

# Summary of Meshtalk WIP

Jay  
CAU EAI Lab  
Lab address Here  
<https://yukinyaa.github.io>

Data Class	Resolution	Frame rate
Face Video	80 cameras	30 FPS
Face Mesh	1,672 vertices	30 FPS
Audio	16kHz	-
Mel Spectrogram	80 Dimension	10ms(100 FPS)

Table 1. Captured, then processed Datasets

## Abstract

The article summarizes Meshtalk[3]. The goal of this practice is to get familiar with  $\text{\LaTeX}$  and technical writing(hopefully).

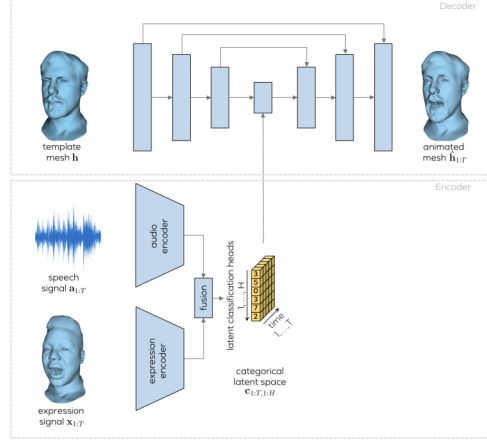


Figure 1. The network diagram[3].

## 1. Introduction to Meshtalk

Meshtalk is a generic method for generating full facial mesh animation from speech. It can generate lip-sync animation from a single frame of generic human facial mesh, and also can mix in emotional information from mesh animation.

### 1.1. Dataset

In total, there are 1.4 million frames of 13 hours equivalent audio-visual dataset from 250 subjects, reading 50 phonetically balanced sentences.

For training dataset, first 40 sentences out of 50 about 200 out of 250 subjects, total  $40 \times 200$  dataset is used for training. For evaluation, remaining 10 sentences of 50 subjects are used.

### Mesh Dataset

Face motion is captured with synchronized cameras, and processed into high-detail mesh, including eyelid, hairstyle, etc.

### Audio Dataset

Audio data is recorded in 16kHz as shown on Table 1. For every mesh, 600ms of audio snippet is processed into Mel spectrogram, every 10ms, in 80 dimension. Hence,  $a_t \in \mathbb{R}^{60 \times 80}$

### 1.2. Network Design

The network resembles Variational Autoencoder with multiple latent space as shown in Figure 1. Target mesh  $\hat{h}$  is estimated from template mesh  $h$  and latent space  $c$  by computing function  $\mathcal{D}$ .

$$\hat{h}_{1:T} = \mathcal{D}(h, c_{1:T,1:H}) \quad (1)$$

Sequence of latent space  $c_{1:T}$  is derived from audio sequence  $a_{1:T}$  and expression signal mesh sequence  $x_{1:T}$ .  $c$  and  $a$  are first mapped to  $T * H * C$  dimensional latent space, then passed through Gumbel-softmax[2] over every classification head.

$$c_{1:T,1:H} = [\text{Gumbel}(\text{enc}_{t,h,1:C})]_{1:T,1:H} \quad (2)$$

$$\text{enc}_{1:T,1:H,1:C} = \tilde{\xi}(x_{1:T}, a_{1:T}) \quad (3)$$

### 1.3. Training

The solution uses a novel cross-modality loss for calculating loss function for the network.

$$\begin{aligned} \mathcal{L}_{xMod} = & \sum_{t=1}^T \sum_{v=1}^V \mathcal{M}_v^{(upper)} (||\hat{h}_{t,v}^{(expr)} - x_{t,v}||) \\ & + \sum_{t=1}^T \sum_{v=1}^V \mathcal{M}_v^{(mouth)} (||\hat{h}_{t,v}^{(audio)} - x_{t,v}||) \end{aligned} \quad (4)$$

Where  $\hat{h}_{t,v}^{(expr)}$  is estimated with correct expression signal and random speech signal, and  $\mathcal{M}^{(upper)}$  is a mask that assigned a higher weight to vertices around the mouth, and low weight to others. Correspondingly,  $\hat{h}_{t,v}^{(audio)}$  is estimated with random expression signal and random speech signal.  $\mathcal{M}^{(mouth)}$  is a mask that assigns a lower weight around the mouth, and vice versa. Loss for eye is defined as following.

$$\mathcal{L}_{eyelid} = \sum_{t=1}^T \sum_{v=1}^V \mathcal{M}_v^{(eyelid)} (||\hat{h}_{t,v}^{(expr)} - x_{t,v}||) \quad (5)$$

$\mathcal{M}^{(eyelid)}$  is a specific eye loss, which was crucial.

Final Loss term is defined as  $\mathcal{L} = \mathcal{L}_{xMod} + \mathcal{L}_{eyelid}$ , which gives both term equal weight.

## 2. Evaluation

To evaluate the network, the network is compared with various other networks with both numerical analysis and user-study.

### 2.1. multi-model input

The network is compared with networks that lacks speech signal and uses  $\ell_2$  loss, and identical network with Meshtalk trained with  $\ell_2$  loss. Evaluation strategy is perplexity  $PP = p(c_{1:T,1:H} | a_{1:T})^{(1 - \frac{1}{T*H})}$  in this case.

### 2.2. Disentanglement evaluation

Latent space is visualized to demonstrate latent configurations that are affected by audio input and expression input are well separated to each other. Also, vertexes affected by each feature are computed. Interesting feature is that not only jaw and lip vertexes are affected by audio, but also eyebrow is significantly affected.

### 2.3. lip-vertex error

Lip-vertex error rate(in maximal  $\ell_2$  error, mm) is computed, and compared with VOCA[1].

### 2.4. Perceptual evaluation

Compared with VOCA and ground truth, user study is conducted, result is shown at 2

	favorability			ours better or equal
	competitor	equal	ours	
ours vs VOCA[1]				
full-face	24.7%	20.9%	54.4%	75.3%
lip sync	24.7%	20.9%	54.4%	75.3%
upper face	24.7%	20.9%	54.4%	75.3%
ours vs ground truth				
full-face	24.7%	20.9%	54.4%	75.3%
lip sync	24.7%	20.9%	54.4%	75.3%
upper face	24.7%	20.9%	54.4%	75.3%

Table 2. Perceptual study

## 2.5. examples

Example with re-targeting, dubbing is demonstrated

## 3. limitation

(1) Look-ahead about 100ms is required for audio input. This is beneficial for lip-sync for sound such as ‘p’. (2) The model cannot be ran in real-time on low-cost devices. (3) If software fails to track certain parts, correlation fails. e.g hair overlaps eyebrow.

## References

- [1] D. Cudeiro, T. Bolkart, C. Laidlaw, A. Ranjan, and M. J. Black. Capture, learning, and synthesis of 3d speaking styles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [2] E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax, 2017. 1
- [3] A. Richard, M. Zollhoefer, Y. Wen, F. de la Torre, and Y. Sheikh. Meshtalk: 3d face animation from speech using cross-modality disentanglement, 2021. 1