

Audio-driven Talking Face Generation: A Brief Survey

Ling-Hao CHEN chenlinghao@idea.edu.cn Outline



Problem Definition

- Demo
- □ Road Map
 - End-to-end, GAN-based
 - 3D model based
 - (2D/3D) landmark based
- **D** Summary
- Challenges



Driven source: audio



Identity Reference



Audio Source







Driven source: audio & video

Input:





Video Source

Identity Reference

Audio Source







D Preliminaries



Driven Source S_k



□ ATVGNet





[1]: Chen, Lele, et al. "Hierarchical cross-modal talking face generation with dynamic pixel-wise loss." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019.



D FACIAL



[1]: Zhang, Chenxu, et al. "Facial: Synthesizing dynamic talking face with implicit attribute learning." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.



D PCAVS



Identity Reference

Output

Identity Reference

Audio Source

[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation." Computer Vision and Pattern Recognition (2021).



□ PCAVS-Eye



Identity Reference

Output

Identity Reference

Audio Source

[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation." Computer Vision and Pattern Recognition (2021).



D End-to-end, GAN-based

- **D** 3D model based
- □ (2D/3D) landmark based



D End-to-end, GAN-based

- **D** 3D model based
- □ (2D/3D) landmark based



Adversarial Talking Face Generation

PROBLEM: How to decouple the voice-related actions in a video from the identity information of the face?



[1]: Zhou, Hang, et al. "Talking face generation by adversarially disentangled audio-visual representation." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. No. 01. 2019.



Adversarial Talking Face Generation

PROBLEM: How to decouple the voice-related actions in a video from the identity information of the face?



 [1]: Zhou, Hang, et al. "Talking face generation by adversarially disentangled audio-visual representation." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. No. 01.

 2019.

 粤港澳大湾区数字经济研究院

 International Digital Economy Academy



Adversarial Talking Face Generation

PROBLEM: How to decouple the voice-related actions in a video from the identity information of the face?

Training adversarially between P-id Encoder and W-id classifier.



[1]: Zhou, Hang, et al. "Talking face generation by adversarially disentangled audio-visual representation." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. No. 01. 2019. 粤港澳大湾区数字经济研究院 International Digital Economy Academy www.idea.edu.cn



➤ Wav2Lip

PROBLEM: How to do accurate lip alignment?



[1]: Prajwal, K. R., et al. "A lip sync expert is all you need for speech to lip generation in the wild." Proceedings of the 28th ACM International Conference on Multimedia. 2020.



➤ Wav2Lip

PROBLEM: How to do accurate lip alignment?

- ✓ Introduce the SyncNet as Lip-Sync Expert.
- $\checkmark\,$ Training with audio and video jointly.
- ✓ Data Augmentation.
 - Image occlusion.
 - Random Sampling.



[1]: Prajwal, K. R., et al. "A lip sync expert is all you need for speech to lip generation in the wild." Proceedings of the 28th ACM International Conference on Multimedia. 2020.



> PCAVS

PROBLEM: How to model the pose of the head?

Pose Controllable:



[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).



> PCAVS

Freamwork Overview:



[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).

idea

> PCAVS

Identity and Non-Identity Encoder (pre-train jointly):

Data augmentation:





✓ Perspective Transformation

- ✓ Color Transfer
- ✓ Centered Random Crop

[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).

idea

> PCAVS

Identity and Non-Identity Encoder (pre-train jointly):

Disentangled representation learning:





 d_p : low-dimension is better

Learning pose

[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).
 [2]: Burkov, Egor, et al. "Neural head reenactment with latent pose descriptors." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
 [3]: Burkov, Egor, et al. "Neural head reenactment with latent pose descriptors." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
 [3]: Burkov, Egor, et al. "Neural head reenactment with latent pose descriptors." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.



> PCAVS



[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).
 [2]: Egor Burkov et al. "Neural Head Reenactment with Latent Pose Descriptors" Computer Vision and Pattern Recognition (2020).
 [3]: Bay System Computer Vision and Pattern Recognition (2020).

nternational Digital Economy Academy

 $\mathbf{F}_{c(k)}^{v-}$

Speech Content



> PCAVS

Final stage:

$$\mathcal{L}_{\text{GAN}} = \underset{G}{\text{min}} \underset{D}{\text{min}} \sum_{n=1}^{N_D} \left(\mathbb{E}_{I_{(k)}} [\log D_n(I_{(k)})] + \mathbb{E}_{f_{cat(k)}} [\log(1 - D_n(G(f_{cat(k)})]), \\ \mathcal{L}_{L_1} = \sum_{i=1}^{N_D} \|D_n(I_{(k)}) - D_n(G(f_{cat(k)}))\|_1, \\ \mathcal{L}_{\text{vgg}} = \sum_{i=1}^{N_D} \|VGG_n(I_{(k)}) - VGG_n(G(f_{cat(k)}))\|_1. \\ \end{pmatrix}$$

$$\mathcal{L}_{total} = \mathcal{L}_{\text{GAN}} + \lambda_1 \mathcal{L}_{L_1} + \lambda_v \mathcal{L}_{\text{vgg}} + \lambda_c \mathcal{L}_c + \lambda_i \mathcal{L}_i$$

[1]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).

