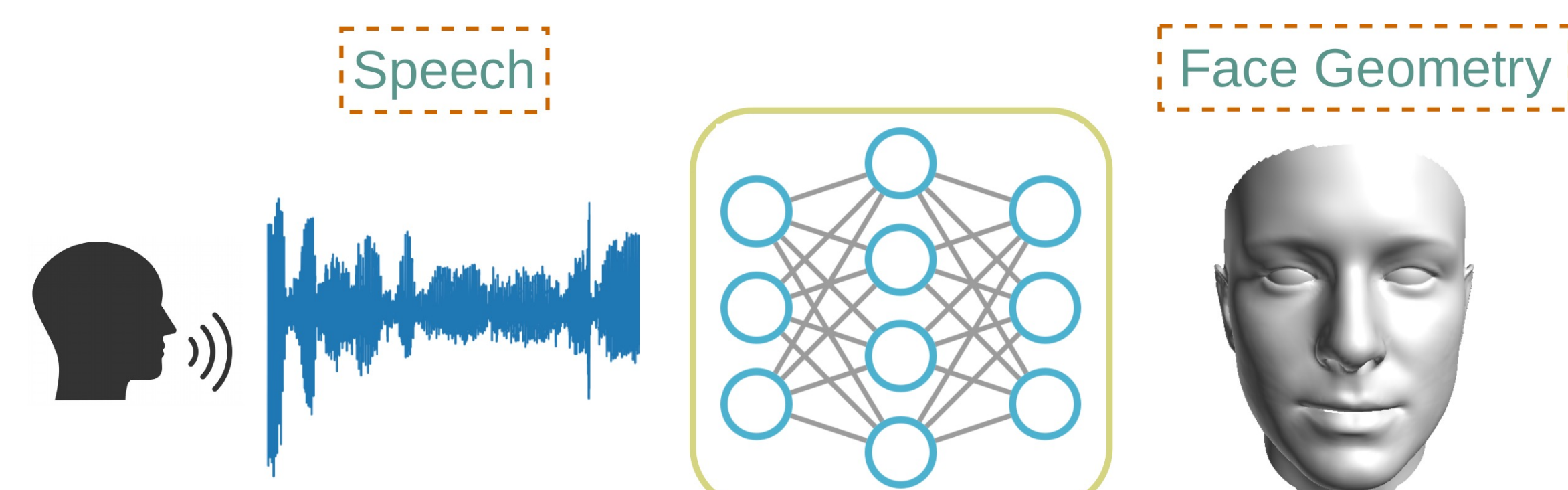
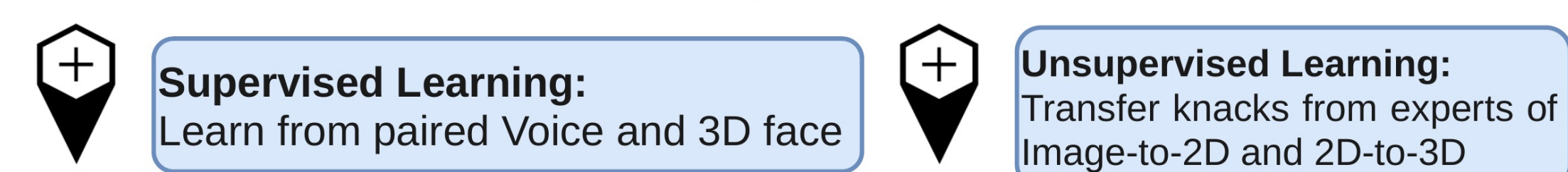


