



CS5647: Sound and Music Computing

# Singing Voice Synthesis with Avatar Generation

Team 6

# Singing Face Generation



NUS  
National University  
of Singapore

School of  
Computing

## 1. Voice Synthesis

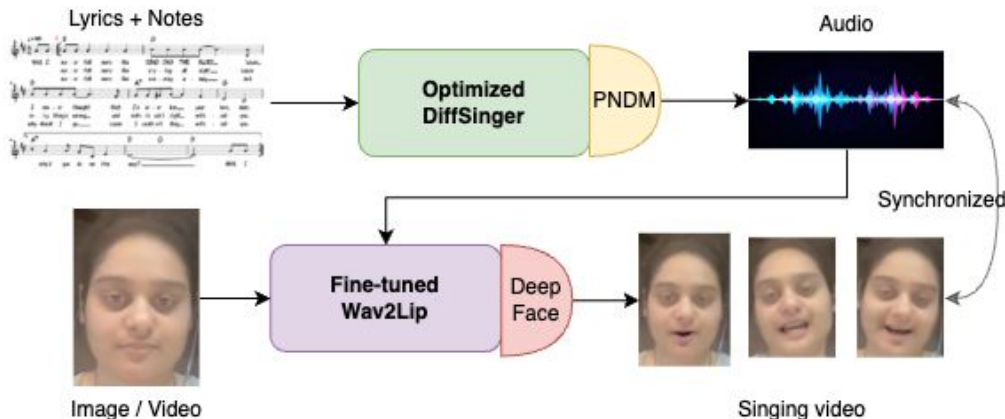
1. Baseline DiffSinger
2. Better Vocoder (BigVGAN)
3. Faster Inference using PNDM

