# JambaTalk: Speech-Driven 3D Talking Head Generation based on a Hybrid Transformer-Mamba Model

Farzaneh Jafari[1], Stefano Berretti[2], Anup Basu[1]
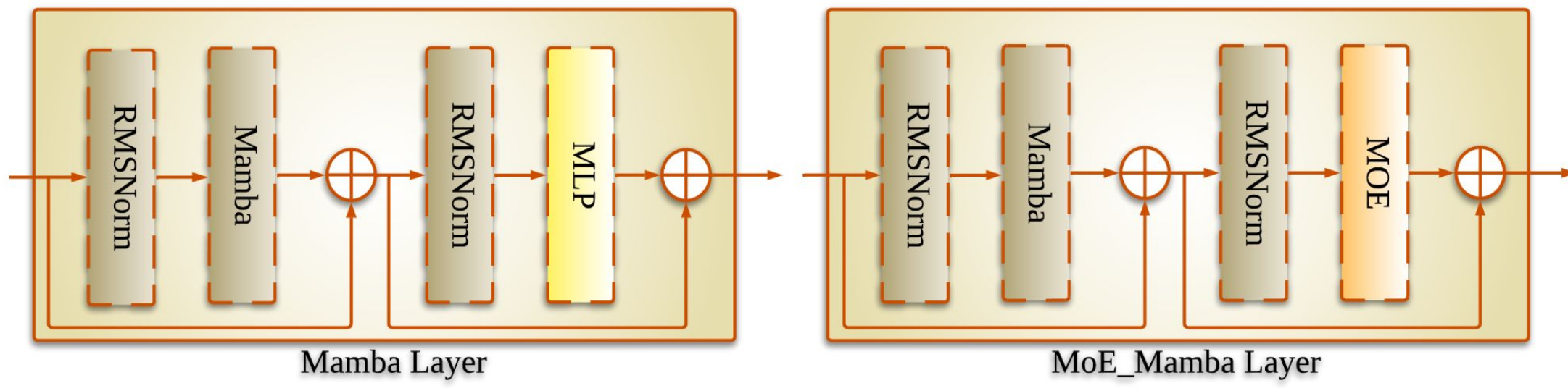
[1]University of Alberta, Canada | [2]University of Florence, Italy
