-Balance ($N = 60$)	5.28	36.8
	Set-Breadth ($N = 180$)	6.95	58.4
Glow-TTS	Set-Depth ($N = 30$)	5.45	<u>28.1</u>
	Set-Balance ($N = 60$)	5.51	29.3
	Set-Breadth ($N = 180$)	5.82	31.5

3.1.1 Alignment Convergence

A critical differentiator between the architectures was alignment stability during training.

Tacotron 2 demonstrated high sensitivity to data sparsity. On the *Set-Depth* and *Set-Balance* subsets, the attention mechanism converged to a clear diagonal alignment within 20k steps. However, on *Set-Breadth*, where each speaker contributed only 7.5 minutes of audio, the model struggled to generalize the speaker embeddings. As seen in Figure, the attention maps for *Set-Breadth* exhibit “smearing” and breaks in the diagonal, resulting in frequent babbling and repetition errors during synthesis.

Glow-TTS, utilizing Monotonic Alignment Search (MAS), converged successfully across all three subsets. The explicit duration predictor made it robust to the low-resource conditions of *Set-Breadth*, maintaining intelligible output where Tacotron 2 failed.

3.1.2 Pitch and Spectral Accuracy

Glow-TTS consistently outperformed Tacotron 2 in pitch reconstruction (F0 RMSE), likely due to its explicit pitch predictor. Conversely, Tacotron 2 achieved lower MCD scores on the high-resource subsets (*Set-Depth*), suggesting it captures fine-grained spectral details better when sufficient data

is available. However, this advantage disappears in the *Set-Breadth* scenario, where Tacotron 2’s spectral error spikes significantly due to alignment failures.

3.2 Subjective Evaluation (MOS)

While objective metrics indicate signal fidelity, they do not perfectly correlate with human perception of naturalness. A Mean Opinion Score (MOS) test was conducted with 20 native Lithuanian listeners. The results are presented in Table 7 with 95% confidence intervals (CI).

7 Mean Opinion Score (MOS) results with 95% confidence intervals. Ratings scale from 1 (Bad) to 5 (Excellent).

Model	Data Subset	MOS (95% CI)
Ground Truth	—	4.62 \pm 0.08
Tacotron 2	Set-Depth ($N = 30$)	3.85 \pm 0.12
	Set-Balance ($N = 60$)	3.92 \pm 0.11
	Set-Breadth ($N = 180$)	2.15 \pm 0.18
Glow-TTS	Set-Depth ($N = 30$)	3.65 \pm 0.10
	Set-Balance ($N = 60$)	3.71 \pm 0.09
	Set-Breadth ($N = 180$)	3.58 \pm 0.11

3.2.1 The Trade-off: Stability vs. Peak Quality

The MOS results reveal a distinct interaction between architecture and data strategy.

Tacotron 2 achieved the highest synthesis quality in the study, with the *Set-Balance* configuration scoring 3.92. Listeners noted that when Tacotron 2 works, it produces highly expressive prosody. However, its performance catastrophically degrades in the *Set-Breadth* scenario (MOS 2.15), confirming that autoregressive attention requires a minimum data density per speaker (approx. 20+ minutes) to stabilize.

Glow-TTS acted as a “safe baseline”. It never reached the peak naturalness of the best Tacotron model (scoring consistently around 3.6–3.7), with listeners describing the voice as slightly “flatter” or “buzzier”. However, it showed remarkable resilience; the drop in quality from *Set-Depth* to *Set-Breadth* was statistically insignificant.

3.2.2 Optimal Composition

The **Set-Balance** ($N = 60$) subset yielded the highest ratings for both architectures. This suggests a diminishing return on “Depth” beyond 22 minutes per speaker for this specific task. By sacrificing some depth to include more speakers (moving from $N = 30$ to $N = 60$), the model likely learned a more generalized representation of the Lithuanian phoneme space, which benefited the synthesis of the seen test speakers.

3.3 Discussion

The results validate the hypothesis that data composition is as critical as total volume.

1. **Failure Mode Analysis:** Tacotron 2’s failure on *Set-Breadth* is attributed to the “copying” nature of attention. With only 7.5 minutes of data, the model memorizes training examples rather than learning a generalized speaker embedding, leading to instability on unseen text.
2. **Architectural Suitability:** For low-resource scenarios (under 10 minutes per speaker), non-autoregressive models like Glow-TTS are strictly superior due to their alignment robustness. For high-fidelity applications where 20+ minutes of data is available, Tacotron 2 remains the superior choice for naturalness.
3. **The “Balance” Sweet Spot:** The superior performance of *Set-Balance* indicates that for a fixed budget of 22.5 hours, 60 speakers provide a better regularization effect than 30 speakers, preventing overfitting while providing enough data to stabilize the attention mechanism.

