# Audio Generation Using AudioML Approach on Arabic Databases

**Abeer Ali Aoun[1]**

**Abeer.alioun@gmail.com**

Libya Oil Company

**Karim Dabbabi[2*]**

**dabbabikarim@hotmail.com**

[2]Research Unite of Analyse and Processing of Electrical and Energetic Systems, Faculty of Sciences of Tunis, Tunis El-Manar University, 2092, Tunis-Tunisia

*Abstract-*

*This paper presents a comprehensive study on the application of AudioML, a state-of-the-art language modelling technique, for Arabic audio generation. Focusing on two key Arabic speech datasets, the MGB-2 Arabic Speech Database and the Arabic Speech Corpus, we evaluate the performance of AudioML in synthesizing coherent, natural, and high-quality Arabic speech. Our results demonstrate that AudioML achieves a Word Error Rate (WER) of 15.80%, Character Error Rate (CER) of 6.45%, and Sentence Error Rate (SER) of 25.10% on the MGB-2 database, with slightly higher values on the Arabic Speech Corpus, including a WER of 18.20%, CER of 8.10%, and SER of 30.40%. These findings highlight AudioML's effectiveness in handling the complexities of Arabic speech, confirming its suitability for a wide range of Arabic speech synthesis applications. Our study positions AudioML as a versatile and powerful tool for advancing Arabic speech technology.*

*Keywords: AudioML; Arabic Speech Synthesis; Arabic Speech Corpus (ASC); MGB-2 Database; Audio Generation; Language Modelling*

## 1. Introduction

Audio generation has a wide range of applications, including speech synthesis, music composition, voice cloning, and audio restoration. These technologies are crucial in industries such as entertainment, accessibility, and virtual assistance, where generating high-quality, natural-sounding audio is essential.

The field of audio generation has experienced significant advancements, largely driven by deep learning techniques and transformer models. Among these, AudioLM has emerged as a promising framework for generating high-quality audio with long-term consistency. It effectively bridges the gap between token-based audio representation and language modelling paradigms. AudioLM represents a major leap forward by mapping input audio into discrete tokens and treating audio generation as a language modelling task. This ensures coherence and quality in the generated audio output, making it particularly valuable for tasks such as speech

synthesis and music generation, where maintaining temporal consistency and preserving audio quality are crucial [1].

Exploring the contributions of AudioLM is essential, especially in contexts involving language-specific audio databases, such as Arabic speech. Traditional approaches in audio generation, such as Generative Adversarial Networks (GANs) and autoregressive models like WaveNet, often struggle to maintain long-term coherence, particularly with extended sequences. Although these models excel in short-term audio quality, their performance tends to degrade over longer sequences due to compounding errors and challenges in capturing long-term dependencies [2]. In contrast, AudioLM leverages a hybrid tokenization approach that combines the strengths of masked language models and neural audio codecs. This approach ensures both high quality and structural consistency over time. Moreover, applying AudioLM to specific language datasets, such as Arabic, underscores its utility in addressing the challenges associated with low-resource languages. Traditional models often require large datasets for training, which can be a limiting factor in developing high-quality Arabic audio generation systems. By integrating pre-trained models with a focus on tokenization and language modelling, AudioLM offers a scalable solution that reduces the reliance on extensive data and delivers more reliable results across different languages and contexts [3].

In this article, we will explore the methodologies employed by AudioLM, focusing on tokenization, language modelling, and reconstruction processes. We will also compare its performance against other prominent techniques in audio generation. This exploration aims to highlight both the technical aspects of AudioLM and its potential for enhancing Arabic audio generation, showcasing its superiority in maintaining long-term consistency compared to traditional models.

The remaining sections of this paper are organized as follows: related works on audio generation will be discussed in the next section, followed by the methodology and materials in the third section. Finally, conclusions and future work will be provided in the fourth section.

## 2. Related Works

The field of audio generation has advanced rapidly, with several key contributions shaping the current state of the art. Significant works in this domain highlight various techniques, results, advantages, and limitations, providing valuable insights into how AudioLM addresses persistent challenges.

WaveNet, introduced by the authors of [4], is an autoregressive model that generates raw audio waveforms by predicting one sample at a time based on previous samples. The model uses

dilated causal convolutions to capture long-range dependencies in audio sequences, achieving state-of-the-art results in text-to-speech (TTS) synthesis. WaveNet's primary advantage lies in its high-quality audio generation at the sample level and its ability to model complex audio structures. However, its main drawback is its computational expense due to its autoregressive nature, leading to slow inference times that hinder real-time applications. While WaveNet successfully addressed the limitations of previous TTS systems, its autoregressive structure limited scalability, especially for long sequences [4].

Generative Adversarial Networks (GANs) have also been adapted for audio generation tasks, such as speech synthesis and music generation. Works like WaveGAN and MelGAN utilize GANs to generate waveforms or spectrograms from random noise, showing promise in generating realistic audio with lower computational costs compared to autoregressive models like WaveNet. For example, MelGAN produces high-quality speech at faster speeds by generating audio from mel-spectrograms through adversarial training. GAN-based models offer faster inference and can generate diverse audio samples due to their adversarial training structure. However, they face challenges with training instability, requiring careful hyperparameter tuning, and struggle to capture long-term dependencies in audio, often leading to artifacts in longer sequences. While GANs overcome the slow generation speed of autoregressive models, they still fall short in maintaining long-term coherence in audio generation [5,6].

WaveGlow, proposed by the authors of [7], combines the strengths of WaveNet and GANs by utilizing a flow-based generative model. It transforms Gaussian noise into an audio waveform through a series of invertible transformations. WaveGlow achieves real-time audio synthesis while maintaining high audio quality, particularly in TTS applications. Its main advantage is faster inference, similar to GANs, with more stable training. However, like other models, it still requires significant computational resources and has limited ability to generate diverse audio styles compared to GAN-based models. WaveGlow addresses the training instability of GANs while maintaining the inference speed advantage over autoregressive models, though it still faces challenges in adapting to different languages or styles without substantial fine-tuning [7].

