



# CNVSRC 2023 Technical Report

Chen Chen, Lantian Li, Dong Wang

Tsinghua University & Beijing University of Posts and Telecommunications

2023.12.09

NCMMSC-CNVSRC 2023 Workshop, Suzhou, China



海天瑞声

Speech home

# OUTLINE

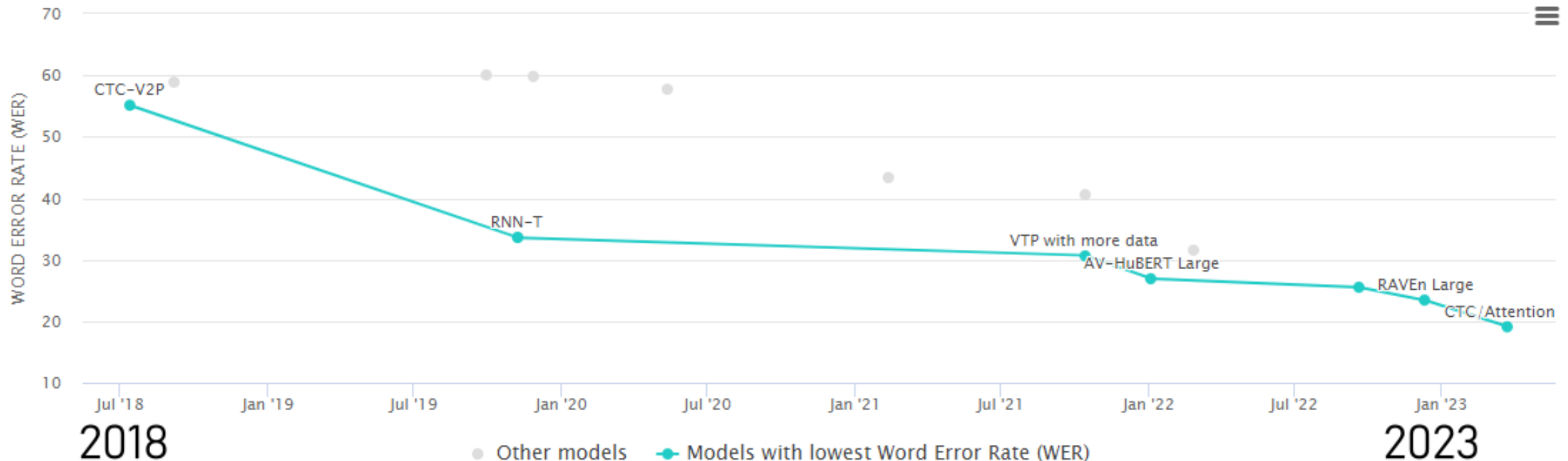
- Data and Tasks
- Baseline
- Technical Summary

# OUTLINE

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

# The Origin of CN-CVS Dataset

- The lack of Chinese Audio-video data constrain the development of Chinese VSR
  - English VSR: Auto-AVSR[1] reach 19.1% WER on LRS3[2] **in the wild**.