## 2. Face Synthesis

1. Baseline Wave2Lip
2. Fine-tune for singing
3. DeepFace augmentation

## 3. Dataset

1. OpenCPop
2. URSing (Audio-Visual Solo Singing)
3. Own dataset (me singing!)





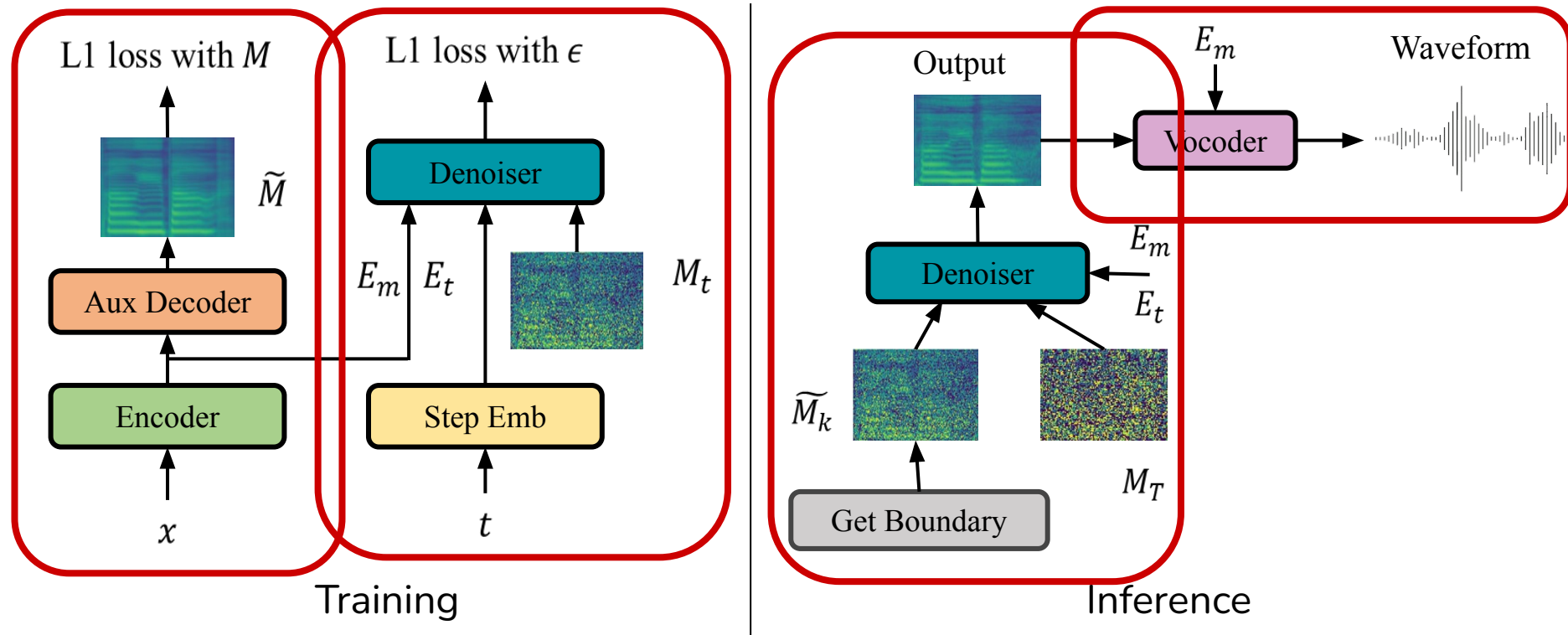
# Voice Synthesis

# DiffSinger



NUS  
National University  
of Singapore

School of  
Computing

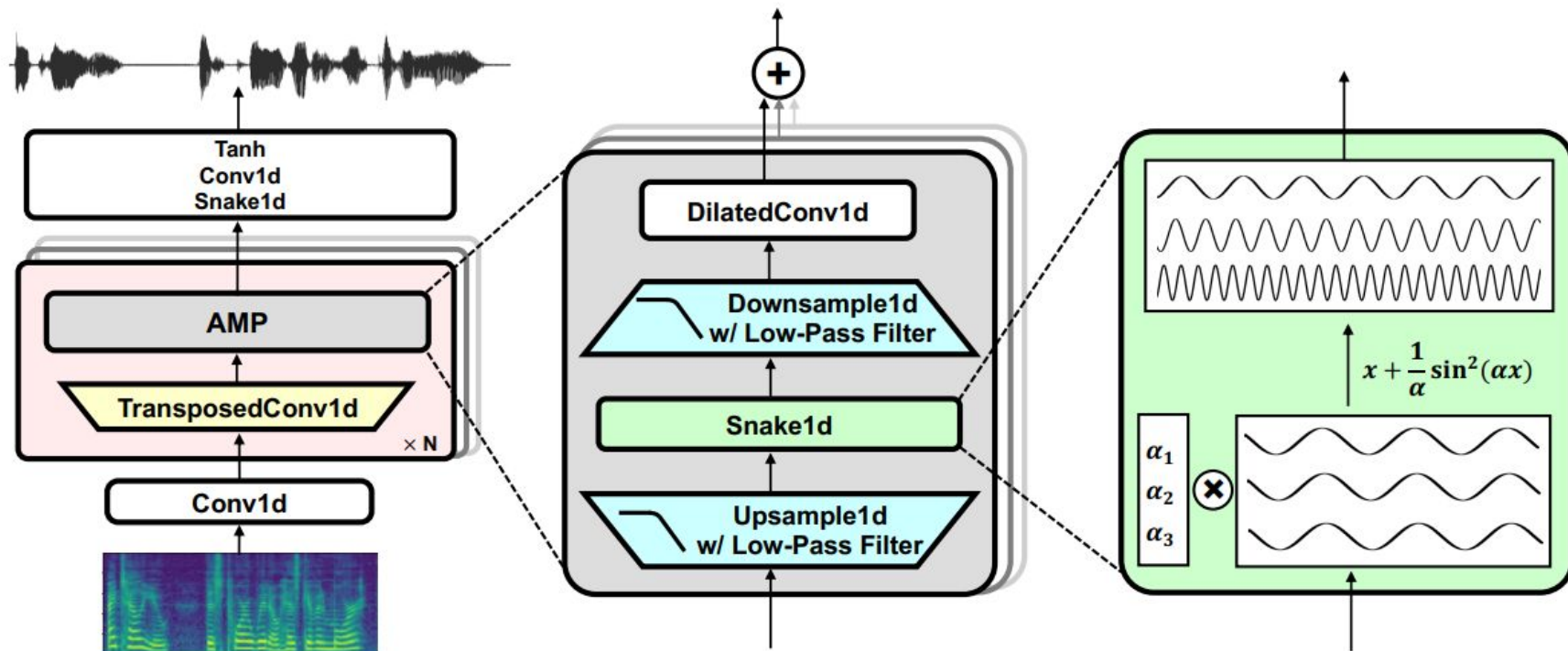


# BigVGAN



NUS  
National University  
of Singapore

School of  
Computing



# Pseudo Normal Differential Methods



NUS  
National University  
of Singapore

School of  
Computing

Noise ratio

Input

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta}(x_t, t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \epsilon_{\theta}(x_t, t) + \sigma_t \epsilon_t.$$
$$\frac{dx}{dt} = -\bar{\alpha}'(t) \left( \frac{x(t)}{2\bar{\alpha}(t)} - \frac{\epsilon_{\theta}(x(t), t)}{2\bar{\alpha}(t)\sqrt{1 - \bar{\alpha}(t)}} \right).$$