## Motivation

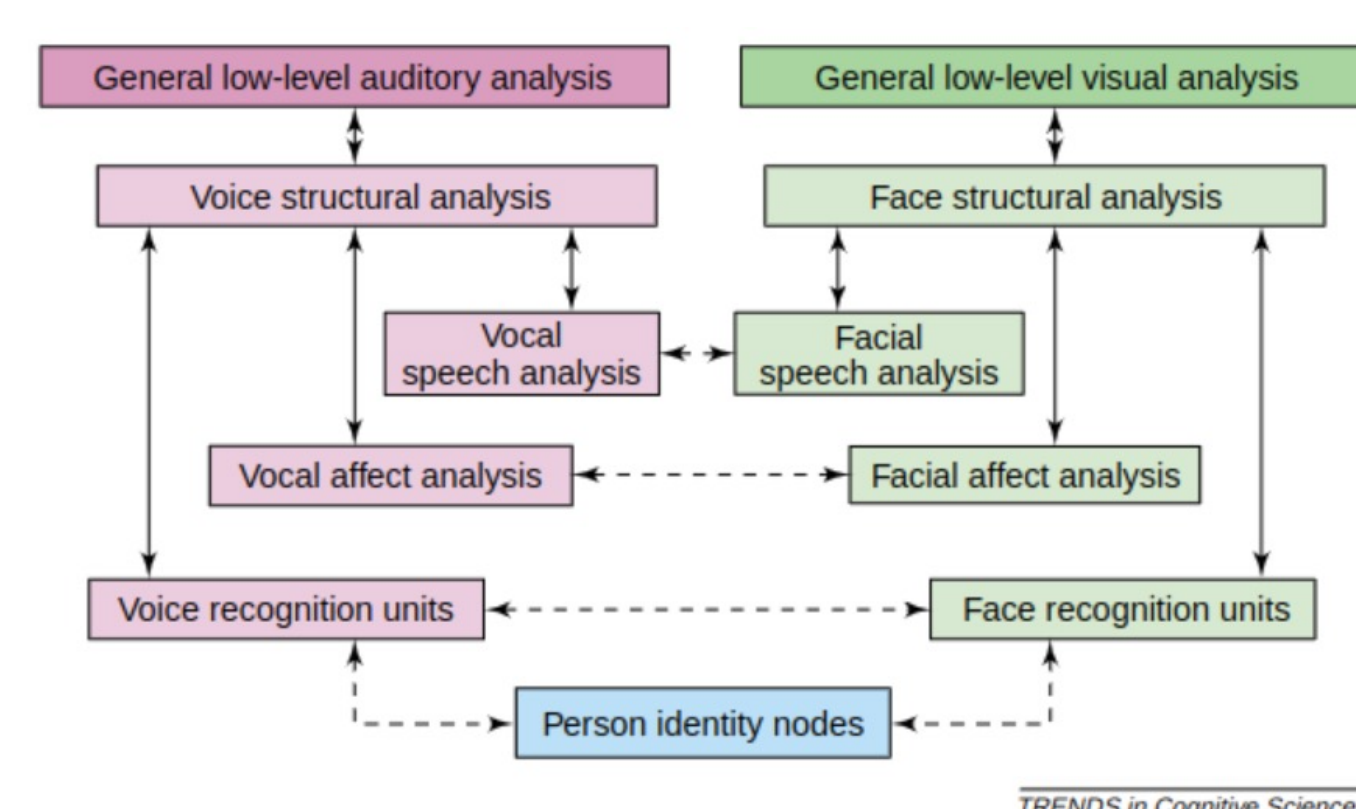
Humans can imagine what a speaker looks like without knowing this person or seeing his face.



A possible task?



Investigate the correlation between the two modality.



Cognitive science finding

[Trends in Cognitive Sciences, 2004, 2020]

## Prior Methods

They use encoder-decoder structure or GAN image synthesis to work on 2D representations.

