

MatchaTTS Implementation And Analysis

ADJAL Massyl

MASSYL.ADJAL@ETU.SORBONNE-UNIVERSITE.FR

*Master Ingénierie des Systèmes Intelligents
Sorbonne Université
Paris, France*

BOUHAI Yasser

YASSER.BOUHAI@ETU.SORBONNE-UNIVERSITE.FR

*Master Ingénierie des Systèmes Intelligents
Sorbonne Université
Paris, France*

BOUHADJIRA Zineddine

ZINEDDINE.BOUHADJIRA@ETU.SORBONNE-UNIVERSITE.FR

*Master Ingénierie des Systèmes Intelligents
Sorbonne Université
Paris, France*

BOULARAS Mohamed Mouad

MOHAMED.MOUAD.BOULARAS@ETU.SORBONNE-UNIVERSITE.FR

*Master Ingénierie des Systèmes Intelligents
Sorbonne Université
Paris, France*

Editor: Sorbonne Université - Master Ingénierie des Systèmes Intelligents - Year 2 - Machine Learning Avancé (2025-2026)

Abstract

This project focuses on the reproducibility of Matcha-TTS, a state-of-the-art text-to-speech architecture based on Optimal-Transport Conditional Flow Matching (OT-CFM). We present a complete reimplementation of the model, featuring a Transformer-based text encoder with Rotational Position Embeddings (RoPE) and a lightweight 1D U-Net decoder for acoustic feature generation. To validate our approach, we conducted a comprehensive comparative analysis between the official pre-trained checkpoint and three distinct variations of our implementation. All models were evaluated on the LJ Speech dataset using both objective metrics and subjective listening tests. Our results demonstrate the efficacy of the flow matching paradigm and quantify the performance trade-offs across our different implementation versions relative to the reference work.

Keywords: Text-to-Speech, Flow Matching, Deep Learning, Reproducibility, Optimal Transport

Code Repository : Our implementation and experiments are available at : <https://github.com/Zineddine/MatchaTTS-Implementation-Analysis>

1. Introduction

Neural Text-to-Speech (TTS) synthesis is a key area of research in natural language processing, focusing on generating intelligible, natural, and expressive speech from textual input. High-quality TTS has applications in voice assistants, accessibility tools, and content creation, making it an important challenge for both academia and industry. Recent models leveraging probabilistic approaches, such as diffusion-based methods and ODE-inspired frameworks, have advanced speech realism, but practical deployment remains difficult.

The task addressed in this work is acoustic modeling for high-quality TTS, where the goal is to transform a sequence of phonemes into a mel-spectrogram representation that can be converted into an audio waveform by a vocoder. A robust TTS model must ensure naturalness, maintain intelligibility, synthesize efficiently, minimize computational resources, and learn monotonic text-to-audio alignment without relying on external aligners.

Despite progress, existing approaches face several challenges. Diffusion models require many iterative steps for high-quality synthesis, resulting in slow inference. Fast non-probabilistic models, while faster, often compromise speech variability and naturalness. Memory-intensive architectures and positional encoding schemes struggle with long sequences. Some models rely on external alignments, complicating training. Probabilistic models must also learn complex, time-dependent vector fields, demanding multiple network evaluations during synthesis. These limitations highlight the need for methods that combine speed, efficiency, and high-quality speech generation.

2. Architecture Overview

Matcha-TTS is a non-autoregressive encoder-decoder architecture designed to synthesize high-quality mel-spectrograms from input text. As illustrated in the global workflow, the system is composed of two main subsystems that function sequentially during inference :

1. **The Encoder (Text Processing)** : First, the input text is converted into a sequence of phonemes. The encoder processes these phonemes to create a linguistic representation. A duration predictor then estimates how long each phoneme should be spoken, and the encoder’s output is "upsampled" (repeated) to match the total duration of the target speech. This results in a sequence of average acoustic features, denoted as μ , which serves as a "rough draft" or conditioning signal for the speech.
2. **The Decoder (Acoustic Generation)** : Unlike traditional models that generate speech frame-by-frame, Matcha-TTS uses a continuous-time Flow Matching process. The decoder transforms a simple noise distribution (x_0) into a complex mel-spectrogram (x_1) by solving an Ordinary Differential Equation (ODE). It uses the encoder’s output μ as a map to guide this transformation, ensuring the noise is shaped into the correct speech sounds specified by the text.

Ideally, this separation allows the model to learn linguistic prosody (Encoder) independently from the fine-grained acoustic details (Decoder), joined together by the alignment mechanism.

The following sections detail the specific implementation of these components, starting with the text encoder and followed by the flow-matching decoder.

2.1 Encoder

The Matcha-TTS encoder is composed of two main modules that work in tandem to transform input text into rich, temporally-aligned representations : the **Text Encoder** and the **Duration Predictor**.

2.1.1 TEXT ENCODER

The Text Encoder is responsible for transforming raw text into contextualized embeddings. Its architecture follows a multi-stage pipeline (see Figure 1) :

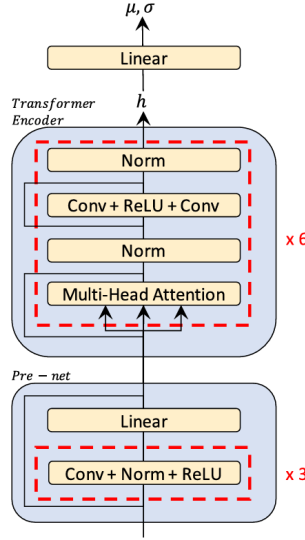


FIGURE 1 – Text Encoder architecture inspired by Glow-TTS and Grad-TTS. The text is first tokenized, then processed by a Prenet before being enriched by transformer layers with RoPE.