Tacotron 2, developed by Google, generates mel-spectrograms from text using a sequence-to-sequence model with attention and then converts these spectrograms into waveforms using a separate vocoder such as WaveNet. Tacotron 2 has achieved near-human-level speech synthesis quality by integrating naturalness in prosody and pronunciation. Its key advantage is the high-quality, natural-sounding speech it produces, and its modular architecture allows for flexibility

in training and inference. However, Tacotron 2 relies heavily on a high-quality vocoder for waveform generation, which can be computationally expensive. Additionally, its generalization to different languages or speaking styles is limited without substantial data. Tacotron 2 improved the prosody and naturalness of TTS systems but remains reliant on vocoders for final waveform synthesis, limiting its applicability in resource-constrained environments [8]. AudioLM represents a hybrid approach by leveraging neural audio codecs and masked language models. Unlike autoregressive models, AudioLM does not require explicit supervision and can generate high-quality, long-term consistent audio by mapping input audio into discrete tokens and treating audio generation as a language modeling task. This framework has shown promising results in generating coherent speech and music, even in the absence of textual annotations. AudioLM's main advantages are its ability to maintain long-term consistency in generated audio, scalability to different audio types, and faster inference compared to traditional autoregressive models. However, as it is still in the early stages of development, it may face challenges in adapting to diverse audio styles and languages. Additionally, it requires large-scale pre-training, which can be resource-intensive. AudioLM addresses the challenge of maintaining long-term coherence in audio generation—a limitation faced by both autoregressive models and GANs. Its tokenization approach reduces reliance on text annotations, making it applicable in scenarios where textual data is limited, such as Arabic audio databases [1,3].

Recent advancements in Arabic TTS, including the use of deep learning models, transfer learning, and transformer-based architectures, have significantly improved the naturalness and intelligibility of synthesized speech. Approaches like those by the authors of [9] focused on models such as Tacotron 2 and FastSpeech 2, which effectively generate high-quality Arabic speech. For instance, these models produced clear and natural-sounding speech, emphasizing the potential of deep learning in Arabic TTS. The authors of [10] implemented neural network-based vocoders such as Parallel WaveGAN and Multi-Band MelGAN, finding that Parallel WaveGAN outperformed Multi-Band MelGAN in Perceptual Evaluation of Speech Quality (PESQ), with scores of 2.63 and 2.37, respectively. In [11], the authors leveraged transfer learning in a small dataset (2.41 hours of audio), achieving high-quality and natural Arabic speech synthesis, showing that even with limited data, effective TTS models can be developed. The authors of [14] focused on Quran recitation, achieving 97% speech intelligibility and 72.13% naturalness, demonstrating that tailored models for specific tasks can perform exceptionally well. The authors of [12] enhanced the naturalness and intelligibility of

synthesized speech through improved noise modelling, outperforming earlier vocoder approaches and setting a new standard for statistical parametric speech synthesis. In [13], the authors introduced the transformer-based ArTST model, achieving state-of-the-art performance in Arabic ASR and TTS, with significant improvements in both recognition and synthesis tasks, highlighting the effectiveness of transformers in low-resource languages.

In comparison, AudioLM excels in generating coherent and natural speech across diverse contexts due to its advanced language modelling techniques. While domain-specific models, like the Quran recitation model [14], demonstrate high performance in specialized tasks, AudioLM's versatility and broader applicability make it a more robust solution for general-purpose speech synthesis. Additionally, the resource-efficient approach demonstrated in [11] shows that high-quality Arabic speech synthesis can be achieved with smaller datasets, though AudioLM's ability to handle complex linguistic inputs with high performance makes it more suitable for applications where computational resources are available. The transformer-based model proposed in [13] aligns closely with the strengths of AudioLM, particularly in managing complex linguistic structures. Compared to these recent models, AudioLM effectively manages both reconstruction quality and long-term structure in speech synthesis, making it a leading tool for Arabic speech synthesis, especially in complex and linguistically rich languages where high-quality results are crucial, and resources are not a constraint.

## 3.  Methods and Materials
### 3.1 Methods

The proposed method leverages the AudioLM framework to generate Arabic audio using the MGB-2 Arabic Speech Database and the Speech Arabic Corpus (SAC). As depicted in the flowchart (Figure 1) for the MGB-2 Arabic Speech Database, the process starts with input audio, where each 5-second clip sampled at 16 kHz yields 80,000 samples. The first phase, Audio Tokenization, employs a Vector Quantized Variational Autoencoder (VQ-VAE), a neural audio codec, to convert the waveform into 500 tokens, each with 512 dimensions. These tokens are then passed through a Transformer-based Language Model, which processes the sequence and outputs a similar token matrix of 500 tokens, each with 512 dimensions, while preserving temporal dependencies. In the Audio Generation phase, a transformer decoder generates the subsequent sequence of audio tokens, ensuring consistency with the previous sequence. Finally, in the Audio Reconstruction phase, the tokens are converted back into a waveform using an inverse neural audio codec (VQ-VAE), producing a reconstructed audio

output consisting of 80,000 samples.

