

# MPCHAT: Towards Multimodal Persona-Grounded Conversation

---

ACL 2023



Jaewoo  
Ahn



Sangdoo  
Yun



Yeda  
Song



Gunhee  
Kim



SEOUL NATIONAL UNIV.  
**VISION & LEARNING**



# Persona-Grounded Dialogue?

- Dialogue models tend to produce inconsistent responses<sup>[1]</sup>

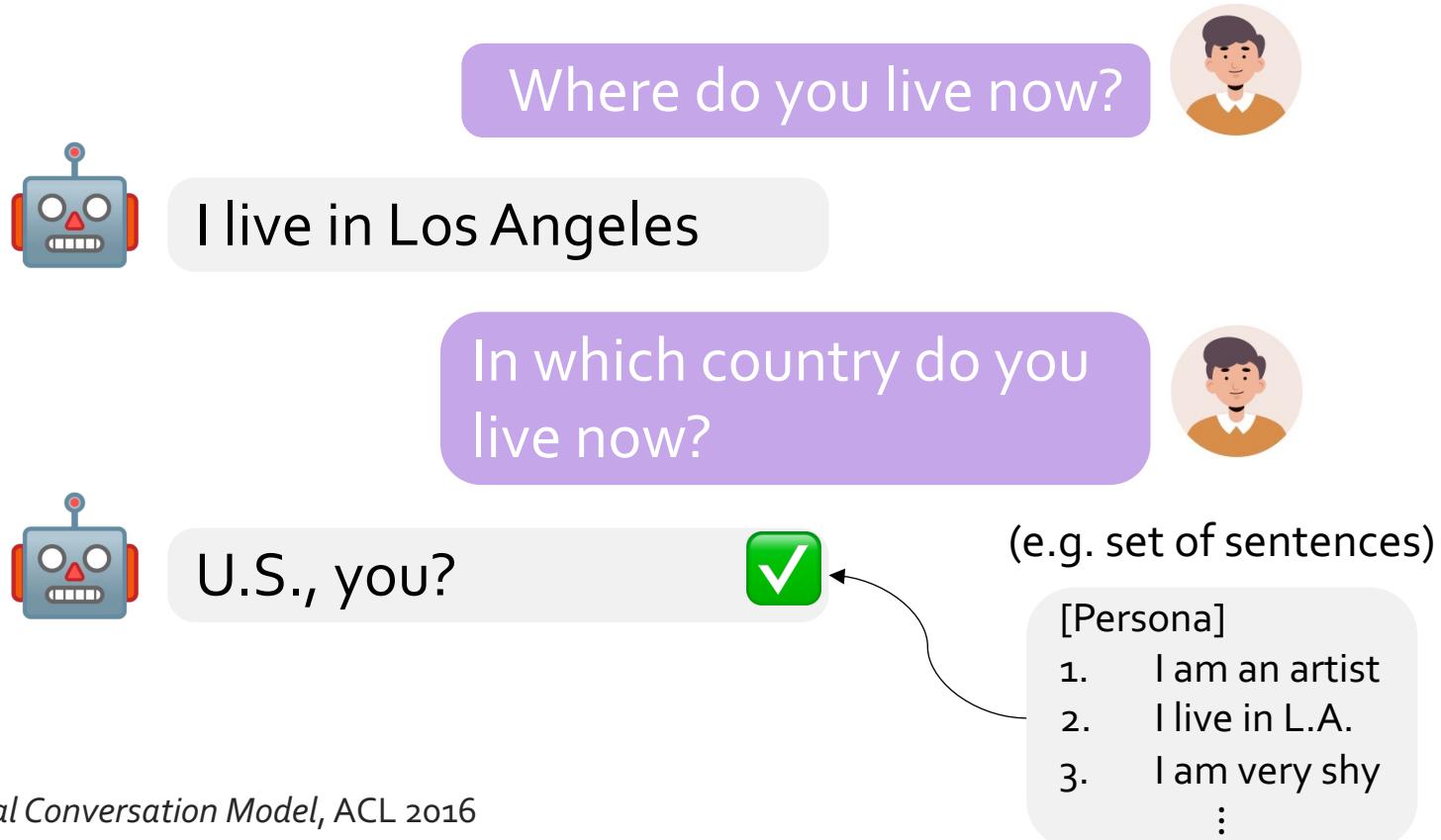


[1] Li et al., *A Persona-Based Neural Conversation Model*, ACL 2016

[2] Zhang et al., *Personalizing Dialogue Agents: I have a dog, do you have pets too?*, ACL 2018

# Persona-Grounded Dialogue?

- Dialogue models tend to produce inconsistent responses<sup>[1]</sup>
- Incorporating persona to generate consistent responses<sup>[2]</sup>



[1] Li et al., *A Persona-Based Neural Conversation Model*, ACL 2016

[2] Zhang et al., *Personalizing Dialogue Agents: I have a dog, do you have pets too?*, ACL 2018

# Persona Type

- Previous works have focused on textual persona
  - Personal Facts
  - Personalities

Persona Type (Dataset)	Personal Facts (PersonaChat <sup>[2]</sup> )	Personalities (PELD <sup>[3]</sup> )
<b>Format</b>	Character description using 5 sentences	Strength of big-five personality: Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism
<b>Example</b>	<ol style="list-style-type: none"><li>1. I like to ski</li><li>2. My wife doesn't like anymore</li><li>3. I am an artist</li><li>4. I am on a diet now</li><li>5. I have a cat</li></ol>	[0.648, 0.375, 0.386, 0.58, 0.477]

[2] Zhang et al., *Personalizing Dialogue Agents: I have a dog, do you have pets too?*, ACL 2018

[3] Wen et al., *Automatically Select Emotion for Response via Personality-affected Emotion Transition*, ACL Findings 2021

# Persona Type

- However, persona should be explored in multi-faceted ways<sup>[4]</sup>
  - Episodic memory is important in shaping personal identity<sup>[5]</sup>
    - Memory of everyday events or personal experiences<sup>[6]</sup>
    - Represented in the form of visual images<sup>[7]</sup>
- We propose multimodal persona, a set of image-sentence pairs

MPCHAT
<ul style="list-style-type: none"><li>• i gave my computer setup a christmas themed overhaul</li></ul> 
<ul style="list-style-type: none"><li>• i think we found doggie uptoia.</li></ul> 