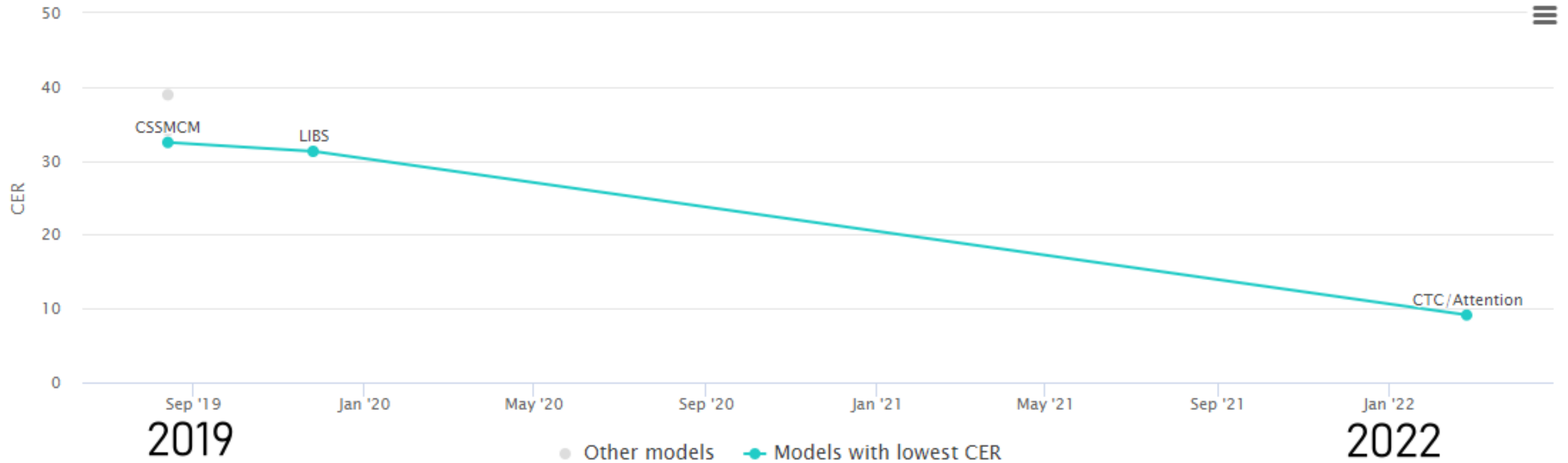
[1] P. Ma, A. Haliassos, A. Fernandez-Lopez, H. Chen, S. Petridis, and M. Pantic, "Auto-AVSR: Audio-Visual Speech Recognition with Automatic Labels," ICASSP 2023

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

<https://paperswithcode.com/sota/lipreading-on-lrs3-ted>

# The Origin of CN-CVS Dataset

- The lack of Chinese Audio-video data constrain the development of Chinese VSR
  - Chinese VSR: [3] reach 9.1% WER on CMLR[4] **in the lab.**



[3] P. Ma, S. Petridis, and M. Pantic, “Visual Speech Recognition for Multiple Languages in the Wild.”

[4] Y. Zhao, R. Xu, and M. Song, “A Cascade Sequence-to-Sequence Model for Chinese Mandarin Lip Reading.”

<https://paperswithcode.com/sota/lipreading-on-cmlr>

# The Origin of CN-CVS Dataset

- The lack of Chinese Audio-video data constrain the development of Chinese VSR

数据集名称	文献	语言	内容类型	数据来源	词汇量	说话人数量	句子数量	时长	摄像机角度
GRID	Cooke et al. (2006)	英语	语法模式	录制朗读	51	33	33000	27h	0°
TCD-TIMIT	Harte et al. (2015)	英语	句	录制朗读	5954	62	6913	≈20h	0°, 30°
Lip2Wav	Prajwal et al. (2020)	英语	句	讲课节目	≈5k/spk	5	-	≈120h	自然角度
LRW	Chung et al. (2017a)	英语	词	电视新闻	500	-	≈539000	173h	自然角度
LRS	Chung et al. (2017b)	英语	句	电视新闻	≈17k	-	118116	75.5h	自然角度
LRS2	Afouras et al. (2019)	英语	句	电视新闻	≈60k	-	≈145000	224.5h	自然角度
LRS3	Afouras et al. (2019)	英语	句	演讲节目	≈70k	≈10k	≈165000	475h	自然角度
VoxCeleb1	Nagrani et al. (2017)	英语	句	网络视频	-	1251	153516	352h	自然角度
VoxCeleb2	Chung et al. (2018)	多语言	句	网络视频	-	6112	1128246	2442h	自然角度
AVSpeech	Ephrat et al. (2018)	多语言	句	演讲授课	-	≈150k	-	≈4700h	自然角度
CAS-VSR-W1k	Yang et al. (2019)	汉语	词	电视节目	1k	>2000	718018	≈140h	自然角度
CMLR	Zhao et al. (2019)	汉语	句	电视新闻	3517	11	102076	≈88h	0°
CN-CVS/News	Chen et al. (2023)	汉语	句	电视新闻	-	28	13016	34.6h	0°
CN-CVS/Speech	Chen et al. (2023)	汉语	句	演讲节目	-	2529	193245	273.4h	自然角度

# Data Collection Pipeline

## 1 Programs selection

新闻类节目  
 单人公开演讲  
 科普演讲  
 ...