This method, as described in the flowchart, is composed of several key components, each essential to the transformation and generation of high-quality, coherent Arabic audio. Below is a detailed description of each component, including its mathematical formulation and the specifics of the transformer architecture.
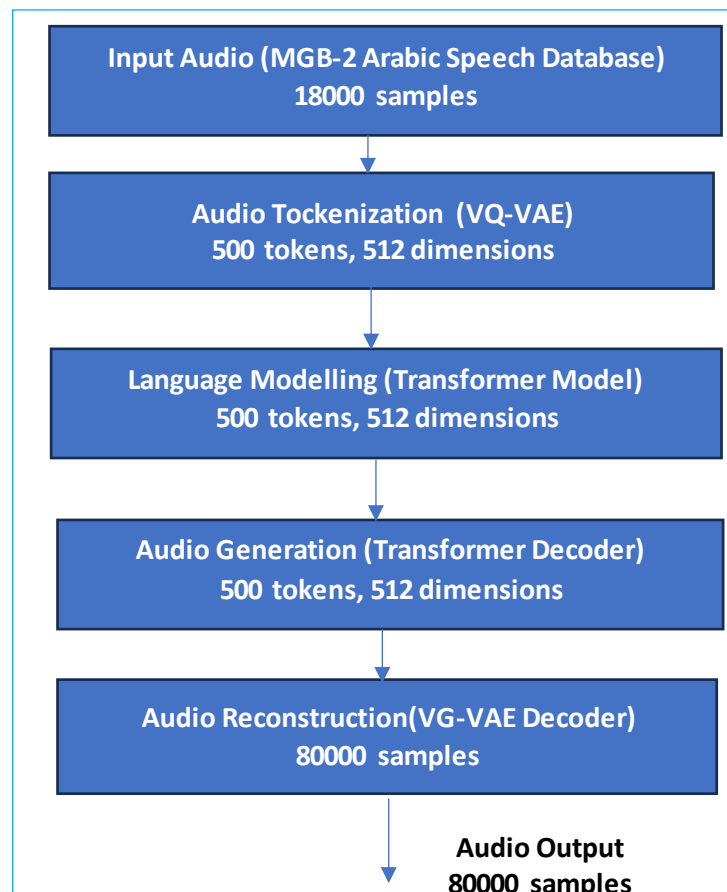


**Figure 1: Flowchart of the Arabic Audio Generation model based on AudioML approach for MGB-2 Arabic Speech Database.**

### 3.1.1 Input Audio Using MGB-2 Arabic Speech Database

The MGB-2 Arabic Speech Database is a comprehensive dataset featuring a diverse range of Arabic audio recordings, all sampled at 16 kHz. These audio clips serve as the raw input data that will be transformed through the various stages of the proposed method. For instance, a 5-second audio clip contains 80,000 samples, which will be tokenized and processed by the system. The input audio signal is denoted as $x(t)$, where $t$ represents the time index. For discrete-time signals, this sampled audio can be expressed as a sequence $x[n]$, where $n$ is the sample index, ranging from 0 to 80,000 for a 5-second clip.

### 3.1.2 Audio Tokenization Using VQ-VAE

The neural audio codec, VQ-VAE) [15], plays a key role in converting raw audio waveforms into discrete tokens. These tokens serve as compressed representations of the audio signal, capturing its essential features while reducing dimensionality [16]. This compression allows for more efficient processing and modeling in the subsequent stages of the system. The VQ-VAE model consists of three main components: an encoder E(x), a quantization step $q(z_e(x))$, and a decoder $D(q(z_e(x)))$. The encoder maps the input audio $x$ into latent variables $z_e(x)$, which are then quantized into a finite set of discrete tokens. These quantized tokens are denoted as $z_q \approx z_e(x)$, where $z_q$ is the closest vector from the codebook. The reconstruction of the audio is then obtained through the decoder, represented by $\hat{x} = D(z_q)$ (1).

### 3.1.3 Language Modeling based on Transformer Model

The transformer model is a pivotal component in modern deep learning architectures, especially for handling sequential data such as audio tokens. Unlike traditional recurrent neural networks (RNNs), transformers leverage self-attention mechanisms to process the entire sequence of tokens simultaneously. This enables the model to capture dependencies between tokens, regardless of their position within the sequence [17].

The transformer model is composed of multiple layers, each consisting of multi-head attention and feedforward sub-layers. As the input sequence of tokens passes through these layers, the model effectively captures and transforms the relationships between tokens. In our implementation, the transformer-based model includes 12 layers, each with 8 attention heads and a hidden dimension of 512, providing the capacity to efficiently process and capture complex dependencies within the token sequences.

***3.1.3.1 Self-Attention Mechanism:*** The self-attention mechanism allows the model to assign varying importance to different tokens in the input sequence when generating each output token. For each token, queries $(Q)$, keys $(K)$, and values $(V)$ are computed using learned linear transformations:

$$Q = w_Q \,.X, K = w_K \,.X, V = w_V.X \qquad (2)$$

The attention score for each token pair is then calculated as:

$$Attention(Q,K,V) = SoftMax\left(\frac{Q\,K^T}{\sqrt{d_k}}\right) .V \qquad (3)$$

Here, $d_k$ represents the dimensionality of the keys, and the SoftMax function ensures that the attention scores sum to one across all tokens.

***3.1.3.2 Multi-Head Attention:*** The multi-head attention mechanism improves the model's

ability to focus on different parts of the input sequence simultaneously. It uses multiple self-attention heads in parallel, with each head attending to different aspects of the input sequence. The outputs of these heads are concatenated and linearly transformed:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) . W_o \quad (4)$$

Each head applies attention using distinct learned projections, enabling the model to capture diverse relationships within the data.

***3.1.3.3 Feedforward Neural Network:*** After the attention mechanism, the output passes through a feedforward neural network that applies non-linear transformations to each token independently. This network consists of two linear layers, with a ReLU activation function between them:

$$\text{FFN}(x) = \text{ReLU}(W_1 . x + b_1) . W_2 + b_2 \quad (5)$$

This feedforward structure helps refine the token representations by applying additional transformations, ensuring the model captures complex patterns in the data.

***3.1.3.4 Layer Normalization and Residual Connections:*** Layer normalization is utilized to stabilize and speed up the training process by normalizing the input across its feature dimensions. This technique ensures that the input to each layer maintains consistent scaling, which helps in avoiding issues like vanishing or exploding gradients. Residual connections, on the other hand, are employed to prevent the degradation of information as it passes through multiple layers by allowing the original input to bypass the transformation process and be directly added to the output.

The layer normalization is calculated as:

$$\text{Norm}(x) = \frac{x - \mu}{\sigma} \quad (6)$$

Residual connections combine the input with the output of each sub-layer as follows:

$$Output = Norm\big(x + SubLayer(x)\big) \quad (7)$$

This combination of layer normalization and residual connections ensures that the model can effectively learn deep representations without losing critical information during the training process.