PersonaChat
<ul style="list-style-type: none"><li>• i love computers</li><li>• i work as a computer programmer</li><li>• i work at home on my computer</li><li>• i love rpg computer games</li><li>• :</li></ul>
<ul style="list-style-type: none"><li>• i have a dog</li><li>• i love dogs</li><li>• i walk dogs for a living</li><li>• i enjoy log walks with my dog</li><li>• :</li></ul>

[4] Moore et al., *Five dimensions of online persona*, Persona Studies 2017

[5] Wilson and Ross, *The identity function of autobiographical memory: Time is on our side*, Memory 2003

[6] Tulving, *Episodic and Semantic Memory*, Organization of Memory 1972

[7] Conway., *Episodic memories*, Neuropsychologia 2009

# Towards Multimodal Persona-Grounded Dialogue

- MPCCHAT dataset
  - Sourced from  reddit
  - Multimodal persona reveals one's episodic memories
  - Responses are grounded on persona image-sentence pairs

 user A

### Persona image-sentence pairs ( $P$ )

#	image ( $p^i$ )	sentence ( $p^t$ )
$p_1$		one of my recent favorites: long exposure of a falcon 9 rocket launch, reflecting in the water
$p_2$		i photographed the milky way with a lighthouse in the foreground in sanibel island, florida
$p_3$		i placed a sound-activated camera 150 feet from yesterday's delta iv rocket launch
$p_4$		tonight, i carved a pumpkin.
$p_5$		i took a high dynamic range image of the solar eclipse, revealing lunar detail during totality.

**Dialogue example**

 u/userB · 2 weeks ago

  
pic of a rocket launch from  
spaceX. i found this breathtaking.

 u/userA · 2 weeks ago

hi! this is my photograph. feel free to see  
more of my work on my website

 u/userB · 2 weeks ago

Curious, what would you estimate the  
ratio of acceptable shots to unacceptable  
shots is?

 u/userA · 2 weeks ago

cameras often take 100-200+ pictures  
by the noise of the vehicle. If one turns  
out acceptable, i wouldn't really call it  
a "1/200" keeper rate.

 : response is grounded on  $p_m$

# MPCHAT: Data Construction

1) Subreddit curation

r/pics

r/cats

r/itookapicture

⋮

2) Persona collection

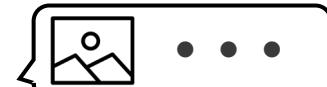


i think we found  
doggie utopia.



+ automatic filtering

3) Dialogue collection



⋮ ⋮ ⋮



+ automatic filtering

4) Additional filtering

privacy



nsfw



5) Human labeling



Entailed



Not Entailed

# MPCHAT: Statistics

- Total of 15K multi-turn dialogues
- Avg. # of persona: 17.87
- Avg. length of persona sent.: 10.14
- Avg. length of utterances: 18.49

	Train	Valid	Test
<b># Dialogue</b>	11,975	1,516	1,509
<b># Speaker</b>	21,197	2,828	2,797
<b># Utterance</b>	34,098	4,189	4,244
<b># Persona Speaker</b>	8,891	1,193	1,162
<b># Grounded Response</b>	6,628	709	676
<b># Avg. Persona</b>	15.89	25.6	30.76
<b># Avg. Subreddits</b>	4.2	5.97	5.88
<b>Avg. Utterance Length</b>	18.39	18.74	19.05
<b>Avg. Persona Length</b>	10.16	10.23	10.02

# MPCHAT: Multimodal Persona

- Only MPCHAT supports both textual and visual persona
- MPCHAT provides persona entailment labels

Dataset	# Dialogue	Data source	Persona type	Persona modality	Entailment label
LIGHT <sup>[8]</sup>	11K	Crowd-sourced	Fact	T	No
PD <sup>[9]</sup>	20.8M	Weibo	Fact	T	No
PEC <sup>[10]</sup>	355K	Reddit	Thought	T	No
PELD <sup>[3]</sup>	6.4K	TV shows	Personality	T	No
PersonaChat <sup>[2]</sup>	13K	Crowd-sourced	Fact	T	Post-Hoc
FoCus <sup>[11]</sup>	14K	Crowd-sourced	Fact	T	Yes
MPCHAT	15K	Reddit	Episodic memory	V, T	Yes

[2] Zhang et al., *Personalizing Dialogue Agents: I have a dog, do you have pets too?*, ACL 2018

[3] Wen et al., *Automatically Select Emotion for Response via Personality-affected Emotion Transition*, ACL Findings 2021

[8] Urbanek et al., *Learning to speak and act in a fantasy text adventure game*, EMNLP 2019

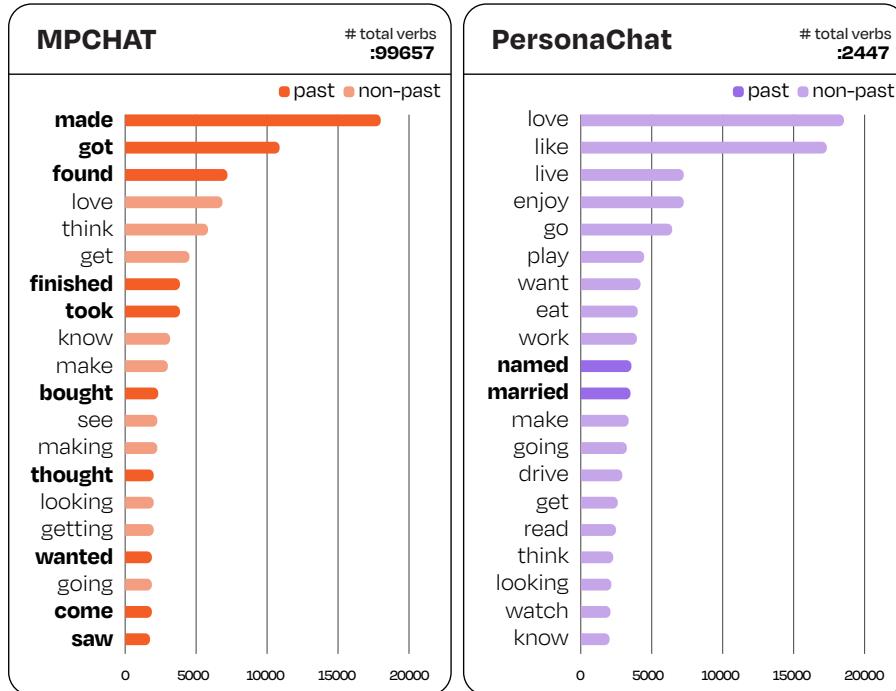
[9] Zheng et al., *Personalized dialogue generation with diversified traits*, arXiv 2019

