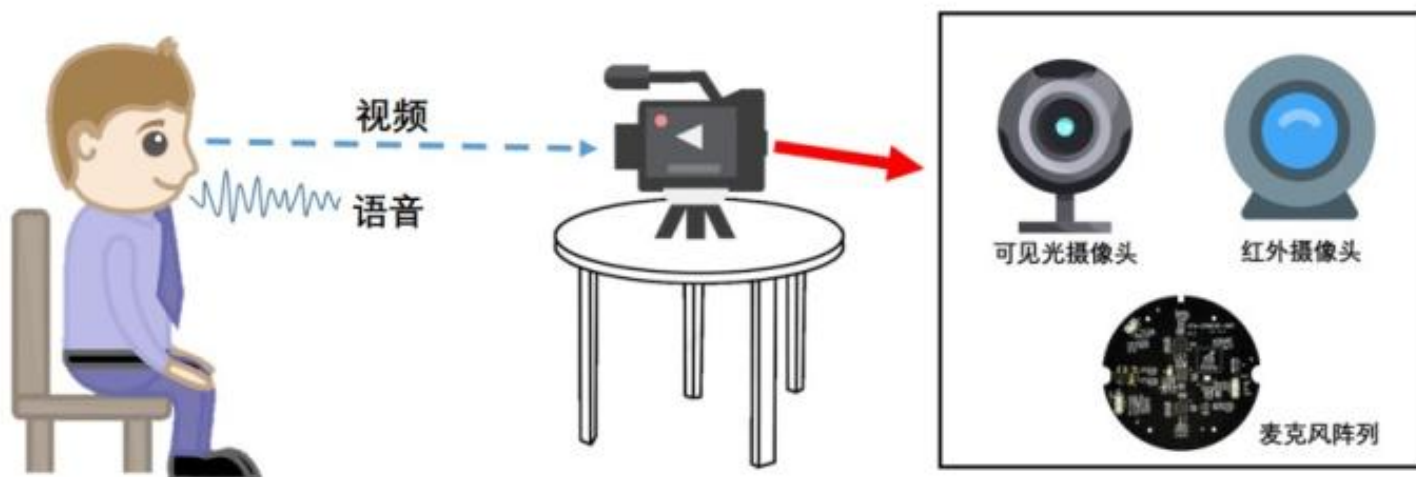


# Thermal\_Speaking\_2021

Shenjingxin  
2021.09

# Collection process



# Randomly generated text(200utt)

≡

1	9	5	4	4	41	6	0	0	2	1	5	81	A	N	K	O	111	O	V	Q	R	A	C	141	搜	索	扬	州	搬	压	机	床		
2	5	6	9	9	42	8	9	9	9	9	2	82	O	N	K	N	112	E	E	X	S	T	H	M	142	我	不	喜	欢	我	小	的		
3	8	2	6	8	43	8	9	2	8	0	7	83	A	J	D	X	113	C	X	Q	Y	B	M	U	143	近	中	国	营	业	厅			
4	5	7	6	7	44	9	9	4	6	7	8	84	P	I	D	H	114	K	E	S	L	R	U	A	144	我	明	大	不	收	拾	你	好	
5	3	2	9	8	45	2	7	6	1	9	5	85	B	Q	M	I	115	E	V	O	L	O	A	W	145	借	八	戒	你	老	公	放	弃	
6	5	7	3	0	46	4	7	4	4	5	4	86	Q	Y	R	P	116	V	J	A	K	W	T	V	146	所	以	我	选	了	放	弃	啊	
7	7	7	4	3	47	3	9	5	1	0	3	87	D	X	Z	O	117	L	R	L	Q	T	V	W	147	家	路	上	小	心	点	啊		
8	2	2	0	6	48	8	8	0	5	0	0	88	I	C	U	K	W	118	S	I	D	O	T	W	148	难	道	不	是	你	的	吗	身	
9	9	5	1	0	49	1	6	3	7	4	5	89	S	J	X	W	F	119	D	O	S	L	L	G	149	现	在	的	我	依	旧	单	身	
					50	5	5	4	7	0	9	90	V	U	C	F	A	120	E	G	Z	L	B	F	150	元	旦	祝	福	有	哪	些	办	
					51	3	9	5	2	9	3	91	V	M	Y	A			121	S	Y	D	D	N	X	151	现	在	想	听	歌	怎	么	办
					52	3	6	2	5	9	8								122	M	B	N	B	R	T	152	我	们	就	这	样	算	了	吧
					53	3	4	6	5	1	4								123	C	J	O	R	W	T	153	我	终	于	失	去	你	了	你
					54	5	6	7	1	2	0															154	其	实	我	也	和	你	一	样
																										155	现	在	油	价	多	少	钱	...