## 2 Videos download

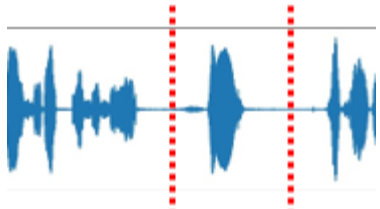


## 3 Shot detection, VAD and Face Tracking

Shot detection  
*Ffmpeg*



VAD  
*pydub*



Face tracking  
*dlib*



## 4 Mouth-Speech Sync

Synchronization  
*SyncNet*



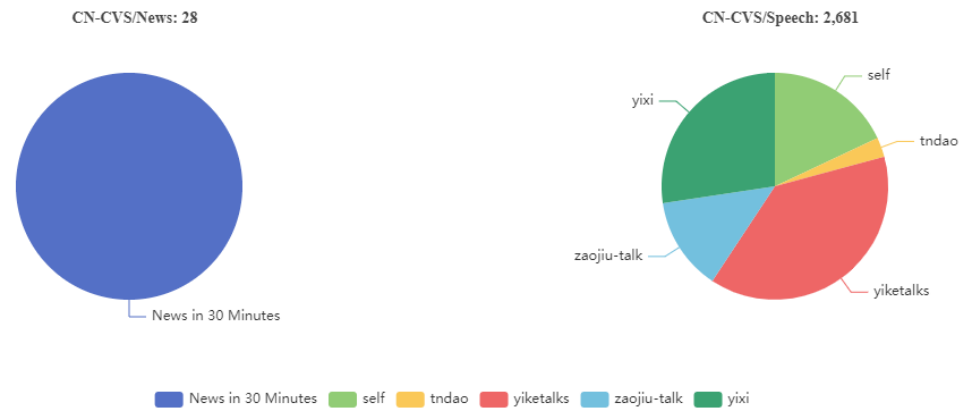
## 5 Human check



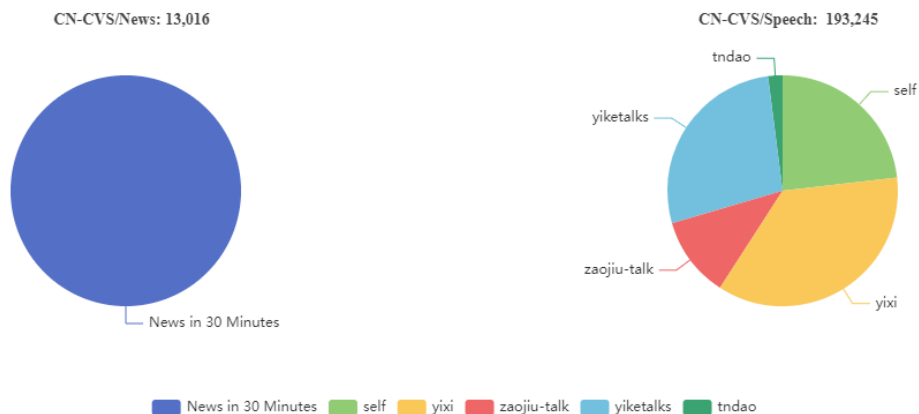
# Data Profile



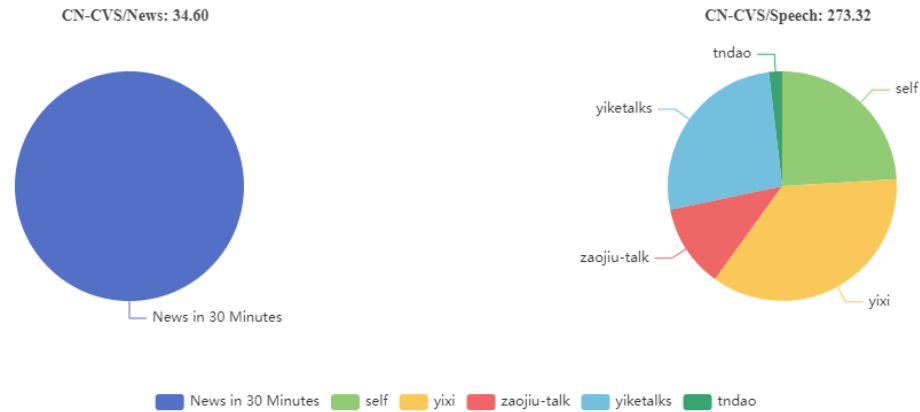
## 2,557 Speakers



## 206,261 Utters



## ~300 Hours





# Data Profile

## □ Additional datasets for CNVSRC 2023

### ● CNVSRC-Single

- 1 Speaker
- 100 Hours Audio-video paired data

### ● CNVSRC-Multi

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

	CNVSRC-Single		CNVSRC-Multi	
DataSet	Dev	Eval	Dev	Eval
# Videos	25,947	2,881	20,450	10,269
# Hours	94.00	8.41	29.24	14.49