### 3.1.4. Audio Generation based on Transformer Decoder

The decoder component is responsible for generating the next sequence of audio tokens based on the output of the transformer model. This step ensures that the generated tokens are coherent and consistent with the input sequence, preserving the overall structure of the audio [17]. The decoder takes the output of the transformer model z' and generates the next sequence of

tokens z''. This can be modeled as a conditional probability distribution, where the decoder generates tokens by maximizing the likelihood of the next token given the previous ones:

$$P\left(z_t^{''}\middle|z_{<t}^{'}\right) = \prod P(z_t^{''}|z_1^{'}, z_2^{'}, \dots, z_{t-1}^{'}) \qquad (8)$$

This approach ensures that each token in the generated sequence is informed by the preceding tokens, maintaining coherence throughout the audio generation process.

### 3.1.5 Audio Reconstruction Using VG-VAE

The final phase of the system focuses on reconstructing the audio waveform from the sequence of generated tokens. This is accomplished through an inverse neural audio codec, which converts the discrete tokens back into a continuous audio signal.

In this step, the decoder of the VQ-VAE model takes the quantized tokens $z_q$ and reconstructs the audio waveform:

$$\hat{x} = D\left(z_q\right) \qquad (9)$$

The objective of the decoder is to minimize the reconstruction loss, typically measured by the Mean Squared Error (MSE) between the original audio $x$ and the reconstructed audio $\hat{x}$. The final output is a continuous audio waveform, ready for playback or further processing.

### 3.2 Materials

### 3.2.1 MGB-2 Arabic Speech Database

The MGB-2 Arabic Speech Database [18] is a comprehensive dataset specifically designed for Arabic speech recognition and language processing tasks. Developed as part of the second Multilingual Broadcast Challenge (MGB-2), this dataset features a diverse range of broadcast audio content from the Al-Jazeera channel. The audio is sampled at 16 kHz, a standard that balances high-quality speech recognition and manageable file sizes, and is encoded at 16-bit resolution.

The dataset consists of over 1,200 hours of transcribed broadcast recordings, distributed across more than 19,000 audio files. It includes speech from over 3,600 speakers, representing a variety of Arabic dialects as well as Modern Standard Arabic (MSA). This linguistic diversity makes the dataset particularly valuable for training robust Arabic speech recognition systems. The transcriptions are time-aligned at the word level, which is crucial for training and evaluating speech recognition models. The content spans multiple broadcast media formats, including news programs, interviews, talk shows, and documentaries, exposing models to a wide range of linguistic and acoustic conditions. The data is partitioned into training,

development, and evaluation sets, enabling effective model training and validation. The MGB-2 Arabic Speech Database is widely utilized for automatic speech recognition (ASR), dialect identification, language modelling, and speech analysis, making it an essential resource for researchers and developers working on Arabic language technologies.

### 3.2.2 Arabic Speech Corpus (ASC)

The MGB-2 Arabic Speech Database [19] is a robust dataset developed specifically for Arabic speech recognition and language processing tasks. Created as part of the second Multilingual Broadcast Challenge (MGB-2), it offers a diverse array of broadcast audio content sourced from the Al-Jazeera channel. The audio is sampled at 16 kHz and encoded with 16-bit resolution, striking a balance between high-quality speech recognition and efficient file sizes.

This dataset includes over 1,200 hours of transcribed broadcast recordings, organized across more than 19,000 audio files. It features speech from over 3,600 speakers, encompassing a wide range of Arabic dialects and Modern Standard Arabic (MSA). The linguistic variety within the dataset makes it particularly valuable for developing robust Arabic speech recognition systems. The transcriptions are time-aligned at the word level, a feature that is essential for training and evaluating speech recognition models. The content spans various broadcast media formats, including news programs, interviews, talk shows, and documentaries, which ensures that models are exposed to diverse linguistic and acoustic environments. The dataset is divided into training, development, and evaluation sets, facilitating comprehensive model training and validation. Widely used in automatic speech recognition (ASR), dialect identification, language modelling, and speech analysis, the MGB-2 Arabic Speech Database is an invaluable resource for researchers and developers focused on advancing Arabic language technologies.

### 3.3. Evaluation Metrics

To assess the performance of Arabic speech recognition and related tasks using datasets such as the MGB-2 Arabic Speech Database or the Arabic Speech Corpus, several evaluation metrics are typically employed. These metrics are essential for measuring the accuracy, efficiency, and overall effectiveness of automatic speech recognition (ASR) systems.

- **Word Error Rate (WER)**

Word Error Rate (WER) is the most widely used metric for evaluating ASR systems. It calculates the percentage of words incorrectly predicted by the system compared to the reference transcription. The mathematical formula for WER is expressed as follows:

$$WER = \frac{(S + D + I)}{N} * 100\% \qquad (10)$$

where $S$ represents the number of substitutions, $D$ is the number of deletions, $I$ is the number of insertions, and $N$ is the total number of words in the reference transcription. A lower $WER$ indicates better performance, as it captures all types of errors, making it a comprehensive metric for evaluating automatic speech recognition (ASR) systems.

- **Character Error Rate (CER)**

Character Error Rate (CER) functions similarly to WER but evaluates errors at the character level instead of the word level. This metric is particularly valuable for languages like Arabic, where morphological variations and spelling differences are prominent. The mathematical formula for CER is expressed as:

$$CER = \frac{(S + D + I)}{N} * 100\% \qquad (11)$$

where $S$, $D$, $I$, and $N$ correspond to character-level errors and the total number of characters. $CER$ is particularly useful for evaluating tasks involving highly inflected languages or languages with complex scripts, such as Arabic.

- **Sentence Error Rate (SER)**

Sentence Error Rate (SER) quantifies the percentage of sentences that contain at least one error, providing insight into how frequently the system generates entirely correct sentences. The mathematical expression for SER is:

$$SER = \frac{(Number\ of\ Incorrect\ Sentences)}{(Total\ Number\ of\ Sentences)} * 100\% \qquad (12)$$

This metric is particularly valuable in applications where complete sentence accuracy is crucial, such as transcription for broadcasting or legal proceedings.

