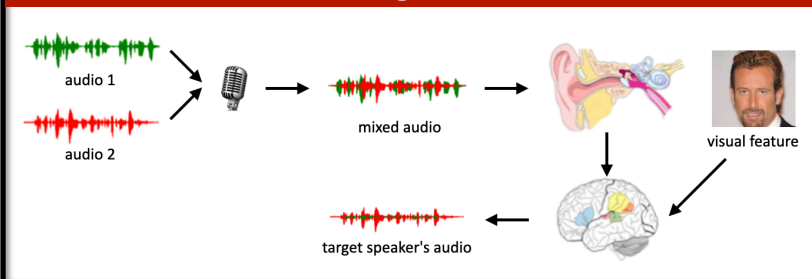


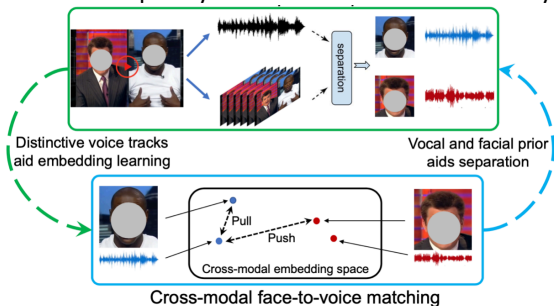


## Background

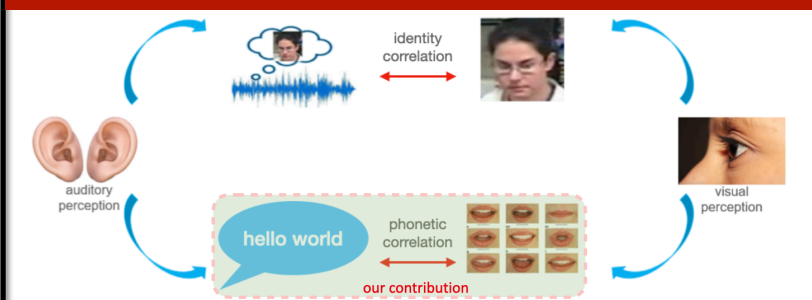


## Previous Work

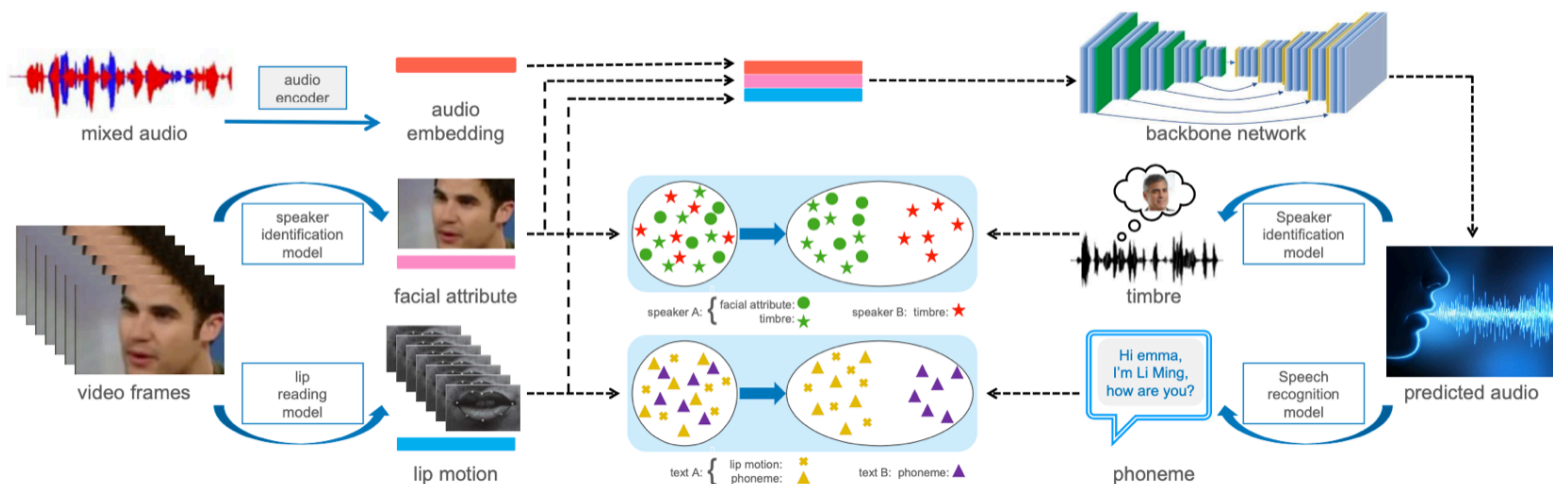
VisualVoice<sup>[1]</sup>: explicitly model the audio-visual identity correlation