# Task Description – 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
T2: Multi-speaker VSR	CN-CVS, CNVSRC-Multi.Dev	No constraint

# Task Description – 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.

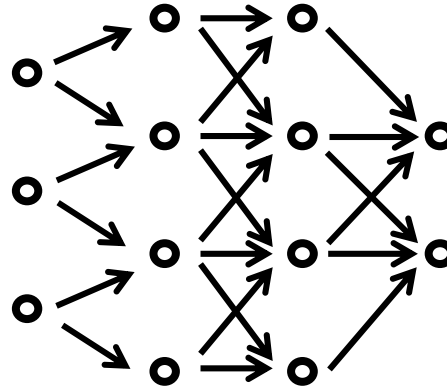
	Fixed Track	Open Track
T1: Single-speaker VSR	CN-CVS, CNVSRC-Single.Dev	No constraint
T2: Multi-speaker VSR	CN-CVS, CNVSRC-Multi.Dev	No constraint

# Task Description – VSR

## □ Definition



Silent Video



VSR System



从嘴唇读出内容 ...  
cong zui chun ...  
from lip movements ...

Text

## □ Performance measurement

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

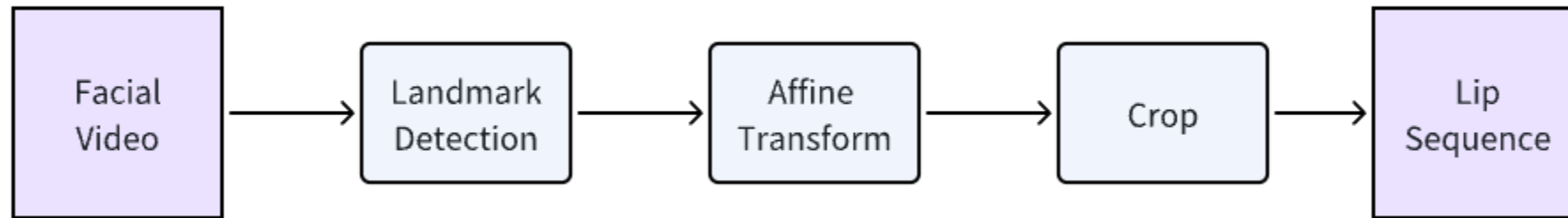
# OUTLINE

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

# Baseline

## □ Data processing

### ● Video Data



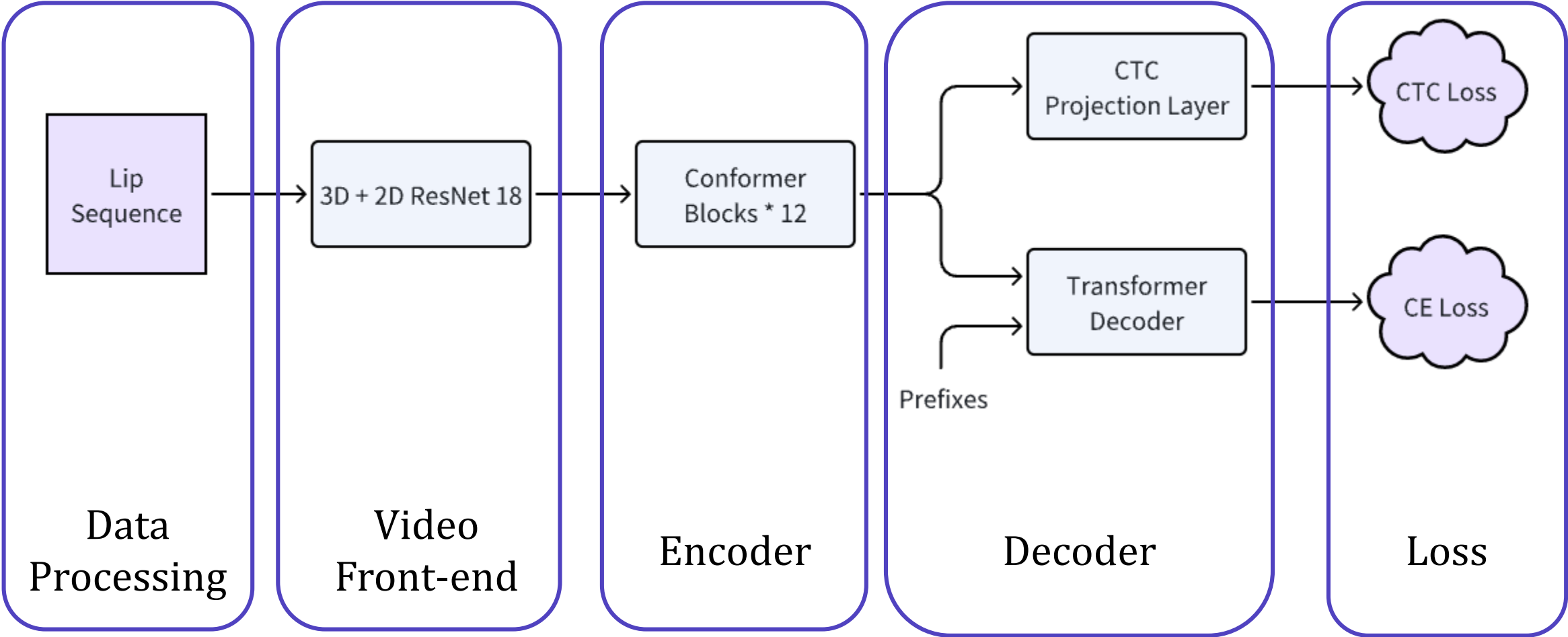
### ● Text Data

- SentencePiece \*

\*<https://pypi.org/project/sentencepiece/>

# Baseline

## Model Structure



# OUTLINE

- Data and Tasks
- Baseline
- Technical Summary**



# Representative Techniques (6 Teams)

Components	Methods
Data processing	Face Detection, Face Alignment, Multi-scale Lip Region Extraction
Data Augmentation	Speed perturbation, Adaptive Time Masking, Random Erase, Flip, Generated Facial Video