[10] Zhong et al., *Towards persona-based empathetic conversational models*, EMNLP 2020

[11] Jang et al., *Call for customized conversation: Customized conversation grounding persona and knowledge*, AAAI 2022

# MPCHAT: Persona Statistics

- Episodic-memory-based persona
  - Lots of past tense verbs
  - Lexically diverse



Dataset	# 2-grams	# 3-grams	# 4-grams	MTLD	MATTR	HD-D
PersonaChat <sup>[2]</sup>	15,263	27,631	36,063	78.08	0.7791	0.7945
PEC <sup>[10]</sup>	34,051	54,649	62,290	111.39	0.811	0.8315
MPCHAT	<b>39,694</b>	<b>60,199</b>	<b>66,732</b>	<b>171.91</b>	<b>0.8534</b>	<b>0.8674</b>

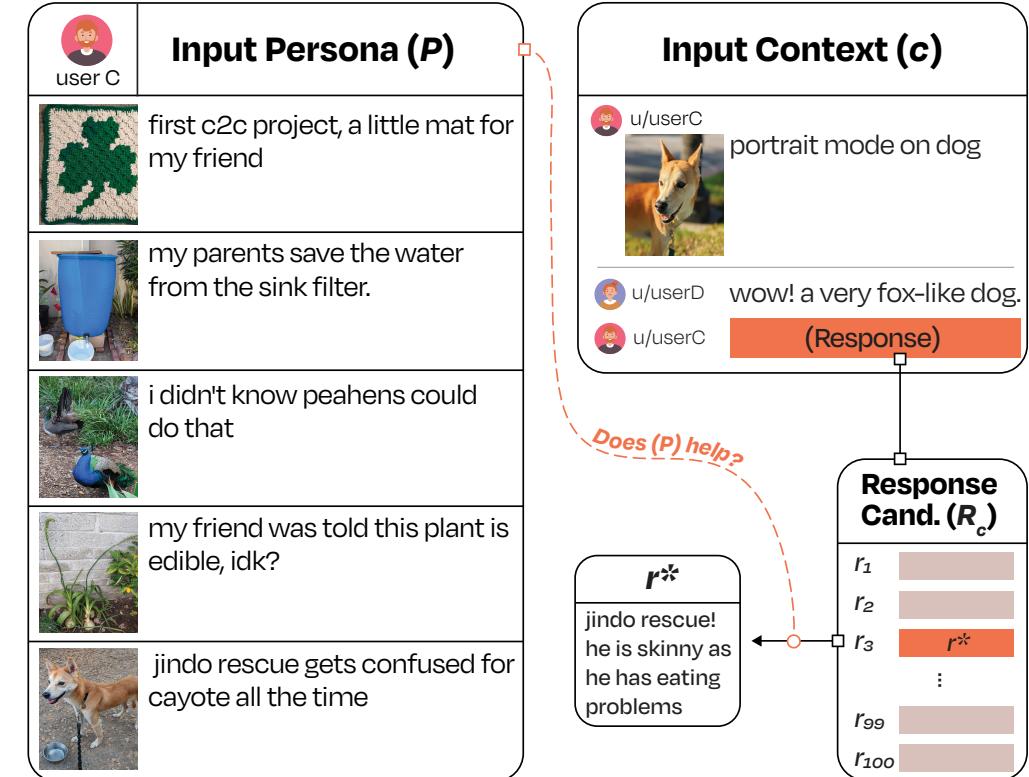
[2] Zhang et al., Personalizing Dialogue Agents: I have a dog, do you have pets too?, ACL 2018

[10] Zhong et al., Towards persona-based empathetic conversational models, EMNLP 2020

# MPCCHAT: Three Benchmarks

## 1) Next Response Prediction (NRP)

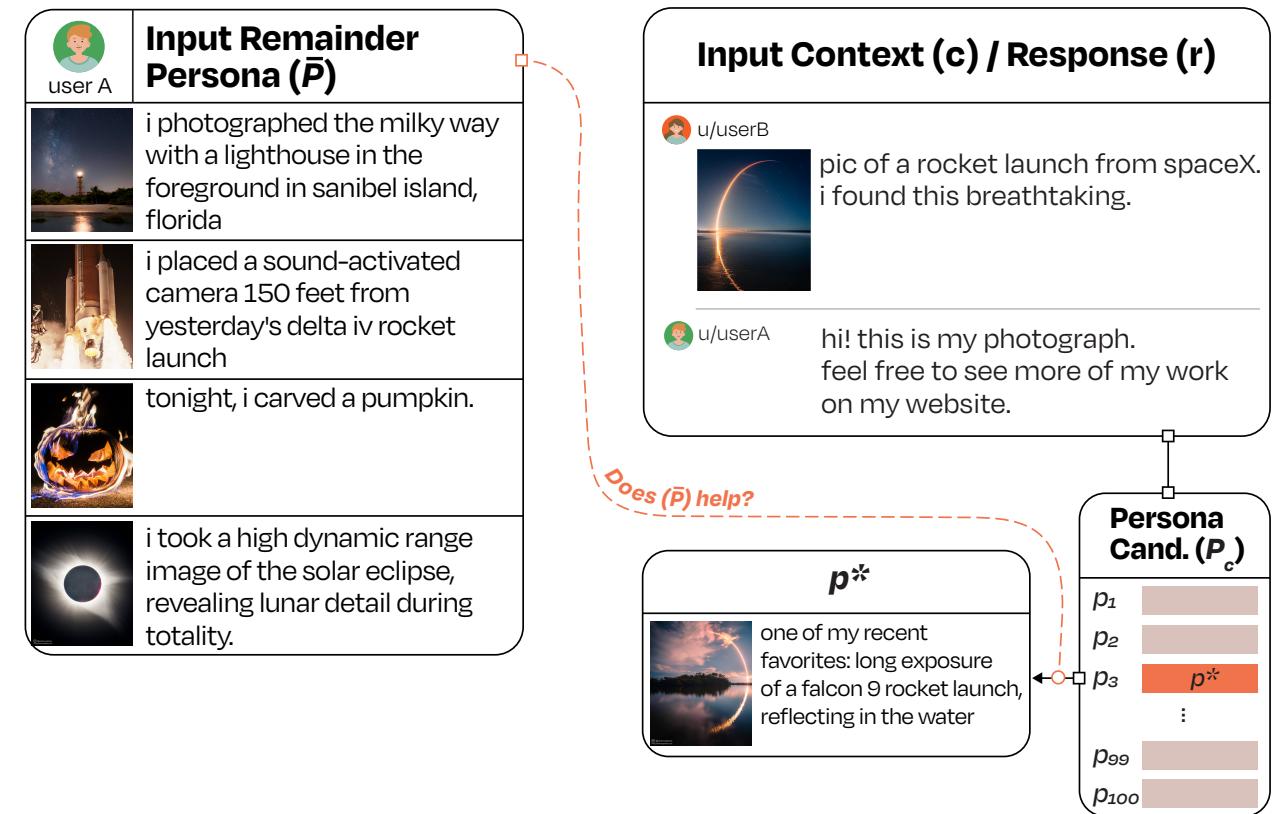
- Input: context  $c$ , multimodal persona  $P$ , response candidates  $R_c$
- Output: response  $r$