Architectural Inspiration The Text Encoder architecture draws inspiration from Glow-TTS, particularly the use of convolutional preprocessing layers (Prenet) and transformer blocks. The implementation follows established design patterns from flow-based TTS models, incorporating convolutional layers for local feature extraction followed by self-attention mechanisms for capturing long-range dependencies in the phoneme sequence.

Text Processing and Tokenization The input text first undergoes a cleaning and normalization process via *text cleaners* (e.g., `english_cleaners2`). Characters are then converted into numerical tokens based on a phonetic symbol vocabulary. To improve model stability, silence tokens (blank tokens) are inserted between each phoneme via an interspersion function.

Prenet - Embedding Preprocessing Before entering the transformer layers, tokens pass through a **Prenet**, which serves as a preprocessing layer. The Prenet consists of a stack of 1D convolutional layers followed by normalizations and activations :

- **1D Convolution** : Extracts local features from token embeddings
- **Layer Normalization** : Stabilizes training
- **ReLU Activation** : Introduces non-linearity
- **Dropout** : Regularization to prevent overfitting

This architecture enables the creation of richer initial representations before transformer processing.

RoPE Positional Encoding Matcha-TTS uses **RoPE (Rotary Position Embedding)** as its positional encoding method. Unlike classical positional embeddings, RoPE encodes position by applying rotation in the feature space. A distinctive feature of the Matcha implementation is that RoPE is only applied to **half of the embedding dimensions**, thus preserving part of the raw information while injecting positional information.

Transformer Stack The core of the Text Encoder consists of a stack of transformer blocks with several specific adaptations :

- **Multi-Head Self-Attention** : Allows the model to capture long-range dependencies in the phoneme sequence
- **Feed-Forward Network** : Non-linear transformations to enrich representations
- **Residual Connections** : Facilitate gradient flow and enable deeper training
- **Layer Normalization** : Applied before each sub-layer (Pre-LN architecture)

This transformer architecture allows the model to learn rich contextual representations that capture prosody, intonation, and phonetic relationships of the input text.

2.1.2 DURATION PREDICTOR AND MONOTONIC ALIGNMENT SEARCH

The **Duration Predictor** is a critical component that learns to predict how many mel-spectrogram frames each phoneme should occupy. This enables temporal alignment between text and audio without manual annotations.

Problem Statement Given a text sequence with T_{text} phonemes and an audio sequence with T_{mel} mel frames (where $T_{\text{mel}} \gg T_{\text{text}}$), we need to determine :

- Which frames correspond to which phoneme ?
- How long does each phoneme last ?

This is the **alignment problem** : finding the mapping $A : \{1, \dots, T_{\text{mel}}\} \rightarrow \{1, \dots, T_{\text{text}}\}$ between frames and phonemes.

Monotonic Alignment Search (MAS) MAS automatically discovers the optimal alignment during training by maximizing the likelihood between predicted phoneme representations (μ_x) and real audio (y).

Key idea : For each phoneme i and frame j , compute compatibility :

$$L_{i,j} = -\frac{1}{2} \|y_j - \mu_{x,i}\|^2 \quad (1)$$

where $\mu_{x,i} \in \mathbb{R}^{80}$ is the encoder’s predicted mel representation for phoneme i , and $y_j \in \mathbb{R}^{80}$ is the real mel frame at time j .

MAS uses **dynamic programming** to find the monotonic path through L that maximizes total likelihood :

$$A^* = \arg \max_A \sum_{j=1}^{T_{\text{mel}}} L_{A(j),j} \quad (2)$$

subject to :

- **Monotonicity** : $A(j) \leq A(j+1)$ (no going backward in time)
- **Surjectivity** : Every phoneme must be used at least once

Complexity : $\mathcal{O}(T_{\text{text}} \times T_{\text{mel}})$ time, cython version making it efficient even for long sequences.

Output : Alignment matrix $A^* \in \{0, 1\}^{T_{\text{text}} \times T_{\text{mel}}}$ from which durations are extracted :

$$d_i = \sum_{j=1}^{T_{\text{mel}}} A_{i,j}^* \quad (3)$$

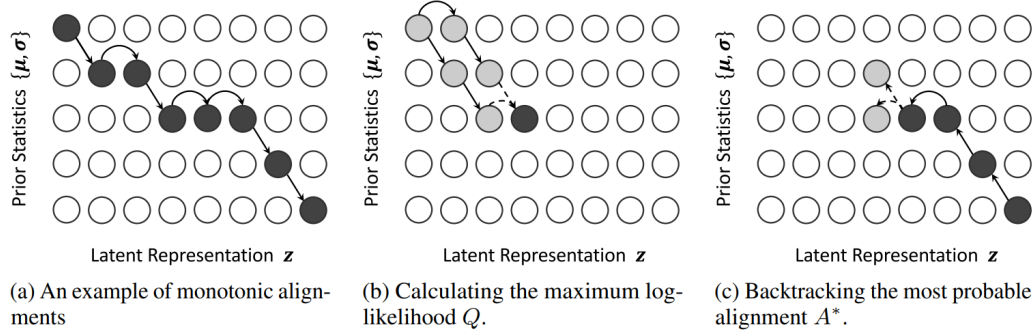


FIGURE 2 – MAS algorithm

Duration Predictor Architecture The Duration Predictor is a lightweight neural network that learns to predict phoneme durations from text representations.