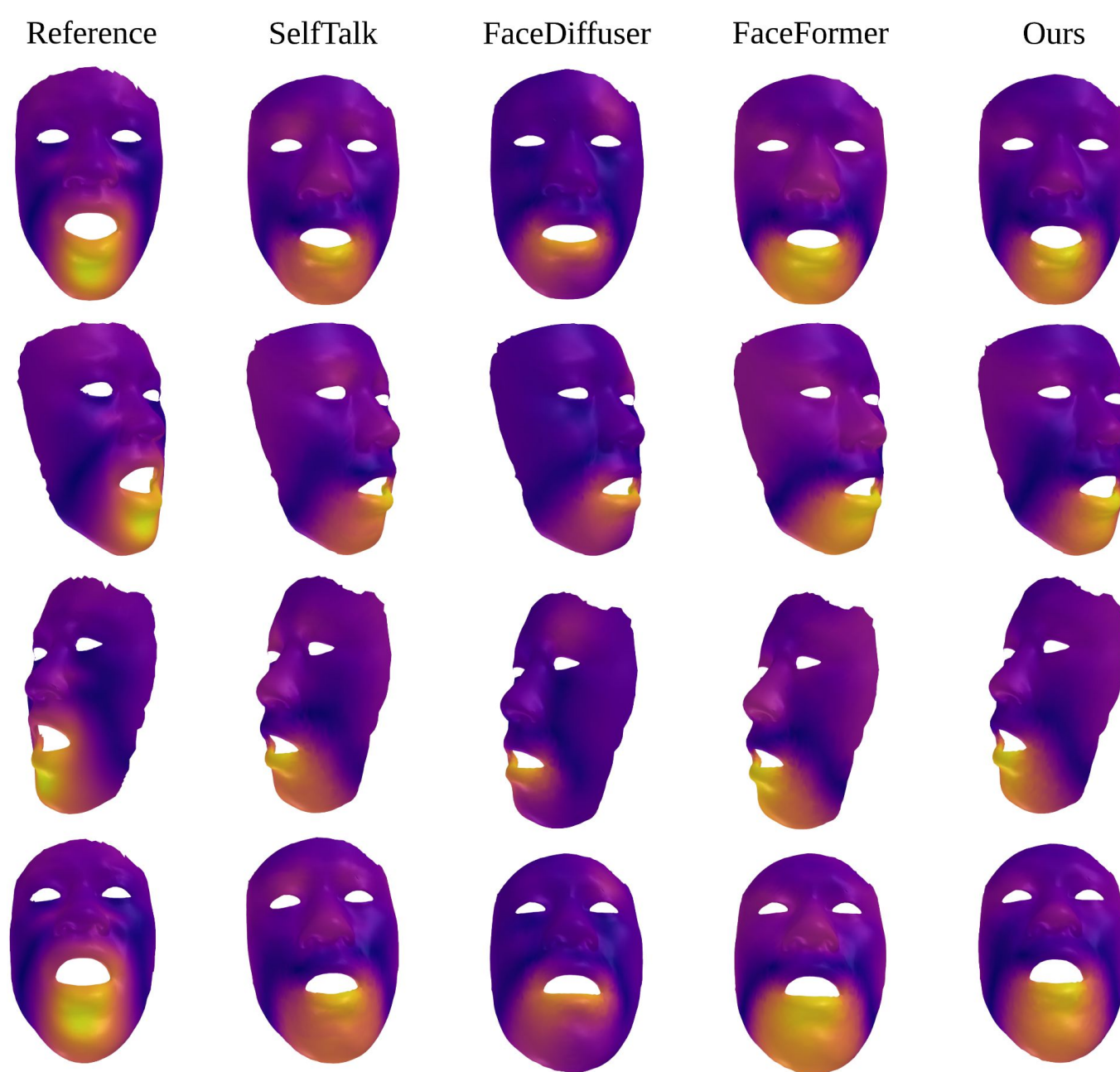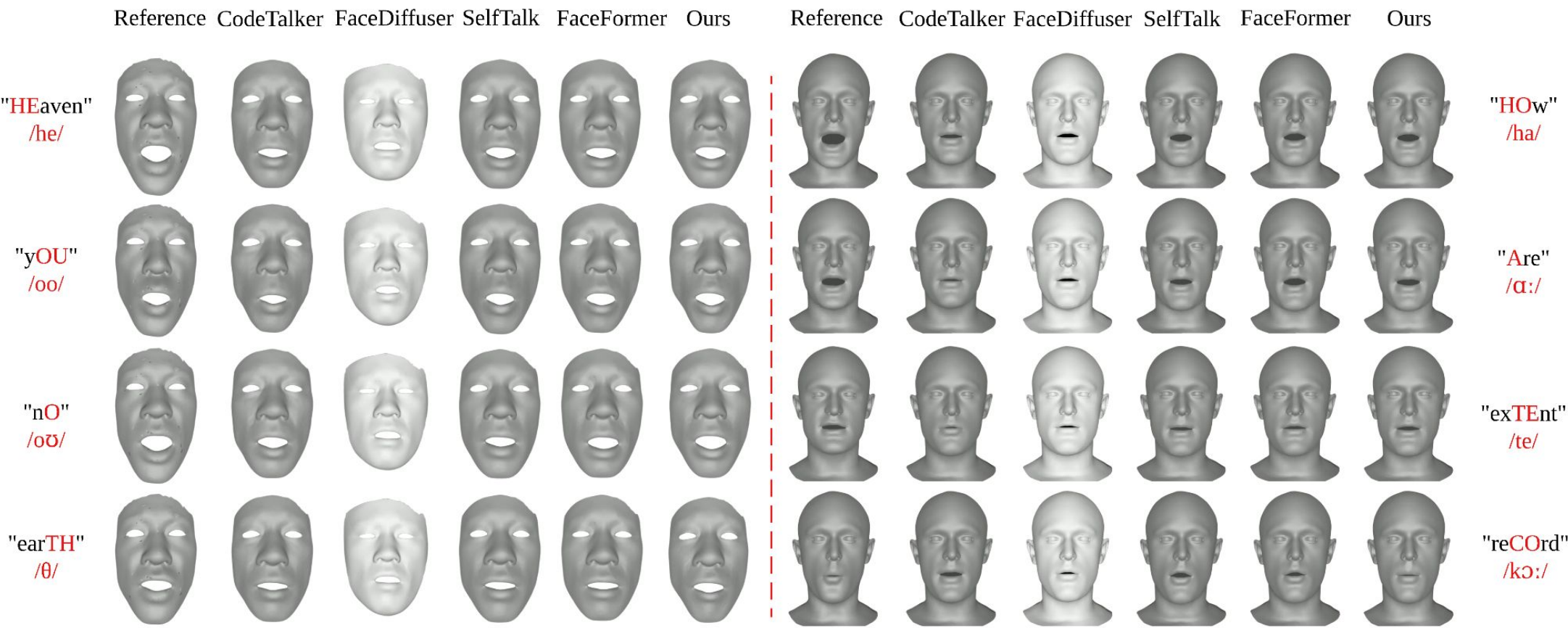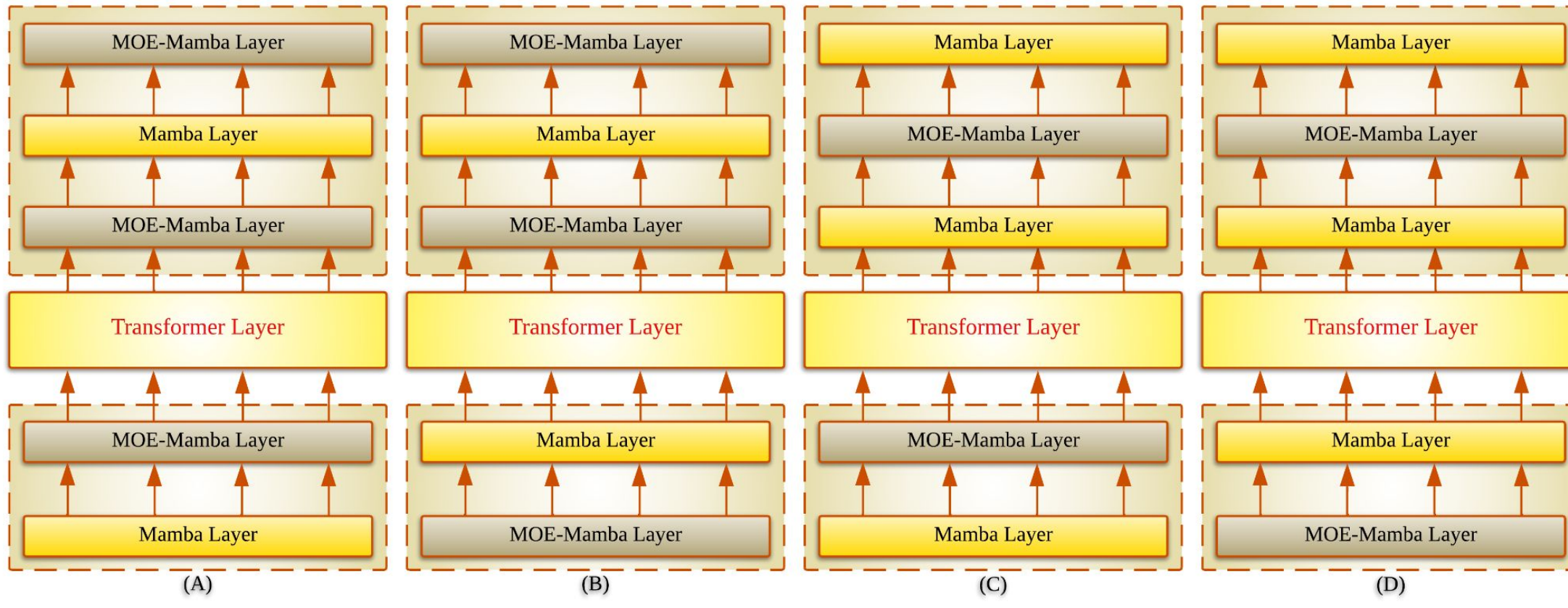
## Problem and Motivation