Video Front-end	3D+2D ResNet18, ResNet3D
Encoder	Conformer, Branchformer, E-Branchformer
Decoder	Transformer Decoder, Bi-Transformer Decoder
Auxiliary design	BPE/Character as Modeling Unit, Phoneme-level prediction task
Loss function	CTC/Attention Loss, CTC Loss in shallow layers, Cross Modality Similarity Loss
Training strategy	Pretrain + Fine-tune, ASR Knowledge Distillation
Language Model	RNNLM, Transformer LM
System fusion	Score-level average

# Technical Highlights

Components	Methods
Data processing	Face Detection, Face Alignment, <b>Multi-scale Lip Region Extraction</b>
Data Augmentation	<b>Speed perturbation</b> , Adaptive Time Masking, Random Erase, Flip, Generated Facial Video
Video Front-end	3D+2D ResNet18, ResNet3D
Encoder	Conformer, Branchformer, E-Branchformer
Decoder	Transformer Decoder, Bi-Transformer Decoder
Auxiliary design	<b>Character as Modeling Unit</b> , Phoneme-level prediction task
Loss function	CTC/Attention Loss, <b>CTC Loss in shallow layers</b> , Cross Modality Similarity Loss
Training strategy	Pretrain + Fine-tune, <b>ASR Knowledge Distillation</b>
Language Model	RNNLM, Transformer LM
System fusion	Score-level average

