

理解非言语交互

陈熙霖

xlchen@ict.ac.cn

中国科学院计算技术研究所

*Look at the person, understand the person for serve the person better
-- A Vision for HCI and Biometrics in Next Decade*

非言语交互 - Jurassic World 2



2023/11/5

2

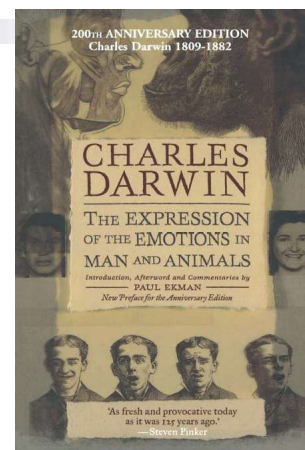
非言语交互

■ 非言语交互的广泛性

- 比言语交互更加稳定而广泛, 甚至在不同类别动物之间也广泛存在
- 人类语言数以百计



Smiling is an universal pass for human society



2023/11/5

3

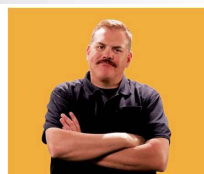
非言语交互传达的信息



矛盾(Conflicting)



强调(Accenting)



替代(Substituting)

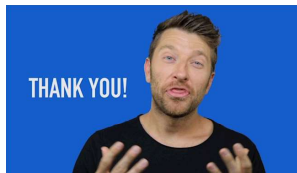
对言语交互的支持(一致)程度

非言语交互通道

- 面部
 - 表情
 - 视线
 - 唇动
- 身体
 - 体势
 - 手势
- ...



重复(Repeating)



补充(Complementing)



调控(Regulating)

2023/11/5

4

人-机无缝交互--AI



2023/11/5

5

面部表情理解

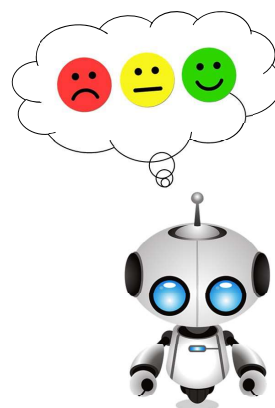
■ 面部表情识别



Basic emotions

Upper Face Action Units					
AU1	AU2	AU4	AU5	AU6	AU7
Inner Brow Raiser	Outer Brow Raiser	Brow Lowerer	Upper Lid Raiser	Cheek Raiser	Lid Tightener
*AU41	*AU42	*AU43	AU44	AU45	AU46
Lip Droop	Slit	Eyes Closed	Squint	Blink	Wink
Lower Face Action Units					
AU9	AU10	AU11	AU12	AU13	AU14
Nose Wrinkler	Upper Lip Raiser	Nasolabial Deepener	Lip Corner Puller	Cheek Puffer	Dimpler
AU15	AU16	AU17	AU18	AU20	AU22
Lip Corner Depressor	Lower Lip Depressor	Chin Raiser	Lip Pucker	Lip Stretcher	Lip Funneler
AU23	AU24	*AU25	*AU26	*AU27	AU28
Lip Tightener	Lip Pressor	Lips Parts	Jaw Drop	Mouth Stretch	Lip Suck

Facial Action Coding System



2023/11/5

6

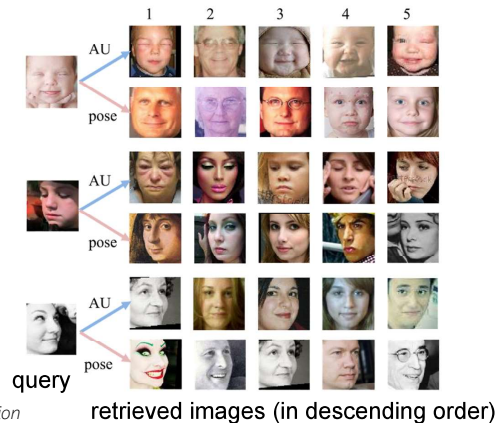
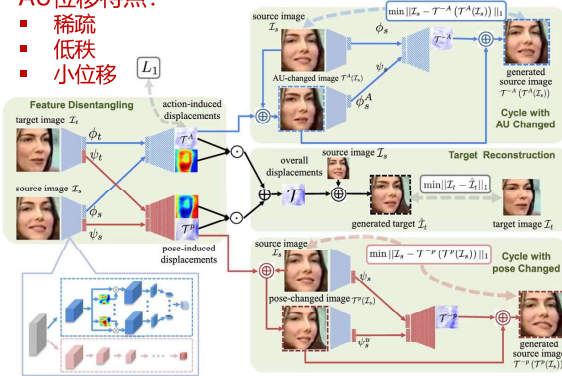
面部表情理解

■ 无监督面部动作表示学习

- Disentangle facial actions and head poses

AU位移特点:

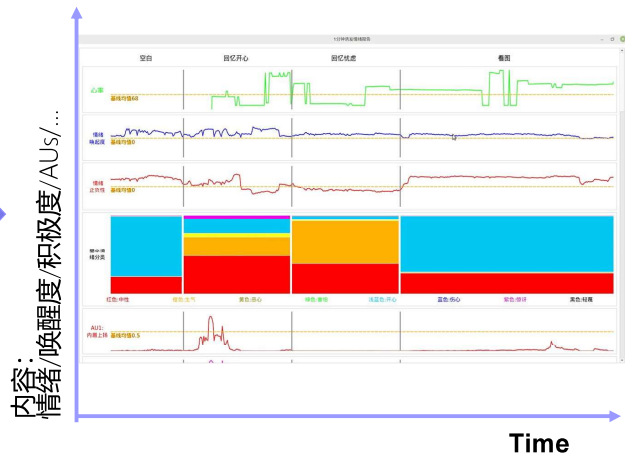
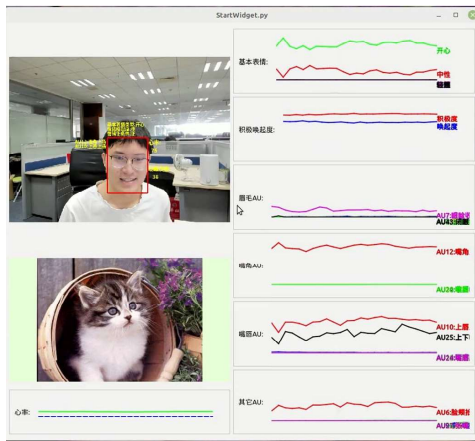
- 稀疏
- 低秩
- 小位移



Yong Li, Jiabei Zeng, Shiguang Shan, Xilin Chen. Self-supervised Representation Learning from Videos for Facial Action Unit Detection. CVPR 2019

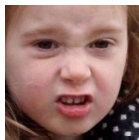
面部表情理解

■ 通过面部表情统计来理解



2023/11/5 PsychoFace V1.0 developed by ICT, CAS.

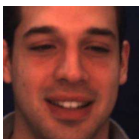
从视觉-语言模型理解面部表情



The corners of the lips pulled upwards, wrinkles near the eyes, raised eyebrows, slightly widened or protruding eyes, and slight upwards or outward pushing of the lips indicate a combination of **anger** and **disgust**. Additionally, the raised upper lip or chin, slightly protruding lower lip, and flared nostrils all further intensify the emotion of **anger** that is evident in this face.