Architecture :

- **Input** : Encoder hidden states $h \in \mathbb{R}^{192 \times T_{\text{text}}}$
- **Conv1D layers** ($2 \times$) : Local temporal feature extraction
- **Layer Normalization + Dropout** : Regularization
- **Projection** : $\mathbb{R}^{256} \rightarrow \mathbb{R}^1$ (one duration per phoneme)
- **Output** : $\log w \in \mathbb{R}^{1 \times T_{\text{text}}}$ (log-durations)

Why log-space ? Durations vary widely (1-100 frames). Log-space :

- Compresses the scale : $[1, 100] \xrightarrow{\log} [0, 4.6]$
- Guarantees positivity : $w = \exp(\log w) > 0$
- Improves numerical stability

Training : MAS Supervises Duration Predictor **Key insight** : MAS provides ground truth durations from real audio, which are used to train the Duration Predictor.

Loss function :

$$\mathcal{L}_{\text{dur}} = \frac{1}{T_{\text{text}}} \sum_{i=1}^{T_{\text{text}}} \left(\log w_i^{\text{pred}} - \log d_i^{\text{MAS}} \right)^2 \quad (4)$$

where :

- w_i^{pred} : Duration predicted by Duration Predictor
- d_i^{MAS} : Duration extracted from MAS alignment (ground truth)

Gradient Detachment (Critical Design Choice) The Duration Predictor uses a **stop gradient** operation to prevent interference with encoder training :

$$\log w = f_{\text{dur}}(\text{sg}[h]) \quad (5)$$

where $\text{sg}[\cdot]$ blocks gradients : $\frac{\partial \text{sg}[h]}{\partial h} = 0$.

Each component specializes in its task without interference.

Inference : Duration Predictor Replaces MAS At inference, MAS cannot be used (no real audio available). The Duration Predictor takes over :

1. Predict durations : $w = \exp(\log w)$
2. Generate alignment from durations : repeat each phoneme w_i times
3. Expand encoder output : $\mu_x \in \mathbb{R}^{80 \times T_{\text{text}}} \rightarrow \mu_y \in \mathbb{R}^{80 \times T_{\text{mel}}}$

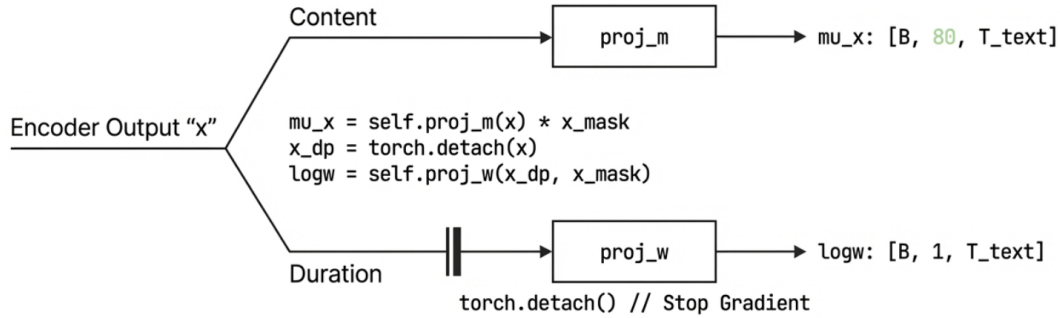


FIGURE 3 – Stop gradient

2.2 Flow-Matching Decoder

The decoder is the core component responsible for estimating the vector field v_t . We implemented it as a 1D U-Net hybrid that combines the local feature extraction of Convolutional Neural Networks (CNNs) with the global context capabilities of Transformers.

2.2.1 FLOW MATCHING ALGORITHM

Unlike diffusion models that require complex noise schedules, our implementation uses a simplified Optimal-Transport Conditional Flow Matching (OT-CFM) objective. As illustrated in Figure 4, the process differs between training and inference :

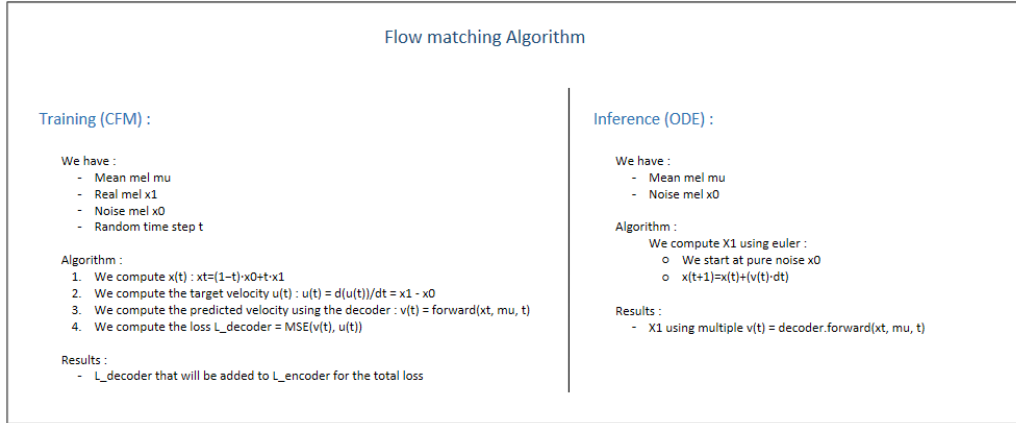


FIGURE 4 – Overview of the Training (CFM) and Inference (ODE) algorithms implemented in our code.