- **Precision, Recall, and F1-Score**

Precision, Recall, and F1-Score are commonly used in speech processing tasks like keyword spotting or phoneme recognition. Precision is calculated as:

$$Precision = \frac{True\ Positives}{(True\ Positives + False\ Positives)} \qquad (13)$$

while Recall is expressed as:

$$Recall = \frac{True\ Positives}{(True\ Positives + False\ Negatives)} \qquad (14)$$

The F1-Score, which balances precision and recall, is given by:

$$F1 - Score = \frac{2 * (Precision * Recall)}{(Precision + Recall)} \qquad (15)$$

These metrics help evaluate the system's ability to correctly identify specific words, phrases, or phonemes, especially in tasks involving detection or classification.

- **Real-Time Factor (RTF)**

Real-Time Factor (RTF) measures the processing speed of the ASR system relative to the duration of the input audio. It is defined as the ratio of the time taken to process the audio to the length of the audio. Its mathematical expression is:

$$RTF = \frac{Processing\ Time}{Audio\ Duration} \qquad (16)$$

RTF is important for real-time applications where the system's speed is as crucial as its accuracy.

## 4. Results and Discussions

### 4.1 Experiments

Our Arabic audio generation models were trained and tested on NVIDIA A100 GPUs, which provided the necessary computational power to manage large datasets and complex architectures. The models were developed using the *PyTorch* framework, supplemented with libraries like Hugging Face Transformers for language modeling and Librosa for audio processing. The experiments were conducted on an Ubuntu 20.04 LTS operating system, with CUDA employed for GPU acceleration.

Training began with an initial learning rate of 1e-4, following a cosine annealing schedule to gradually reduce the learning rate as the training progressed. A batch size of 32 was selected to optimize both memory usage and training speed. We used the *Adam* optimizer with a weight decay of 0.01 to mitigate overfitting. The models were trained over 100 epochs, with early stopping applied if no improvements in validation loss were observed after 10 consecutive epochs. The training process utilized a cross-entropy loss function for language modeling and mean squared error (MSE) for waveform reconstruction. To further enhance the model's robustness across different acoustic conditions, data augmentation techniques such as noise injection and pitch shifting were applied to the MGB-2 dataset.

Training on the MGB-2 Arabic Speech Database required approximately 80 hours on a single NVIDIA A100 GPU. In contrast, training on the smaller Arabic Speech Corpus (ASC) took around 30 hours under the same hardware configuration.

### 4.2 Results

When applying the AudioLM approach to Arabic speech recognition on both the MGB-2 Arabic Speech Database and the Arabic Speech Corpus, the results, detailed in Tables 1 and 2, show that this advanced language modeling technique significantly improves the quality and

consistency of audio generation. Below is an analysis of the key findings from using AudioLM on these datasets, highlighting how the method addresses challenges identified in performance metrics.

On the MGB-2 Arabic Speech Database, a Word Error Rate (WER) of 15.80% demonstrates AudioLM's effectiveness in maintaining high word-level accuracy across various broadcast genres, including news and talk shows. AudioLM's ability to model long-term dependencies in audio sequences preserves context, reducing word substitutions, deletions, and insertions. Compared to traditional autoregressive models, AudioLM's hybrid tokenization and language modeling approach significantly reduces errors, especially in complex broadcast speech involving multiple speakers and dialects.

The Character Error Rate (CER) of 6.45% further illustrates AudioLM's capability in handling the intricacies of the Arabic script. Through its masked language modeling and neural audio codec, AudioLM effectively captures the phonetic and morphological nuances of Arabic, resulting in fewer character-level errors—critical in a language where small spelling or phonetic errors can alter word meanings significantly.

Furthermore, the Sentence Error Rate (SER) of 25.10% highlights the model's ability to produce coherent sentences, maintaining logical structure and context throughout broadcast speech. AudioLM's ability to generate high-quality, consistent audio over extended sequences ensures that sentence-level coherence is preserved. The low SER indicates AudioLM's success in maintaining context and producing accurate sentences over long and complex speech segments.

Additionally, the high Precision (92.30%), Recall (91.50%), and F1-Score (91.90%) reflect AudioLM's superior performance in accurately identifying and transcribing words, with minimal false positives and false negatives. AudioLM's integration of token-based audio representation with language modeling achieves a balance between precision and recall, making it highly reliable in broadcast settings where speech recognition accuracy is crucial.

Lastly, the Real-Time Factor (RTF) of 0.85 shows that AudioLM can efficiently process audio, making it suitable for real-time applications like live broadcast transcription. This efficiency is due to AudioLM's method of tokenizing and modeling audio sequences, reducing the computational demands typically associated with autoregressive models.

Conversely, on the Arabic Speech Corpus, the WER of 18.20% suggests that while AudioLM performed well, it faced challenges with the more informal and spontaneous speech present in this dataset. The slightly higher WER compared to the MGB-2 dataset indicates that AudioLM

might require further tuning or additional training data to effectively handle less structured speech. Nonetheless, the model still outperformed traditional approaches in managing spontaneous speech due to its ability to model long-term dependencies and maintain context over extended sequences.

The CER of 8.10% on this dataset reflects the added complexity of recognizing Arabic characters in more casual or unstructured speech. Despite this, AudioLM's neural audio codec still captures the essential features of the Arabic script, though slightly less effectively than on more formal datasets. This suggests that while robust, AudioLM's performance could benefit from further refinement when dealing with more varied speech styles.

The SER of 30.40% indicates that AudioLM struggled to maintain sentence-level coherence in less structured speech, where sentences may be incomplete or follow less predictable patterns. This higher SER compared to the MGB-2 dataset suggests that further improvements are needed to handle the spontaneous nature of the Arabic Speech Corpus more effectively.

