

# Technical Report



Zehua Liu  
刘泽华



# CNVSRC 2024

Chinese Continuous Visual Speech Recognition Challenge

# CNVSRC2024 Technical Report

**Zehua Liu, Lantian Li, Dong Wang**

Tsinghua University & Beijing University of Posts and Telecommunications

2024.08.16

NCMMSC-CNVSRC2024 Workshop, Urumqi, China



海天瑞声

Speech home

# Outline

- Data and Tasks
- Baseline
- Technical Summary

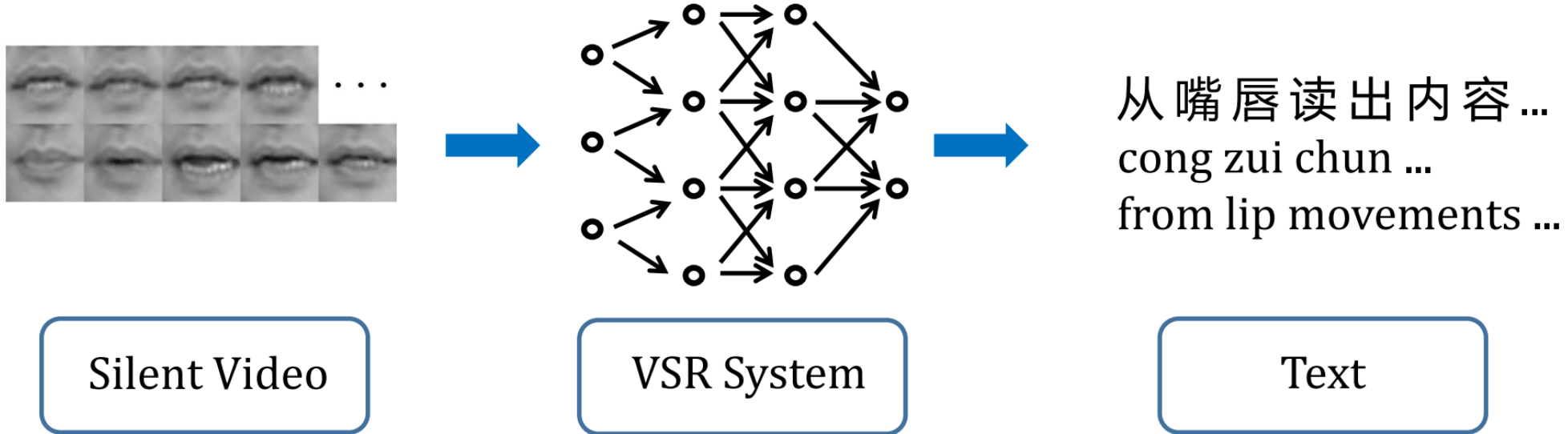
# Outline

- **Data and Tasks**
- **Baseline**
- **Technical Summary**

# Task Description

## □ VSR(Visual Speech Recognition)

### □ Definition



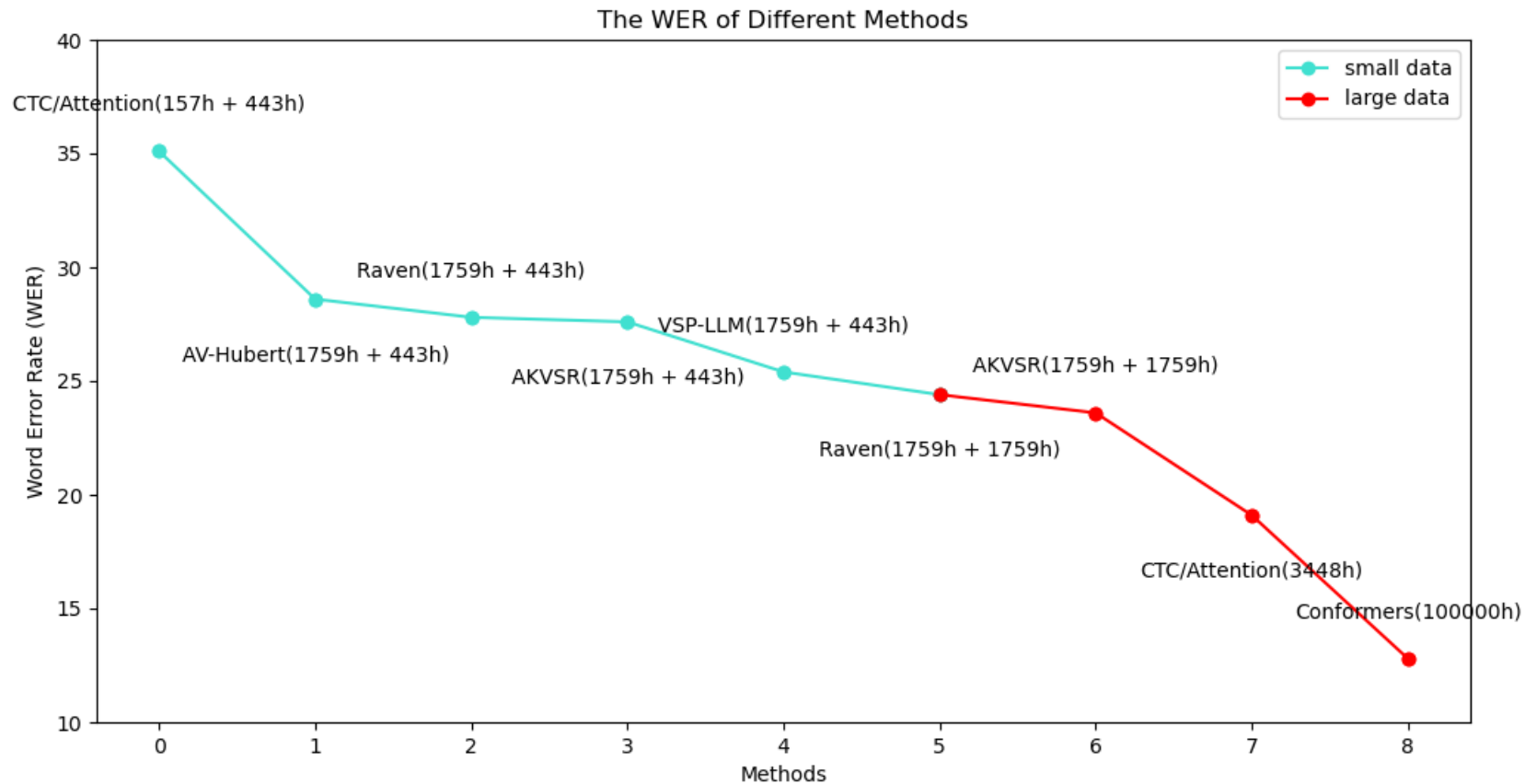
### □ Performance measurement

$$\text{CER} = \frac{\mathcal{N}_{\text{Ins}} + \mathcal{N}_{\text{Subs}} + \mathcal{N}_{\text{Del}}}{\mathcal{N}_{\text{Total}}} \times 100\%$$