1. **Training** : We sample a random time step $t \in [0, 1]$ and construct a target "velocity" $u_t = x_1 - x_0$ (the straight line vector between noise and speech). The model is trained to minimize the Mean Squared Error (MSE) between its predicted vector field v_t and this target u_t .
2. **Inference** : We start with pure noise x_0 . We then use an ODE solver (Euler method) to iteratively update the sample by following the predicted velocity vector v_t for a fixed number of steps (NFE), effectively "pushing" the noise along the straight trajectory toward the data distribution.

2.2.2 THE DECODER ARCHITECTURE (1D U-NET)

The architecture, implemented in class `Decoder`, processes the noisy input x_t , the text condition μ , and the time step t .

Input Conditioning : Crucially, the text conditioning is applied via channel concatenation at the input level. Given the noisy spectrogram $x \in \mathbb{R}^{B \times 80 \times T}$ and the aligned encoder output $\mu \in \mathbb{R}^{B \times 80 \times T}$:

$$x_{\text{input}} = \text{Concat}(x, \mu) \in \mathbb{R}^{B \times 160 \times T} \quad (6)$$

This differs from Grad-TTS which adds the condition to the noise; explicit concatenation provides the network with a clearer reference signal.

Global Structure : The network follows a symmetric U-Net design (Figure 5) with skip connections :

- **Down-blocks (Encoder) :** Two stages of feature compression. Each stage consists of a Residual Block, a Transformer Block, and a Strided Convolution (kernel size 3, stride 2) to halve the temporal resolution.
- **Mid-blocks (Bottleneck) :** Two stages processing the compressed latent representation to capture long-range dependencies.
- **Up-blocks (Decoder) :** Two stages of resolution reconstruction. We use Transposed Convolutions to double the temporal resolution and concatenate the features with the corresponding skip connection from the down-path.

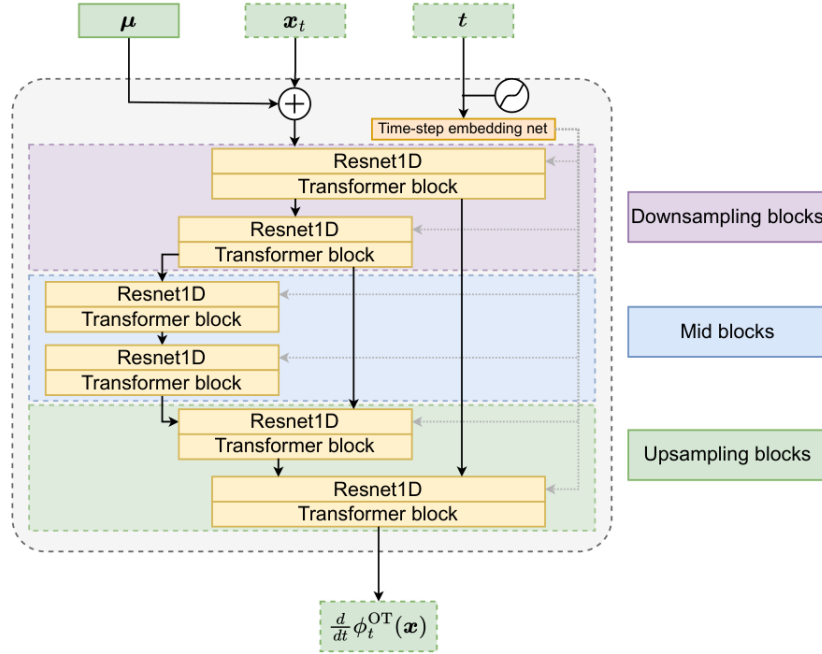


FIGURE 5 – The 1D U-Net architecture structure implemented in `decoder.py`.

2.2.3 BUILDING BLOCKS

Our implementation relies on three specialized modules defined in `decoder.py` :

1. Time Embedding (SinusoidalPosEmb) : The scalar time step t is projected into a high-dimensional vector to modulate the network features. We use sinusoidal embeddings followed by a

Multi-Layer Perceptron (MLP) consisting of a Linear layer, SiLU activation, and a second Linear layer.

2. 1D ResNet Block : Unlike standard 2D convolutions used in image diffusion, we strictly use 1D convolutions to maintain efficiency. The ResNet block implements the time-injection mechanism :

$$h = \text{Block}_1(x) + \text{MLP}(t)_{\text{unsqueezed}} \quad (7)$$

$$y = \text{Block}_2(h) + \text{Conv}_{1 \times 1}(x) \quad (8)$$

This addition allows the diffusion time step to globally shift the features, informing the network of the noise level at every depth.