Differential Equation



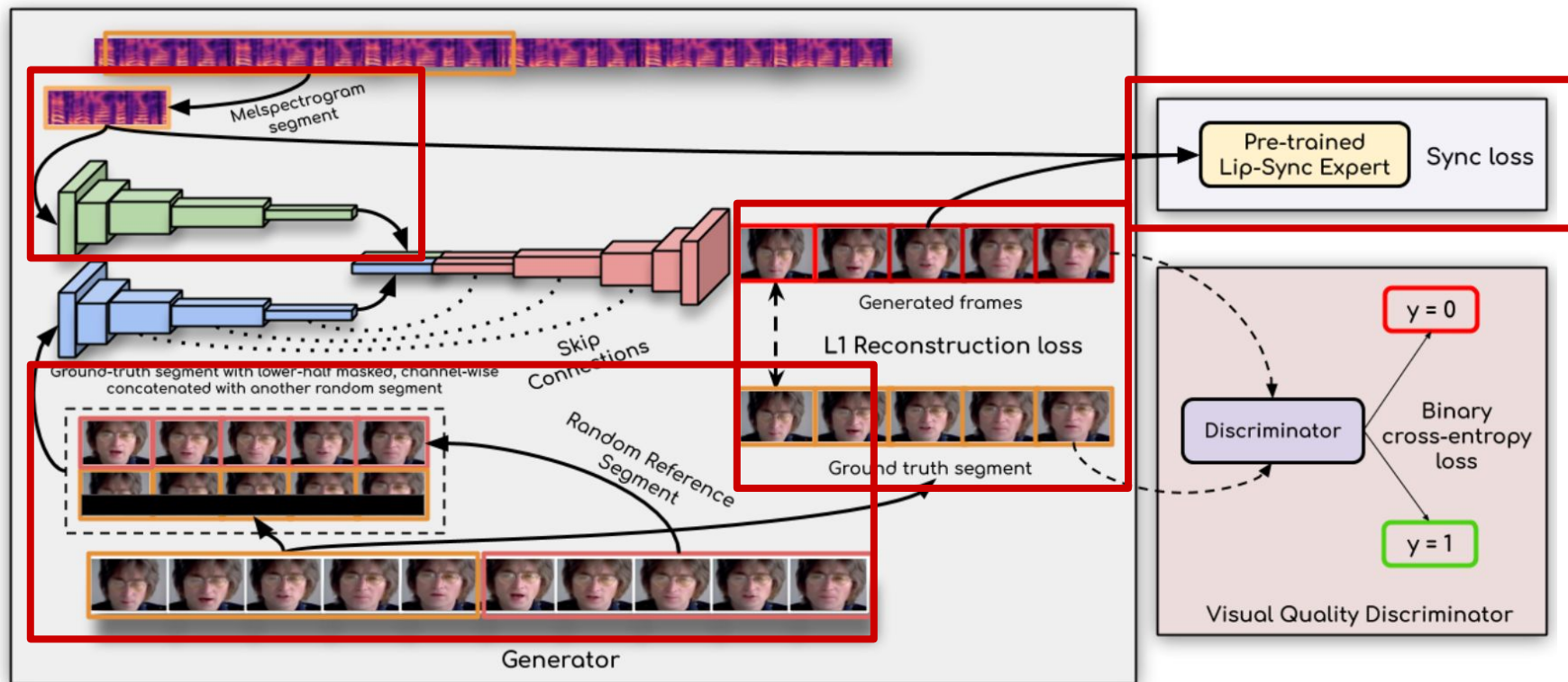
# Face Synthesis

# Wav2Lip



NUS  
National University  
of Singapore

School of  
Computing

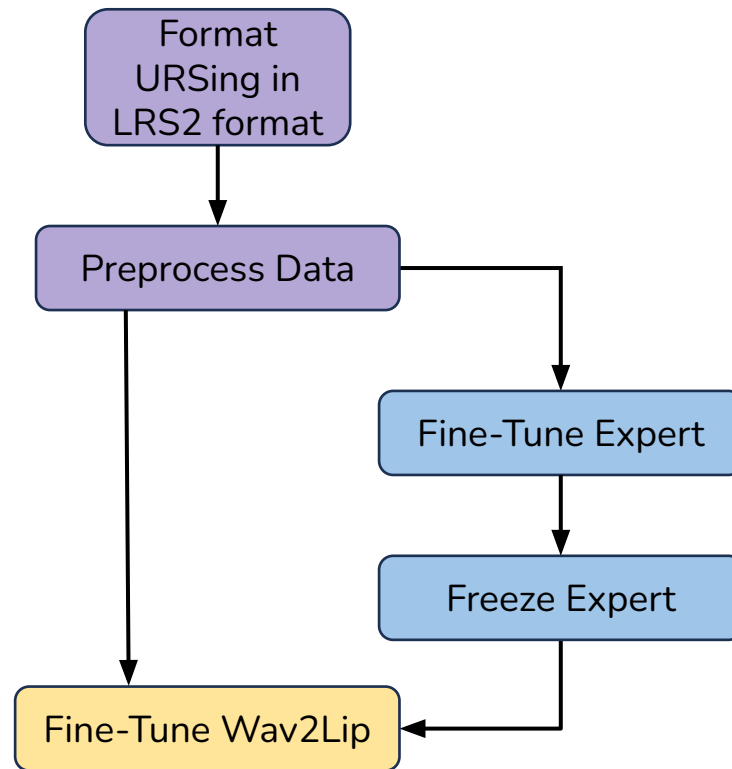




# Fine-Tuned Wav2Lip



- URSing (Audio-Visual Solo Singing)
  - 10 audio-video songs (total length ~ 3mins each)
  - $\Rightarrow$  60 clips (30s each)
- Freeze all but last layer of Expert Discriminator.
- Fine-tune Wav2Lip
  - Freeze Expert Discriminator
  - Freeze all but last layer of Wav2Lip
    - Encoders + Visual Quality Discriminator



## Goals:

- Improve Wav2Lip by introducing demographic information to the pipeline such as age, ethnicity and gender

## Model:

- Nine layer neural network to achieve near human levels of accuracy
- Capable of Face Detection, Analysis and Verification