# Data processing

Components	Methods																																																																																								
Data processing	Face Detection, Face Alignment, <b>Multi-scale Lip Region Extraction</b>																																																																																								
Data Augmentation	<b>Speed perturbation</b> , Adaptive Time Masking, Random Erase, Flip, Generated Facial Video																																																																																								
Video Front-end	3D+2D ResNet																																																																																								
Encoder	<p style="text-align: center;"><b>T237</b></p> <pre> graph LR     A[Raw Video Size: (224 x 224)] -- x1 --&gt; B[Lip Extractor]     B -- x1 --&gt; C[Lip Video Size: (N x N)]           </pre>																																																																																								
Decoder																																																																																									
Auxiliary																																																																																									
Loss func	<table border="1"> <thead> <tr> <th>System</th> <th>Encoder</th> <th>Crop</th> <th>SP</th> <th>T1.Dev</th> <th>T1.Eval</th> <th>T2.Dev</th> <th>T2.Eval</th> </tr> </thead> <tbody> <tr> <td>Baseline<sup>6</sup></td> <td>Conf</td> <td>96</td> <td>✗</td> <td>48.57</td> <td>48.60</td> <td>58.77</td> <td>58.37</td> </tr> <tr> <td>M1</td> <td>Conf</td> <td>96</td> <td>✓</td> <td>39.43</td> <td>39.99</td> <td>46.08</td> <td>45.73</td> </tr> <tr> <td>M2</td> <td>Branch</td> <td>96</td> <td>✓</td> <td>39.00</td> <td>39.36</td> <td>46.63</td> <td>46.37</td> </tr> <tr> <td>M3</td> <td>E-Branch</td> <td>96</td> <td>✓</td> <td>38.59</td> <td>38.61</td> <td>46.26</td> <td>45.80</td> </tr> <tr> <td>M4</td> <td>E-Branch</td> <td>48</td> <td>✗</td> <td>46.88</td> <td>45.81</td> <td>55.58</td> <td>55.51</td> </tr> <tr> <td>M5</td> <td>E-Branch</td> <td>64</td> <td>✗</td> <td>44.40</td> <td>43.59</td> <td>53.64</td> <td>52.98</td> </tr> <tr> <td>M6</td> <td>E-Branch</td> <td>80</td> <td>✗</td> <td>42.95</td> <td>42.26</td> <td>50.77</td> <td>50.38</td> </tr> <tr> <td>M7</td> <td>E-Branch</td> <td>96</td> <td>✗</td> <td>40.56</td> <td>40.42</td> <td>47.16</td> <td>46.53</td> </tr> <tr> <td>M8</td> <td>E-Branch</td> <td>112</td> <td>✗</td> <td>38.46</td> <td>38.95</td> <td>45.17</td> <td>44.87</td> </tr> <tr> <td>ROVER</td> <td>-</td> <td>-</td> <td>-</td> <td><b>34.47</b></td> <td><b>34.76</b></td> <td><b>41.39</b></td> <td><b>41.06</b></td> </tr> </tbody> </table>	System	Encoder	Crop	SP	T1.Dev	T1.Eval	T2.Dev	T2.Eval	Baseline <sup>6</sup>	Conf	96	✗	48.57	48.60	58.77	58.37	M1	Conf	96	✓	39.43	39.99	46.08	45.73	M2	Branch	96	✓	39.00	39.36	46.63	46.37	M3	E-Branch	96	✓	38.59	38.61	46.26	45.80	M4	E-Branch	48	✗	46.88	45.81	55.58	55.51	M5	E-Branch	64	✗	44.40	43.59	53.64	52.98	M6	E-Branch	80	✗	42.95	42.26	50.77	50.38	M7	E-Branch	96	✗	40.56	40.42	47.16	46.53	M8	E-Branch	112	✗	38.46	38.95	45.17	44.87	ROVER	-	-	-	<b>34.47</b>	<b>34.76</b>	<b>41.39</b>	<b>41.06</b>
System	Encoder	Crop	SP	T1.Dev	T1.Eval	T2.Dev	T2.Eval																																																																																		
Baseline <sup>6</sup>	Conf	96	✗	48.57	48.60	58.77	58.37																																																																																		
M1	Conf	96	✓	39.43	39.99	46.08	45.73																																																																																		
M2	Branch	96	✓	39.00	39.36	46.63	46.37																																																																																		
M3	E-Branch	96	✓	38.59	38.61	46.26	45.80																																																																																		
M4	E-Branch	48	✗	46.88	45.81	55.58	55.51																																																																																		
M5	E-Branch	64	✗	44.40	43.59	53.64	52.98																																																																																		
M6	E-Branch	80	✗	42.95	42.26	50.77	50.38																																																																																		
M7	E-Branch	96	✗	40.56	40.42	47.16	46.53																																																																																		
M8	E-Branch	112	✗	38.46	38.95	45.17	44.87																																																																																		
ROVER	-	-	-	<b>34.47</b>	<b>34.76</b>	<b>41.39</b>	<b>41.06</b>																																																																																		
Training s																																																																																									
Language Mo																																																																																									
System fusion																																																																																									