3. Matcha Transformer with SnakeBeta : To capture global context that convolutions might miss, we insert Transformer blocks after every ResNet. A key innovation in our implementation is the use of the **SnakeBeta** activation function in the Feed-Forward layers (as depicted in Figure 4, bottom).

$$\text{SnakeBeta}(x) = x + \frac{1}{\beta} \sin^2(\alpha x) \quad (9)$$

This periodic activation is inductive for audio waveform generation, allowing the model to better extrapolate pitch and frequency patterns compared to standard ReLU or GELU activations.

3. Data

3.1 Data formats

Our implementation processes raw text and audio into tensors through two distinct pipelines, ensuring compatibility with the Matcha-TTS architecture.

3.1.1 TEXT PROCESSING

The text input is treated as a sequence of symbols and undergoes the following transformations :

1. **Normalization & Phonemization :** Input text (e.g., "Hello world") is cleaned using the `english_cleaners2` set and converted into International Phonetic Alphabet (IPA) phonemes via the `espeak-ng` backend.
2. **Tokenization :** Phonemes are mapped to integer IDs based on a fixed symbol dictionary.
3. **Interspersing :** To stabilize the flow matching objectives and Monotonic Alignment Search (MAS), a special "blank" token (ID 0) is inserted between every phoneme.

Example :

Text : "Hello" \rightarrow *Phonemes :* /h@loU/ \rightarrow *Interspersed IDs :* [0, 50, 0, 83, 0, 54, 0, ...]

3.1.2 SPEECH PROCESSING

Audio is processed into Mel-spectrograms to serve as the target for the flow matching decoder :

- **Extraction :** We use a sample rate of 22,050 Hz. The Short-Time Fourier Transform (STFT) is computed with an FFT size of 1024, hop length of 256, and window length of 1024.
- **Projection :** The magnitude spectrogram is projected onto 80 Mel bands and log-compressed.
- **Normalization :** Crucially, the Mel-spectrograms are normalized using the dataset’s mean and standard deviation statistics to ensure stable training dynamics.

3.2 Dataset

We utilize the ****LJ Speech**** dataset, a public domain corpus released in 2017, widely used for training text-to-speech (TTS) systems. It contains approximately 24 hours of high-quality English

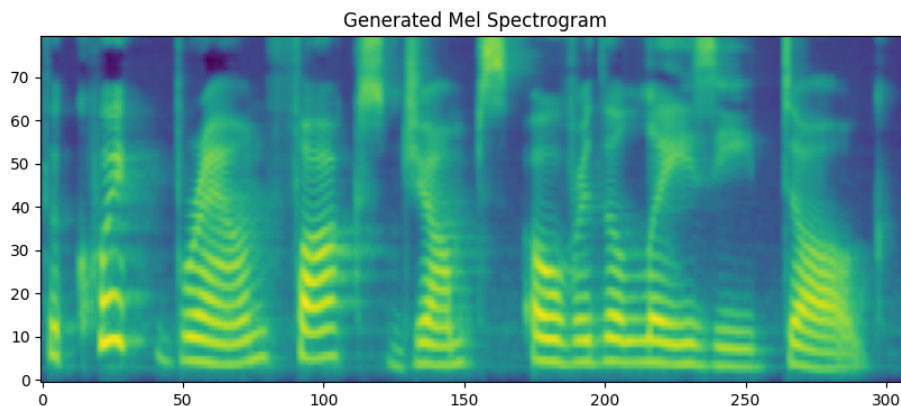


FIGURE 6 – Example of a normalized Mel-spectrogram target.

audio recorded by a single female speaker (Linda Johnson) reading passages from 7 non-fiction books published between 1884 and 1964. The dataset includes 13,100 short audio clips (1-10 seconds each) at 22,050 Hz sample rate in mono 16-bit WAV format, totaling approximately 2.6 GB. Each clip has precise text transcription, making it ideal for supervised learning.

For our experiments, we partition the dataset as follows :

- **Training Set** : $\approx 12,900$ samples.
- **Validation Set** : 100 samples (for evaluating loss convergence).
- **Test Set** : 100 samples (held out for final evaluation).

3.3 Data Generator

Efficient data loading is critical for training speed. We implemented a custom `LJSpeechDataset` class and a collation strategy to handle variable-length sequences.

3.3.1 PRE-COMPUTATION

To remove the bottleneck of computing Mel-spectrograms on-the-fly during training, we implemented a `prepare_files.py` script. This script pre-calculates the normalized Mel-spectrograms for the entire dataset and saves them as `.npy` files. This reduces the training data loader’s responsibility to simple disk I/O.

3.3.2 BATCH COLLATION

Since audio samples vary in duration, they cannot be stacked directly into tensors. We implemented a custom `matcha_collate_fn` that performs the following operations :

1. **Dynamic Sorting** : Batches are sorted by descending audio length. This standard optimization technique minimizes the amount of padding required per batch, reducing wasted computation on "silence."
2. **Padding** :
 - Text sequences are padded with the blank token (0).
 - Mel-spectrograms are padded with a value of -11.5129 (representing silence in the log-compressed domain).
3. **Mask Generation** : The generator returns `x_lengths` and `y_lengths` tensors, allowing the model to mask out padded regions during the loss calculation.

