# DRVC: A FRAMEWORK OF ANY-TO-ANY VOICE CONVERSION WITH SELF-SUPERVISED LEARNING

Authors: Qiqi Wang, Xulong Zhang, Jianzong Wang, Ning Cheng, Jing Xiao
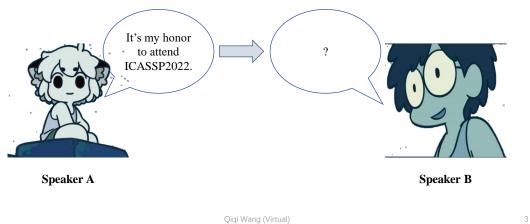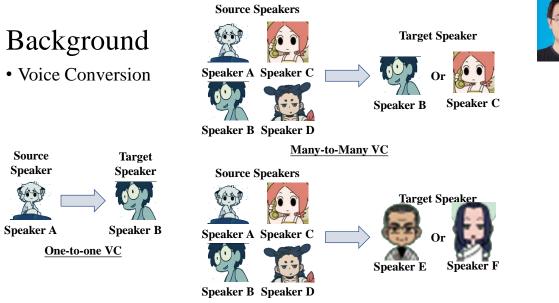
Speaker: Qiqi Wang (virtual)

10th April 2022

## Outline

- Background
- DRVC
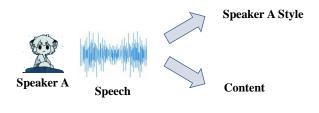- Experiments
- Conclusion

# Background

• Voice Conversion



**Speaker A**

**Speaker B**

# Background

• Voice Conversion

**Source Speakers**



**Target Speaker**

**Speaker A  Speaker C**

**Speaker B  Speaker D**

**Speaker B      Speaker C**

Or

**Many-to-Many VC**

**Source**
**Speaker**

**Target**
**Speaker**

**Speaker A**        **Speaker B**

**One-to-one VC**

**Source Speakers**

**Speaker A  Speaker C**

**Speaker B  Speaker D**

**Target Speaker**

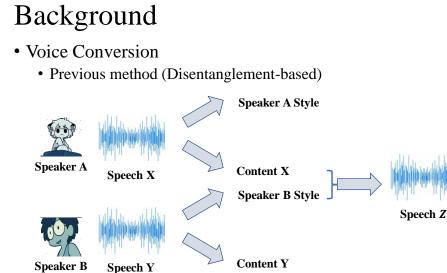**Speaker E      Speaker F**

Or

**Any-to-Any VC**

# Background

- Voice Conversion
  - Previous method (Disentanglement-based)

**Assumption:** *Speech information consists of speaker style and content information.*



Speaker A Style

Content

Speaker A    Speech

# Background

- Voice Conversion
  - Previous method (Disentanglement-based)



Speaker A Style

Speaker A    Speech X

Content X

Speaker B Style

Speech *Z*

Speaker B    Speech Y

Content Y

# Background

- Voice Conversion
  - Shortages



Fixed Size of Content information

Distanglement is incomplete

# DRVC

- Speech Distanglement
  - Two encoders
    - Speaker Style Encoder: $E_S$
    - Content Encoder: $E_{Con}$



$E_S$ → Speaker A Style

$E_{Con}$ → Content X

$$\{x_c, x_s\} = \{E_{Con}(X), E_S(X)\}$$

# DRVC

- Speech Distanglement
  - Generator $G$



$$\tilde{x} = G(y_c, x_s) = G(\{E_{Con}(Y), E_S(X)\})$$

# DRVC

- Two Stage Conversion
  - First Conversion
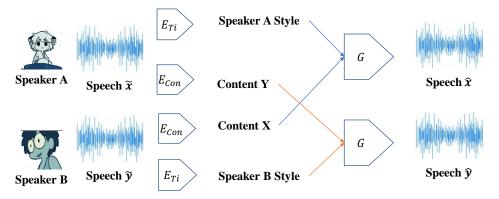
# DRVC

- Two Stage Conversion
  - Second Conversion

# DRVC

- Loss Function
  - Cycle Loss



$$\mathcal{L}_{cycle} = E_{x,y}[||\hat{x} - X|| + ||\hat{y} - Y||]$$

# DRVC

- Loss Function
  - Same Loss



$$\mathcal{L}_{same} = E[|\tilde{y}_c - x_c| + |\tilde{x}_c - y_c|] + E[|\tilde{x}_s - x_s| + |\tilde{y}_s - y_s|]$$

