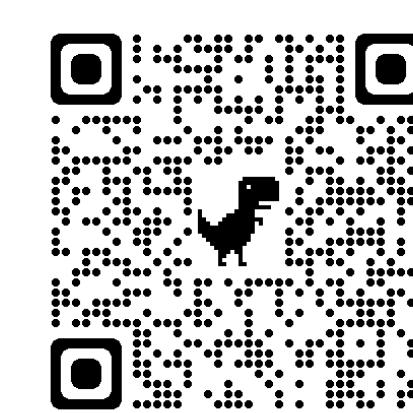


PVChat: Personalized Video Chat with One-Shot Learning

Yufei Shi^{1*} Weilong Yan^{2*} Gang Xu³ Yumeng Li⁴ Yucheng Chen¹ Zhenxi Li¹Fei Yu³ Ming Li^{3†} Si Yong Yeo^{1†}¹Lee Kong Chian School of Medicine, NTU Singapore ²National University of Singapore³Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ) ⁴Nankai University

* Equal Contribution † Corresponding author

Code: <https://github.com/DavidYan2001/PVChat>Email: yufei005@e.ntu.edu.sg

Problem Definition and Contributions

Goal: Enable one-shot, identity-aware video understanding and QA from a single reference video, generalizing across scenes and multi-subject settings

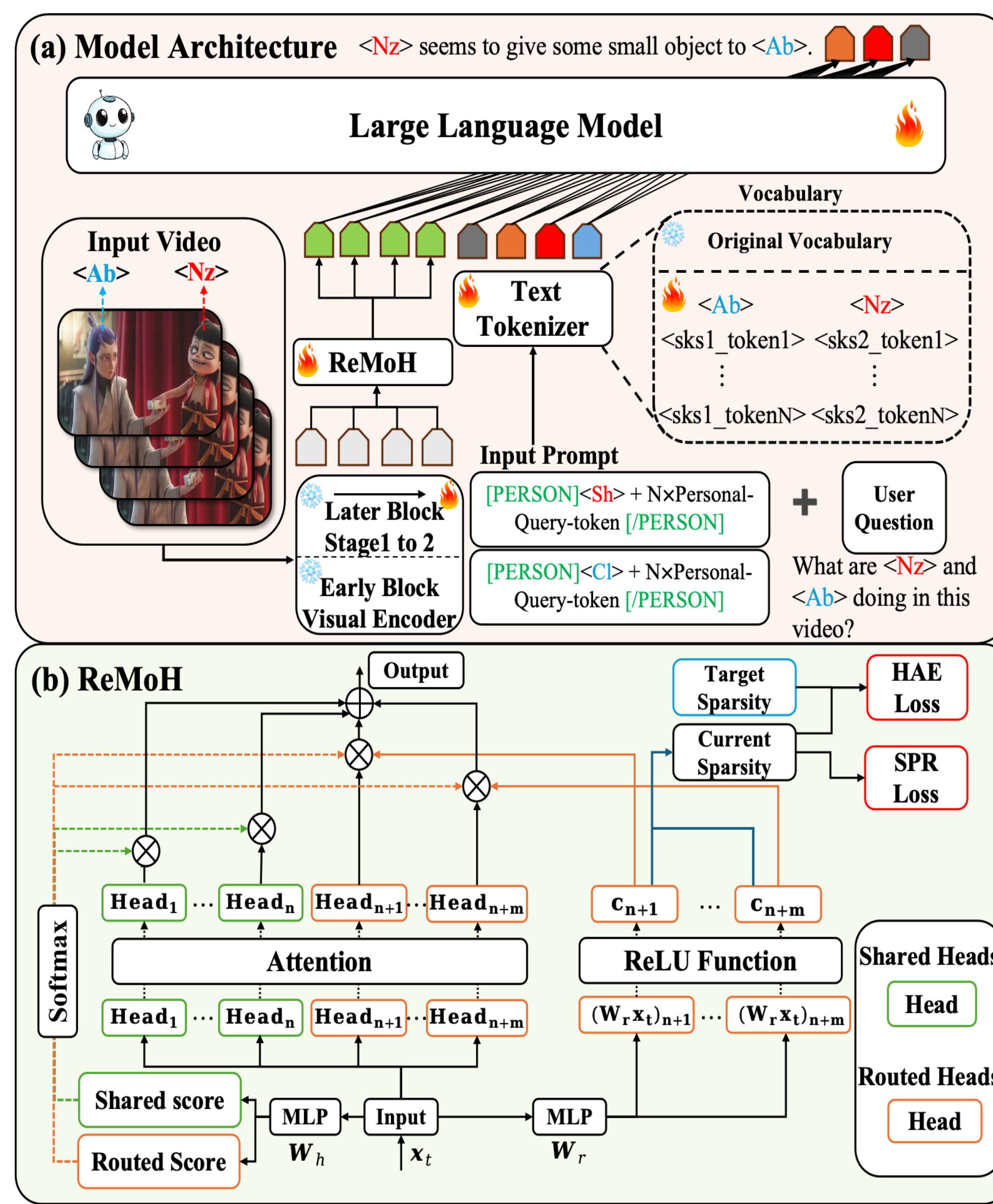
Motivation:

- General-purpose ViLMMs lack reliable person-specific reasoning due to non-personalized training.
- Prior personalization is image-centric, missing motion/temporal cues essential for video.

Key Contributions:

- PVChat: First ViLMM to personalize from a single reference video, enabling one-shot, identity-aware video understanding and user-specific QA.
- Data pipeline: Single-clip augmentation and QA generation using off-the-shelf tools to create identity-preserving positives and look-alike hard negatives; we release a diverse dataset (6 scenarios, 304 source videos, 2,304 generated videos, 30k+ QA pairs).
- Modeling: ReLU-Routing Mixture-of-Heads plus Smooth Proximity Regularization and Head Activation Enhancement for stable training and stronger extraction of individual-specific cues.

Model and ReMoH (ReLU Routing Mixture-of-Heads Attention)



SPA and HAE Loss

SPA (Smooth Proximity Regularization)

- To encourage the sparsity of activation, which helps to reduce computational cost and boost specific-domain learning. T_s and R_s is the Target and current sparsity.

$$\mathcal{L}_{SPR} = \beta_p \cdot \mathcal{L}_{Reg}$$

$$\beta_{p+1} = \beta_p \cdot e^{k \cdot (T_s - R_s)}$$

$$\mathcal{L}_{Reg} = \left\| \frac{1}{n} (\mathbf{W}_r \mathbf{x}_t) \right\|$$

HAE (Head Activation Enhancement)

- L_{Reg} just focuses on increasing the sparsity (usually makes all output prone to 0). we design a HAE Loss to activate more heads to avoid some experts always being asleep

$$\mathcal{L}_{HAE} = e^{2 \cdot (R_s - T_s)} - 1 \quad \text{if } R_s > T_s.$$

- L_{HAE} helps some heads to be more active when they fall into a state of zero.

Final Loss

$$\mathcal{L} = \mathcal{L}_{LM} + \mathcal{L}_{SPR} + \mathcal{L}_{HAE}.$$

\mathcal{L}_{LM} is the cross-entropy loss.