Face image synthesis:



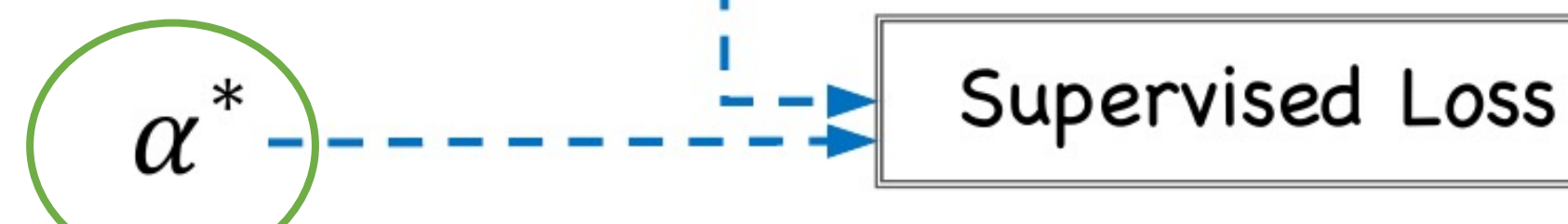
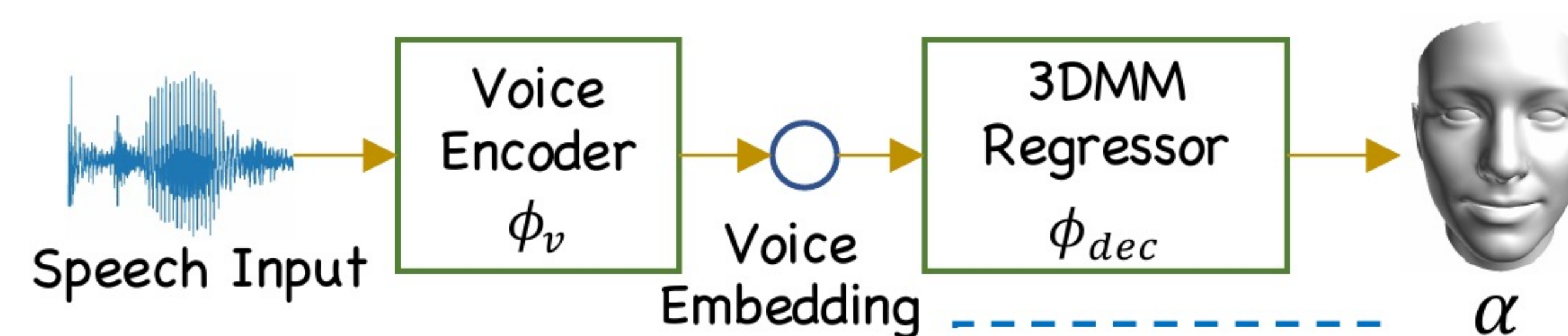
- Irrelevant factors: hair, beards, background, accessories
- Controversial Factor: Race and Ethnicity
- Hard to quantify the reconstruction preciseness