# DeepFace Examples



NUS  
National University  
of Singapore

School of  
Computing



{'age': 24,  
'dominant\_gender': 'Man',  
'dominant\_race': 'asian'}



{'age': 35,  
'dominant\_gender': 'Woman',  
'dominant\_race': 'white'}



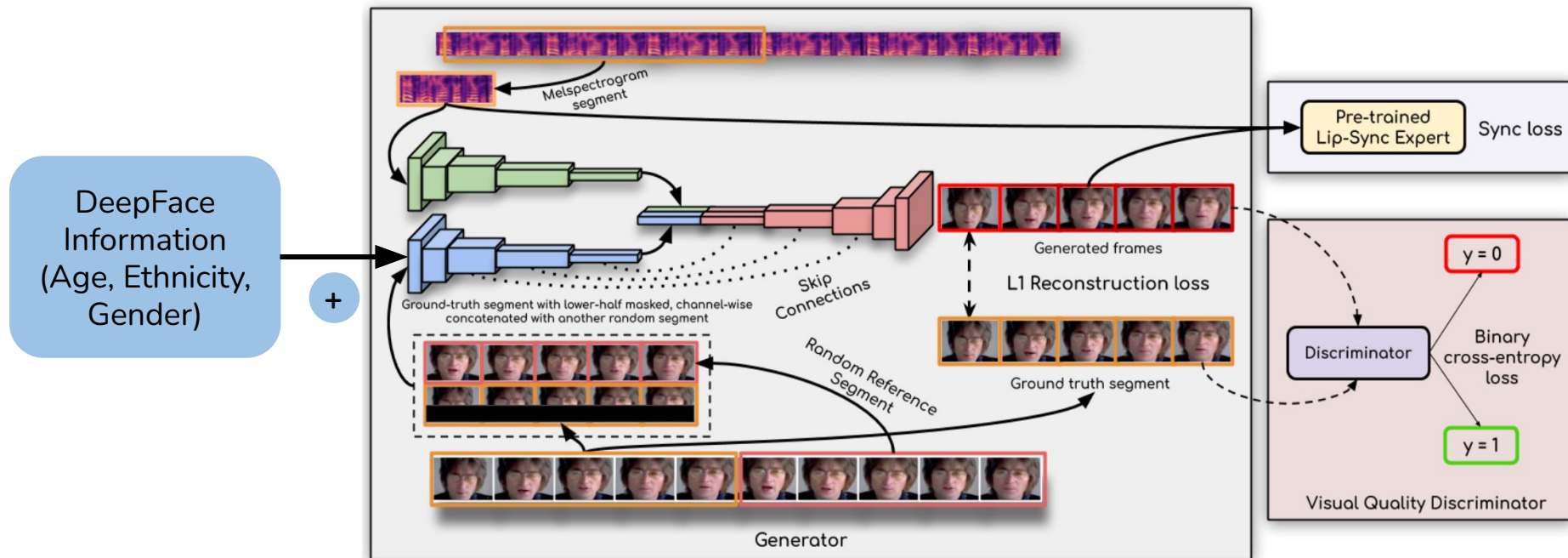
{'age': 28,  
'dominant\_gender': 'Woman',  
'dominant\_race': 'black'}

# DeepFace Augmentation



NUS  
National University  
of Singapore

School of  
Computing





# Results



# Results - DiffSinger

---



NUS  
National University  
of Singapore

School of  
Computing

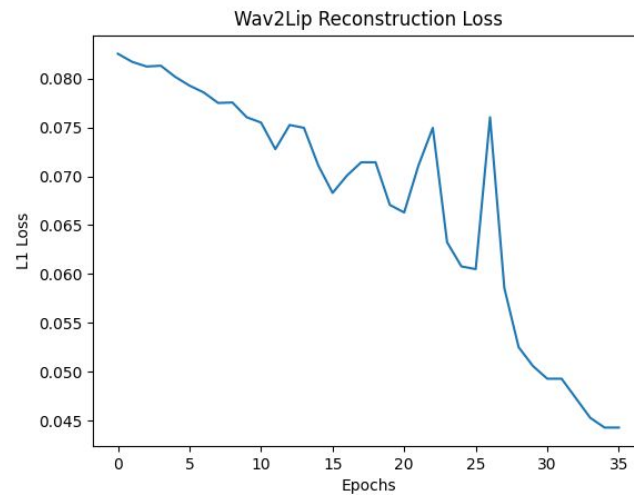
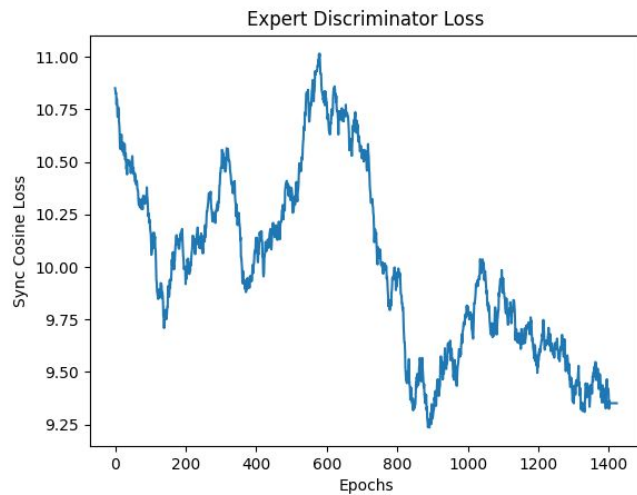
- Signal to Noise Ratio for Vocoder
  - HifiGAN: -0.00022 
  - BigVGAN: -0.00044 
- PNDM inference speedup
  - 40x in 25 steps vs 1000 steps