Current talking head generation faces key challenges:

- **Limited sequence length:** Transformers struggle with long sequences due to quadratic complexity
- **Memory constraints:** High computational cost restricts real-time applications
- **Quality trade-offs:** No model excels across all metrics

**Our Solution:** Hybrid Transformer-Mamba architecture combining attention mechanisms with State Space Models for efficient long-sequence processing.
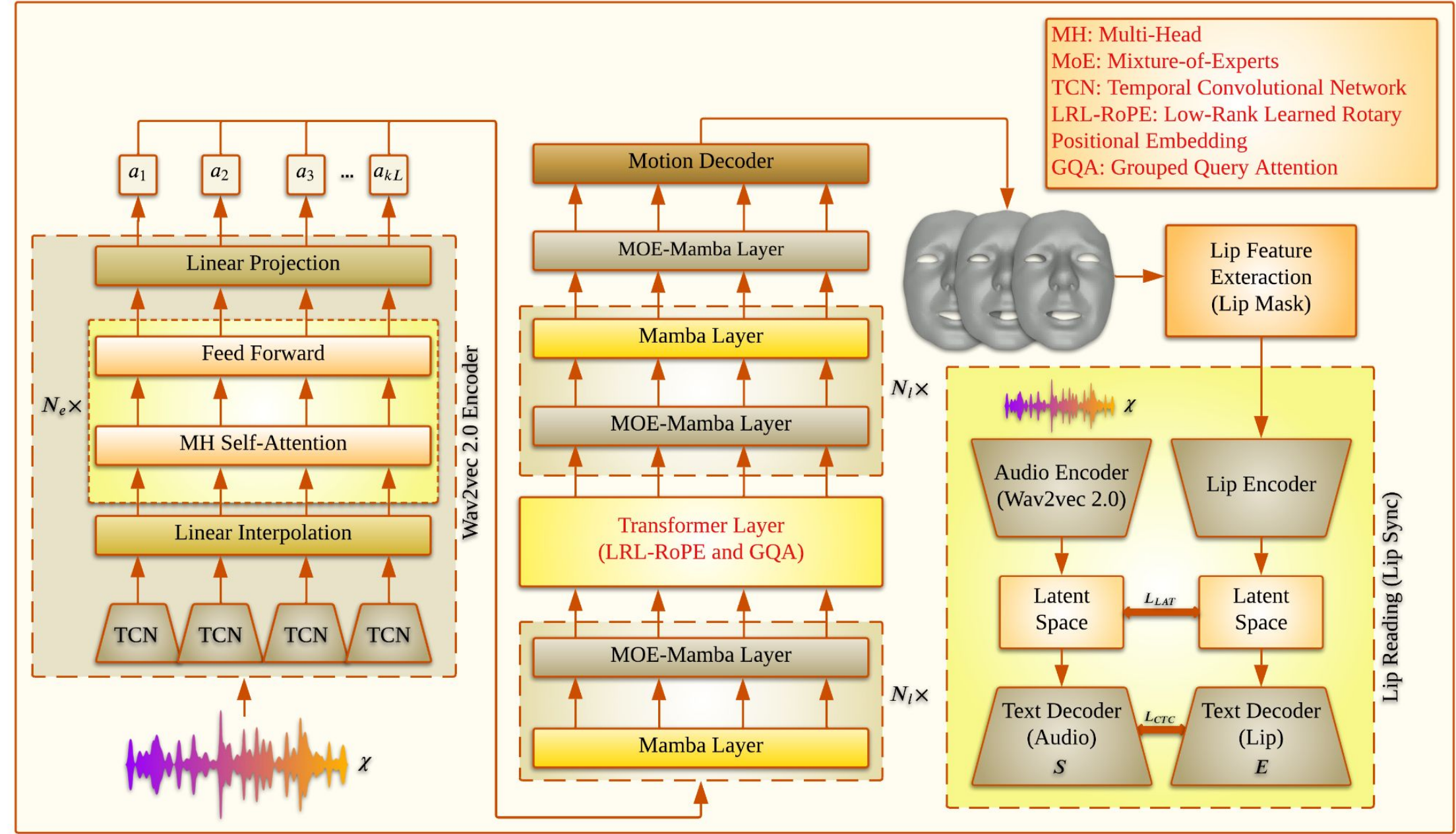