# DRVC

- Loss Function
  - Domain Loss



$$\mathcal{L}_{domain} = -\frac{1}{2}(\sum_i A(i)C(x_s) + \sum_i B(i)C(y_s))$$

# DRVC

- Loss Function
  - Adversarial Loss



| Speaker A | Speech X | Speech $\tilde{x}$ | Speech $\hat{x}$ | | Real Or Synthetic ? |
|---|---|---|---|---|---|

Discriminator

| Speaker B | Speech Y | Speech $\tilde{y}$ | Speech $\hat{y}$ |
|---|---|---|---|

# Experiments

- Data
  - VCC2018

| Sources Speakers | |
|---|---|
| VCC2SF1 | VCC2SM1 |
| VCC2SF2 | VCC2SM2 |
| VCC2SF4 | VCC2SM4 |
| VCC2TF2 | VCC2TM2 |

| Target Speakers | |
|---|---|
| VCC2SF4 | VCC2SM4 |
| VCC2TF2 | VCC2TM2 |

**Many-to-Many VC**

| Target Speakers | |
|---|---|
| VCC2SF3 | VCC2SM3 |
| VCC2TF1 | VCC2TM1 |

**Any-to-Any VC**

# Experiments
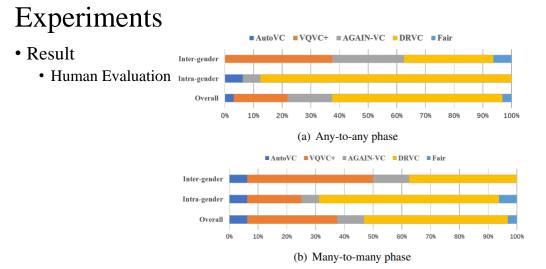
- Result
  - MCD & MOS

**Table 1**. Comparison of different models in any-to-any and many-to-many. ⇓ means lower score is better, and ⇑ means bigger score is better.

| Methods | Any-to-Any | | Many-to-Many | |
|---|---|---|---|---|
| | MCD ⇓ | MOS ⇑ | MCD ⇓ | MOS ⇑ |
| Real | - | $4.65 \pm 0.12$ | - | $4.66 \pm 0.21$ |
| VQVC+ | $7.47 \pm 0.07$ | $2.52 \pm 0.42$ | $7.78 \pm 0.07$ | $2.62 \pm 0.22$ |
| AutoVC | $7.69 \pm 0.21$ | $2.95 \pm 0.56$ | $7.61 \pm 0.17$ | $3.17 \pm 0.65$ |
| AGAIN-VC | $7.42 \pm 0.19$ | $2.45 \pm 0.34$ | $7.64 \pm 0.21$ | $2.47 \pm 0.58$ |
| **DRVC** | $\mathbf{7.39 \pm 0.05}$ | $\mathbf{3.32 \pm 0.36}$ | $\mathbf{7.59 \pm 0.04}$ | $\mathbf{3.51 \pm 0.52}$ |

# Experiments

- Result
  - Human Evaluation



(a) Any-to-any phase



(b) Many-to-many phase

# Experiments

- Result
  - Ablation experiments

**Table 2**. Ablation experiments on the proposed model. ⇓ means lower score is better.

| Model | MCD⇓ |
|---|---|
| DRVC *w/o* Cycle Loss | 7.68 ± 0.26 |
| DRVC *w/o* Identity Loss | 7.63 ± 0.14 |
| DRVC *w/o* Domain Loss | 7.72 ± 0.12 |
| DRVC *w/o* Voice Same Loss | 7.75 ± 0.32 |
| DRVC *w/o* Content Same Loss | 7.50 ± 0.32 |
| DRVC *w/o* Adversarial Loss | 7.72 ± 0.35 |
| **DRVC** | **7.39 ± 0.05** |

# Conclusion

- Contribution

  - We propose a end-to-end framework, DRVC, to address the untangle overlapping problem without circumspection choose the content sizes.

  - Both the subjective and objective results show our model has better performance.

# Thanks for you listening

Acknowledge & Notes:
- All anime character images are from the '*The Legend of LUOXIAOHEI*'.
- The presentation speech video, including the voice and personal video, is auto synthesis by PingAn Technology Co. Ltd.

Qiqi Wang (Virtual)                                        21