粤港澳大湾区数字经济研究院 International Digital Economy Academy



> AD-NeRF

Target: Generate the talking face of a specific person in any pose.



[1]: Yudong Guo et al. "AD-NeRF: Audio Driven Neural Radiance Fields for Talking Head Synthesis." arXiv: Computer Vision and Pattern Recognition (2021): n. pag.



> AD-NeRF

Target: Generate the talking face of a specific person in any pose.



[1]: Yudong Guo et al. "AD-NeRF: Audio Driven Neural Radiance Fields for Talking Head Synthesis." arXiv: Computer Vision and Pattern Recognition (2021): n. pag.

粵港澳大湾区数字经济研究院 International Digital Economy Academy



Speech2Talking-Face

Contribution: The discrimination of identity features is strengthened; for the first time, face reconstruction and speech-driven face speech generation are completed in the same framework.



 [1]: Yasheng Sun et al. "Speech2Talking-Face: Inferring and Driving a Face with Synchronized Audio-Visual Representation" International Joint Conference on Artificial Intelligence (2021).

 粤港澳大湾区数字经济研究院 International Digital Economy Academy
 25



Speech2Talking-Face

Identity space:



[1]: Yasheng Sun et al. "Speech2Talking-Face: Inferring and Driving a Face with Synchronized Audio-Visual Representation" International Joint Conference on Artificial Intelligence (2021). 粤港澳大湾区数字经济研究院



Speech2Talking-Face

Identity-irrelevant space:

Final stage:

$$\mathcal{L}_{all} = \mathcal{L}_{adv} + \lambda_r \mathcal{L}_{rec} + \lambda_c \mathcal{L}_c + \lambda_{sync}^{id} \mathcal{L}_{sync}^{id} + \lambda_{sync}^s \mathcal{L}_{sync}^s$$

Adversarial loss, Reconstruction loss, Several contrastive loss



[1]: Yasheng Sun et al. "Speech2Talking-Face: Inferring and Driving a Face with Synchronized Audio-Visual Representation" International Joint Conference on Artificial Intelligence (2021). 粤港澳大湾区数字经济研究院 27 International Digital Economy Academy



D End-to-end, GAN-based

- **D** 3D model based
- □ (2D/3D) landmark based



D End-to-end, GAN-based

- **D** 3D model based
- □ (2D/3D) landmark based



> ADEVP

PROBLEM: How to decouple the emotion and the speech content?

Q: How to length-align speech with uneven length?A: DTW Algorithm.

We will train E_c , E_e , D in this stage.

Loss:

$$L_{cross} = \|D(E_c(x_{i,m}), E_e(x_{j,n})) - x_{i,n}\|_2 + \|D(E_c(x_{j,n}), E_e(x_{i,m})) - x_{j,m}\|_2$$
$$L_{self} = \|D(E_c(x_{i,m}), E_e(x_{i,m})) - x_{i,m}\|_2 + \|D(E_c(x_{j,n}), E_e(x_{j,n})) - x_{j,n}\|_2.$$



Cross-reconstruction for disentanglemen

[1]: Ji, Xinya, et al. "Audio-driven emotional video portraits." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.



> ADEVP

PROBLEM: How to decouple the emotion and the speech content?

Loss:

$$\begin{split} L_{cross} &= \|D(E_{c}(x_{i,m}), E_{e}(x_{j,n})) - x_{i,n}\|_{2} \\ &+ \|D(E_{c}(x_{j,n}), E_{e}(x_{i,m})) - x_{j,m}\|_{2} \\ L_{self} &= \|D(E_{c}(x_{i,m}), E_{e}(x_{i,m})) - x_{i,m}\|_{2} \\ &+ \|D(E_{c}(x_{j,n}), E_{e}(x_{j,n})) - x_{j,n}\|_{2}. \\ L_{cla} &= -\sum_{k=1}^{N} (p_{k} * \log q_{k}) \\ L_{con} &= \sum_{k=i,j} \|E_{c}(x_{k,m}) - E_{c}(x_{k,n})\|_{1} \\ L_{dis} &= L_{cross} + L_{self} + \lambda_{cla}L_{cla} + \lambda_{con}L_{con} \end{split}$$



Cross-reconstruction for disentanglemen

[1]: Ji, Xinya, et al. "Audio-driven emotional video portraits." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.

粤港澳大湾区数字经济研究院 International Digital Economy Academy



> ADEVP

PROBLEM: How to decouple the emotion and the speech content?



Overview of Emotional Video Portrait algorithm

[1]: Ji, Xinya, et al. "Audio-driven emotional video portraits." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.



➢ FACIAL

(1) How do explicit and implicit attributes potentially affect each other? (2) How to model implicit properties?



[1]: Chenxu Zhang et al. "FACIAL: Synthesizing Dynamic Talking Face with Implicit Attribute Learning" Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021. 澳大湾区数字经济研究院



➢ FACIAL

The generated videos have not only synchronized lip movements, but also natural head movements and eye blink information.