MH: Multi-Head
MoE: Mixture-of-Experts
TCN: Temporal Convolutional Network
LRL-RoPE: Low-Rank Learned Rotary Positional Embedding
GQA: Grouped Query Attention

Details of the Mamba and MoE_Mamba layers in the JambaTalk Decoder. Both layers begin with an RMSNorm normalization followed by a Mamba block for sequence modeling and include residual connections to preserve gradient flow. In the standard Mamba Layer (left), the Mamba output is followed by another RMSNorm and a feedforward MLP block. In contrast, the MoE_Mamba Layer (right) replaces the MLP with a Mixture-of-Experts (MoE) module, enabling dynamic expert routing per token and enhancing model capacity while maintaining computational efficiency.
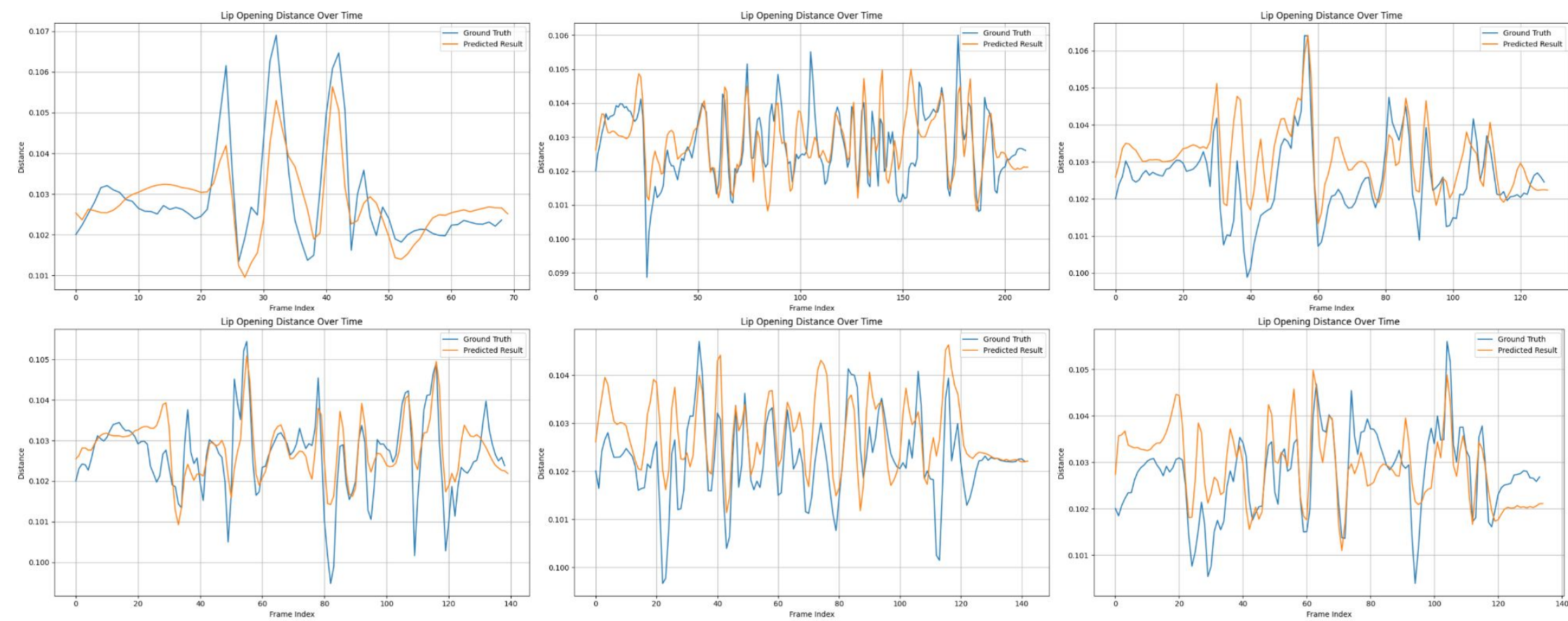