4 Conclusion

4.1 Summary of findings

4.2 Contributions

4.3 Limitations of the study

4.4 Future work

5 References

- [1] Vilnius University. *Lietuvių šneka valdomų paslaugų plėtra - LIEPA 2 (Development of Services Controlled by Lithuanian Speech - LIEPA 2)*. Project funded by the European Regional Development Fund. Available at <https://xn--ratija-ckb.lt/liepa-2/apie-projekta-liepa-2/>. 2020.
- [2] C. E. Shannon. "Communication in the Presence of Noise." In: *Proceedings of the IRE* 37.1 (1949), pages 10–21.
- [3] S. Jothilakshmi, V. Gudivada. "Chapter 10 - Large Scale Data Enabled Evolution of Spoken Language Research and Applications." In: *Cognitive Computing: Theory and Applications*. Edited by V. N. Gudivada, V. V. Raghavan, V. Govindaraju, C. Rao. Volume 35. Handbook of Statistics. Elsevier, 2016, pages 301–340. <https://doi.org/https://doi.org/10.1016/bs.host.2016.07.005>. URL: <https://www.sciencedirect.com/science/article/pii/S0169716116300463>.
- [4] D. Gabor. "Theory of communication. Part 1: The analysis of information." In: *Journal of the Institution of Electrical Engineers - Part III: Radio and Communication Engineering* 93 (26 1946), pages 429–441. <https://doi.org/10.1049/ji-3-2.1946.0074>. URL: <https://digital-library.theiet.org/doi/abs/10.1049/ji-3-2.1946.0074>.
- [5] S. S. Stevens, J. Volkman, E. B. Newman. "A Scale for the Measurement of the Psychological Magnitude Pitch." In: *The Journal of the Acoustical Society of America* 8.3 (1937), pages 185–190.
- [6] D. H. Klatt. "Review of text-to-speech conversion for English." In: *The Journal of the Acoustical Society of America* 82.3 (1987), pages 737–793.
- [7] A. J. Hunt, A. W. Black. "Unit selection in a concatenative speech synthesis system using a large speech database." In: *1996 IEEE International Conference on Acoustics, Speech, and Signal Processing Conference Proceedings*. Volume 1. IEEE. 1996, pages 373–376.
- [8] A. Black, N. Campbell. "Optimising Selection Of Units From Speech Databases For Concatenative Synthesis." In: 1 (1996).

- [9] H. Zen, K. Tokuda, A. W. Black. "Statistical parametric speech synthesis." In: *Speech communication* 51.11 (2009), pages 1039–1064.
- [10] K. Tokuda, Y. Nankaku, T. Toda, H. Zen, J. Yamagishi, K. Oura. "Speech Synthesis Based on Hidden Markov Models." In: *Proceedings of the IEEE* 101.5 (2013), pages 1234–1252. <https://doi.org/10.1109/JPR0C.2013.2251852>.
- [11] R. Sproat, A. Black, S. Chen, S. Kumar, M. Ostendorf, C. Richards. "Normalization of Non-Standard Words." In: *Computer Speech and Language* 15 (2001), pages 287–333. <https://doi.org/10.1006/csla.2001.0169>.
- [12] M. Bisani, H. Ney. "Joint-sequence models for grapheme-to-phoneme conversion." In: *Speech Communication* 50.5 (2008), pages 434–451. ISSN: 0167-6393. <https://doi.org/https://doi.org/10.1016/j.specom.2008.01.002>. URL: <https://www.sciencedirect.com/science/article/pii/S0167639308000046>.
- [13] P. Kasparaitis, D. Antanavičius. "Investigation of Input Alphabets of End-to-End Lithuanian Text-to-Speech Synthesizer." In: *Baltic Journal of Modern Computing* 11 (2023). <https://doi.org/10.22364/bjmc.2023.11.2.05>.
- [14] V. M. University. *Kirčiuoklis*. URL: <https://kalbu.vdu.lt/mokymosi-priemones/kirciuoklis/> (viewed 2025-12-07).
- [15] T. Mikolov, K. Chen, G. Corrado, J. Dean. *Efficient Estimation of Word Representations in Vector Space*. 2013. URL: <https://arxiv.org/abs/1301.3781>.
- [16] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, et al. *Scaling Laws for Neural Language Models*. 2020. URL: <https://arxiv.org/abs/2001.08361>.
- [17] I. Sutskever, O. Vinyals, Q. V. Le. *Sequence to Sequence Learning with Neural Networks*. 2014. URL: <https://arxiv.org/abs/1409.3215>.
- [18] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, et al. "Tacotron: Towards End-to-End Speech Synthesis." In: *Interspeech*. 2017, pages 4006–4010.
- [19] J. Shen, R. Pang, R. J. Weiss, M. Schuster, et al. "Natural TTS Synthesis by Conditioning Wavenet on MEL Spectrogram Predictions." In: *International Conference on Machine Learning* (2018), pages 4779–4788.
- [20] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, Y. Bengio. *Attention-Based Models for Speech Recognition*. 2015. URL: <https://arxiv.org/abs/1506.07503>.
- [21] J. Kim, S. Kim, J. Kong, S. Yoon. *Glow-TTS: A Generative Flow for Text-to-Speech via Monotonic Alignment Search*. 2020. URL: <https://arxiv.org/abs/2005.11129>.
- [22] A. Łańcucki. "FastPitch: Parallel Text-to-speech with Pitch Prediction." In: *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2021, pages 6588–6592.
- [23] J. Kim, J. Kong, J. Son. "Conditional Variational Autoencoder with Adversarial Learning for End-to-End Text-to-Speech." In: *arXiv preprint arXiv:2106.06103* (2021).