## Our study focus: 3D geometry



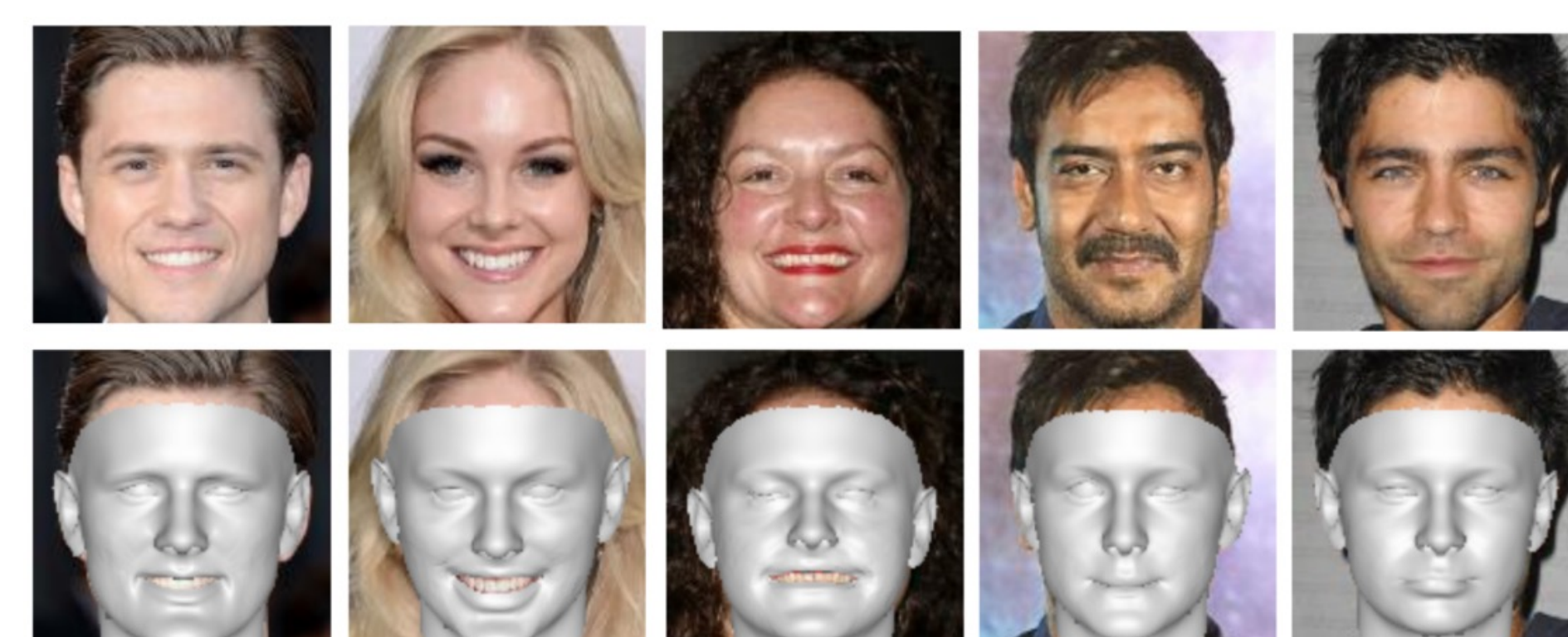
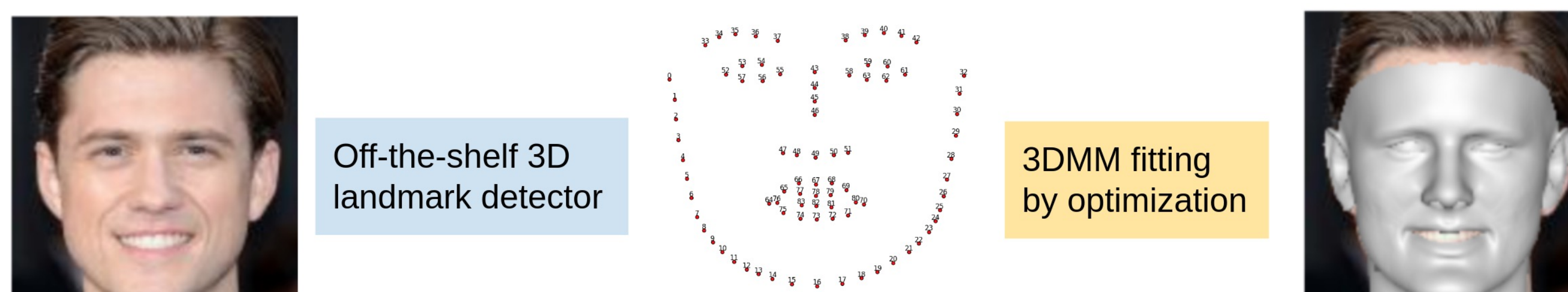
Mesh representation,  
using BFM Face