**Overview of JambaTalk:** The Wav2Vec 2.0 model is used to extract features from the input speech, with the encoder initialized using pre-trained weights from the original model. These encoded features are passed to the JambaTalk decoder, which generates a sequence of animated 3D face meshes. The Transformer layer incorporates Low-Rank Learned Rotary Positional Embedding (LRL-RoPE) and Grouped-Query Attention (GQA), providing a computation-efficient alternative to traditional Transformers. The lip feature extraction block then converts motion decoder outputs into lip deformation features by selecting lip vertices with a lip mask, which are processed by a Transformer-based lip encoder in the lip reader module to synchronize lip shapes.
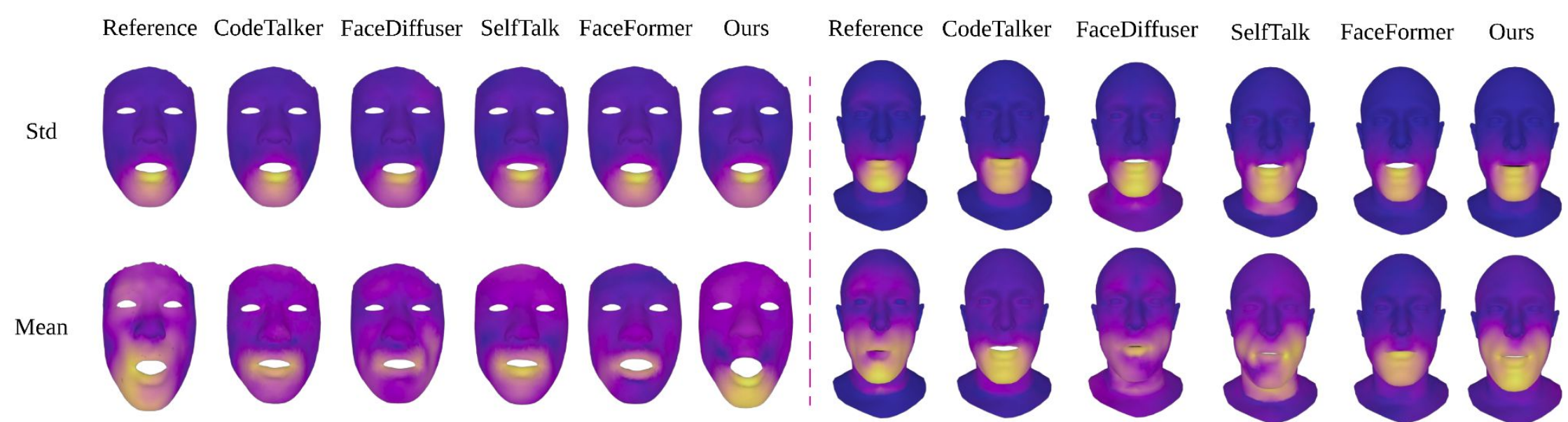