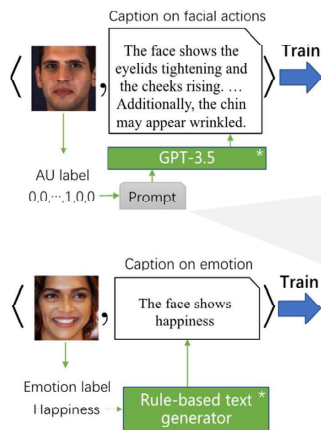
The slightly raised and furrowed eyebrows, along with the creased skin around the eyes, indicate a state of **alertness** or **concentration**. These actions also suggest that the individual is feeling **anxious** or **frustrated**. Additionally, the slight squint or furrow of the forehead and wrinkles on the skin around the eyes, as well as the raised eyebrows, further infer a feeling of **sadness**.



The corners of the lips pulled upwards and outward, the lower eyelids lifted, and the cheeks slightly elevated indicate **happiness**, while the tightened lower eyelid and horizontally stretched lips suggest a genuine sense of **joy**.

从视觉-语言模型理解面部表情

Step 1: Synthesize the training image-text pairs



Step 2: Train AU-BLIP and Emot-BLIP to describe on single aspect

Label: AU10, AU14, AU17, AU24

one-shot prompt with GPT-3.5

<List all 27 Action Units (AUs) defined by the Facial Action Coding System (FACS)>
You are currently acting as an AU description expert, and you should answer questions according to this example:

Question: "Describe briefly the facial actions or status of a face which contains AU2, AU4, AU9, AU20, AU25 in 1-9 sentences. Do not describe the emotion or facial expressions, and do not mention 'AU' in the reply."

Answer: "The face shows the outer eyebrows raised, forehead wrinkles, a wrinkled nose, a stretched lip, and either parted lips or a dropped jaw."

Question: Describe briefly the facial actions or status of a face which contains AU10, AU14, AU17, AU24 in 1-9 sentences. Do not describe the emotion or facial expressions, and do not mention 'AU' in the reply.

Answer:

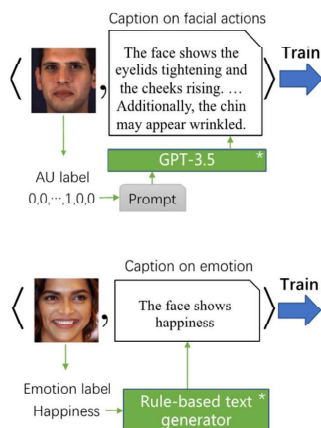
(output) The face shows an upward movement of the upper lip, dimpling at the corners of the mouth, a slight upward pull of the chin, and a pressing together of the lips.

2023/11/5 Yujian Yuan, Jiabei, Shiguang Shan. Describe Your Facial Expressions by Linking Image Encoders and Large Language Models. BMVC2023.

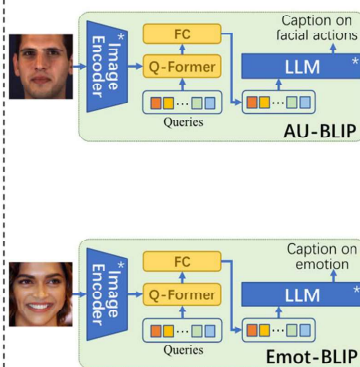
10

从视觉-语言模型理解面部表情

Step 1: Synthesize the training image-text pairs



Step 2: Train AU-BLIP and Emot-BLIP to describe on single aspect



Step 3: Train Exp-BLIP to describe on both aspects

2023/11/5 Yujian Yuan, Jiabei, Shiguang Shan. Describe Your Facial Expressions by Linking Image Encoders and Large Language Models. BMVC2023.

11

从视觉-语言模型理解面部表情

zero-shot prompt

Combine the facial action sentence with the emotion sentence. This combined sentence should describe all the facial actions and emotions mentioned in these two sentences by pointing out how each emotion is inferred from the corresponding facial actions.

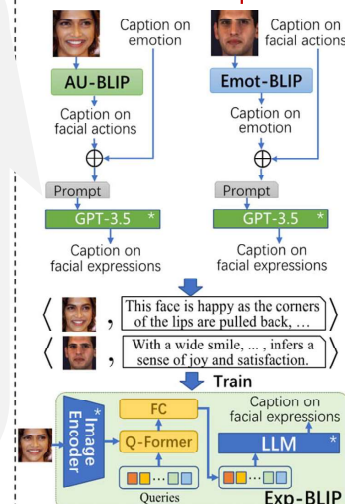
Facial Action Sentence: "The corners of the lips are pulled towards the ears, resulting in a widening of the mouth. The lower lip may also be slightly pushed forward or pouted."

Emotion Sentence: "The actions on the face show relief, contentment and happiness."

output by GPT-3.5

The facial actions, including the pulling of the corners of the lips towards the ears resulting in a wider mouth and the slight pushing forward or pouting of the lower lip, suggest a combination of relief, contentment, and happiness emotions. The pulling of the corners of the lips upwards towards the ears indicates a sense of happiness or joy, while the slight pushing forward or pouting of the lower lip suggests a degree of contentment.

Step 3: Train Exp-BLIP to describe on both aspects



2023/11/5

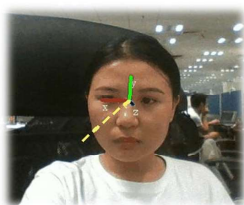
12

视线(Gaze)

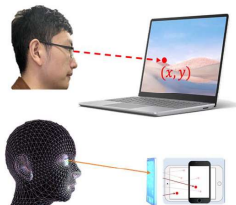


The eyes are the windows to the soul.

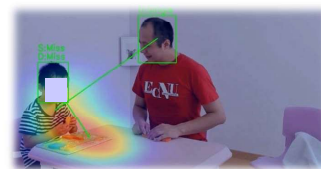
---An old proverb



3D视线指向



屏幕上的关注点



交流中的视线跟踪

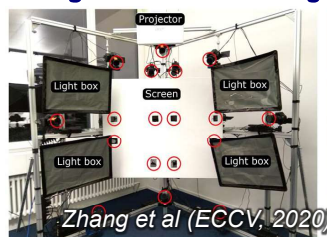
2023/11/5

13

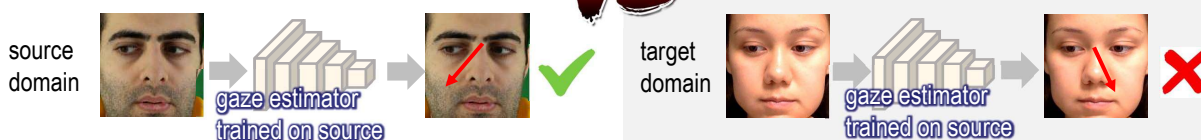
视线估计

挑战: 不同环境的影响

During the data collecting (source domain)