# Data Augmentation

Components	Methods
Data processing	Face Detection, Face Alignment, Lip Region Extraction
Data Augmentation	<b>Speed perturbation</b> , <b>T237, T238</b> , <b>ASR Knowledge Distillation</b> , Random Erase, Flip, Generated Facial Video
Video Front-end	3D+2D ResNet18, ResNet3D
Encoder	Conformer, Branchformer, E-Branchformer
Decoder	Transformer Decoder, Bi-Transformer Decoder
Auxiliary design	Character as Modeling Unit, Phoneme-level prediction task
Loss function	CTC/Attention Loss, CTC Loss in shallow layers, Cross Modality Similarity Loss
Training strategy	Pretrain + Fine-tune, ASR Knowledge Distillation
Language Model	RNNLM, Transformer LM
System fusion	Score-level average

# Auxiliary design

Components	Methods
Data processing	Face Detection, Face Alignment,
Data Augmentation	Speed perturbation, Adaptive Ti
Video Front-end	3D+2D ResNet101, ResNet3D
Encoder	Conformer <b>T266, T267</b> Bi-Bra
Decoder	Transformer Decoder, Bi-Transformer Decoder
Auxiliary design	<b>Character as Modeling Unit</b> , Phoneme-level prediction task
Loss function	
Training strategy	
Language Model	
System fusion	

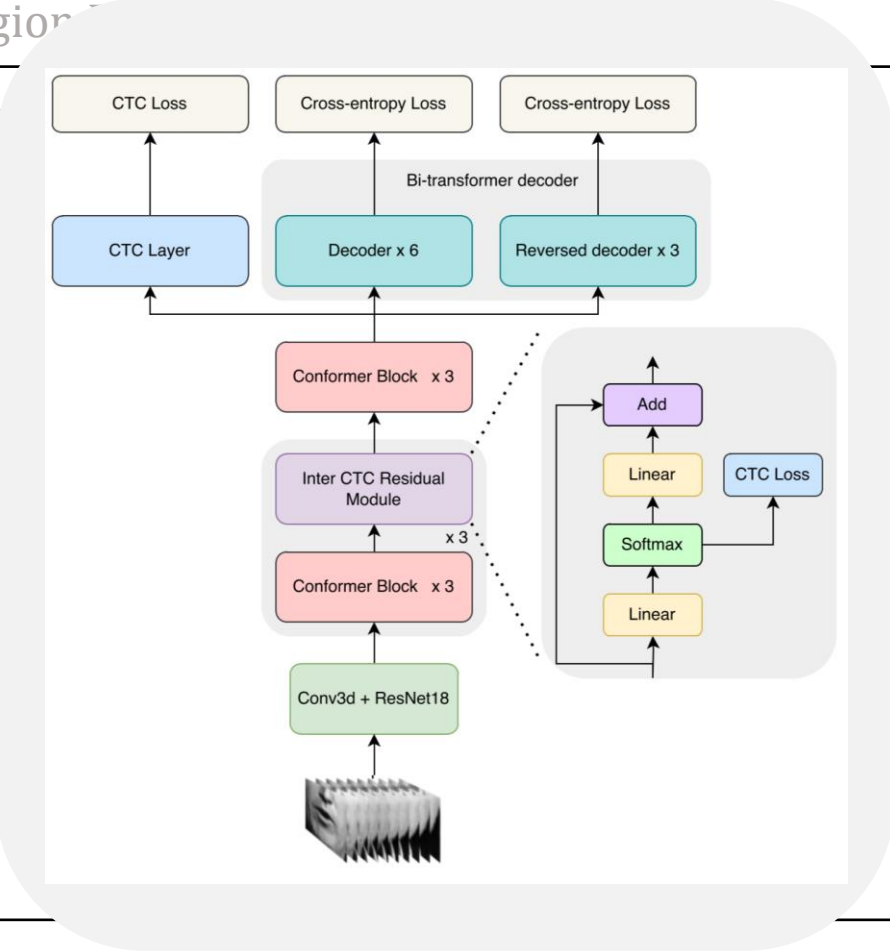
System	Model	CER (%)
M1	Proposed system	40.46
M2	M1 - RNNLM	40.62
M3	M2 - char unit	42.36
M4	M3 - Bi-transformer decoder	43.19
M5	M4 - Inter CTC residual module	48.57