**Algorithm 1** Low-Rank Learned RoPE (LRL-RoPE)

**Require:** $x \in \mathbb{R}^{B \times T \times d}$, $pos\_idx \in \mathbb{R}^{T}$, $W_1 \in \mathbb{R}^{d \times r}$, $W_2 \in \mathbb{R}^{r \times d}$

1: **function** LowRankLearnedRoPE($x$, $pos\_idx$, $W_1$, $W_2$)
2:    $\theta \leftarrow 1/(10000^{\text{arange}(0,d,2)/d})$     ▷ predefined frequencies, size $[d/2]$
3:    $\text{angles} \leftarrow pos\_idx[:, None] \cdot \theta[None, :]$     ▷ $[T, d/2]$
4:    $\sin\_pos \leftarrow \sin(\text{angles})$
5:    $\cos\_pos \leftarrow \cos(\text{angles})$
6:    $\text{base\_emb} \leftarrow \text{concat}(\sin\_pos, \cos\_pos, \dim=-1)$     ▷ $[T, d]$
7:    $\text{learned\_emb} \leftarrow \text{base\_emb} \cdot W_1^{\top} \cdot W_2^{\top}$     ▷ low-rank learnable correction with transposes
8:    $\Delta\sin, \Delta\cos \leftarrow \text{split}(\text{learned\_emb}, 2)$
9:    $\sin\_final \leftarrow \sin\_pos + \Delta\sin$
10:   $\cos\_final \leftarrow \cos\_pos + \Delta\cos$
11:   $x_{\text{even}} \leftarrow x[..., 0 :: 2]$
12:   $x_{\text{odd}} \leftarrow x[..., 1 :: 2]$
13:   $x_{\text{rot\_even}} \leftarrow x_{\text{even}} * \cos\_final[None, :, :] - x_{\text{odd}} * \sin\_final[None, :, :]$
14:   $x_{\text{rot\_odd}} \leftarrow x_{\text{even}} * \sin\_final[None, :, :] + x_{\text{odd}} * \cos\_final[None, :, :]$
15:   **return** interleave($x_{\text{rot\_even}}, x_{\text{rot\_odd}}$)
16: **end function**





A visual comparison of frames from synthesized facial animation sequences produced by various methods, alongside reference frames from the ground-truth sequence. The red utterances are depicted in the visual frames. Our approach generates lip shapes that closely resemble the reference frames. Left: $BIWI\_6$ Test-B. Right: Vocaset Test.



Lip opening distance over time, showing the variation in 3D Euclidean distance between the upper and lower lip landmarks for each video frame. Peaks indicate moments when the mouth is open wider, while valleys correspond to smaller openings or closed lips on $BIWI\_6$ Test-B dataset.

A visual comparison of frames from synthesized facial animation sequences produced by various methods, alongside reference frames from the ground-truth sequence. The red utterances are depicted in the visual frames. Our approach generates lip shapes that closely resemble the reference frames. Left: $BIWI\_6$ Test-B. Right: Vocaset Test.

Qualitative comparison of mouth dynamics between the reference and different models (SelfTalk, FaceDiffuser, FaceFormer, and our proposed JambaTalk). The mean heatmaps visualize the lip and jaw motion intensity during speech.



The temporal statistics (mean and standard deviation) of motion variations between adjacent frames in the sequence on Vocaset Test and $BIWI\_6$ Test-B datasets.

Project Page     arXiv     GitHub