# Results - Wav2Lip Fine-tuning



NUS  
National University  
of Singapore

School of  
Computing



Trained on **three RTX3090**



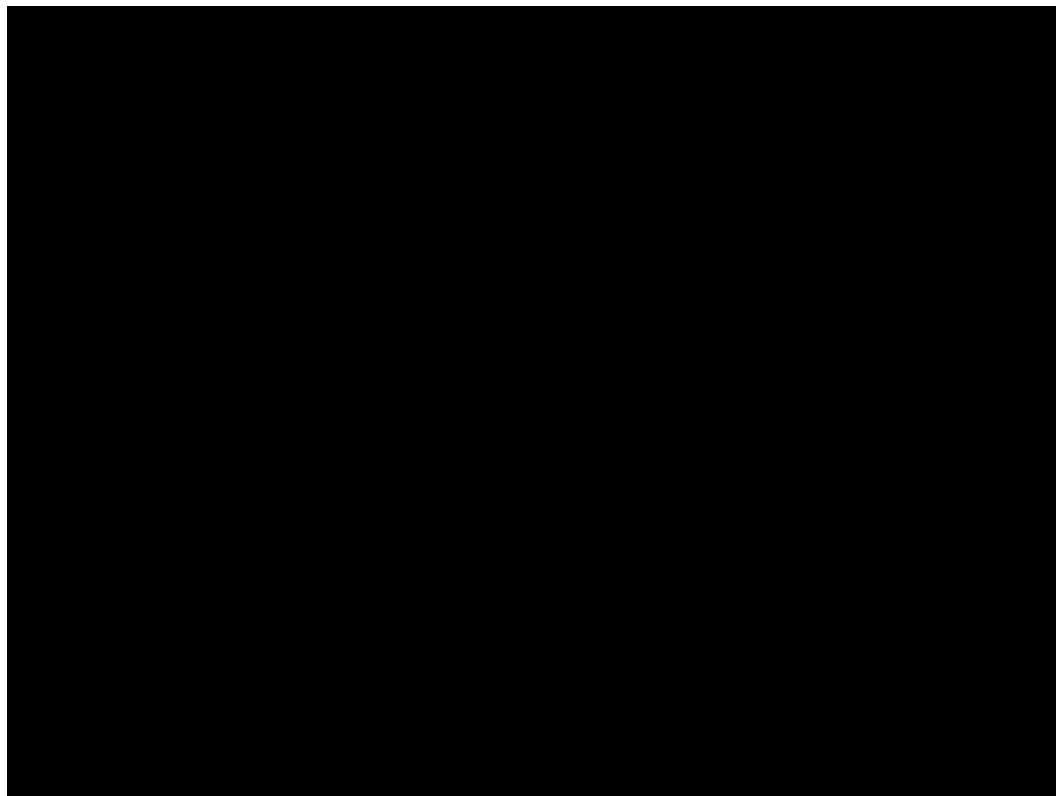
# Results - Wav2Lip DeepFace Augmentation

---



**NUS**  
National University  
of Singapore

School of  
Computing







# Demo

- [1] K R Prajwal, Rudrabha Mukhopadhyay, Vinay Namboodiri, C V Jawahar. (2020). **A Lip Sync Expert Is All You Need for Speech to Lip Generation In The Wild.**
- [2] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, Fei Wang. (2022). **SadTalker: Learning Realistic 3D Motion Coefficients for Stylized Audio-Driven Single Image Talking Face Animation**
- [3] Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, & Zhou Zhao. (2022). **DiffSinger: Singing Voice Synthesis via Shallow Diffusion Mechanism.**
- [4] Luping Liu, Yi Ren, Zhijie Lin, & Zhou Zhao. (2022). **Pseudo Numerical Methods for Diffusion Models on Manifolds.**
- [5] Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catanzaro, & Sungroh Yoon. (2023). **BigVGAN: A Universal Neural Vocoder with Large-Scale Training.**



# Thank You!



# Future Work

# Possible Improvements

---



**NUS**  
National University  
of Singapore

School of  
Computing

- Improving the DeepFace augmentation to seamlessly integrate the information into the input face data, allowing us to further improve the system and remove the artifacts.