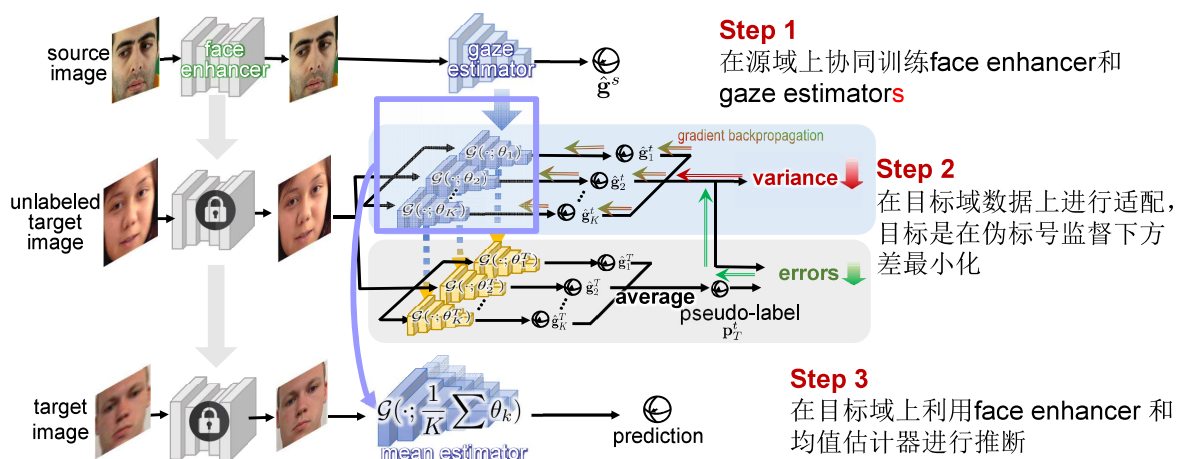
In real world application (target domain)



2023/11/5

14

视线估计



实验结果对比

Method	$\mathcal{D}_E \rightarrow \mathcal{D}_M$	$\mathcal{D}_E \rightarrow \mathcal{D}_D$	$\mathcal{D}_G \rightarrow \mathcal{D}_M$	$\mathcal{D}_G \rightarrow \mathcal{D}_D$
Only Source	7.50	7.88	7.23	8.02
w/o source				
PureGaze [2]	7.08	7.48	9.28	9.32
PnP-GA(oma) [4]	5.65	-	6.86	-
CSA [6]	5.37	6.77	7.30	7.73
RUDA [1]	5.70	6.29	6.20	5.86
w/ source				
Gaze360 [3]	5.97	7.84	7.38	9.61
GazeAdv [5]	6.75	8.10	8.19	12.27
PnP-GA [4]	5.53	5.87	6.18	7.92
CRGA [6]	5.68	5.72	6.09	6.68
UnReGA ⁻	5.35	6.06	5.58	5.84
UnReGA	5.11	5.70	5.42	5.80

[1]. Bao et al., CVPR 2022
 [2]. Cheng et al., AAAI 2022
 [3]. Kellnhofer, et al., CVPR 2019
 [4]. Liu et al., ICCV 2021
 [5]. Wang et al., CVPR 2019
 [6]. Wang et al., CVPR 2022
 UnReGA: proposed method
 UnReGA⁻: UnReGA w/o face enhancer

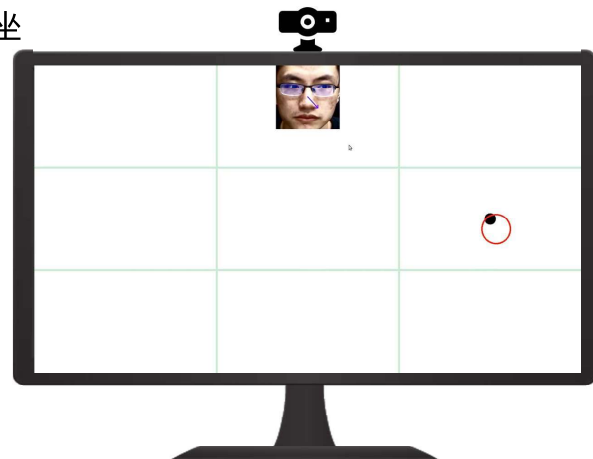
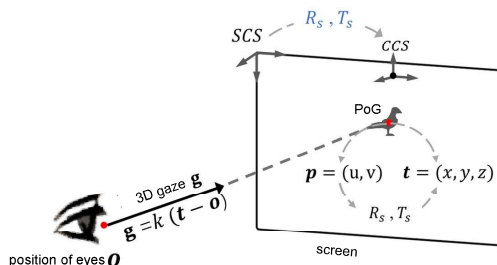
2023/11/5 \mathcal{D}_E - ETH-XGaze, \mathcal{D}_G - Gaze360, \mathcal{D}_M - MPIIGaze, \mathcal{D}_D - EyeDiap

16

屏幕上视线点(Point of Gaze, PoG)

■ 估计PoG

- 需要从摄像机坐标系(CCS)到屏幕坐标系(SCS)的变换
- CCS下的空间视线方向
- CCS下的眼睛的位置



Realtime PoG estimation demo system developed by ICT, CAS.

2023/11/5

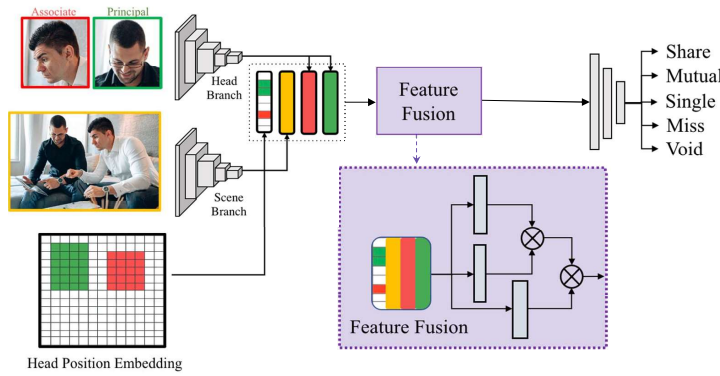
17

双向视线交流模式

■ 覆盖所有可能的个体与其伙伴的可能视线状态



双向视线交流模式



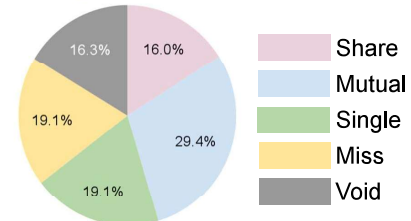
数据集

数据源

UCO-LAEO和
VACATION中的部分数据

370个视频片段
每个片段3-15秒

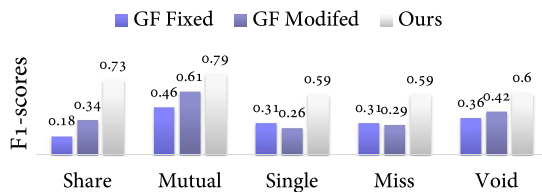
新的标注: 5 patterns



2023/11/5

19

双向视线交流模式



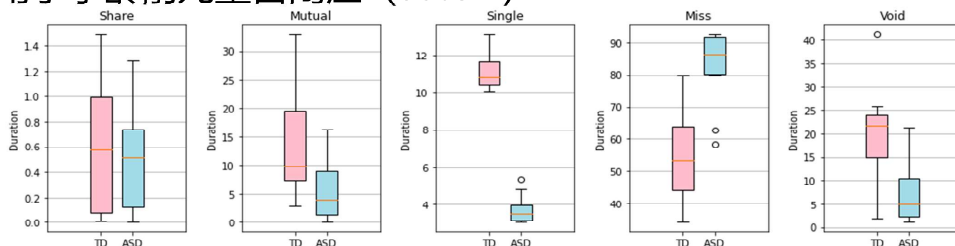
Method	Looking-At-Each-Other (AP)			Share (Acc.) VideoCoAtt
	UCO-LAEO	AVA-LAEO	OF-MG	
LAEO-Net[Marin-Jiménez et al. 2022]	79.5	50.6	-	-
AAAI'21[Doosti et al. 2021]	65.1	72.2	70.1	-
CVPR'18[Fan et al. 2018]	-	-	-	71.4
Ours	80.3	82.5	72.1	73.9