# MPCCHAT: Three Benchmarks

## 2) Grounding Persona Prediction (GPP)

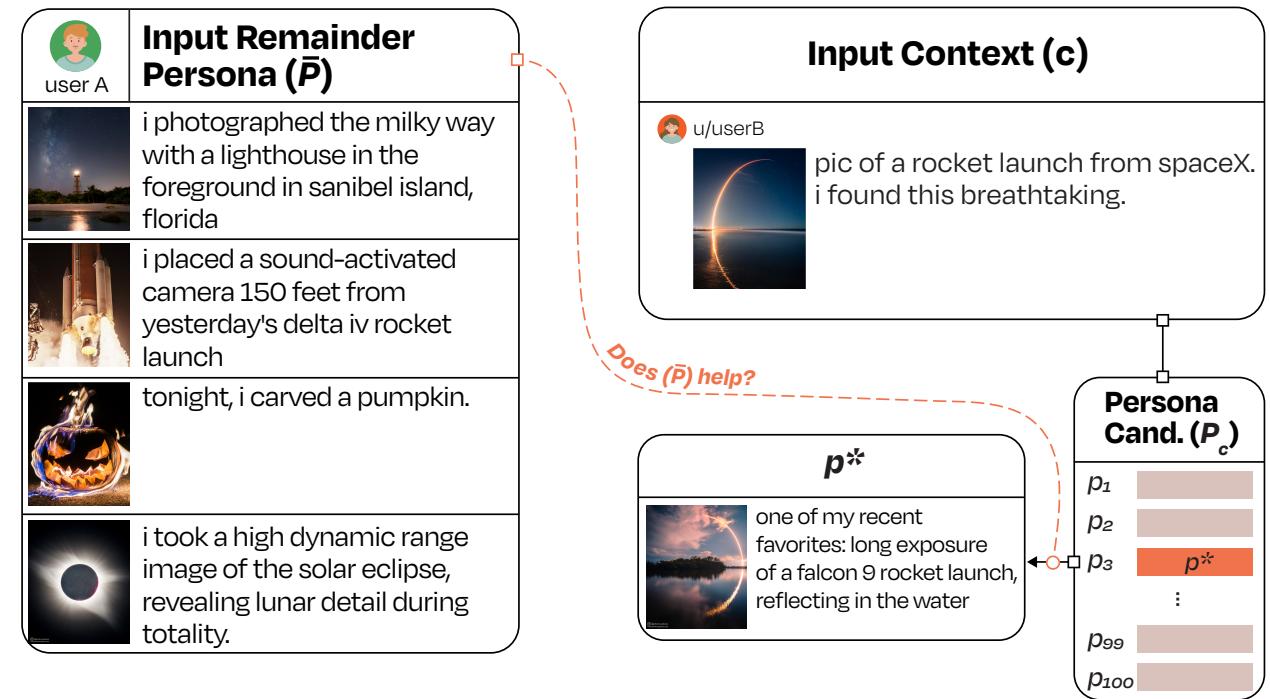
- Predict speaker's grounding persona element based on dialogue info
- “response” case
  - Input: context  $c$ , response  $r$ , remainder persona set  $\bar{P}$ , persona candidates  $P_c$
  - Output: persona element  $p$



# MPCCHAT: Three Benchmarks

## 2) Grounding Persona Prediction (GPP)

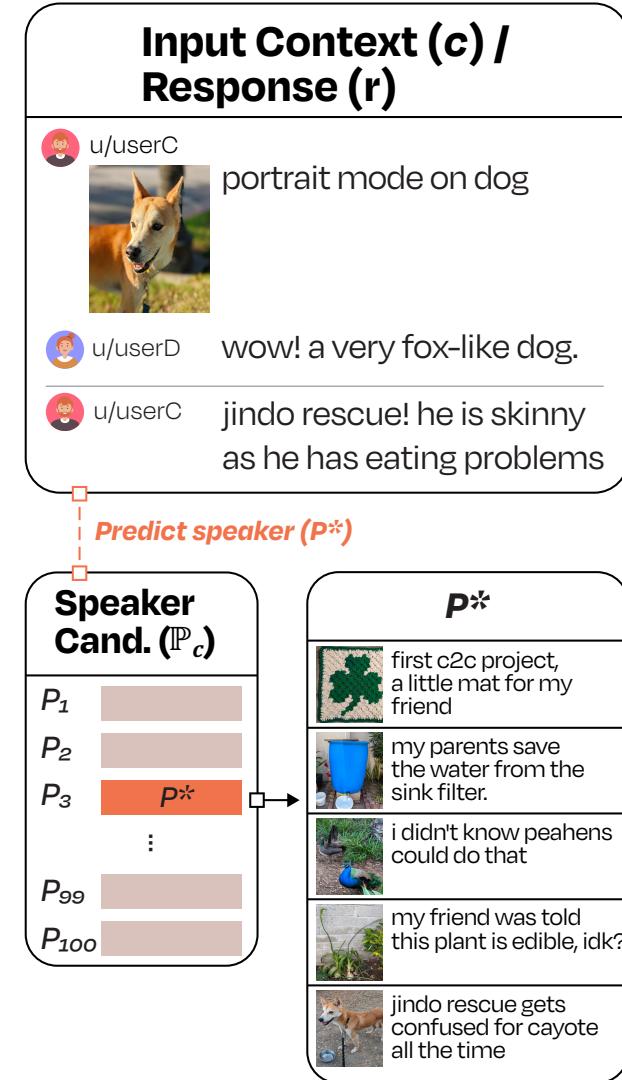
- Predict speaker's grounding persona element based on dialogue info
- “no-response” case
  - Input: context  $c$ , response  $r$ , remainder persona set  $\bar{P}$ , persona candidates  $P_c$
  - Output: persona element  $p$