utt\_1-40(40)  
随机的4位数字

utt\_41-80(40)  
随机的6位数字

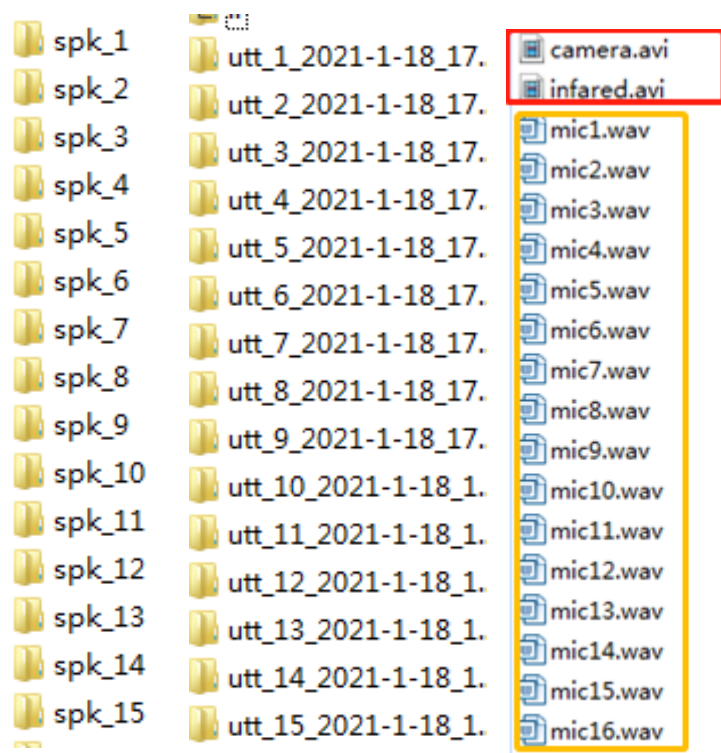
utt\_81-110(30)  
4位随机的英文字母

utt\_111-140(30)  
6位随机的英文字母

Utt\_141-200(60)  
随机的一句中文句子

# Structure of dataset

一共包含247个id: spk\_1,spk\_2,spk\_3.....spk\_252,其中缺少spk\_17、 spk\_73、  
spk\_96、 spk\_242、 spk\_251文件夹



# Preliminary inspection work

红外视频里人脸过大或者过小的问题

1 spk\_8 utt\_195-200  
2 spk\_55 utt\_19  
3 spk\_55 utt\_45

问题一：音频没有录制成功

1 spk\_55 utt\_78  
2 spk\_18 utt\_1-200  
3 spk\_19 utt\_1-200  
4 spk\_21 utt\_135  
5 spk\_30 utt\_80  
6 spk\_38 utt\_2

问题二：出现了其他人的说话声

1 spk\_2 utt\_200  
2 spk\_3 utt\_200  
3 spk\_4 utt\_200  
3 spk\_5 utt\_200  
4 spk\_6 utt\_12  
5 spk\_6 utt\_200  
6 spk\_7 utt\_200  
7 spk\_8 utt\_200  
1 spk\_9 utt\_200  
2 spk\_10 utt\_200  
3 spk\_12 utt\_200  
4 spk\_13 utt\_200  
5 spk\_14 utt\_200  
6 spk\_15 utt\_200  
7 spk\_16 utt\_200  
8 spk\_20 utt\_200  
9 spk\_21 utt\_14  
10 spk\_21 utt\_145-147  
11 spk\_22 utt\_149  
12 spk\_22 utt\_150  
13 spk\_22 utt\_200  
14 spk\_23 utt\_7  
15 spk\_23 utt\_134-140  
16 spk\_23 utt\_147

问题三：非人声的背景噪声过大

1 spk\_2 utt\_174  
1 spk\_13 utt\_108  
2 spk\_21 utt\_126  
3 spk\_21 utt\_127  
4 spk\_22 utt\_81  
5 spk\_22 utt\_92  
6 spk\_22 utt\_112  
7 spk\_22 utt\_168  
8 spk\_23 utt\_6  
9 spk\_23 utt\_9-12  
10 spk\_23 utt\_14-15  
11 spk\_23 utt\_59  
12 spk\_23 utt\_64  
13 spk\_23 utt\_66-70  
14 spk\_23 utt\_94  
15 spk\_23 utt\_96-98  
16 spk\_23 utt\_107  
17 spk\_23 utt\_125-126  
18 spk\_23 utt\_133  
19 spk\_23 utt\_141-145  
20 spk\_23 utt\_151  
21 spk\_23 utt\_155-157  
22 spk\_23 utt\_175-182  
23 spk\_23 utt\_186  
24 spk\_23 utt\_190-191  
25 spk\_23 utt\_198  
26 spk\_24 utt\_60-62  
27 spk\_24 utt\_76-79  
28 spk\_24 utt\_149-150  
29 spk\_24 utt\_159  
30 spk\_24 utt\_165  
31 spk\_24 utt\_167  
32 spk\_25 utt\_16-19

问题四：红外数据中出现了其他非当前说话人的图像

1 spk\_12 utt\_200  
2 spk\_15 utt\_1  
3 spk\_18 utt\_198  
4 spk\_18 utt\_199  
5 spk\_19 utt\_28-32  
6 spk\_19 utt\_35-44  
7 spk\_28 utt\_61-64  
8 spk\_29 utt\_28-30