## Supervised Learning Model



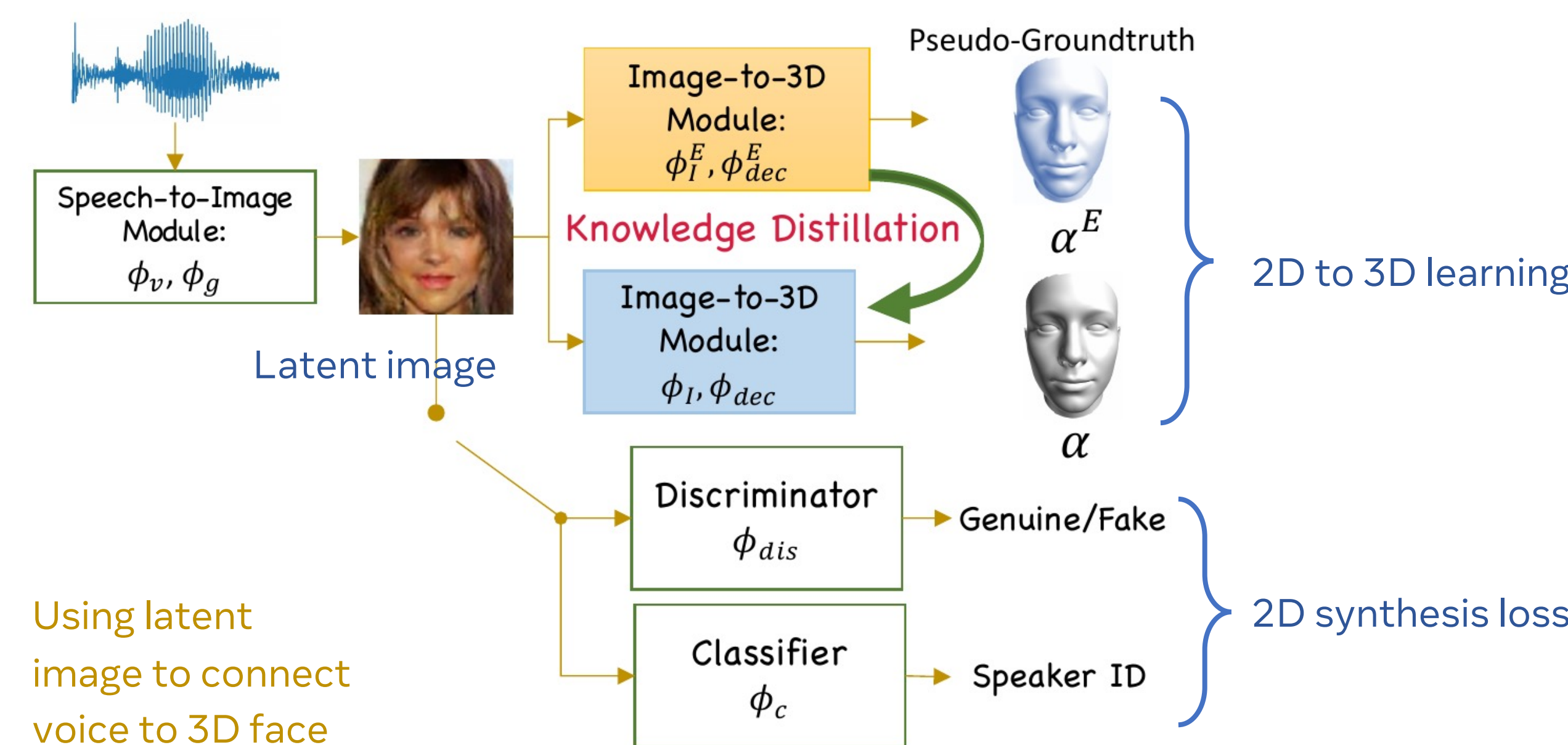
3DMM parameters: controlling face deformation

How to get paired 3D face and voice data?