2023/11/5

20

双向视线交流模式

用于学龄前儿童自闭症 (autism)



Null hypothesis

H_{1_0} : the duration of 'Share' pattern is the same between children with and without autism
 H_{2_0} : the duration of 'Mutual' pattern is the same between children with and without autism
 H_{3_0} : the duration of 'Single' pattern is the same between children with and without autism
 H_{4_0} : the duration of 'Miss' pattern is the same between children with and without autism
 H_{5_0} : the duration of 'Void' pattern is the same between children with and without autism

t-statistics p-value

-0.46 0.66
-2.12 0.048 (*)
-19.00 0.0000 (***)
4.54 0.0000 (***)
-3.07 0.006 (**)

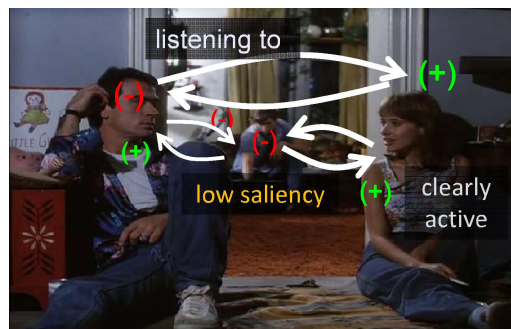
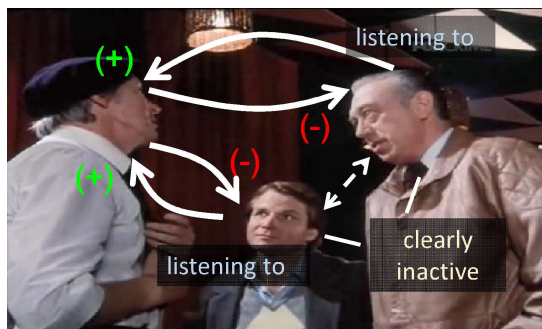
2023/11/5

21

理解视线之外的交互

■ 多方交互

- 谁是交流中的主角?
- 谁正在谈话?
-



2023/11/5

22

理解视线之外的交互

■ 主动说话人检测

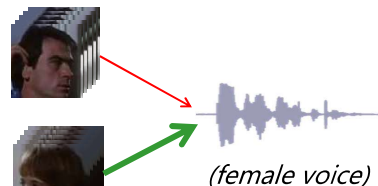
- 从上下文线索进行判别



场景布局与
说话人的位置
(空间上下文)



视觉注意力
(关系上下文)

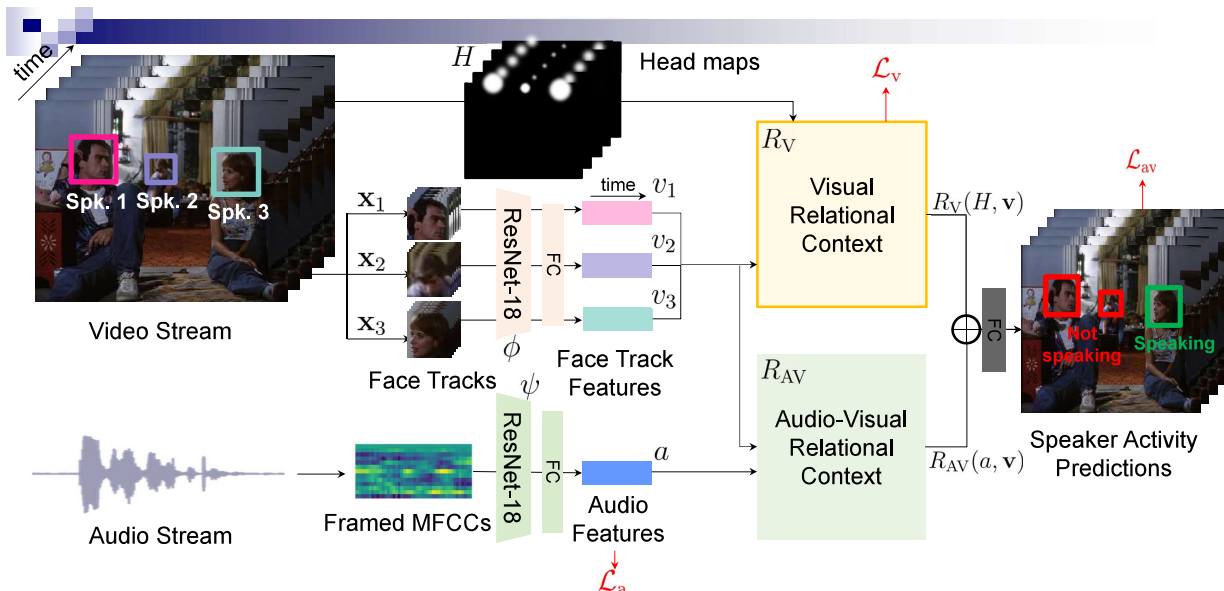


模态匹配对比
(关系上下文)