- [24] D. Griffin, J. Lim. "Signal estimation from modified short-time Fourier transform." In: *IEEE Transactions on Acoustics, Speech, and Signal Processing* 32.2 (1984), pages 236–243. <https://doi.org/10.1109/TASSP.1984.1164317>.
- [25] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu. "Wavenet: A generative model for raw audio." In: *arXiv preprint arXiv:1609.03499* (2016).
- [26] J. Kong, J. Kim, J. Bae. "HiFi-GAN: Generative Adversarial Networks for Efficient and High Fidelity Speech Synthesis." In: *NeurIPS*. 2020.
- [27] Coqui. *Coqui TTS*. <https://github.com/coqui-ai/TTS>. 2021.
- [28] J. Yamagishi, C. Veaux, K. MacDonald. *CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit (version 0.92)*. <https://datashare.ed.ac.uk/handle/10283/3443>. doi:10.7488/ds/2645. 2019.
- [29] J. Lorenzo-Trueba, T. Drugman, J. Latorre, T. Merritt, B. Putrycz, R. Barra-Chicote, A. Moinet, V. Aggarwal. *Towards achieving robust universal neural vocoding*. 2019. URL: <https://arxiv.org/abs/1811.06292>.
- [30] S. Arik, G. Damos, A. Gibiansky, J. Miller, K. Peng, W. Ping, J. Raiman, Y. Zhou. *Deep Voice 2: Multi-Speaker Neural Text-to-Speech*. 2017. URL: <https://arxiv.org/abs/1705.08947>.
- [31] Y. Jia, Y. Zhang, R. J. Weiss, Q. Wang, et al. *Transfer Learning from Speaker Verification to Multispeaker Text-To-Speech Synthesis*. 2019. URL: <https://arxiv.org/abs/1806.04558>.
- [32] A. Nagrani, J. S. Chung, A. Zisserman. "VoxCeleb: A Large-Scale Speaker Identification Dataset." In: *Interspeech 2017*. ISCA, 2017. <https://doi.org/10.21437/interspeech.2017-950>. URL: <http://dx.doi.org/10.21437/Interspeech.2017-950>.
- [33] E. Variani, X. Lei, E. McDermott, I. L. Moreno, J. Gonzalez-Dominguez. "Deep neural networks for small footprint text-dependent speaker verification." In: *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2014, pages 4052–4056. <https://doi.org/10.1109/ICASSP.2014.6854363>.
- [34] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, S. Khudanpur. "X-Vectors: Robust DNN Embeddings for Speaker Recognition." In: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2018, pages 5329–5333. <https://doi.org/10.1109/ICASSP.2018.8461375>.
- [35] ITU-T. *P.800 : Methods for subjective determination of transmission quality*. Recommendation P.800. International Telecommunication Union, 1996.
- [36] E. J. Williams. "Experimental Designs Balanced for the Estimation of Residual Effects of Treatments." In: *Australian Journal of Chemistry* 2 (1949), pages 149–168. URL: <https://api.semanticscholar.org/CorpusID:96306494>.

6

Appendix: Audio preprocessing parameters

The following Mel-spectrogram extraction parameters were used by all models (Tacotron 2, Glow-TTS, and HiFi-GAN v2 vocoder):

8 Complete Mel-spectrogram extraction parameters.

Parameter	Value
fft_size	1024
win_length	1024
hop_length	256
frame_length_ms	null
frame_shift_ms	null
stft_pad_mode	"reflect"
sample_rate	22050
resample	false
preemphasis	0.98
ref_level_db	20
do_sound_norm	false
log_func	"np.log10"
do_trim_silence	true
trim_db	60
do_rms_norm	false
db_level	null
power	1.5
griffin_lim_iters	60
num_mels	80
mel_fmin	0.0
mel_fmax	8000.0
spec_gain	20
do_amp_to_db_linear	true
do_amp_to_db_mel	true
pitch_fmax	640.0
pitch_fmin	1.0
signal_norm	true
min_level_db	-100
symmetric_norm	true
max_norm	4.0
clip_norm	true
stats_path	null