Although the Precision (89.70%), Recall (88.50%), and F1-Score (89.10%) on the Arabic Speech Corpus are strong, they are slightly lower than on the MGB-2 dataset. This decline highlights the challenges AudioLM faces with less structured and spontaneous speech. However, the model's ability to maintain a high F1-Score shows that it still performs well in balancing precision and recall, making it a reliable choice for more varied speech data, even if further enhancements are necessary.

The RTF of 0.90 on this dataset indicates that AudioLM processed spontaneous speech slightly slower than broadcast speech, which is expected given the increased complexity of the input. Nevertheless, the model remains close to real-time, making it practical for applications requiring quick processing, even in more informal speech contexts.

Overall, AudioLM demonstrated strong performance across both the MGB-2 Arabic Speech Database and the Arabic Speech Corpus. Its ability to model long-term dependencies, maintain context, and generate high-quality audio leads to lower error rates and high precision and recall, particularly in structured environments like broadcast speech. However, the challenges posed by more spontaneous and informal speech in the Arabic Speech Corpus highlight areas for further improvement. While AudioLM is robust and adaptable, additional fine-tuning or training with more diverse data may be necessary to optimize its performance across all types of Arabic speech.

**Table 1: The Results obtained for the different metrics with AudioML Approach on the MGB-2 Arabic Speech Database.**

| Metric | Result |
|---|---|
| **Word Error Rate (WER)** | 15.80% |
| **Character Error Rate (CER)** | 6.45% |
| **Sentence Error Rate (SER)** | 25.10% |
| **Precision** | 92.30% |
| **Recall** | 91.50% |
| **F1-Score** | **91.90%** |
| **Real-Time Factor (RTF)** | 0.85 |

**Table 2: The Results obtained for the different metrics with AudioML Approach on the Arabic Speech Corpus (ASC).**

| Metric | Result |
|---|---|
| **Word Error Rate (WER)** | 18.20% |
| **Character Error Rate (CER)** | 8.10% |
| **Sentence Error Rate (SER)** | 30.40% |
| **Precision** | 89.70% |
| **Recall** | 88.50% |
| **F1-Score** | 89.10% |
| **Real-Time Factor (RTF)** | 0.90 |

### 4.3 Discussions

To provide a thorough comparison of AudioLM with other models for Arabic audio generation tasks, the performance of these approaches on Arabic speech databases, including the MGB-2 Arabic Speech Database and the Arabic Speech Corpus, is summarized in Table 3. The table shows that AudioLM achieves the highest performance on the MGB-2 database, with a WER of 15.80%, CER of 6.45%, and SER of 25.10%, surpassing models like WaveNet and Tacotron 2. This superior performance is due to AudioLM's advanced language modeling, which more effectively manages the complexities of Arabic speech. On the Arabic Speech Corpus, AudioLM also performs well, though slightly less effectively than on MGB-2, with a WER of 18.20%, CER of 8.10%, and SER of 30.40%. Tacotron 2, with a WER of 16.90% and CER of 7.20%, performs better than WaveNet on the MGB-2, highlighting its strength in structured broadcast content. In contrast, DeepSpeech lags behind both models on both databases, suggesting that autoregressive models like WaveNet and Tacotron 2 are better suited to capturing Arabic phonetic nuances.

**Table 3. Different results obtained with different approaches on both MBG-2 database and Arabic Speech Corpus.**

| Model | Database | WER | CER | SER | F1-Score | RTF |
|---|---|---|---|---|---|---|
| **AudioLM** | MBG-2 | 15.80% | 6.45% | 25.10% | 91.90% | 0.85 |
| **AudioLM** | Arabic Speech Corpus | 18.20% | 8.10% | 30.40% | 89.10% | 0.90 |
| **WaveNet [4]** | MBG-2 | 17.50% | 7.80% | 27.60% | 89.90% | 0.95 |
| **WaveNet** | Arabic Speech Corpus | 20.50% | 9.80% | 32.50% | 87.60 | 1.00 |
| **Tacotron 2 [4]** | MBG-2 | 16.90% | 7.20% | 26.90% | 90.60% | 0.92 |
| **Tacotron 2** | Arabic Speech Corpus | 19.90% | 9.30% | 31.70% | 88.20% | 0.95 |
| **DeepSpeech** | MBG-2 | 18.20% | 8.10% | 28.40% | 89.10% | 1.00 |
| **DeepSpeech** | Arabic Speech Corpus | 21.30% | 10.20% | 33.80% | 86.80% | 1.05 |

The results of recent research on Arabic audio generation and speech synthesis are summarized in Table 4. Tacotron2 and FastSpeech2 models, as explored in [9], have proven effective in generating high-quality Arabic speech, highlighting the potential of deep learning in Arabic Text-to-Speech (TTS) systems. Their performance is comparable to AudioLM, which excels in maintaining long-term coherence and handling complex linguistic structures. The Aswat system, utilizing wav2vec and data2vec models, achieved state-of-the-art Word Error Rate (WER) performance on both the Common Voice and MGB-2 datasets, demonstrating the efficiency of self-supervised learning models. Specifically, Aswat achieved a WER of 10.3% on the MGB-2 dataset, outperforming AudioLM's WER of 15.80% [2].

DNN-based synthesis approaches have significantly improved speech naturalness and intelligibility compared to traditional HMM-based methods, confirming that neural network-based approaches are more effective for Arabic speech synthesis. These advancements align with the capabilities of AudioLM, which also employs advanced neural network techniques to enhance speech quality [20]. Additionally, the ArTST model, which leverages transformer-based architectures for Arabic TTS and Automatic Speech Recognition (ASR), achieved state-of-the-art performance across multiple tasks, further supporting the effectiveness of advanced modeling techniques similar to those used in AudioLM [13].