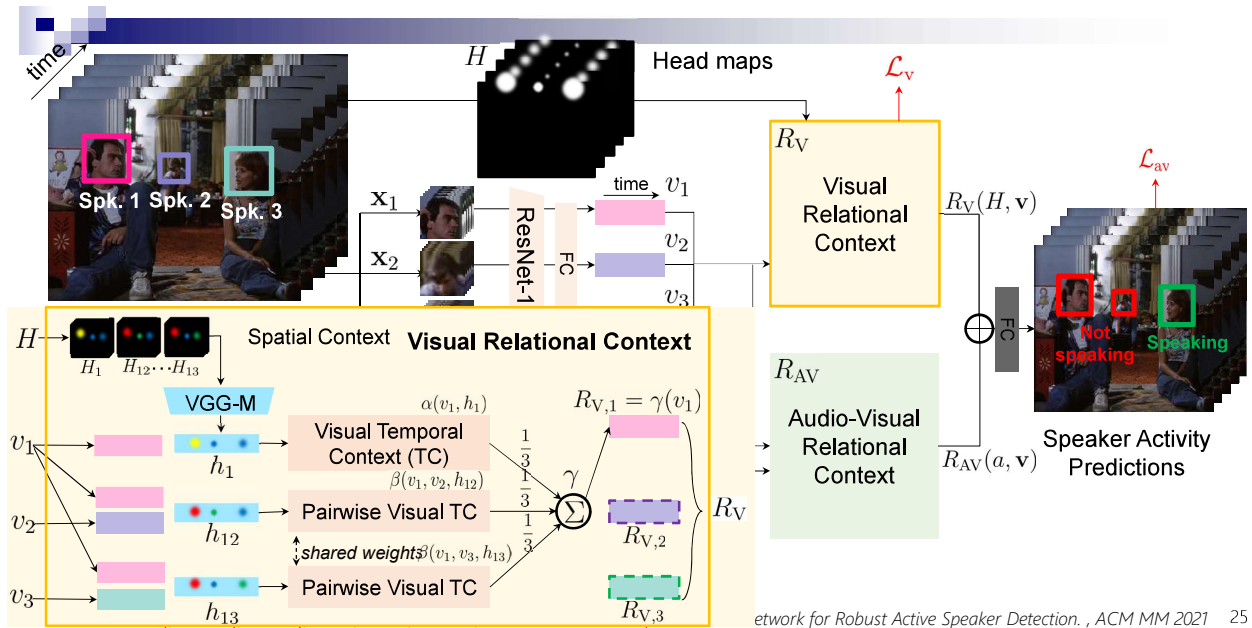
2023/11/5

23

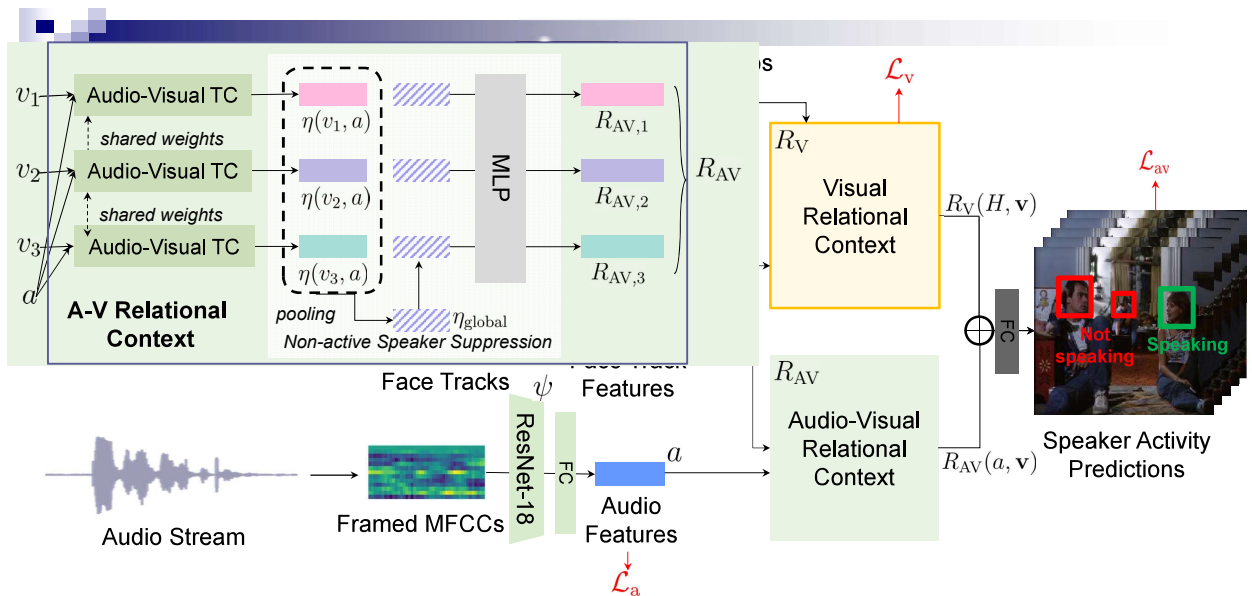
理解视线之外的交互



理解视线之外的交互



理解视线之外的交互



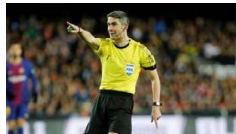
2023/11/5 Y. Zhang, S. Liang, S. Yang, X. Liu, Z. Wu, S. Shan, X. Chen. UniCon: Unified Context Network for Robust Active Speaker Detection. , ACM MM 2021 26

理解视线之外的交互



手势/姿势作为交互手段

- 手势/姿态传达了丰富的关于指令、情绪、态度等信息



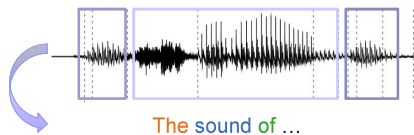
Passionate Sign Language Interpreter At Rock Gig
https://www.youtube.com/watch?v=DYoB_A8GZ08

2023/11/5

28

手势/手语识别的挑战

- 自由手势的任意性与差别性
- 手语受限于细粒度标注的困难和数据的稀缺



The sound of ...

ASR dataset scales
LibraSpeech (1000 hours)
Whisper (680,000 hours)

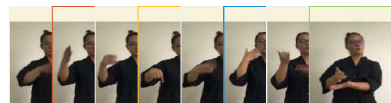


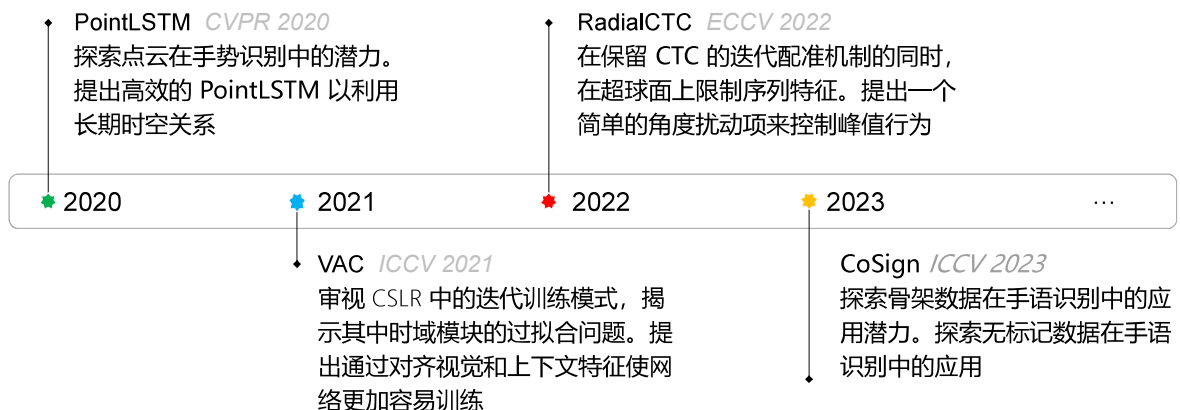
Table under cat is

SLR dataset scales
Phoenix14 (12.5 hours)
CSL-Daily (23.3 hours)

2023/11/5

29

近期的一些工作



2023/11/5

30

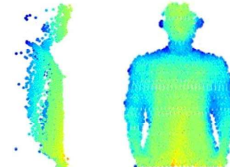
选择点云的原因

■ 选择手势识别的输入模式

- 视频
 - 计算成本高
 - 复杂的背景和照明影响
- 点云
 - + 高效且稳定
 - 难以找到对应点



RGB Video



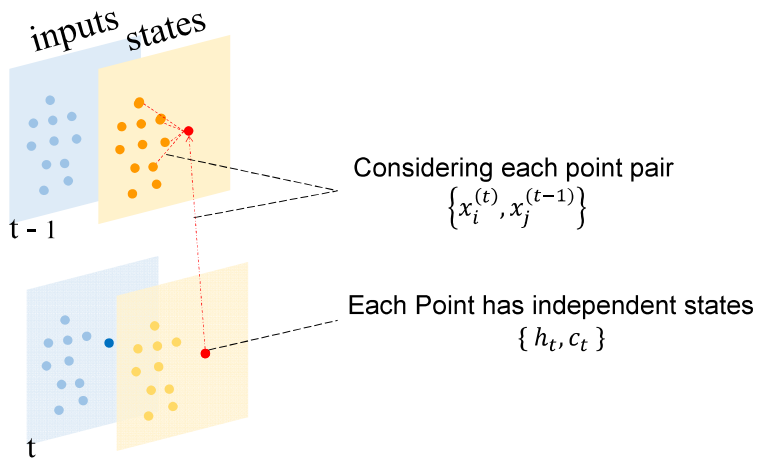
Point Cloud

2023/11/5 Yuecong Min, Yanxiao Zhang, Xiujuan Chai, Xilin Chen, An Efficient PointLSTM for Point Clouds based Gesture Recognition, CVPR 2020