method	training set(hours)	finetuning set(hours)		single-speaker dev CER	multi-speakers dev CER
		T1	T2		
baseline	287	83.7	18	48.57%	58.77%
char model unit	287	83.7	18	43.59%	56.77%
+video aug	287	166	36	42.30%	54.74%
+RNN LM				42.18%	
+transformer LM				42.16%	
+model fusion				<b>41.50%</b>	

# Loss function

Components	Methods
Data processing	T266
Data Augmentation	
Video Frame	
Encoder	
Decoder	
Auxiliary	
Loss function	CTC/Attention Loss, <b>CTC Loss in shallow layers</b> , CTC Loss
Training strategy	Pretrain + Fine-tune, ASR Knowledge Distillation
Language Model	RNNLM, Transformer LM
System fusion	Score-level average

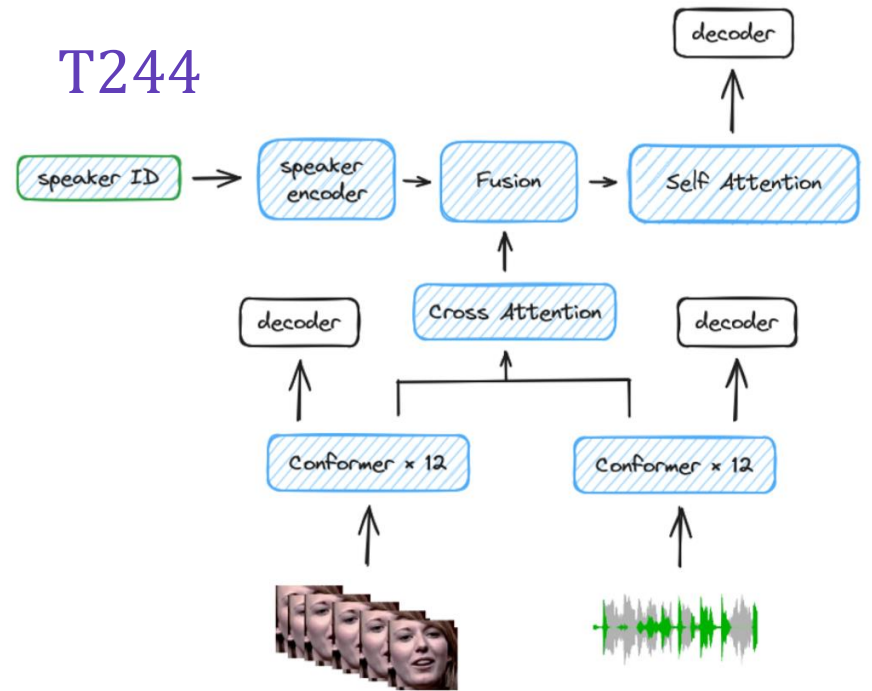
System	Model	CER (%)
M1	Proposed system	40.46
M2	M1 - RNNLM	40.62
M3	M2 - char unit	42.36
M4	M3 - Bi-transformer decoder	43.19
M5	M4 - Inter CTC residual module	48.57



# Training strategy

Components	Methods
<p>Dat: <b>T290</b></p> <p>Dat:</p> <p>Vid:</p> <p>Enc:</p> <p>Dec:</p> <p>Aux:</p> <p>Los:</p>	<p>Ground Truth Transcriptions → CTC/Attention Loss</p> <p>scale Lip Regi</p> <p>sking, Rando</p> <p>mer</p> <p>Decoder</p> <p>level predicti</p> <p>w layers, Cro</p>
<p>Training strategy</p>	<p>Pretrain + Fine-tune, <b>ASR Knowledge Distillation</b></p>
<p>Language Model</p>	<p>RNNLM, Transformer LM</p>
<p>System fusion</p>	<p>Score-level average</p>

**T244**





# CNVSRC 2023

Chinese Continuous Visual Speech Recognition Challenge

# Many Thanks !



海天瑞声

Speechhome