# Task Description

## □ VSR(Visual Speech Recognition) in English

### □ LRS3[1]

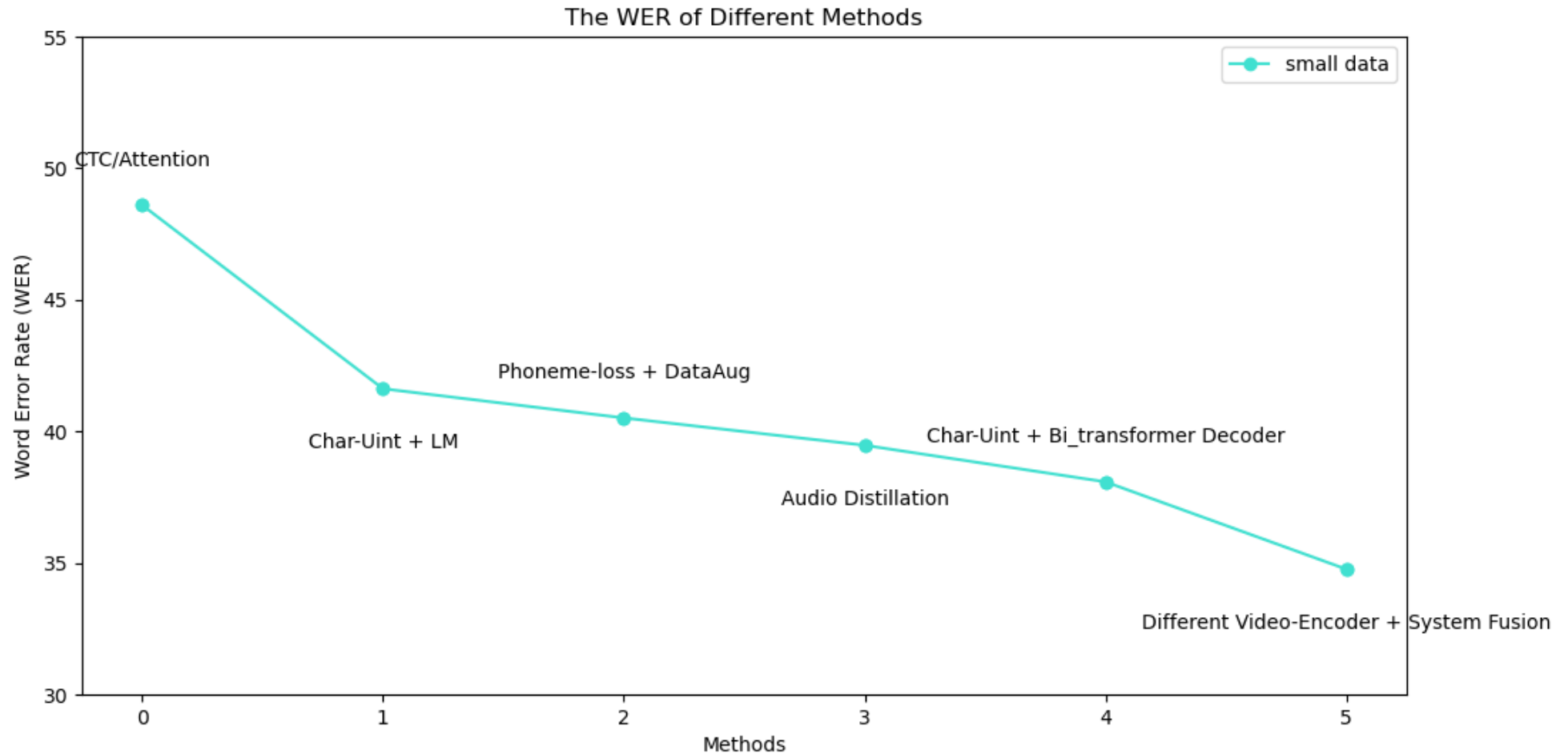


[1]T. Afouras, J. S. Chung, and A. Zisserman, “LRS3-TED: a large-scale dataset for visual speech recognition.”

# Task Description

## □ VSR(Visual Speech Recognition) in Chinese

### □ CNVSRC 2023 Single-Speaker Fixed Track



# Data

## □ Dataset for CNVSRC 2024

### ● CNCVS

- 2557 Speakers
- 300 Hours Audio-video paired data

### ● CN-CVS2-P1

- Additional Dataset for open track
- 200 Hours Audio-video paired data

### ● CNVSRC – Single

- 1 Speaker
- 100 Hours Audio-video paired data ,collect from Internet

### ● CNVSRC – Multi

- 43 Speakers
- 1 Hour per Speaker Audio-video paired Data



# Task Description

## □ T1 Single-Speaker VSR

### ● Fixed Track

- **ONLY** CN-CVS and CNVSRC-Single.Dev is allowed for training/tuning **ALL** the components of the system.
- This track is designed to compare different techniques under the **SAME** data resource.

### ● Open Track

- **ANY** data sources can be used for developing **ALL** the components of the system.
- This track is designed to examine the performance **Frontier** of the present technologies.

	Fixed Track	Open Track
T1: Single-speaker VSR	CN-CVS, CNVSRC-Single.Dev	No constraint (e.g. CN-CVS2-P1)
T2: Multi-speaker VSR	CN-CVS, CNVSRC-Multi.Dev	No constraint (e.g. CN-CVS2-P1)

# Task Description

## □ T2 Multi-Speaker VSR

### ● Fixed Track

- **ONLY** CN-CVS and CNVSRC-Multi.Dev is allowed for training/tuning **ALL** the components of the system.
- This track is designed to compare different techniques under the **SAME** data resource.

### ● Open Track

- **ANY** data sources can be used for developing **ALL** the components of the system.
- This track is designed to examine the performance **Frontier** of the present technologies.

	Fixed Track	Open Track
T1: Single-speaker VSR	CN-CVS, CNVSRC-Single.Dev	No constraint (e.g. CN-CVS2-P1)
T2: Multi-speaker VSR	CN-CVS, CNVSRC-Multi.Dev	No constraint (e.g. CN-CVS2-P1)