31

PointLSTM

■ 查找过往邻点的对应点

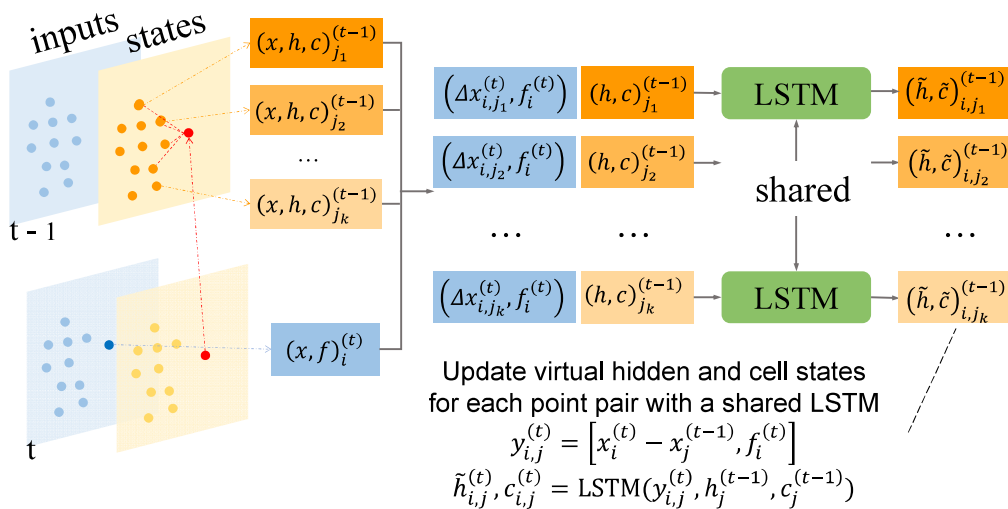


2023/11/5

32

PointLSTM

■ 查找过往邻点的对应点

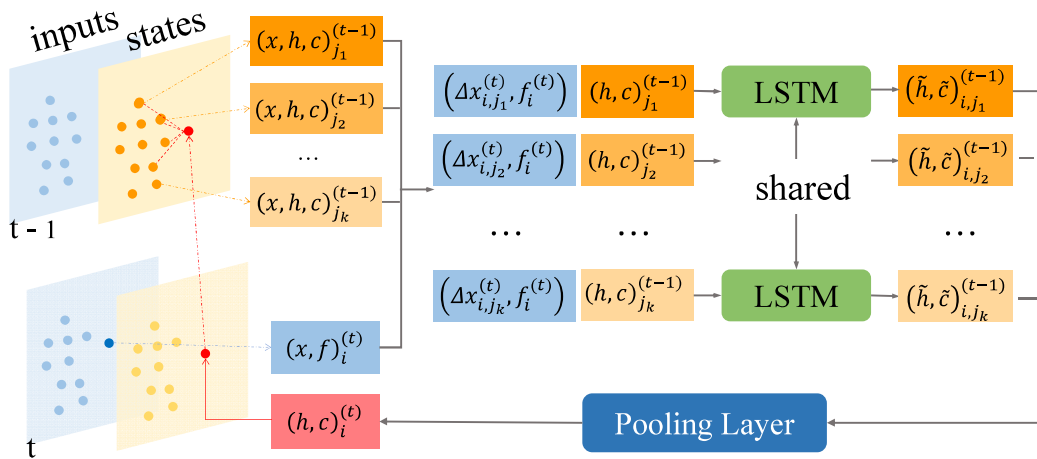


2023/11/5

33

PointLSTM

- 查找过往邻点的对应点



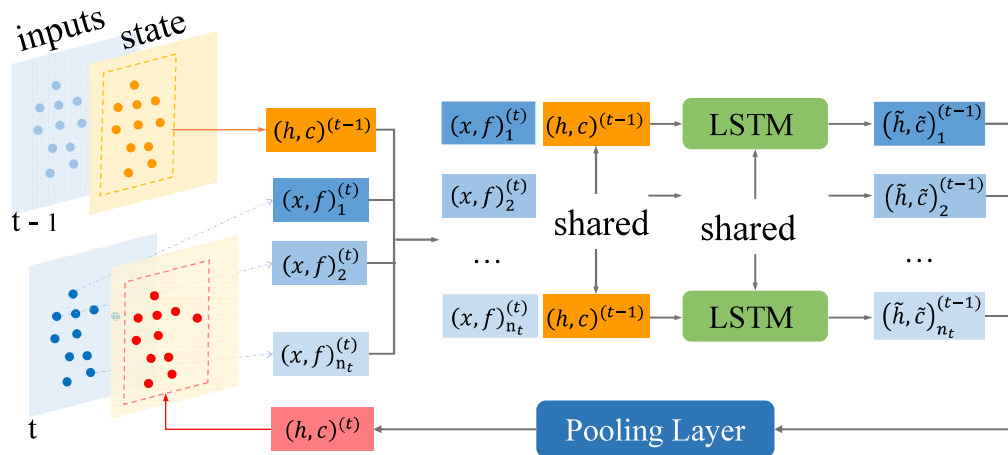
2023/11/5

The final states are obtained through a pooling layer

34

PointLSTM-PSS

- 点共享状态 (h_t, c_t) 简化版

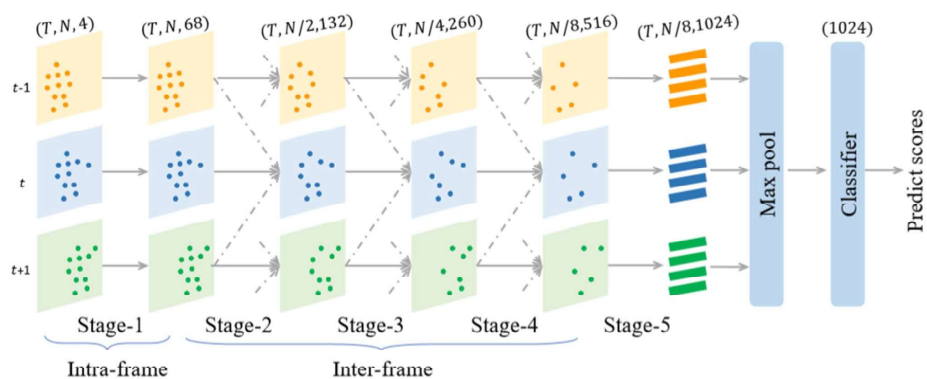


2023/11/5

35

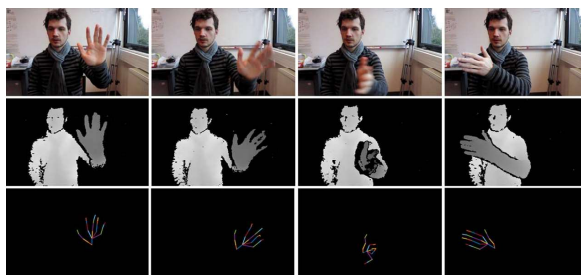
网络结构

- Baseline



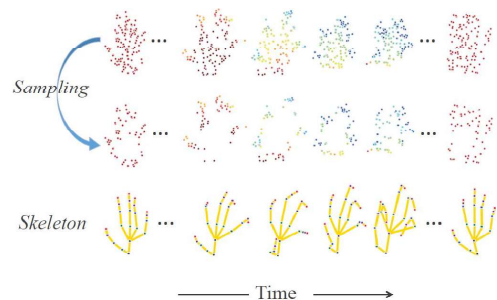
SHREC'17上的实验

- 14 gestures × one finger or whole hand × 28 participants



<http://www-rech.telecom-lille.fr/shrec2017-hand/#gestures>

Sampling 64 points from 128 points for each frame

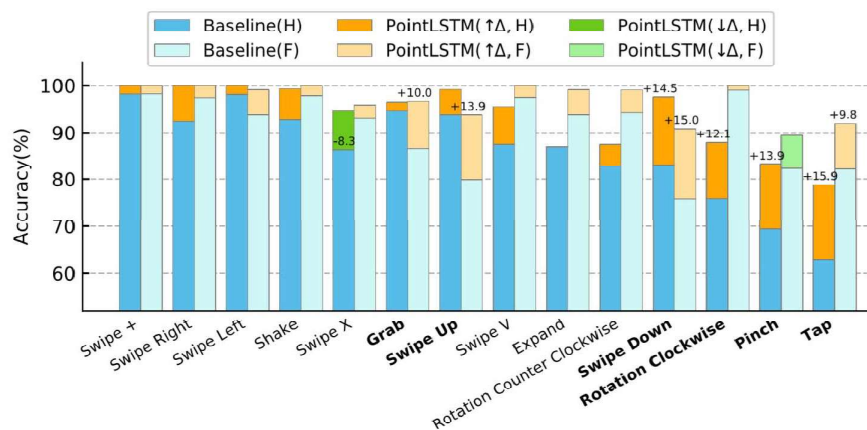


2023/11/5

37

SHREC'17上的实验

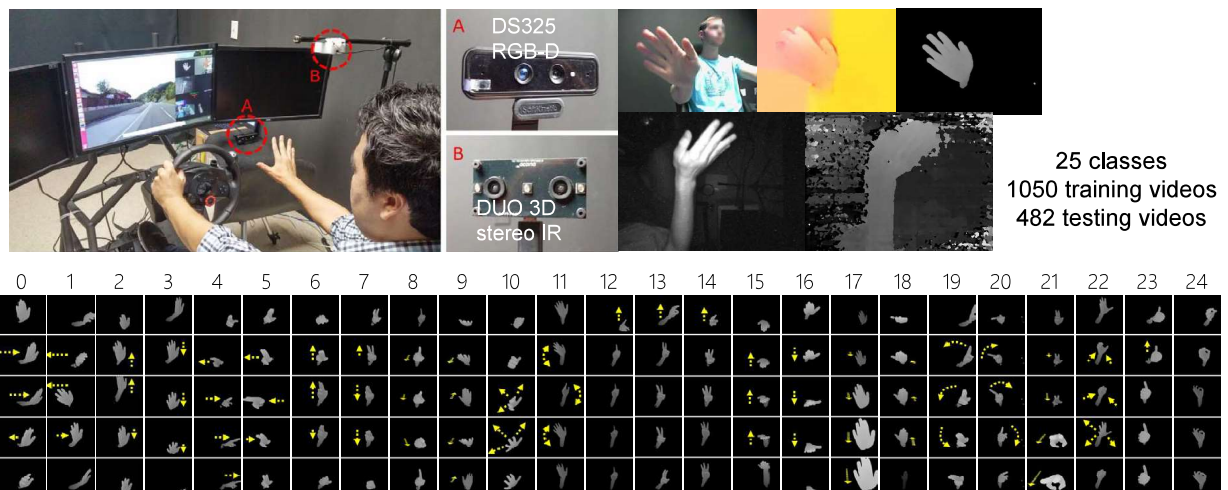
- PointLSTM-middle (94.70%) vs. Baseline (88.90%)



2023/11/5

38

NvGesture上的实验



P. Molchanov, X. Yang, S. Gupta, K. Kim, S. Tyree, & J. Kautz. Online detection and classification of dynamic hand gestures with recurrent 3d convolutional neural network. CVPR 2016.

2023/11/5

39

NvGesture上的实验结果

Method	Modality	Accuracy
R3DCNN [1]	IR image	63.5%
R3DCNN [1]	Optical Flow	77.8%
R3DCNN [1]	Depth Video	80.3%
PreRNN [2]		84.4%
MTUT [3]		84.9%
R3DCNN [1]	RGB Video	74.1%
PreRNN [2]		76.5%
MTUT [3]		81.3%
PointNet++ [4]	Point Clouds	63.9%
FlickerNet [5]		86.3%
Baseline		85.9(± 0.5)%
PointLSTM-early		87.9(± 0.7)%
PointLSTM-PSS		87.3(± 0.4)%
PointLSTM-middle		86.9(± 0.6)%