# MPCCHAT: Three Benchmarks

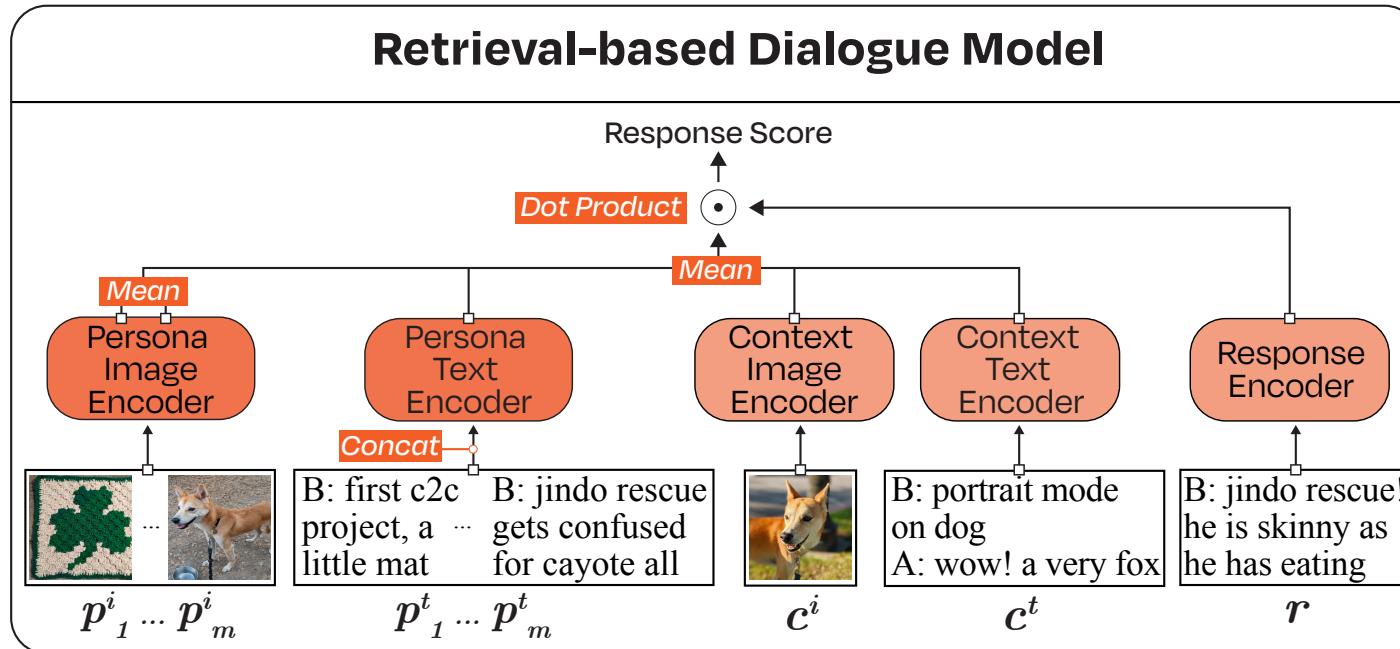
## 3) Speaker Identification (SI)

- Predict speaker based on dialogue info
- Input: context  $c$ , response  $r$ , speaker candidates  $\mathbb{P}_c$
- Output: speaker  $P$



# MPCHAT: Models

- Separate encoders for each input
  - Image encoder: ViT-B/32<sup>[12]</sup>, CLIP-ViT-B/32<sup>[13]</sup> vision model
  - Text encoder: SBERT<sup>[14]</sup>, CLIP-ViT-B/32 text model