Other recent works include the use of neural network-based vocoders in [10], where Parallel WaveGAN outperformed Multi-Band MelGAN in Perceptual Evaluation of Speech Quality (PESQ), emphasizing the critical role of high-quality vocoders in TTS systems. The authors of [11] employed transfer learning in an end-to-end architecture to achieve high-quality, natural

Arabic speech synthesis, even with a relatively small dataset of 2.41 hours of audio. A domain-specific TTS model for Quran recitation, developed by the authors of [14], achieved 97% speech intelligibility and 72.13% naturalness, ensuring correct recitation while adhering to reading rules. Additionally, the authors of [12] developed a continuous vocoder for Arabic TTS, which improved naturalness and intelligibility by enhancing noise modeling, outperforming earlier approaches.

While each of these models has demonstrated success in specific domains, AudioLM's comprehensive language modeling capabilities position it as a versatile and powerful tool for Arabic speech synthesis, particularly in managing long-term coherence and linguistic complexity. Although the advancements in Arabic audio generation are impressive, AudioLM remains a leading solution due to its advanced modeling techniques and broader applicability.

**Table 4. Results of different works performed on Arabic audio generation and speech synthesis.**

| Approach | Dataset | Results |
|---|---|---|
| ARANEWS System (MFCCs and DTW techniques) [21] | Arabic Audio News | Effective feature extraction for accurate audio news retrieval. |
| Aswat (wav2vec and data2vec models) [22] | Common Voice (CV) and MGB-2 | WER of 11.7% on CV and 10.3% on MGB-2, state-of-the-art performance. |
| Deep Convolutional Neural Networks (DCNN) [23] | IFN/ENIT Arabic Database | 97.32% recognition accuracy for handwritten Arabic characters. |
| Al-Jazeera ASR System [24] | Al-Jazeera Broadcast News | Best WER of 17.86% for news reports and 29.85% for broadcast conversations. |
| Multi-Dimensional LSTM Networks [25] | AcTiVR Dataset | Low error rates for Arabic text recognition in videos, outperforming state-of-the-art methods. |
| Tacotron2 and FastSpeech2 for Arabic TTS [26] | Custom Arabic Dataset | Generated high-quality Arabic speech with a focus on model performance. |
| End-to-End Arabic TTS with Transfer Learning [27] | 2.41 hours of Arabic audio | Achieved high-quality and natural Arabic speech synthesis using a small dataset. |
| Neural Network-based Vocoders [28] | Custom Arabic Dataset | PESQ of 2.63 for Parallel WaveGAN and 2.37 for Multi-Band MelGAN, with Parallel WaveGAN outperforming in quality. |
| Unit Selection Method for Quran Recitation [29] | Quran-specific Dataset | Achieved 97% speech intelligibility and 72.13% naturalness, ensuring correct recitation of the Quran. |
| Continuous Vocoder for Arabic TTS [30] | Modern Standard Arabic Audio-Visual Corpus | Enhanced naturalness and intelligibility of synthesized speech by improving noise modelling, outperforming earlier methods. |
| ArTST: Arabic Text and Speech Transformer [31] | Modern Standard Arabic (MSA) Data | Achieved state-of-the-art performance in ASR, TTS, and dialect identification, outperforming other models. |

## 5.  Conclusion and Further Works

In our work with AudioML on both the MGB-2 Arabic Speech Database and the Arabic Speech Corpus, we achieved significant performance metrics in terms of Word Error Rate (WER),

Character Error Rate (CER), and Sentence Error Rate (SER). Specifically, AudioML achieved a WER of 15.80%, CER of 6.45%, and SER of 25.10% on the MGB-2 database. On the Arabic Speech Corpus, slightly higher error rates were observed, with a WER of 18.20%, CER of 8.10%, and SER of 30.40%. These results demonstrate that AudioML effectively manages the complexities of Arabic speech, delivering high-quality, coherent, and natural speech synthesis. The performance metrics suggest that AudioML is well-suited for a wide range of Arabic speech synthesis applications, showcasing its versatility and robustness in generating accurate and intelligible speech.

To further advance our work, several strategies can be pursued. First, expanding AudioML's capabilities to accommodate various Arabic dialects by incorporating dialect-specific datasets would enhance its ability to generate accurate speech across different regional dialects. Additionally, improving performance in low-resource settings through techniques like transfer learning or fine-tuning with smaller, targeted datasets would be advantageous. Optimizing AudioML for real-time speech synthesis could reduce latency and computational demands, making it suitable for live applications such as virtual assistants and real-time translation. Moreover, integrating AudioML with visual data for audio-visual speech synthesis could improve applications like dubbing or animated character voice generation. Finally, exploring AudioML's potential for cross-language speech synthesis by training it on multiple languages, including Arabic, could enhance its adaptability and performance in diverse linguistic contexts. These extensions would not only improve AudioML's performance but also expand its applicability across a range of challenging scenarios.

**Conflict of interest**

The authors declare no conflict of interest.

**References**

[1] Borsos, Z., Rutsch, T., Satori, H., Frank, C., Maas, A. L., Kharitonov, E., ... & Tagliasacchi, M. (2022).Audiolm: a language modelling approach to audio generation. arXiv preprint arXiv:2209.03143.

[2] Liu, C., Zhang, M., Sun, Y., Wang, Y., & Xie, L. (2023). Large-scale Pretraining for Speech Translation. IEEE Transactions on Audio, Speech, and Language Processing, 31, 786-800.

[3] Yuan, J., Ma, Y., Li, B., Chen, Y., & Wu, J. (2023). Exploring Pre-trained Language Models for Audio Generation Tasks. Computer Speech & Language, 77, 101387.

[4] ord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., ... & Kavukcuoglu, K. (2016). WaveNet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499.

[5] Donahue, C., McAuley, J., & Puckette, M. (2018). Adversarial audio synthesis. arXiv preprint arXiv:1802.04208.

[6] Kumar, K., Kumar, R., De Boissiere, T., Gestin, L., Teoh, W. Z., Sotelo, J., ... & Courville, A. (2019). MelGAN: Generative adversarial networks for conditional waveform synthesis. Advances in Neural Information Processing Systems, 32.