PointLSTM-late		87.5(± 1.0)%
Human	RGB Video	88.4%

- [1] P. Molchanov, X. Yang, S. Gupta, K. Kim, S. Tyree, & J. Kautz. Online detection and classification of dynamic hand gestures with recurrent 3d convolutional neural network. CVPR 2016.
- [2] X. Yang, P. Molchanov, & J. Kautz. Making convolutional networks recurrent for visual sequence learning. CVPR 2018.
- [3] M. Abavisani, H. Reza, V. Joze, & V. Patel. Improving the performance of unimodal dynamic handgesture recognition with multimodal training. ECCV 2019.
- [4] C. Rui, Z. Qi, L. Yi, H. Su, & L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. NeurIPS 2017.
- [5] Y. Min, X. Chai, L. Zhao, & X. Chen. Flickernet: Adaptive 3d gesture recognition from sparse point clouds. BMVC 2019.

2023/11/5

40

从手势到行为动作 – MSR Action 3D



MSR Action3D*

- 20 classes \times 10 subjects
- Cover various movements of arms, legs, torso and their combinations
- 567 sequences
- 20 joint locations

方法	输入模态	精度
Gram Handkel [1]	skeleton	94.74%
MeteorNet [2]	Point Clouds	88.50%
Baseline		87.62(± 1.48)%
PointLSTM-early		91.78(± 3.10)%
PointLSTM-PSS		90.79(± 3.14)%
PointLSTM-middle		91.08(± 3.43)%
PointLSTM-late		92.29(± 3.09)%

- [1] X. Zhang, Y. Wang, M. Gou, M. Sznai, & O. Camps. Efficient temporal sequence comparison and classification using gram matrix embeddings on a Riemannian manifold. CVPR 2016.
- [2] X. Liu, M. Yan, & J. Bohg. MeteorNet: Deep learning on dynamic 3d point cloud sequences. CVPR 2019.