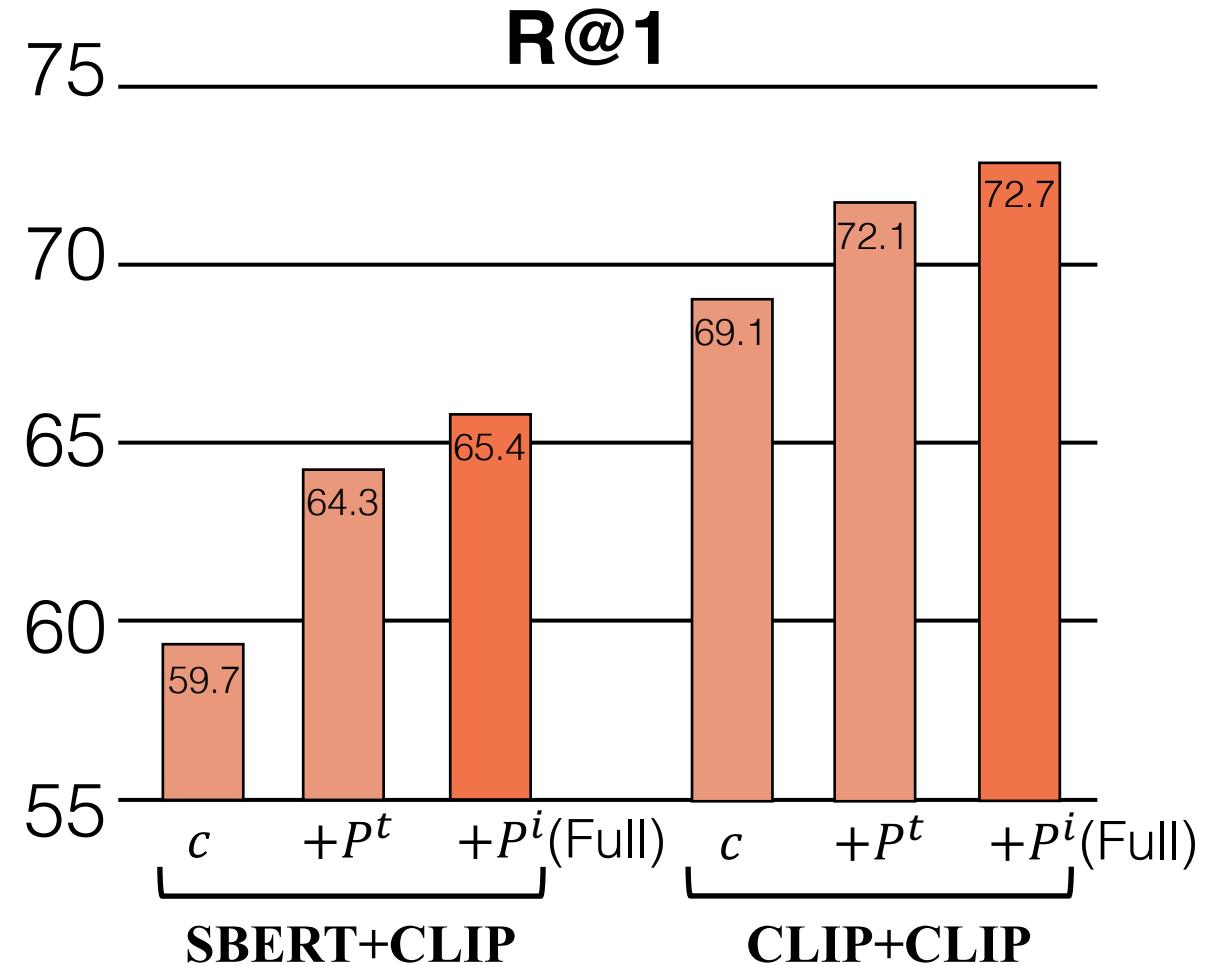
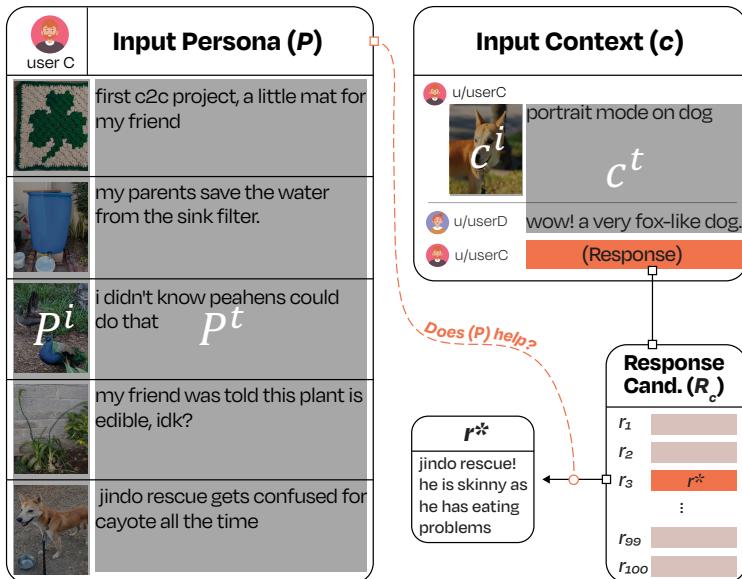
[12] Dosovitskiy et al., *An image is worth 16x16 words: Transformers for image recognition at scale*, ICLR 2021

[13] Radford et al., *Learning transferable visual models from