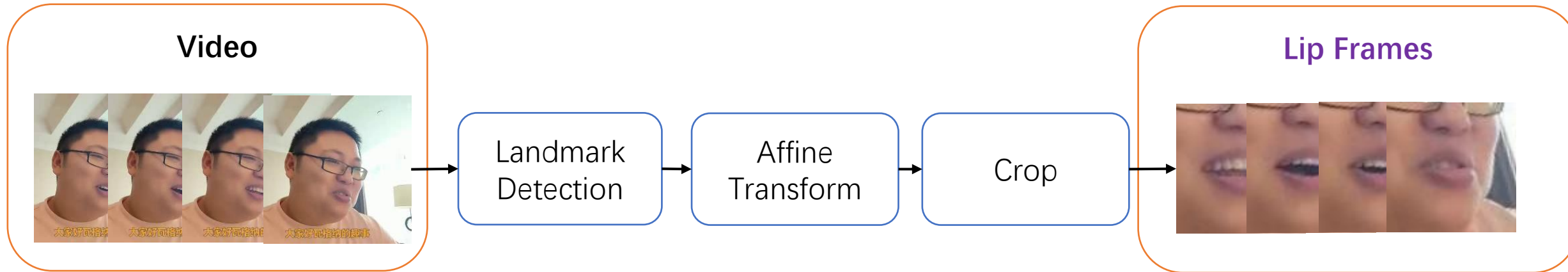
# Outline

- Data and Tasks
- **Baseline**
- Technical Summary

# Baseline

## □ Data processing

### □ Video Preprocess

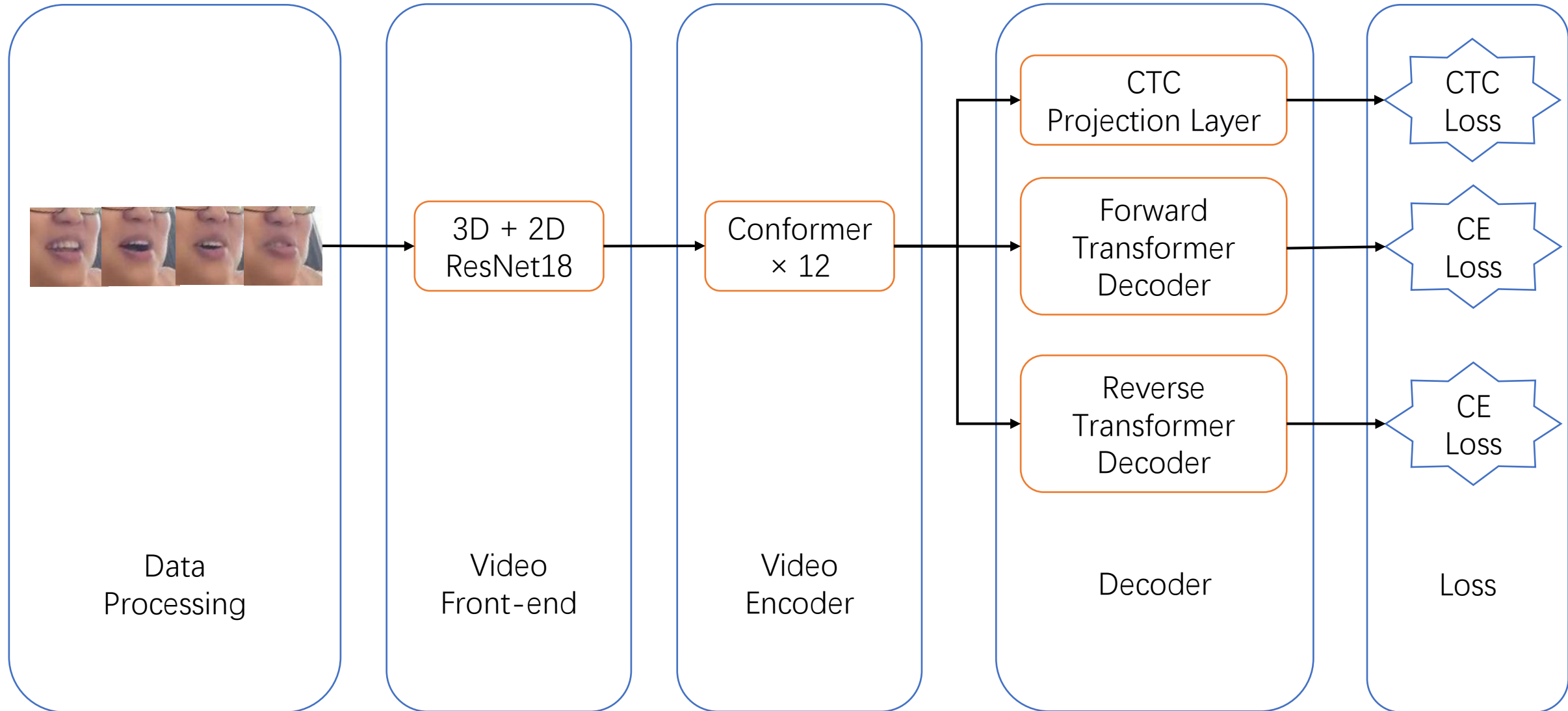


### □ Model Input

The input of the VSR model is **Lip Frames**

# Baseline

## □ Model Structure



# Baseline

## □ CNVSR2024 vs CNVSR 2023

### □ CNVSR 2024 Improvement

- Reverse Transformer Decoder
- Model units from subword units (5,904 units) into char units (4,468 units)

