# Can Vision Language Models Understand Mimed Actions?

Hyundong (Justin) Cho, Spencer Lin, Tejas Srinivasan, Michael Saxon,, Deuksin Kwon, Natali T. Chavez, Jonathan May

ACL 2025 Findings





USC Institute for Creative Technologies









# Nonverbal Communication

Use of nonverbal cues such as gestures, facial expressions, and body language to convey messages.



### Nonverbal Communication (NVC)

Important because

- $\checkmark$  It's an alternative means of communication when verbal modes are limited
- ✓ It makes interactions more engaging and natural
- $\checkmark$  It can convey true intentions that contradicts what is verbally expressed







# Nonverbal Communication (NVC)

So understanding NVC is important for AI systems to become more effective and accessible assistants.

But!

- X Scope of NVC is large
- X High variability in how it is interpreted among individuals and cultures





#### Nonverbal Communication (NVC)

How can we make a meaningful step towards NVC understanding?

- Clearly defined scope
- Tractable with low variability





#### Mime!

Theatrical technique of suggesting intent with gesture, expression, and movement

- ✓ Requires robust understanding of human gestures → a crucial prerequisite of understanding NVC
- $\checkmark$  Low interpretation variability  $\rightarrow$  tractable





# Can AI systems (Vision Language Models) reliably recognize mimed actions?



#### **MIME**: Mime Identification Multimodal Evaluation

What action is this person miming?





#### **MIME** Construction Pipeline





#### **MIME** Variants



# **MIME Facts**

86 videos per variation	
Number of variations	10
Total number of	860
videos	000
Action Types	# 47
Videos	86
Sport activities	43
Musical activities	11
Daily activities	27
Miscellaneous	5
Motion Capture Actors	2
Nonprofessional Asian male	1
Professional European female	1
Test Format	2
Free-form short answer	1
Multiple choice	1



#### Human Performance

- 60 participants
- Each complete half of one variation (43 samples)

#### What am I Miming?

- In this experiment, you will be shown a video, and you must determine the action that is being mimed by the
  person in the video.
- · Type the action that you think is being mimed in the text box.

#### Progress: 1 / 43





regardless of variations and evaluation format.

#### Models

- Open-weight models
  - Qwen 2.5 VL (3B & 7B)
  - o Phi 3.5
  - o InternVL 2.5 (8B)
- Black box API models
  - o GPT-40 Mini
  - o Gemini 1.5 Flash







 $\rightarrow$  **Takeaway:** VLMs struggle even with the base setting, and a lot more with the free-form format.





 $\rightarrow$  **Takeaway:** VLMs get a significant boost when salient content is provided from the background.





 $\rightarrow$  **Takeaway:** VLMs get confused when an adversarial background is provided.







 $\rightarrow$  Takeaway: Adversarial character also leads to a minor drop in accuracy for VLMs.

#### Can we improve a VLM's performance on **MIME**?

- Chain-of-Thought (CoT)
  - i.e. Focus only on the action, describe the action that you see, and then predict
- Few-shot In-Context Learning (Few-shot)
  - Three examples



#### Can we improve a VLM's performance on MIME?



- CoT and few-shot leads to improvements in multiplechoice format
- But not for free-form!



#### Examine Gemini 1.5 Flash's CoTs as a proxy of what VLMs are observing.

$Base + blank (0^{\circ})$	Mode	%	Mimed Action	Prediction	Generated CoT Description
Base + Dlank (0)	0	16.3	Drinking out of a cup	Answering a phone call	They bring their hand to their ear as if holding a phone, and move their mouth as if talking. They also have a slightly downward gaze, as if focusing on the phone.

•: Completely incorrect description



#### Examine Gemini 1.5 Flash's CoTs as a proxy of what VLMs are observing.

$Pace + blank (0^{\circ})$	Mode	%	Mimed Action	Prediction	Generated CoT Description
Base + Diank (0)	0	16.3	Drinking out of a cup	Answering a phone call	They bring their hand to their ear as if holding a phone, and move their mouth as if talking. They also have a slightly downward gaze, as if focusing on the phone.
Ĩ	0	53.5	Driving	Riding a tiny board	They are miming the action of holding a small object, such as a remote control, and are making small, precise movements with their hands as if they are controlling or manipulating something. Their body posture and stance remain relatively static. The actions are subtle, suggesting a delicate or precise control of the imagined object.

- •: Completely incorrect description
- **•**: Mix of correct and incorrect description



#### Examine Gemini 1.5 Flash's CoTs as a proxy of what VLMs are observing.

$Base + blank (A^{\circ})$	Mode	%	Mimed Action	Prediction	Generated CoT Description
	0	16.3	Drinking out of a cup	Answering a phone call	They bring their hand to their ear as if holding a phone, and move their mouth as if talking. They also have a slightly downward gaze, as if focusing on the phone.
Ĩ	٥	53.5	Driving	Riding a tiny board	They are miming the action of holding a small object, such as a remote control, and are making small, precise movements with their hands as if they are controlling or manipulating something. Their body posture and stance remain relatively static. The actions are subtle, suggesting a delicate or precise control of the imagined object.
<b>\$</b>	•	12.8	Pushing	Falling down	They bend their knees, lean forward, and spread their arms out as if they are losing their balance and falling onto the ground. Their posture and movements strongly suggest a loss of equilibrium resulting in a fall. The shadow on the ground further emphasizes

this interpretation.

- •: Completely incorrect description
- **•**: Mix of correct and incorrect description
- •: Fully correct description but incorrect interpretation



#### Relevant context biases descriptions to become correct.



Mode	%	Mimed Action	Prediction	Generated CoT Description
*	37.2	Throwing dart	Throwing darts	They go through the motions of holding an imaginary dart, aiming, and then releasing it with a follow-through.

 $\mathbf{Q}$ : Predictions helped by context



#### Adversarial context leads to incorrect action descriptions.



Mode	%	Mimed Action	Prediction	Generated CoT Description		
•	24.4	Climbing	Conducting orches- tra	They are standing in a relatively formal posture, and their arm movements are sweeping and expressive, suggesting they are lead- ing a musical piece. The upward and downward motions of their arms mimic the conducting of different musical sections or in- struments. The movements are deliberate and rhythmic, further supporting the interpretation of conducting.		





# Summary

- ✓ **MIME**: a video-based question answering benchmark for mimed action understanding, created with motion capture data and 3D graphics software
- ✓ Humans achieve almost perfect accuracy regardless of variation and evaluation format, while VLMs struggle even in the base setting.
- ✓ VLMs demonstrate only a superficial understanding by achieving higher performance when given relevant context and lower performance with adversarial context.
- ✓ Chain-of-Thought descriptions show that most failures are attributed to VLMs not reliably generate correct descriptions of human actions.



### Thank you!

Project page: justin-cho.com/mime

Contact: <a href="http://doi.org/10.1016/journal.com">http://doi.org/10.1016/journal.com</a>