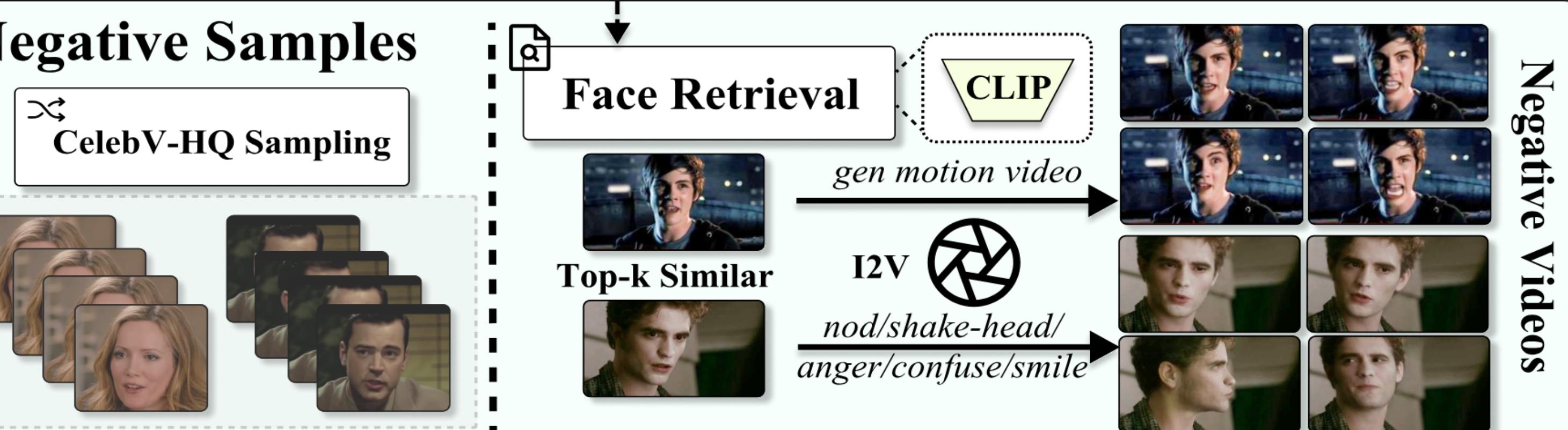
Data Pipeline and Experimental Results

Video Expansion Pipeline

Positive Samples



Negative Samples

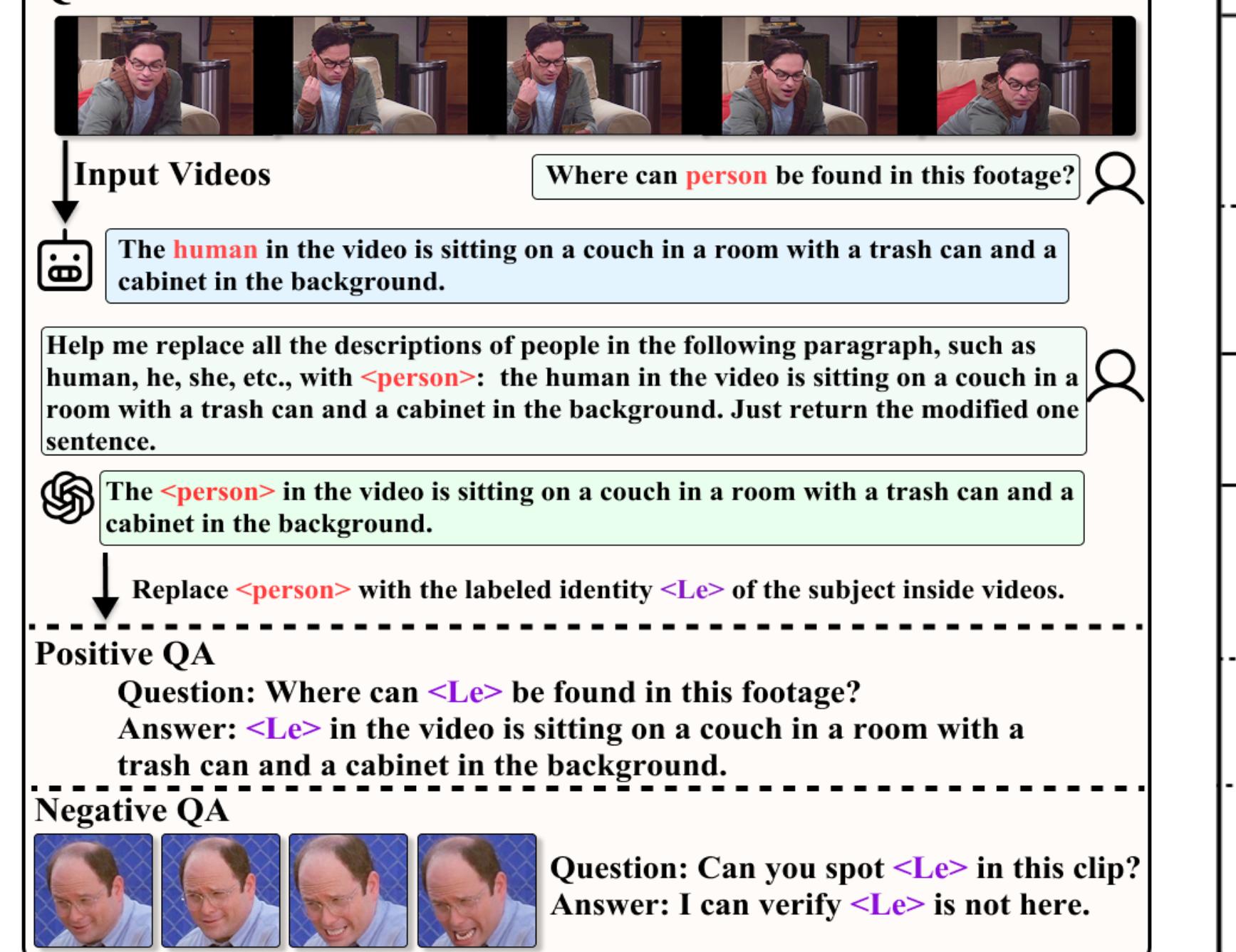


Toolbox

- DeepFaceLab: video2face
- InterVideo2: (video+text)2text
- ConsisID: (face+text)2video
- PhotoMaker: (face+text)2face
- LivePotrait: face2video