问题五：少读漏读或重复

1 spk\_1 utt\_42  
2 spk\_4 utt\_167  
3 spk\_4 utt\_177  
1 spk\_10 utt\_134  
2 spk\_10 utt\_150  
3 spk\_13 utt\_3  
4 spk\_15 utt\_160  
5 spk\_16 utt\_164  
6 spk\_16 utt\_165  
7 spk\_16 utt\_168  
8 spk\_16 utt\_170  
9 spk\_16 utt\_179  
10 spk\_20 utt\_150  
11 spk\_20 utt\_152  
12 spk\_20 utt\_180  
13 spk\_21 utt\_159  
14 spk\_21 utt\_126  
15 spk\_21 utt\_173  
16 spk\_22 utt\_128  
17 spk\_22 utt\_133  
18 spk\_22 utt\_160  
19 spk\_23 utt\_48  
20 spk\_23 utt\_62

# First step

## 1. 查找缺失的spk 缺少5个spk\_id

## 2. 音频缺失 无声、打不开、文件缺失

id	缺失utt	缺失utt条数
spk_88	utt_200	1
spk_98	utt_200	1
spk_118	utt_187	1
spk_135	utt_102(没声)	1
spk_209	utt_169-200	32
spk_218	utt_96-200	105
spk_220	utt_186-200	15
spk_223	utt_87(没声)-200	114
spk_227	utt_140-200	61
spk_229	utt_121-200	80
spk_233	utt_151-200	50
spk_243	utt_121-200	80

## 3. 视频缺失

Q1: 文件打不开或者不存在

视频有问题的id	infrared ( 问题红外视频 )	camera ( 问题可见光视频 )
spk_22	utt_31-111	
spk_23	utt_154-200	
spk_26	utt_100-169	
spk_35	utt_140-184	
spk_41	utt_43-109	
spk_47	utt_68-69 ; utt_132-139	
spk_57	utt_22-104	
spk_58	utt_16-26	
spk_63	utt_116-139	
spk_86		utt_1-129
spk_88	utt_200	utt_200
spk_112	utt_74-86	
spk_118	utt_187	utt_187

Q2: 视频过短不足一秒

id	infrared	camera	备注
spk_193	utt_26		转图片只有3张
spk_234	utt_664		转了3张
spk_150	utt_16;utt_67;		只转了2张图片
spk_178	utt_8		转了2张

问题音频: 1117utt

问题视频: 1030utt

# Check the video



camera0



camera1



camera2



camera3



camera4



camera5



camera6



camera7



camera8



camera9



infared0



infared1



infared2



infared3



infared4



infared5



infared6



infared7



infared8



infared9

# Scene classification



scene1



scene2



scene3



scene4



scene5



scene6



scene7



scene8



scene9



# Problem of video

- 1.视频文件无法打开
- 2.视频文件缺失
- 3.视频过短与音频对不上（转成图片仅仅几帧）
- 4.视频中出现其他人

# Check the audio-GMM-HMM

247spk语音总的 utt : 48652, 其中中文 (数字和句子) 34105utt; 英文: 14178utt

1

英文的语音识别模型

英文模型词错误率 (wer) : 2.9%  
英文模型句错误率 (ser) : 4.78%

英文字母合计: 233896个wav文件, 解码的有233892个wav (14960个utt), 针对能够解码的wav, 解码后与参考文本不一致的有11191个wav文件 (1021个utt)。  
解码文本不一致率: 6.8%

2

中文的语音识别模型

中文数字模型词错误率 (wer) : 2.01%  
中文数字模型句错误率 (ser) : 4.03%

中文数字音频合计315904个wav, 其中解码的有315904个wav (19744个utt), 解码后与参考文本不一致的有9824个wav文件 (614个utt)。  
解码文本不一致率: 3.1%

中文汉字模型词错误率 (wer) : 19.17%  
中文汉字模型句错误率 (ser) : 27.90%

中文汉字音频合计229764个wav, 其中解码的有229744个wav (14361个utt), 解码后与参考文本不一致的有62576个wav文件 (3913个utt)。  
解码文本不一致率: 27.2%

结论:

汇总	总数utt	读错的utt	正确的utt	读错率
英文部分	14547	291	14256	2%
中文部分	34105	1297	32808	3.8%

# Problem of the audio

- 1.16路音频不全缺失的(部分缺失; 全部缺失)
- 2.音频存在但打不开或者没声音
- 3.发音与文本是否相符, 是否读错, 重复, 漏读等问题
4. 因噪音过大导致识别错



## Conclusion

可见光视频：247人，48590个视频

红外视频：247人，47863个视频

英文字母音频：244人，14256个utt

中文音频：238人，13661个utt

中文数字：245人，19175个utt

# Experiments of THS2021(DNF)

Dataset : 每个spk的每个utt中取5张图片, 每个人 $200*5=1000$ 张, 随机选50个spk作为测试集, 197个spk作为训练集

# Tsne: 训练过程：上可见光下红外