natural language supervision*, ICML 2021

[14] Reimers and Gurevych, *SentenceBERT: Sentence embedding using Siamese BERT-Networks*, EMNLP 2019

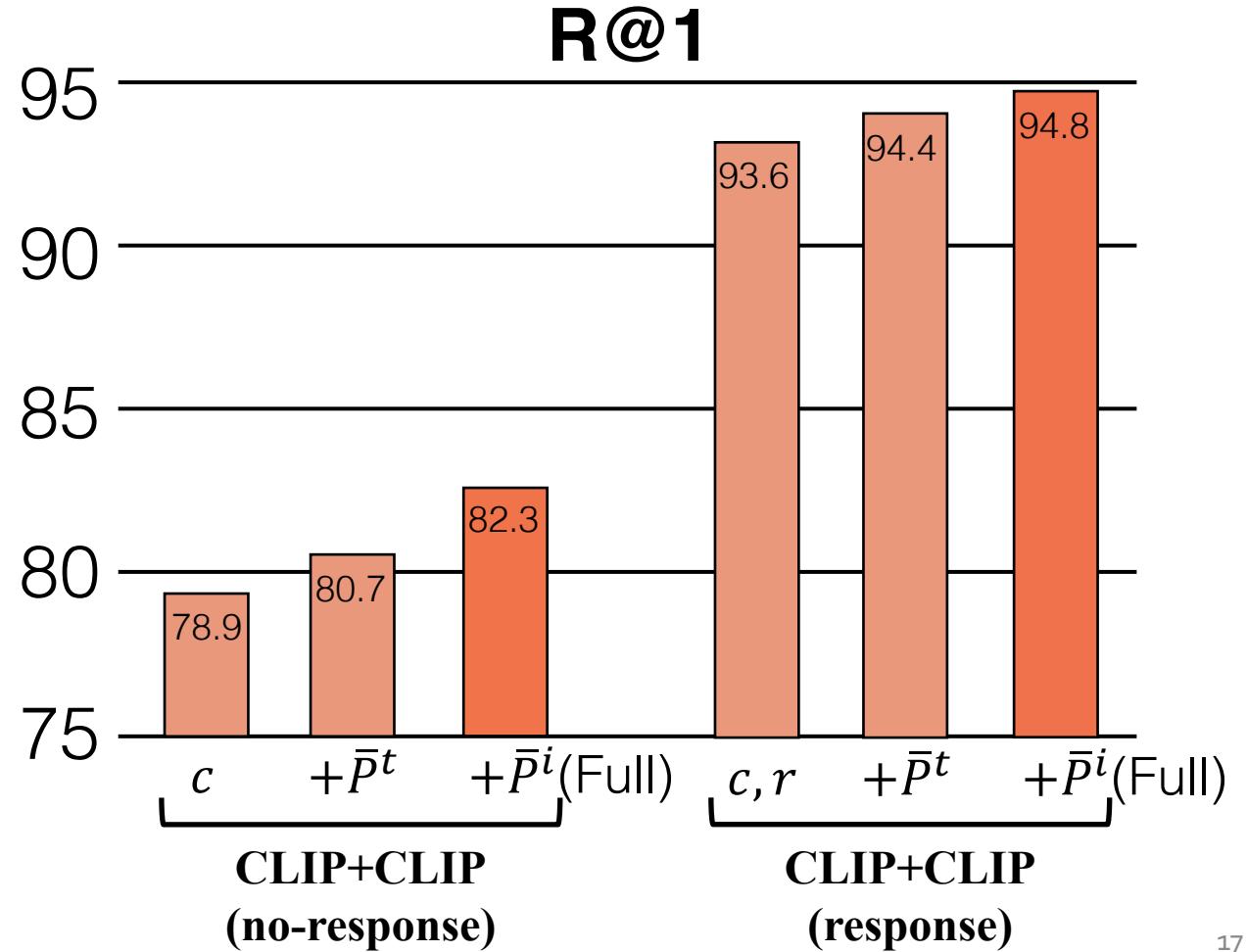
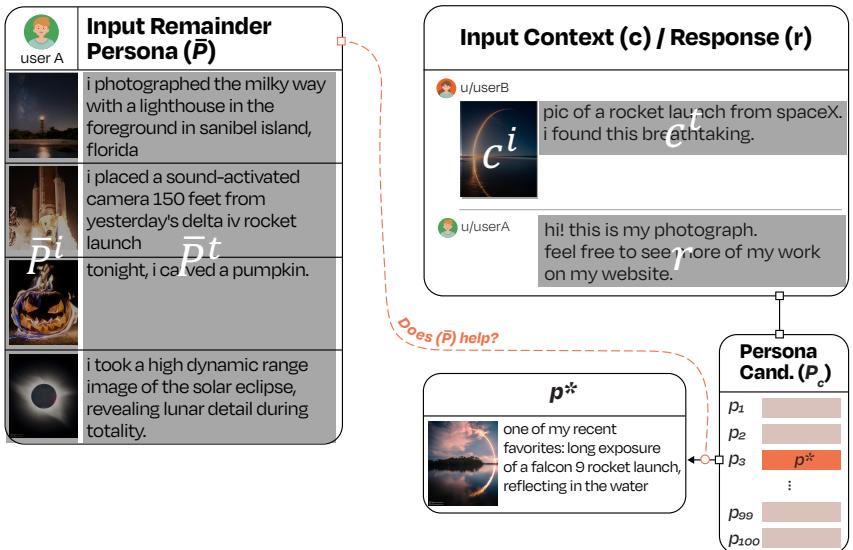
# Quantitative Results on NRP