Voxceleb-3D Dataset: Paired voice and 3DMM parameter datasets

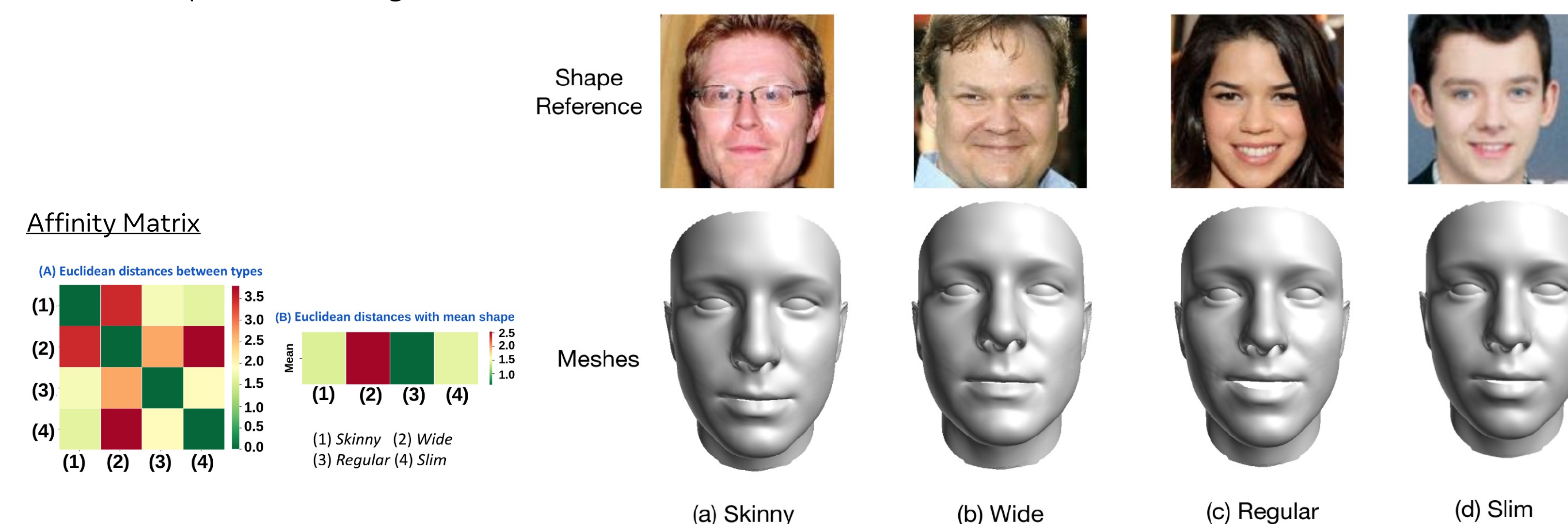
## Unsupervised Learning Model



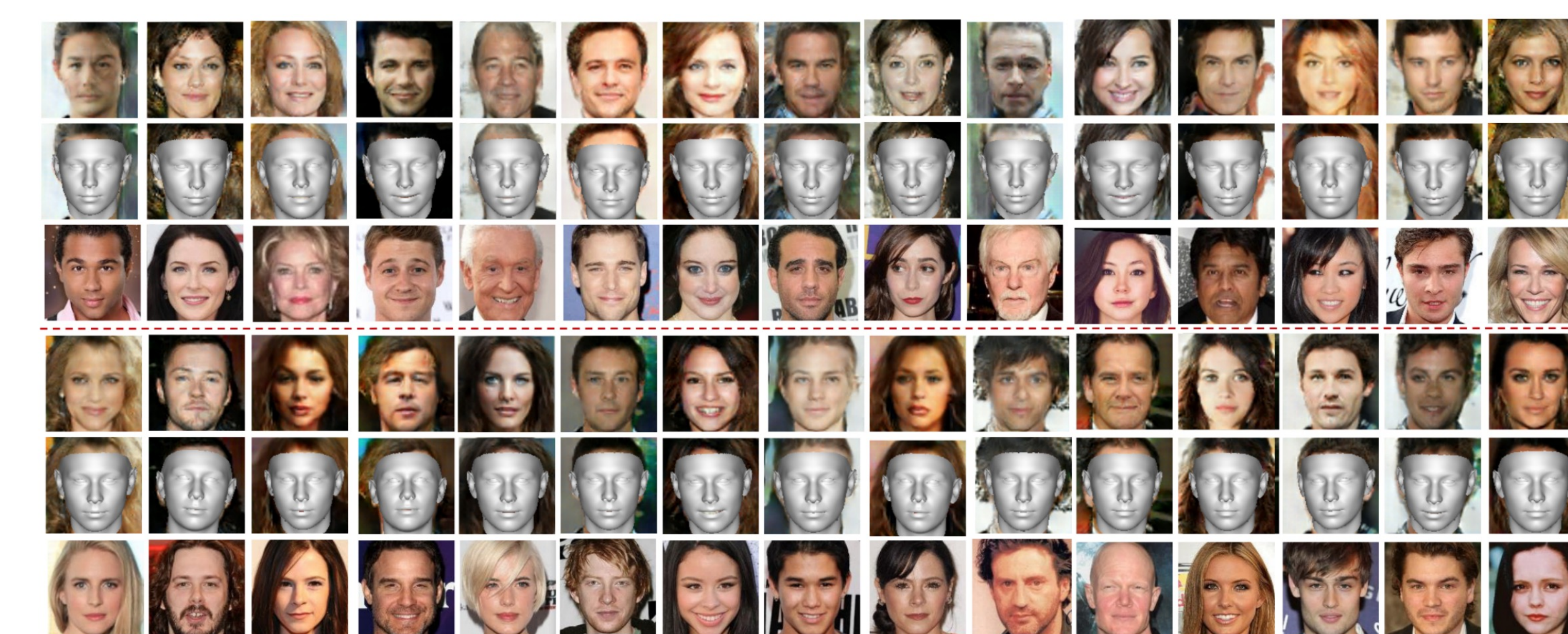
Using latent image to connect voice to 3D face