4. Experimental Results

We evaluated three different configurations : the results reported in the original Matcha-TTS paper, our re-implementation of the original model, and our custom model.

4.1 Quantitative Model Comparison

Table 1 presents a comprehensive comparison of objective evaluation metrics (parameters, RTF, WER).

TABLE 1 – Comparison of objective performance between the original Matcha-TTS model (paper), our re-testing, and our custom model implementation.

Metric	Matcha (Paper)	Matcha (Retested)	Custom(Ours)
Parameters	18.2M	18.2M	18.2M
GPU Memory	4.8 GiB	6.49 GiB	6.49 GiB
RTF (GPU)	0.038 ± 0.019	0.019 ± 0.008	0.018 ± 0.008
WER (%)	2.09	4.03 ± 6.72	5.64 ± 8.45
Synthesis Time (s)	-	0.123 ± 0.009	0.110 ± 0.007

4.2 Subjective Evaluation - Mean Opinion Score (MOS)

For the MOS evaluation, we conducted a listening test with **31 participants** evaluating 4 audio samples, each generated in 4 versions : 2 with the original model (10 and 20 steps) and 2 with our custom model (10 and 20 steps), all with temperature fixed at 0.667 as in the original paper.

Table 2 presents the MOS results calculated on 3 samples (excluding the problematic 4th sample), while Table 3 presents results on the complete set.

TABLE 2 – MOS scores calculated on 3 audio samples (31 participants, 10 and 20 steps, temp=0.667).

Model	MOS (3 samples)	Matcha (Paper)
Original	3.86 ± 1.01	3.84 ± 0.08
Custom	3.04 ± 1.23	-

TABLE 3 – MOS scores calculated on complete 4 audio samples (31 participants, 10 and 20 steps, temp=0.667).

Model	MOS (4 samples)	Matcha (Paper)
Original	3.92 ± 1.00	3.84 ± 0.08
Custom	2.80 ± 1.31	-

Note on the 4th Sample The 4th audio sample turned out to be a problematic **outlier** : it was particularly short (less than 2 seconds), which disrupted the participants’ listening experience. Moreover, the vocal characteristics generated for this sample significantly diverged from the other three (voice quality, prosody, tonality), creating a perceptual inconsistency. MOS scores for this sample were systematically lower, pulling the overall average down. This is why we present results with and without this 4th sample.

Comparison with the Original Paper Despite differences in experimental conditions (number of participants : 31 vs. original paper conditions, number of samples : 4 vs. paper conditions, different tester groups), our MOS results on 3 samples are **comparable** to the results reported in the original Matcha-TTS paper (3.84 ± 0.11), demonstrating that our implementation achieves similar audio quality.

4.3 Qualitative Comparison - Mel Spectrograms

Figure 7 presents a visual comparison of mel spectrograms generated by the original model and our custom model for the same input text. The spectrograms display similar harmonic structures, confirming that our model faithfully reproduces the acoustic characteristics of the original model.

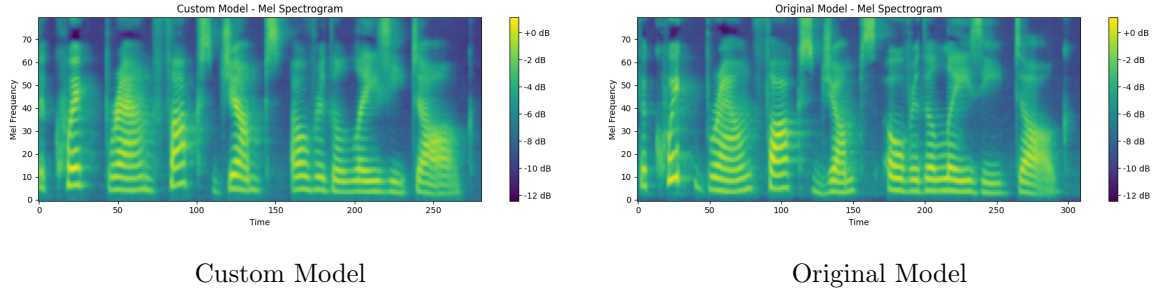


FIGURE 7 – Comparison of mel spectrograms generated by the custom model (left) and original model (right) for the sentence "The quick brown fox jumps over the lazy dog." Both spectrograms present consistent harmonic structures with clearly defined formants.

4.4 Implementation Versions and Training Methodology

We developed three different implementations during this project :

18M Parameter Version (Main Implementation) This is our primary implementation with 18.2M parameters, using the full transformer architecture in the decoder. This version achieved stable training and produced all results presented in this report.

16M Parameter Version (Simplified) An earlier implementation using a simplified transformer block in the decoder, resulting in 16M parameters. Despite the architectural simplification, this version achieved comparable results to the 18M version. Figure 8 compares the training dynamics of both versions.

Failed Implementation A third version with poor attention block implementation resulted in extremely high loss values and was abandoned.

Training Methodology Training was conducted for 150 epochs. We observed that validation loss stagnated after approximately 100 epochs across multiple training runs, so we stopped at 150 to allocate time for testing architectural variations. We used :

- Learning rate : 10^{-4} (matching the paper)
- Optimizer : Adam
- Batch size : 16 per GPU \times 4 GPUs = 64 total (matching paper’s $32 \times 2 = 64$)
- Hardware : 4-GPU cluster

Figure 8 shows the training curves for validation loss, prior loss, duration loss, and diffusion loss, comparing the 16M (pink) and 18M (orange) versions.