2023/11/5

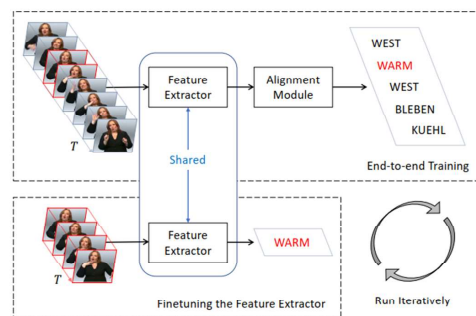
* Wanqing Li, Zhengyou Zhang, and Zicheng Liu. Action recognition based on a bag of 3d points. CVPRW 2010.

41

视频手语识别

■ 捕捉精细的视觉线索

- SLR 模型很难捕捉到正确的视觉线索
 - 迭代训练方案在 SLR 中得到广泛应用



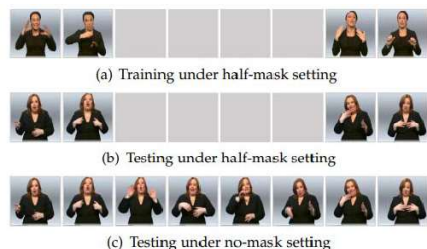
(a) Iterative training scheme.

2023/11/5

42

捕捉精细的视觉线索

- SLR 模型很难捕捉到正确的视觉线索
 - 迭代训练方案在 SLR 中得到广泛应用
- SLR 模型更倾向于根据上下文信息进行预测
 - SLR 模型即使在屏蔽一半序列的情况下也能拟合出训练集



Method	Evaluation setting	WER on Training set (↓)	WER on Dev set (↓)
Baseline	half-mask	17.5	59.5
Baseline	no-mask	31.6	49.6
Baseline+VAC	half-mask	22.2	50.1
Baseline+VAC	no-mask	16.8	29.9

2023/11/5

43

利用 VAC(Visual Alignment Constraints) 实现端到端训练

- 利用视觉特征和上下文特征之间的对应关系
- 对视觉特征引入辅助分类器
- Visual Enhancement Constraint (VEC)

$$L_{VE} = L_{CTC}^v = -\log p(l|x; \theta^v)$$

- Visual Alignment Constraint (VAC)

$$L_{VA} = \text{KL} \left(\text{softmax} \left(\frac{Z}{\tau} \right), \text{softmax} \left(\frac{\hat{Z}}{\tau} \right) \right)$$

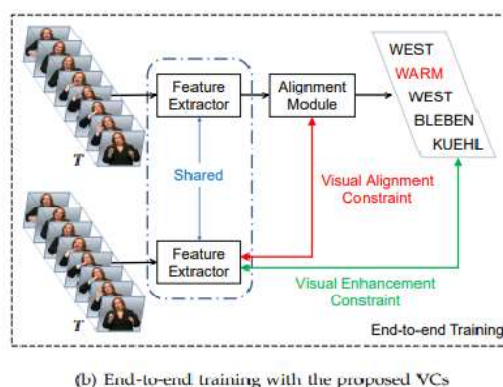
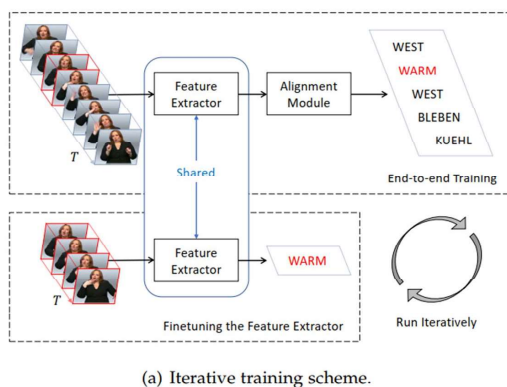
- 更简化版本的 VAC:
共享辅助分类器和主分类器的权重

2023/11/5

44

利用VAC实现端到端训练

- 迭代训练 vs. 使用 VAC 的端到端训练



Phoenix14的实验

■ 电视上的'Real-life'采集

- 9 signers, 5.6K training samples (10.7 hours), 1081 classes



2023/11/5

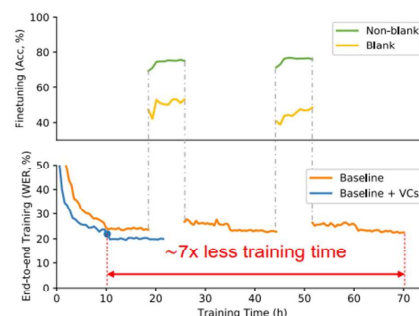
46

Phoenix14上的实验结果

■ 利用VAC的端到端训练 vs. 迭代训练模式

- 训练效率更高, 精度更高

Iterations	Constraint	Dev	Test
1	-	23.4	24.7
2	-	22.6	22.7
3	-	22.3	23.0
-	-	23.1	24.2
-	VEC	20.7	21.0
-	VEC & VAC	19.6	19.7



2023/11/5

47

Phoenix14上的实验结果

Method	Iteration	Dev	Test
SubUNet [1]		40.8	40.7
Re-Sign [2]	✓	27.1	26.8
CNN+LSTM+HMM [3]	✓	26.0	26.0
FCN [4]		23.7	23.9
DNF [5]	✓	23.1	22.9
CMA [6]	✓	21.3	21.9
STMC [7]	✓	21.1	20.7
Baseline		22.7	23.5
Baseline+VAC		19.6	19.7

[1] N. C. Camgoz, S. Hadfield, O. Koller, R. Bowden. "Subunets: End-to-end hand shape and continuous sign language recognition." ICCV, 2017.

[2] O. Koller, S. Zargaran, H. Ney. "Re-sign: Re-aligned end-to-end sequence modelling with deep recurrent CNN-HMMs." CVPR, 2017.

[3] O. Koller, N. C. Camgoz, H. Ney, R. Bowden. "Weakly supervised learning with multi-stream CNN-LSTM-HMMs to discover sequential parallelism in sign language videos." TPAMI, 2019.

[4] K. L. Cheng, Z. Yang, Q. Chen, Y.-W. Tai. "Fully Convolutional Networks for Continuous Sign Language Recognition." ECCV, 2020.

[5] R. Cui, H. Liu, C. Zhang. "A deep neural framework for continuous sign language recognition by iterative training." TMM, 2019.

[6] J. Pu, W. Zhou, H. Hu, H. Li. "Boosting Continuous Sign Language Recognition via Cross Modality Augmentation." ACM MM, 2020.

[7] H. Zhou, W. Zhou, Y. Zhou, H. Li. "Spatial-temporal multi-cue network for continuous sign language recognition, AAAI 2020.

2023/11/5

48

Phoenix14上的实验结果



Pay more attention on facial expression with VAC

2023/11/5

49

几点思考

- 非言语通道的理解对于未来自然的人-机共生环境是极其重要的
- 与言语交流相比，非言语互动的不确定性更大，更具挑战性
- 非言语交流应成为生物特征识别和人机交互领域的关注焦点

*Look at the person, understand the person for serve the person better
 -- A Vision for HCI and Biometrics in Next Decade*

2023/11/5

感谢我的同事和学生

- 山世光
- 曾加贝
- 杨双
- 何明捷
- 柴秀娟
- 闵越聪
- 李勇
- 蔡昕
- 常菲
- 张远航

2023/11/5

51

Thanks for your attention

*Look at the person, understand the person for serve the person better
-- A Vision for HCI and Biometrics in Next Decade*