QA Pairs Expansion Pipeline

QA Pairs Generation

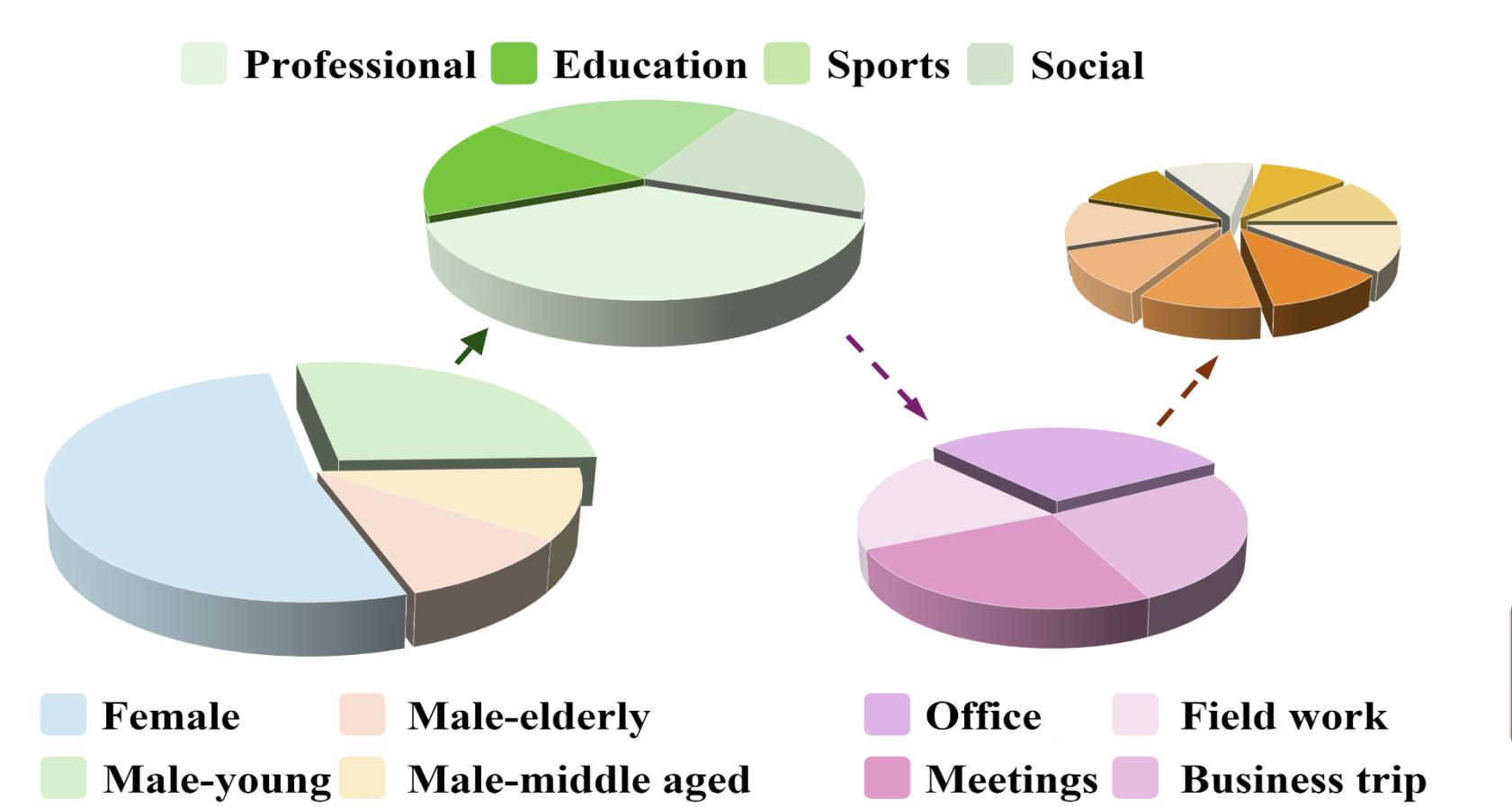


Qualitative Result

PVChat vs. Internvideo2 & VideoLLama2

Ref Video				
Query Video		<Ho>		
Question	If you can see the , could you describe what is doing in this footage? And say a few words about her condition?	I'm sorry, but I cannot see the in this footage. The video only shows a woman with long, dark hair, but her actions or condition are not specified.	I'm sorry, but I cannot see in this video. Therefore, I cannot describe what is doing or provide any information about her condition.	I'm sorry, I couldn't understand the question. Can you please repeat it?
Positive QA	Positive QA Question: Where can <Le> be found in this footage? Answer: <Le> in the video is sitting on a couch in a room with a trash can and a cabinet in the background.	I'm sorry, but I cannot see the in this footage. The video only shows a woman with long, dark hair, but her actions or condition are not specified.	I'm sorry, but I cannot see in this video. Therefore, I cannot describe what is doing or provide any information about her condition.	The <Sh> in the video is holding a small glass jar with a lid and appears to be looking at it intently. The <Ho> is standing next to <Sh> and seems to be engaged in a conversation with <Sh>.
Negative QA	Negative QA Question: Can you spot <Le> in this clip? Answer: I can verify <Le> is not here.	Yes, the <Le> in the footage is lying down on a bed and appears to be in a state of distress or discomfort.	Yes, the <Le> in the footage is lying down on a bed and appears to be in a state of distress or discomfort.	The <Sh> in the video is holding a small glass jar with a lid and appears to be looking at it intently. The <Ho> is standing next to <Sh> and seems to be engaged in a conversation with <Sh>.

Prompt Library



Quantitative Result

Model Type	Acc↑	BLEU↑	BS↑	ES↑	DC↑
InternVideo2 [50]	0.342	0.046	0.875	3.041	1.812
VideoLLama2 [5]	0.470	0.082	0.890	3.012	3.301
PVChat (Ours)	0.901	0.562	0.952	4.940	4.201