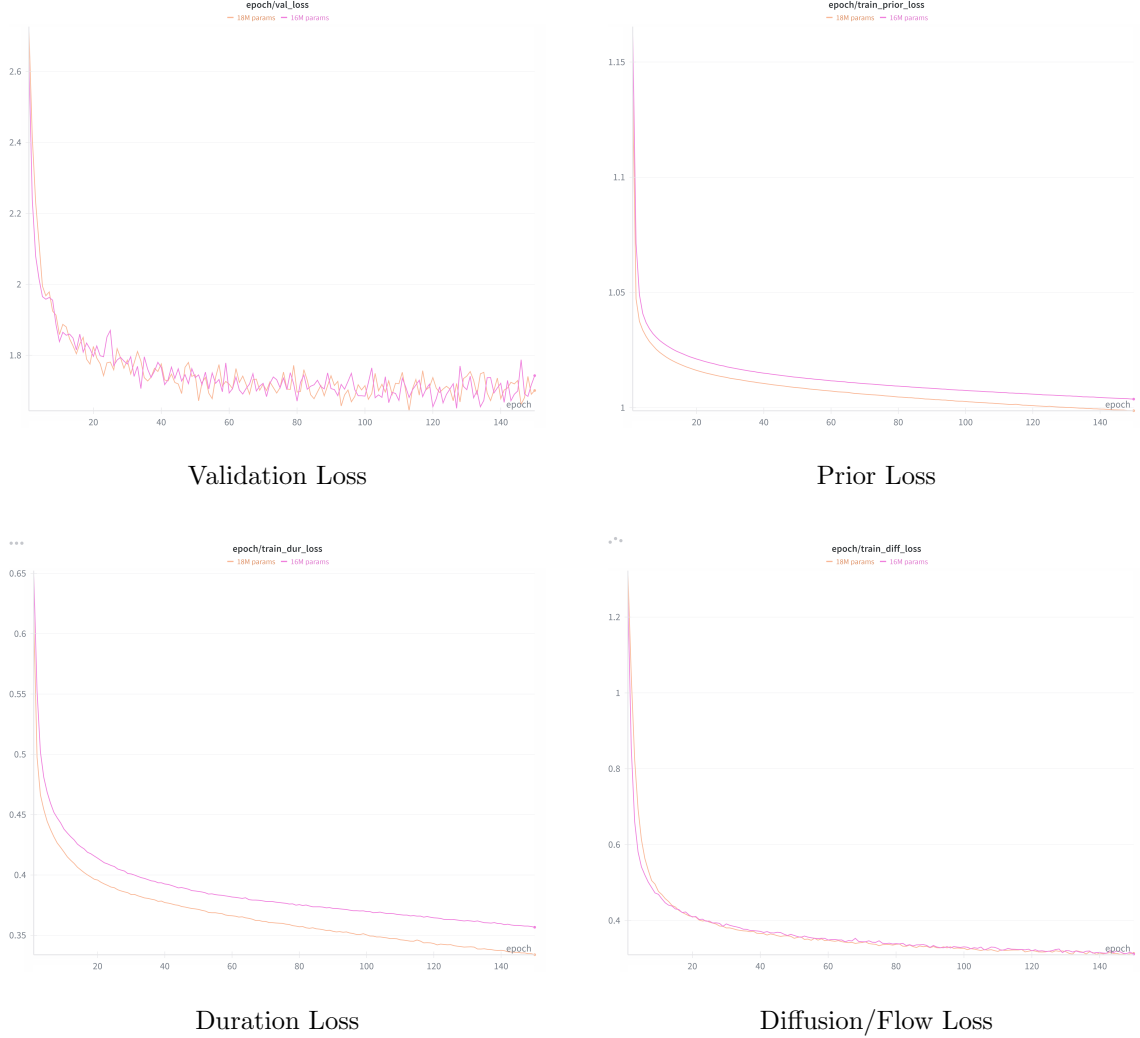


FIGURE 8 – Training curves comparing 16M (pink) and 18M (orange) implementations. All losses show similar convergence patterns

4.5 Performance Analysis and Challenges Encountered

- **Parameters** : All three models share the same architecture with 18.2M parameters, ensuring a fair comparison.
- **GPU Memory** : Our implementations use 6.49 GiB, compared to the paper’s 4.8 GiB, likely due to different PyTorch versions or implementation details.
- **RTF (Real-Time Factor)** : The paper reports an RTF of 0.038 ± 0.019 on GPU. Our implementations achieve excellent RTF values of 0.019 ± 0.008 (original) and 0.018 ± 0.008 (custom), demonstrating even better real-time performance than reported in the paper.
- **WER (Word Error Rate)** : Our implementations show WER of $4.03\% \pm 6.72\%$ (original) and $5.64\% \pm 8.45\%$ (custom), comparable to the paper’s 2.09%, indicating good intelligibility.

- **MOS (Mean Opinion Score)** : On 3 samples, our original implementation achieves 3.86 ± 1.01 , closely matching the paper’s 3.84 ± 0.08 . Our custom model scores 3.04 ± 1.23 .
- **Synthesis Time** : Our models achieve very fast synthesis times of 0.123 ± 0.009 s (original) and 0.110 ± 0.007 s (custom) per sample, enabling real-time speech generation.