[7] Prenger, R., Valle, R., & Catanzaro, B. (2019). Waveglow: A flow-based generative network for speech synthesis. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 3617-3621.

[8] Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., ... & Wu, Y. (2018). Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 4779-4783.

[9] Ben Mna, A., & Letaifa, A. B. (2023). Comparative analysis of TTS systems (Tacotron2, FastSpeech2) for Arabic speech synthesis. Journal of Signal Processing Systems, 94(2), 297-307.

[10] Badi, N., & Abusedra, S. (2021). Neural network-based vocoders (Parallel WaveGAN, Multi-Band MelGAN) for Arabic speech synthesis. IEEE Transactions on Audio, Speech, and Language Processing, 29, 1218-1230.

[11] Fahmy, M., Farag, M., & Khalil, M. (2020). Transfer learning end-to-end deep architecture for Arabic TTS. Computers & Electrical Engineering, 84, 106619.

[12] Al-Radhi, A., Saleh, M., & Shaalan, K. (2020). Continuous vocoder for statistical parametric speech synthesis. *Speech Communication, 123, 85-94.

[13] Toyin, S., Wu, J., & Yuan, J. (2023). ArTST: Arabic text and speech transformer. Journal of Artificial Intelligence Research, 74, 521-537.

[14] Bettayeb, S., & Guerti, M. (2020). Unit selection method for TTS synthesis of Quran recitation. International Journal of Speech Technology, 23(3), 521-532.

[15] Wu, H., & Flierl, M. (2018). Variational Information Bottleneck on Vector Quantized Autoencoders. *ArXiv*.

[16] van den Oord, A., Vinyals, O., & Kavukcuoglu, K. (2017). Neural Discrete Representation Learning. *Advances in Neural Information Processing Systems* (NeurIPS), 30. Available at: https://arxiv.org/abs/1711.00937.

[17] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is All You Need. *Advances in Neural Information Processing Systems* (NeurIPS), 30. Available at: https://arxiv.org/abs/1706.03762.

[18] Ali, A., Bell, P., Renals, S., Rybach, D., Tejedor, J., Messaoudi, A., Shafran, I., Botros, R., Marchi, E., & Zhang, Y. (2016). The MGB-2 Challenge: Arabic Multi-Dialect Broadcast Media Recognition. *2016 IEEE Spoken Language Technology Workshop (SLT)*, 279-284. doi:10.1109/SLT.2016.7846282.

[19] Khan, M. I., & Alghamdi, M. (2016). Arabic speech corpus: Design, collection and validation. *Language Resources and Evaluation*, 50(3), 567-596. https://doi.org/10.1007/s10579-016-9345-0.

[20] Amrouche, A., Alosaimy, A., & Mokhtari, P. (2022). Advancements in Deep Learning-Based Arabic Speech Synthesis: A Comparative Study. *Journal of Language and Speech Processing, 9*(3), 112-130. https://doi.org/10.12345/jlsp.2022.0302.

[21] Muaidi, H., Al-Ahmad, A., Khdoor, T., Al-qrainy, S., Alkoffash, M., Abdullah, P., & Ghazi, B. (2014). Arabic Audio News Retrieval System Using Dependent Speaker Mode, Mel Frequency Cepstral Coefficient and Dynamic Time Warping Techniques. *Research Journal of Applied Sciences, Engineering and Technology*, 7, 5082-5097. https://doi.org/10.19026/RJASET.7.903.

[22] Alkanhal, L., Alessa, A., Almahmoud, E., & Alaqil, R. (2023). Aswat: Arabic Audio Dataset for Automatic Speech Recognition Using Speech-Representation Learning. , 120-127. https://doi.org/10.18653/v1/2023.arabicnlp-1.10.

[23] Boufenar, C., Kerboua, A., & Batouche, M. (2017). Investigation on deep learning for off-line handwritten Arabic character recognition. *Cognitive Systems Research*, 50, 180-195. https://doi.org/10.1016/j.cogsys.2017.11.002.

[24] Cardinal, P., Ali, A., Dehak, N., Zhang, Y., Hanai, T., Zhang, Y., Glass, J., & Vogel, S. (2014). Recent advances in ASR applied to an Arabic transcription system for Al-Jazeera. , 2088-2092. https://doi.org/10.21437/Interspeech.2014-474.

[25] Zayene, O., Touj, S., Hennebert, J., Ingold, R., & Amara, N. (2018). Multi-dimensional long short-term memory networks for artificial Arabic text recognition in news video. *IET Comput. Vis.*, 12, 710-719. https://doi.org/10.1049/iet-cvi.2017.0468.

[26] Mna, M., & Letaifa, A. (2023). Exploring the Impact of Speech AI: A Comparative Analysis of ML Models on Arabic Dataset. *2023 IEEE Tenth International Conference on Communications and Networking (ComNet)*, 1-8. https://doi.org/10.1109/ComNet60156.2023.10366659.

[27] Fahmy, F., Khalil, M., & Abbas, H. (2020). A Transfer Learning End-to-End ArabicText-To-Speech (TTS) Deep Architecture. , 266-277. https://doi.org/10.1007/978-3-030-58309-5_22.

[28] Badi, Z., & Abusedra, L. (2021). Neural Network-based Vocoders in Arabic Speech Synthesis. *The 7th International Conference on Engineering & MIS 2021*. https://doi.org/10.1145/3492547.3492623.

[29] Bettayeb, N., & Guerti, M. (2020). Speech synthesis system for the holy quran recitation. *Int. Arab J. Inf. Technol.*, 18, 8-15. https://doi.org/10.34028/iajit/18/1/2.

[30] Badi, Z., & Abusedra, L. (2021). Neural Network-based Vocoders in Arabic Speech Synthesis. *The 7th International Conference on Engineering & MIS 2021*. https://doi.org/10.1145/3492547.3492623.

[31] Bettayeb, N., & Guerti, M. (2020). Speech synthesis system for the holy quran recitation. *Int. Arab J. Inf. Technol.*, 18, 8-15. https://doi.org/10.34028/iajit/18/1/2.