### □ Result

	T1: Single-speaker VSR		T2: Multi-speaker VSR	
	CNVSR 2024	CNVSR 2023	CNVSR 2024	CNVSR 2023
Valid	41.22%	48.57%	52.42%	58.77%
Eval	39.66%	48.60%	52.20%	58.37%

# Outline

- Data and Tasks
- Baseline
- **Technical Summary**

# Technical Summary

Components	Methods
Data processing	Face Detection, Face Alignment, Multi-scale Lip Region Extraction
Data Augmentation	Speed Perturbation, Adaptive Time Masking, Random Crop, Flip, Color Transformation, Image Size Crop
Visual Front-end	Resnet3D, Enhanced Resnet3D
Encoder	Conformer, Branchformer, E-Branchformer
Decoder	Transformer Decoder, Reverse Transformer Decoder, S4D Decoder, Reverse S4D Decoder
Loss function	CTC/Attention Loss, KLDivLoss
Training strategy	Pretrain, Fine-tune
System fusion	Score-level Average



# Technical Highlights

Components	Methods
Data processing	Face Detection, Face Alignment, Multi-scale Lip Region Extraction
Data Augmentation	Speed Perturbation, Adaptive Time Masking, Random Crop, Flip, Color Transformation, <b>Image Size Crop</b>
Visual Front-end	Resnet3D, <b>Enhanced Resnet3D</b>
Encoder	Conformer, Branchformer, E-Branchformer
Decoder	Transformer Decoder, Reverse Transformer Decoder, <b>S4D Decoder, Reverse S4D Decoder</b>
Loss function	CTC/Attention Loss, KLDivLoss
Training strategy	<b>Pretrain</b> , Fine-tune
System fusion	<b>Score-level Average</b>

# Data Augmentation

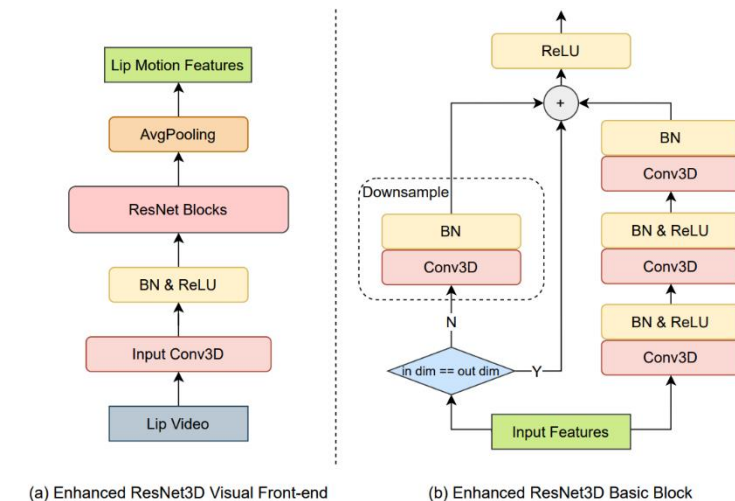
Components	Methods
Data processing	Face Detection, Face Alignment, Multi-scale Lip Region Extraction
Data Augmentation	Speed Perturbation, Adaptive Time Masking, Random Crop, Flip, Color Transformation, <b>Image Size Crop</b>
Visual Front-end	Resnet3D, <b>Enhanced Resnet3D</b>
Encoder	Conformer, Branchformer, E-Branchformer
Decoder	Transformer Decoder, Reverse Transformer Decoder, <b>S4D Decoder</b>
Loss function	CTC/Attention Loss, KLDivLoss
Training strategy	<b>Pretrain</b> , Fine-tune
System fusion	<b>Score-level Average</b>

System	Visual Frontend	Encoder	Decoder	Crop	T1.Dev
M1	ResNet3D	E-Branchformer	Transformer	80	41.76
M2	ResNet3D	Conformer	Transformer	96	39.27
M3	ResNet3D	Branchformer	Transformer	96	38.86
M4	ResNet3D	E-Branchformer	Transformer	96	38.39
M5	Enhanced ResNet3D	E-Branchformer	Transformer	96	38.30
M6	ResNet3D	E-Branchformer	Transformer	112	38.13
M7	Enhanced ResNet3D	E-Branchformer	Bi-Transformer	112	36.60
M8	Enhanced ResNet3D	E-Branchformer	Bi-Transformer	128	36.48
R1	ROVER Eval of M1~M8				
R2	ROVER R1 and Eval.FT of M1~M8				