[1]: Chenxu Zhang et al. "FACIAL: Synthesizing Dynamic Talking Face with Implicit Attribute Learning" Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021. 澳大湾区数字经济研究院 34



D End-to-end, GAN-based

- **D** 3D model based
- □ (2D/3D) landmark based



D End-to-end, GAN-based

- **D** 3D model based
- □ (2D/3D) landmark based



ATVG-Net

The generated videos have not only synchronized lip movements, but also natural head movements and eye blink information.



[1]: Chen, Lele, et al. "Hierarchical cross-modal talking face generation with dynamic pixel-wise loss." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019. International Digital Economy Academy



MakeltTalk

Alignment of lip movements, facial expressions, head movements.



[1]: Zhou, Yang, et al. "MakeltTalk: speaker-aware talking-head animation." ACM Transactions on Graphics (TOG) 39.6 (2020): 1-15.



MakeltTalk

Alignment of lip movements, facial expressions, head movements.

- **\Box** Learning active dynamics: Δp .
- Get dynamic landmarks of driven object.
- The GLSL-based method is used to stretch the landmarks for cartoon images to achieve the purpose; the images are generated using the Image2Image architecture.



[1]: Zhou, Yang, et al. "MakeltTalk: speaker-aware talking-head animation." ACM Transactions on Graphics (TOG) 39.6 (2020): 1-15.

<u>≜a</u> Methods	🗹 GAN-based	🗹 Landmark-based	☑ 3D-model-based	≡ Input	■ Journal/Conference
Adversarial TFG				Audio Referenced Video	AAAI-19
ATVG-Net				Audio Image	CVPR-19
Wav2Lip				Audio Referenced Video	ACM MM-20
MakeItTalk				Audio Image	ACM TOG-21
ADEVP				Audio Image	CVPR-21
PCAVS				Audio Referenced Video Image	CVPR-21
FICIAL				Audio Image	CVPR-21
Speech2Talking-Face				Audio Referenced Video Image	IJCAI-21
Imitating ADTF-Sync				Audio Referenced Video Image	ACM MM-21
Info-Bottleneck TFG				Audio Image	CIKM-21(short)
AD-NeRF				Audio Image	ICCV-21



- □ Jitter between video frames.
- How to make posture and expression controllable and editable?
- How to better decouple the information between face ID and lip movements, blinks, and facial expression?
- Robustness of the model and the sensitivity of noise.
- Generated faces are not enough yet.

[1]: Zhou, Yang, et al. "MakeltTalk: speaker-aware talking-head animation." ACM Transactions on Graphics (TOG) 39.6 (2020): 1-15.

[2]: Zhou, Hang, et al. "Talking face generation by adversarially disentangled audio-visual representation." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. No. 01. 2019.

[3]: Prajwal, K. R., et al. "A lip sync expert is all you need for speech to lip generation in the wild." Proceedings of the 28th ACM International Conference on Multimedia. 2020.

[4]: Hang Zhou et al. "Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation.." Computer Vision and Pattern Recognition (2021).

[5]: Burkov, Egor, et al. "Neural head reenactment with latent pose descriptors." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.

[6]: Yudong Guo et al. "AD-NeRF: Audio Driven Neural Radiance Fields for Talking Head Synthesis." arXiv: Computer Vision and Pattern Recognition (2021): n. pag.
 [7]: Yasheng Sun et al. "Speech2Talking-Face: Inferring and Driving a Face with Synchronized Audio-Visual Representation" International Joint Conference on Artificial Intelligence (2021).

[8]: Ji, Xinya, et al. "Audio-driven emotional video portraits." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.

[9]: Chenxu Zhang et al. "FACIAL: Synthesizing Dynamic Talking Face with Implicit Attribute Learning" Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.

[10]: Chen, Lele, et al. "Hierarchical cross-modal talking face generation with dynamic pixel-wise loss." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019.

[11]: Ben Mildenhall et al. "NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis" European Conference on Computer Vision (2020).

[12]: Wu, Haozhe, et al. "Imitating Arbitrary Talking Style for Realistic Audio-Driven Talking Face Synthesis." Proceedings of the 29th ACM International Conference on Multimedia. 2021.

[13]: Tang, Jie, et al. "Talking Face Generation Based on Information Bottleneck and Complementary Representations." Proceedings of the 30th ACM International Conference on Information & Knowledge Management. 2021.

[14]: Doukas, Michail Christos, Stefanos Zafeiriou, and Viktoriia Sharmanska. "HeadGAN: One-shot Neural Head Synthesis and Editing." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.

[15]: Huang, Po-Hsiang, Fu-En Yang, and Yu-Chiang Frank Wang. "Learning Identity-Invariant Motion Representations for Cross-ID Face Reenactment." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.

THANK YOU

Ling-Hao CHEN (陈 凌灏) chenlinghao@idea.edu.cn School of Computer Science and Technology Xidian University G-343, 266 Xinglong Section of Xifeng Road, Xi'an, ShaanXi, P.R. China.



www.idea.edu.cn