## Motivation



## The Proposed Framework



Contrastive learning:

$$\mathcal{L}_1 = \max\{d(\mathbf{i}_{A_1}^a, \mathbf{i}_{A_2}^v) - d(\mathbf{i}_{A_1}^a, \mathbf{i}_B^v) + m, 0\}$$

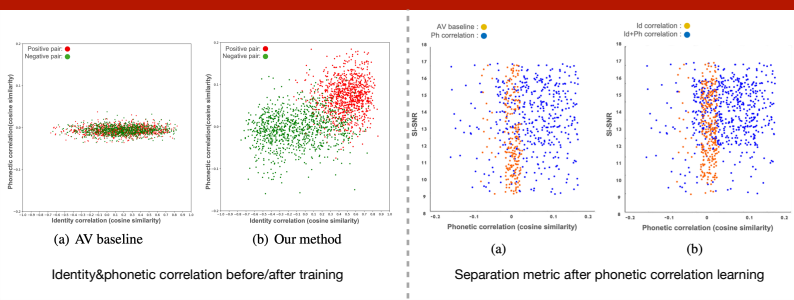
$$\mathcal{L}_2 = \max\{d(\mathbf{p}_A^a, \mathbf{p}_A^v) - d(\mathbf{p}_A^a, \mathbf{p}_B^v) + m, 0\}$$

Adversarial training:

$$\mathcal{L}_G = \min_G \mathbb{E}_{\mathbf{x} \sim i^v} \log(D(\mathbf{x})) + \mathbb{E}_{\mathbf{x} \sim i^a} \log(1 - D(\mathbf{x}))$$

$$\mathcal{L}_D = \max_D \mathbb{E}_{\mathbf{x} \sim i^v} \log(D(\mathbf{x})) + \mathbb{E}_{\mathbf{x} \sim i^a} \log(1 - D(\mathbf{x}))$$

## Visualization



## Experiments



	SDR	PESQ	STOI
[2](AV Baseline)	8.46	2.27	0.843
[2](CMC loss)	8.85	2.39	0.854
Ours(AV baseline)	9.392	2.536	0.851
Ours(triplet)	<b>9.623</b>	<b>2.545</b>	<b>0.855</b>
Ours(adversarial)	<b>9.982</b>	<b>2.584</b>	<b>0.861</b>

	SDR	SIR	SAR	PESQ	STOI	SI-SNR
[1](Reported)	10.2	17.2	11.3	2.83	0.87	-
[1](Released)	7.023	13.708	9.546	2.569	0.792	6.471
[1](Our impl.)	7.692	14.347	10.195	2.579	0.791	7.467
Ours(triplet)	<b>8.178</b>	<b>14.692</b>	<b>10.38</b>	<b>2.6</b>	<b>0.793</b>	<b>7.676</b>
Ours(adversarial)	<b>8.949</b>	<b>16.012</b>	<b>10.79</b>	<b>2.687</b>	<b>0.811</b>	<b>8.477</b>

[1] R. Gao and K. Grauman. "VisualVoice: Audio-Visual Speech Separation with Cross-Modal Consistency". In CVPR 2021.

[2] N. Makishima, M. Ihori, A. Takashima, T. Tanaka, S. Orihashi, and R. Masumura, "Audio-visual speech separation using cross-modal correspondence loss". In ICASSP 2021.

[3] T. Afouras, J. S. Chung, A. Zisserman LRS3-TED: a large-scale dataset for visual speech recognition arXiv preprint arXiv:1809.00496

[4] J. S. Chung\*, A. Nagrani\*, A. Zisserman VoxCeleb2: Deep Speaker Recognition INTERSPEECH, 2018.