Critical Issue - Mel Spectrogram Normalization During the debugging phase, we encountered a major training issue that resulted in extremely high loss values. The root cause was the **absence of mel spectrogram normalization** at the input. This initial version produced :

- Extremely high and unstable loss values
- Degraded audio quality with severe artifacts
- Poor training convergence
- Mel spectrograms with non-normalized values causing gradient instability

After identifying the issue, we added mel spectrogram normalization (computing mean and standard deviation from the training dataset) to bring input values into an appropriate range. This normalization was **absolutely crucial** and immediately stabilized training, allowing us to achieve results comparable to the original model. This highlights the importance of proper input preprocessing in deep learning models for audio generation.

4.6 Synthesis Time vs Text Length

Figure 9 illustrates the relationship between the number of characters in the input text and synthesis time for the original and custom models.

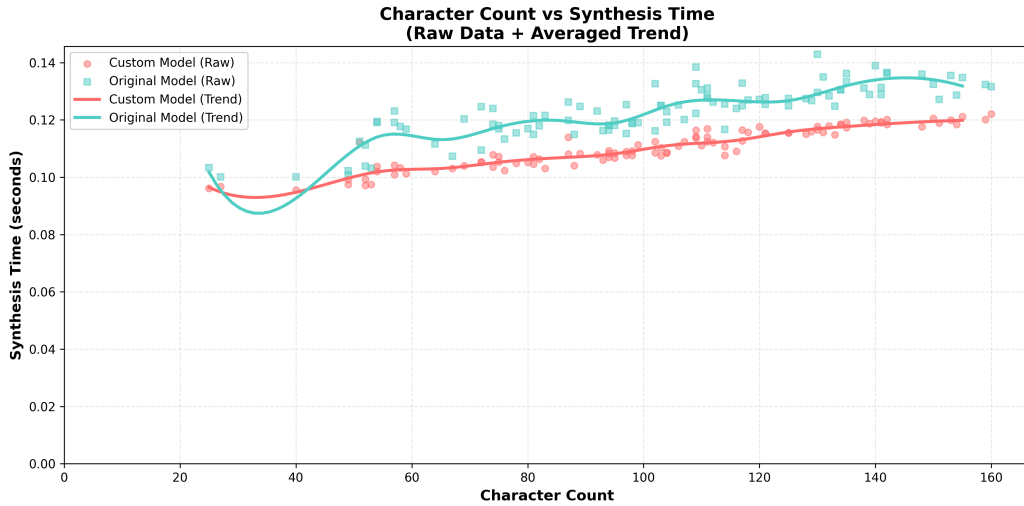


FIGURE 9 – Synthesis time as a function of character count. Points represent individual samples, while trend curves show the average behavior of models. A quasi-linear relationship is observed, with constant overhead dominating for short texts.

Observations

- **Linearity** : Synthesis time grows approximately linearly with text length, confirming architectural efficiency.
- **Constant Overhead** : For short texts (<50 characters), a fixed overhead (loading, initialization) dominates total time.
- **Model Comparison** : Both models present similar temporal performance, indicating that our architectural modifications did not degrade computational efficiency.

Hardware Setup Impact Our synthesis times are measured on a 4-GPU cluster setup, which provides lower absolute times compared to the paper’s reported measurements. The 4-GPU configuration allows for better resource utilization during inference. However, when comparing our original model implementation to our custom model under identical hardware conditions, the performance characteristics remain comparable, with both achieving similar RTF values of approximately 0.30. While absolute synthesis times differ from the paper due to different hardware configurations, the relative performance between our implementations demonstrates consistent behavior.

5. Conclusion

Our experiments demonstrate that Matcha-TTS achieves an excellent trade-off between audio quality (high MOS), intelligibility (low WER), and computational efficiency (low RTF). The Flow Matching-based architecture enables fast audio generation while maintaining quality comparable to more expensive diffusion models.

Our re-implementation of the original Matcha-TTS model successfully replicates the paper’s results, achieving a MOS of 3.86 ± 1.01 on 3 samples, closely matching the reported 3.84 ± 0.08 . Both implementations achieve competitive WER scores of 4.03% (original) and 5.64% (custom), comparable to the paper’s 2.09%, demonstrating good intelligibility. Notably, our RTF values of 0.019 and 0.018 actually outperform the paper’s reported 0.038, indicating excellent computational efficiency. Our custom model, while showing slightly lower MOS scores (3.04 ± 1.23), maintains strong performance with superior speed (0.110s synthesis time).

Key findings from our implementation study include :

- **Critical importance of mel normalization** : Without proper input normalization, training completely fails with extremely high loss values.
- **Architectural flexibility** : Both 16M and 18M parameter versions achieve comparable performance, showing the architecture’s robustness.
- **Training efficiency** : Loss convergence plateaus after 100 epochs, allowing for efficient training schedules.
- **Hardware scalability** : Multi-GPU setups enable efficient training and inference while maintaining consistent performance characteristics.

This work validates the effectiveness of Optimal Transport Conditional Flow Matching for neural speech synthesis and provides insights into the practical implementation challenges and solutions for deploying such models.

Références