Q1: Is it feasible to predict visually reasonable face meshes from voice?

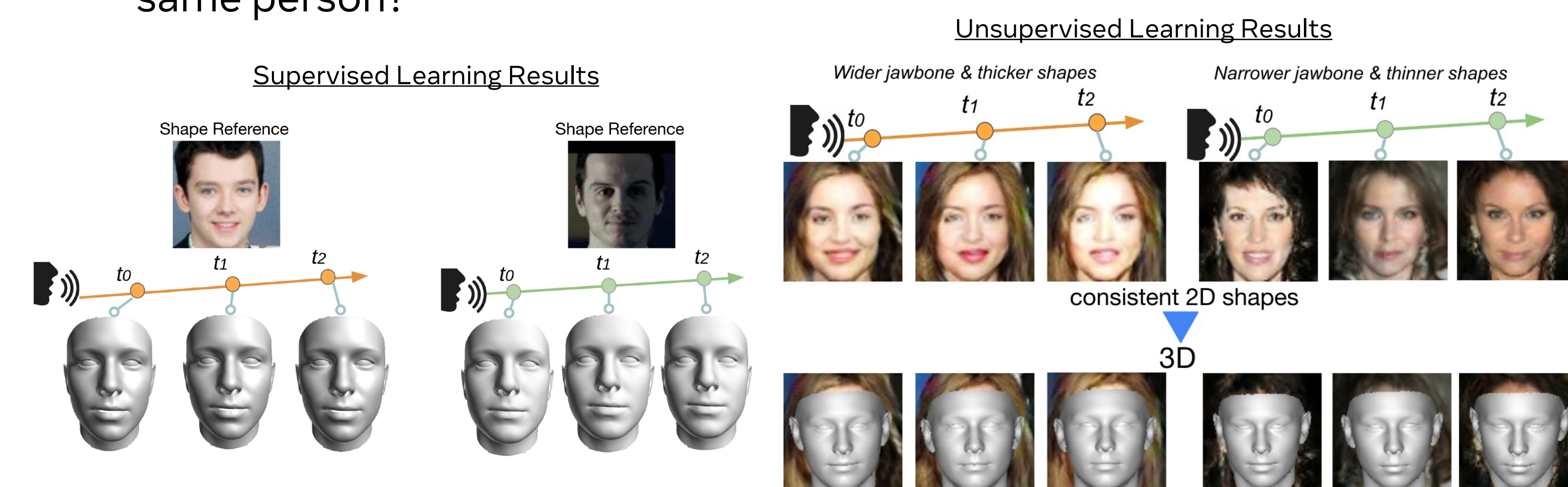
Supervised Learning Results



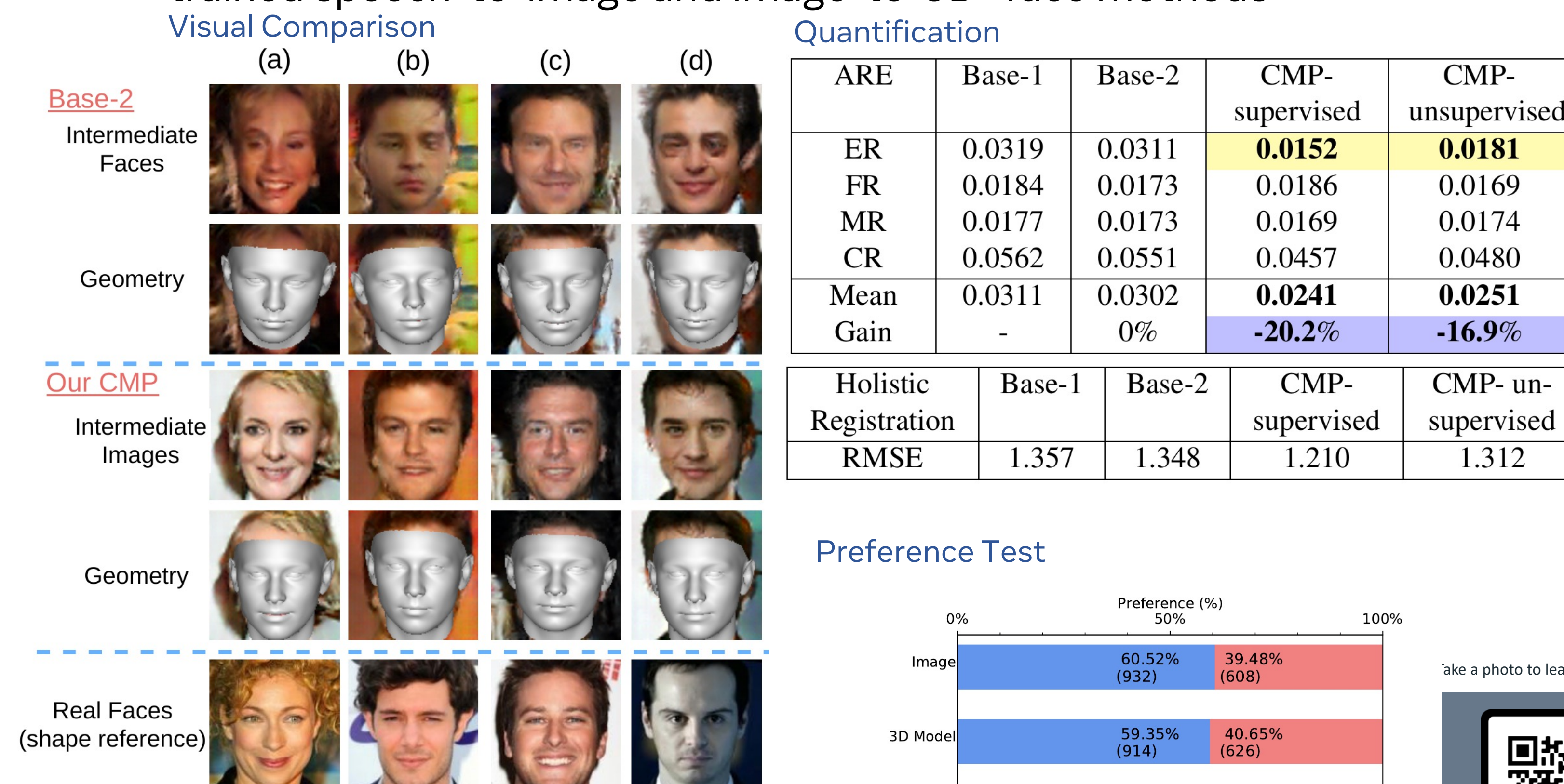
Unsupervised Learning Results



Q2: How stable is the mesh prediction from different utterances of the same person?



Q3: Comparison with face meshes produced by cascading separately trained speech-to-image and image-to-3D- face methods



Preference Test