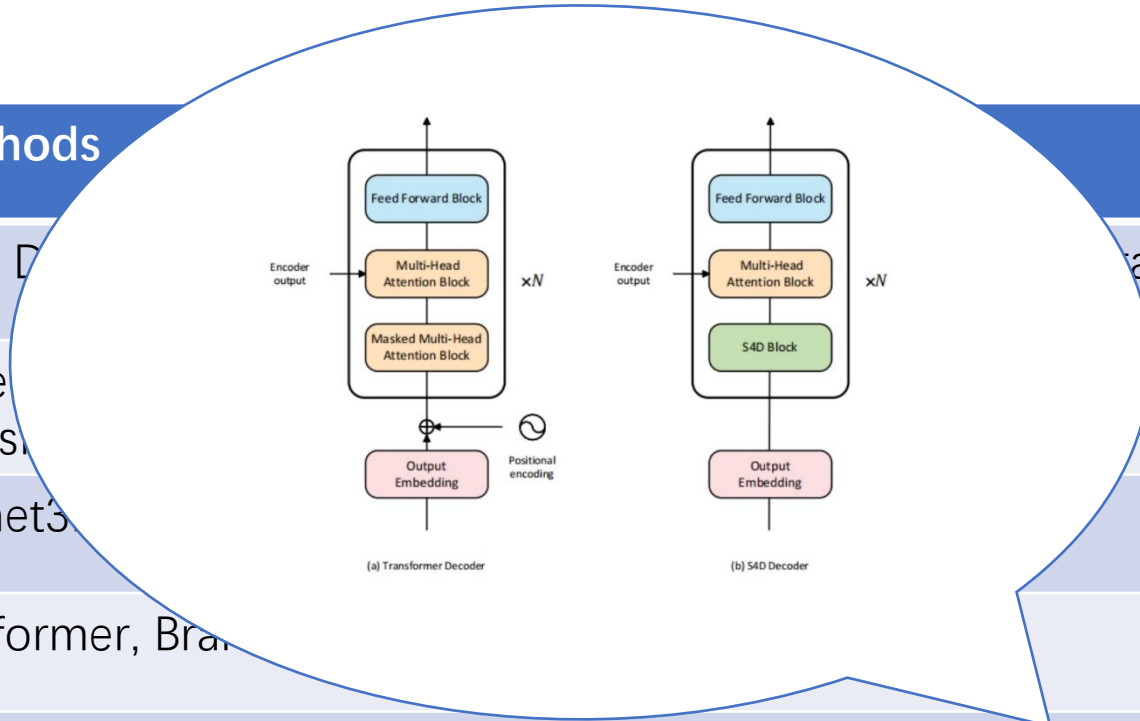
# Visual Front-end

Components	Methods
Data processing	Face Detection, Face Alignment, Multi-scale Lip Region Extraction
Data Augmentation	Speed Perturbation, Adaptive Time Masking, Random Crop, Flip, Color Transformation, <b>Image Size Crop</b>
Visual Front-end	Resnet3D, <b>Enhanced Resnet3D</b>
Encoder	Conformer, Branchformer, E-B
Decoder	Transformer Decoder, Revers <b>S4D Decoder</b>
Loss function	CTC/Attention Loss, KLDivLos
Training strategy	<b>Pretrain</b> , Fine-tune
System fusion	<b>Score-level Average</b>