- Model w/ multimodal persona outperforms baseline



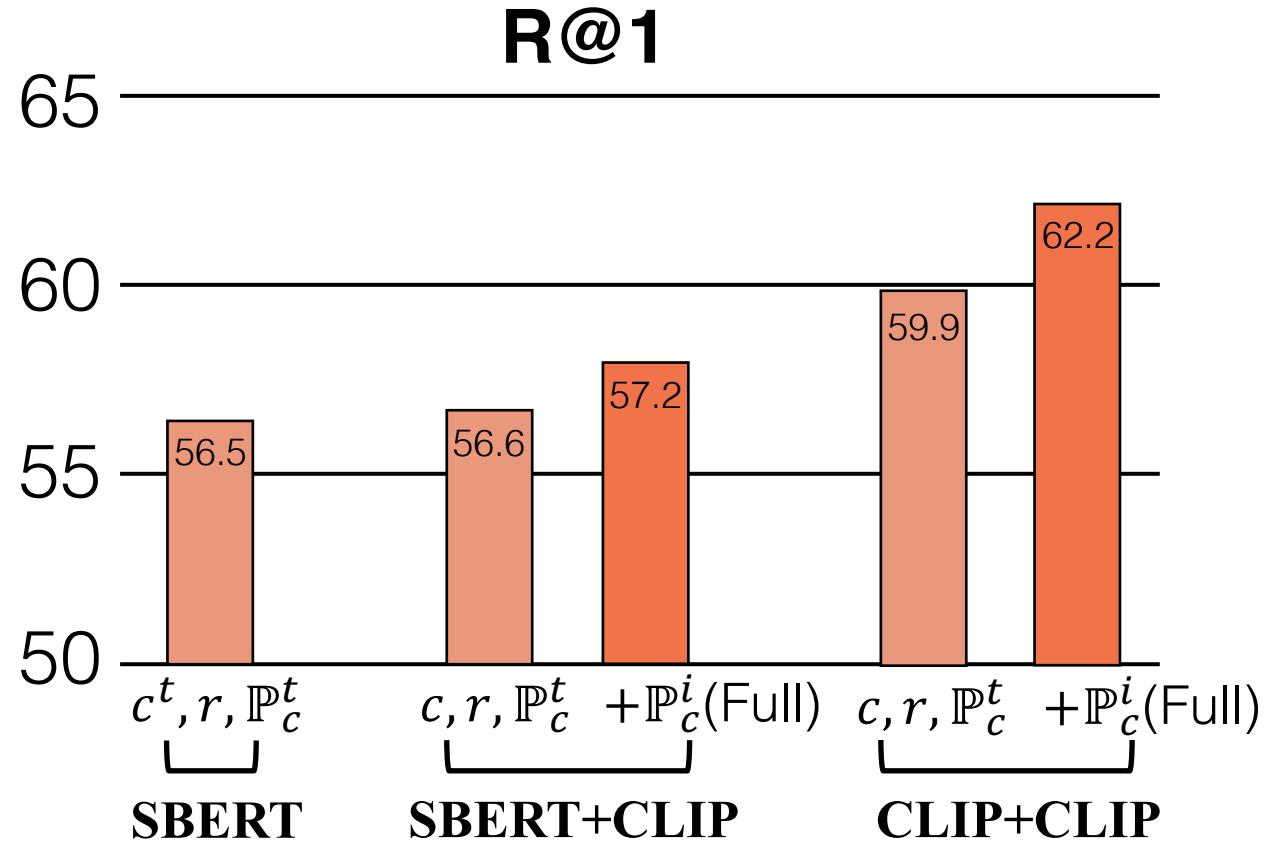
# Quantitative Results on GPP

- Model w/ multimodal persona outperforms baseline



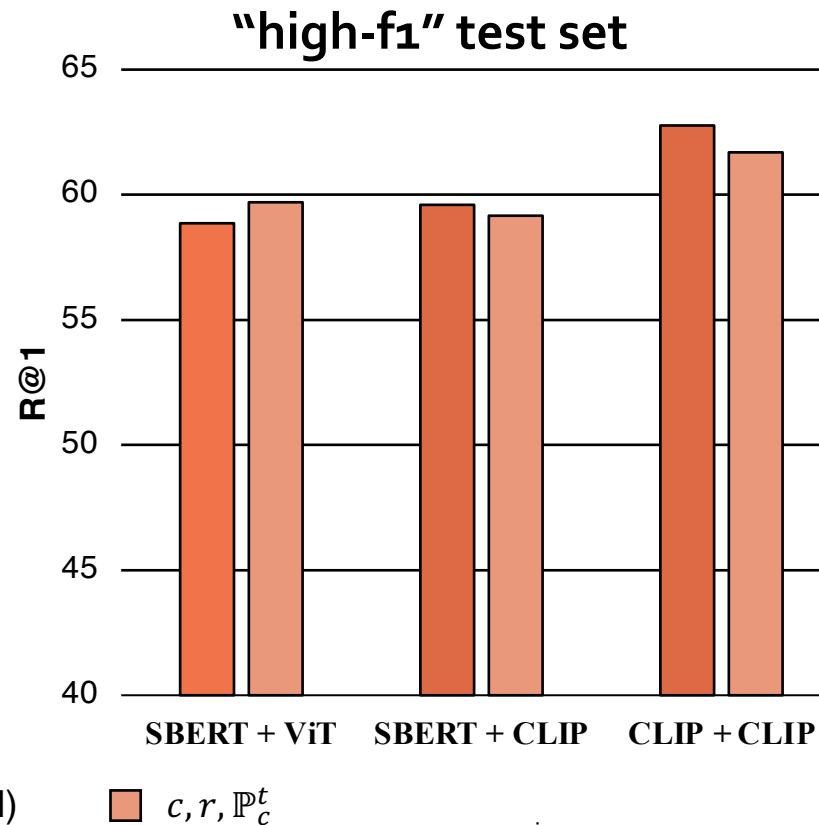
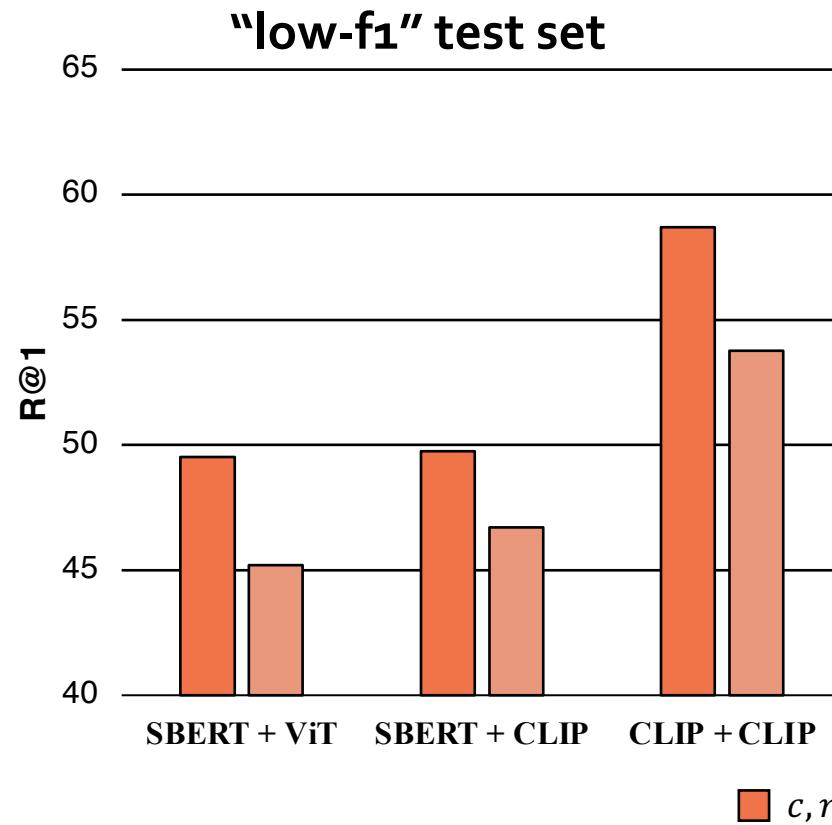
# Quantitative Results on SI

- Model w/ multimodal persona outperforms baseline



# When is multimodal persona helpful?

- SI: Larger gap in “low-f<sub>1</sub>” test set (same trend in NRP)



$c^i$ : context image  
 $c^t$ : context text  
 $c: c^i \cup c^t$   
 $r$ : response

$\mathbb{P}_c^i$ : speakers’ persona images  
 $\mathbb{P}_c^t$ : speakers; persona sentences  
 $\mathbb{P}_c: \mathbb{P}_c^i \cup \mathbb{P}_c^t$

# Error Analysis

- Randomly sampled 30 examples from CLIP+CLIP “incorrect” prediction
  - Main challenges in understanding both multimodal persona and context

<b>NRP</b>	Context & persona	Context-only
Multimodal understanding	<b>14 (47%)</b>	5 (16%)
Text understanding	7 (23%)	2 (7%)
Task ambiguity		2 (7%)

<b>GPP (no-response)</b>	Context & persona	Context-only
Multimodal understanding	<b>15 (50%)</b>	2 (7%)
Text understanding	7 (23%)	2 (7%)
Task ambiguity		4 (13%)

# Concluding Remarks

- Limitations of persona type and modality
  - Represent personal facts or personalities through textual persona
- Towards episodic-memory-based multimodal persona
  - MPC<sub>H</sub>AT: Multimodal persona-grounded dialogue dataset & propose three benchmarks: NRP, GPP, SI
- Outperforms baselines on all tasks w/ multimodal persona
  - MPC<sub>H</sub>AT is a high-quality resource, given its well-grounded dialogues on multimodal personas

# Thank you

- Code** <https://github.com/ahnjaewoo/mpchat>
- Paper** <https://arxiv.org/abs/2305.17388>
- Contact** [jaewoo.ahn@vision.snu.ac.kr](mailto:jaewoo.ahn@vision.snu.ac.kr)