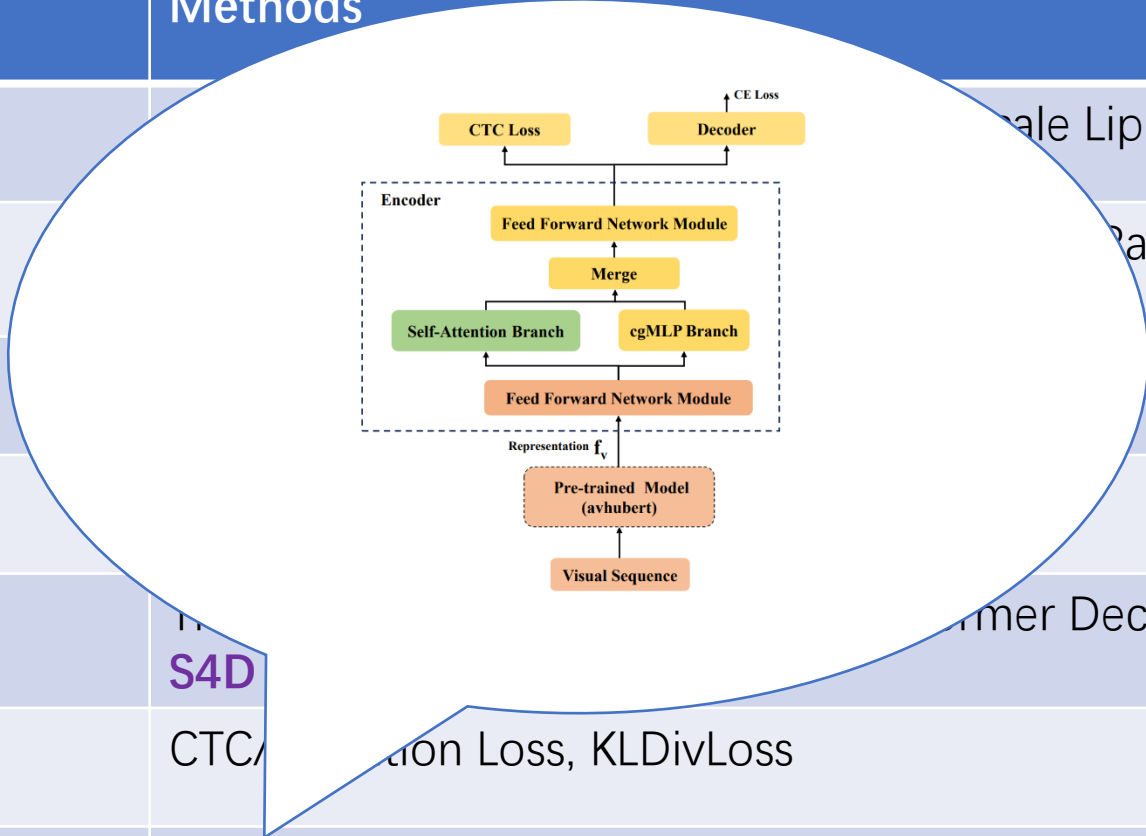
# Decoder

Components	Methods
Data processing	Face D
Data Augmentation	Spee Trans
Visual Front-end	Resnet3
Encoder	Conformer, Bran
Decoder	Transformer Decoder, Reverse Transformer Decoder, <b>S4D Decoder, Reverse S4D Decoder</b>
Loss function	CTC/Attention Loss, KLDivLoss
Training strategy	<b>Pretrain</b> , Fine-tune
System fusion	<b>Score-level Average</b>



# Training Strategy

Components	Methods
Data processing	Scale Lip Region Extraction
Data Augmentation	Random Crop, Flip, Color
Visual Front-end	
Encoder	
Decoder	Former Decoder, <b>S4D Decoder</b> , Reverse
Loss function	CTC/Attention Loss, KLDivLoss
Training strategy	<b>Pretrain</b> , Fine-tune
System fusion	<b>Score-level Average</b>



# System fusion

Components	Methods																																																																																																																									
Data processing	Face Detection, Face Alignment, Multi-scale Lip Region Extraction																																																																																																																									
Data Augmentation	Speed P, Trans, Flip, Color																																																																																																																									
Visual Front-end	<table border="1"> <thead> <tr> <th>System</th> <th>Visual Frontend</th> <th>Encoder</th> <th>Decoder</th> <th>Crop</th> <th>T1.Dev</th> <th>T1.Eval</th> <th>T1.Eval.FT</th> <th>T2.Dev</th> <th>T2.Eval</th> <th>T2.Eval.FT</th> </tr> </thead> <tbody> <tr> <td>M1</td> <td>ResNet3D</td> <td>E-Branchformer</td> <td>Transformer</td> <td>80</td> <td>41.76</td> <td>-</td> <td>-</td> <td>48.50</td> <td>-</td> <td>-</td> </tr> <tr> <td>M2</td> <td>ResNet3D</td> <td>Conformer</td> <td>Transformer</td> <td>96</td> <td>39.27</td> <td>-</td> <td>-</td> <td>46.21</td> <td>-</td> <td>-</td> </tr> <tr> <td>M3</td> <td>ResNet3D</td> <td>Branchformer</td> <td>Transformer</td> <td>96</td> <td>38.86</td> <td>-</td> <td>-</td> <td>46.49</td> <td>-</td> <td>-</td> </tr> <tr> <td>M4</td> <td>ResNet3D</td> <td>E-Branchformer</td> <td>Transformer</td> <td>96</td> <td>38.39</td> <td>-</td> <td>-</td> <td>45.97</td> <td>-</td> <td>-</td> </tr> <tr> <td>M5</td> <td>Enhanced ResNet3D</td> <td>E-Branchformer</td> <td>Transformer</td> <td>96</td> <td>38.30</td> <td>-</td> <td>-</td> <td>45.79</td> <td>45.43</td> <td>41.75</td> </tr> <tr> <td>M6</td> <td>ResNet3D</td> <td>E-Branchformer</td> <td>Transformer</td> <td>112</td> <td>38.13</td> <td>-</td> <td>35.91</td> <td>44.88</td> <td>44.56</td> <td>40.37</td> </tr> <tr> <td>M7</td> <td>Enhanced ResNet3D</td> <td>E-Branchformer</td> <td>Bi-Transformer</td> <td>112</td> <td>36.60</td> <td>36.79</td> <td>34.56</td> <td>43.29</td> <td>43.13</td> <td>38.73</td> </tr> <tr> <td>M8</td> <td>Enhanced ResNet3D</td> <td>E-Branchformer</td> <td>Bi-Transformer</td> <td>128</td> <td>36.48</td> <td>37.04</td> <td><b>34.27</b></td> <td>43.23</td> <td>42.99</td> <td><b>37.93</b></td> </tr> <tr> <td>R1</td> <td colspan="4">ROVER Eval of M1~M8</td> <td></td> <td colspan="2">32.7199</td> <td colspan="3">38.3602</td> </tr> <tr> <td>R2</td> <td colspan="4">ROVER R1 and Eval.FT of M1~M8</td> <td></td> <td colspan="2"><b>30.4679</b></td> <td colspan="3"><b>34.2955</b></td> </tr> </tbody> </table>	System	Visual Frontend	Encoder	Decoder	Crop	T1.Dev	T1.Eval	T1.Eval.FT	T2.Dev	T2.Eval	T2.Eval.FT	M1	ResNet3D	E-Branchformer	Transformer	80	41.76	-	-	48.50	-	-	M2	ResNet3D	Conformer	Transformer	96	39.27	-	-	46.21	-	-	M3	ResNet3D	Branchformer	Transformer	96	38.86	-	-	46.49	-	-	M4	ResNet3D	E-Branchformer	Transformer	96	38.39	-	-	45.97	-	-	M5	Enhanced ResNet3D	E-Branchformer	Transformer	96	38.30	-	-	45.79	45.43	41.75	M6	ResNet3D	E-Branchformer	Transformer	112	38.13	-	35.91	44.88	44.56	40.37	M7	Enhanced ResNet3D	E-Branchformer	Bi-Transformer	112	36.60	36.79	34.56	43.29	43.13	38.73	M8	Enhanced ResNet3D	E-Branchformer	Bi-Transformer	128	36.48	37.04	<b>34.27</b>	43.23	42.99	<b>37.93</b>	R1	ROVER Eval of M1~M8					32.7199		38.3602			R2	ROVER R1 and Eval.FT of M1~M8					<b>30.4679</b>		<b>34.2955</b>		
System		Visual Frontend	Encoder	Decoder	Crop	T1.Dev	T1.Eval	T1.Eval.FT	T2.Dev	T2.Eval	T2.Eval.FT																																																																																																															
M1		ResNet3D	E-Branchformer	Transformer	80	41.76	-	-	48.50	-	-																																																																																																															
M2		ResNet3D	Conformer	Transformer	96	39.27	-	-	46.21	-	-																																																																																																															
M3		ResNet3D	Branchformer	Transformer	96	38.86	-	-	46.49	-	-																																																																																																															
M4		ResNet3D	E-Branchformer	Transformer	96	38.39	-	-	45.97	-	-																																																																																																															
M5		Enhanced ResNet3D	E-Branchformer	Transformer	96	38.30	-	-	45.79	45.43	41.75																																																																																																															
M6		ResNet3D	E-Branchformer	Transformer	112	38.13	-	35.91	44.88	44.56	40.37																																																																																																															
M7		Enhanced ResNet3D	E-Branchformer	Bi-Transformer	112	36.60	36.79	34.56	43.29	43.13	38.73																																																																																																															
M8		Enhanced ResNet3D	E-Branchformer	Bi-Transformer	128	36.48	37.04	<b>34.27</b>	43.23	42.99	<b>37.93</b>																																																																																																															
R1	ROVER Eval of M1~M8					32.7199		38.3602																																																																																																																		
R2	ROVER R1 and Eval.FT of M1~M8					<b>30.4679</b>		<b>34.2955</b>																																																																																																																		
Encoder																																																																																																																										
Decoder																																																																																																																										
Loss function	CTC																																																																																																																									
Training strategy	Pretrain, Fine-tune																																																																																																																									
System fusion	Score-level Average																																																																																																																									

erse



# CNVSRC 2024

Chinese Continuous Visual Speech Recognition Challenge

# Many Thanks!



海天瑞声

Speech home