

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^[1]

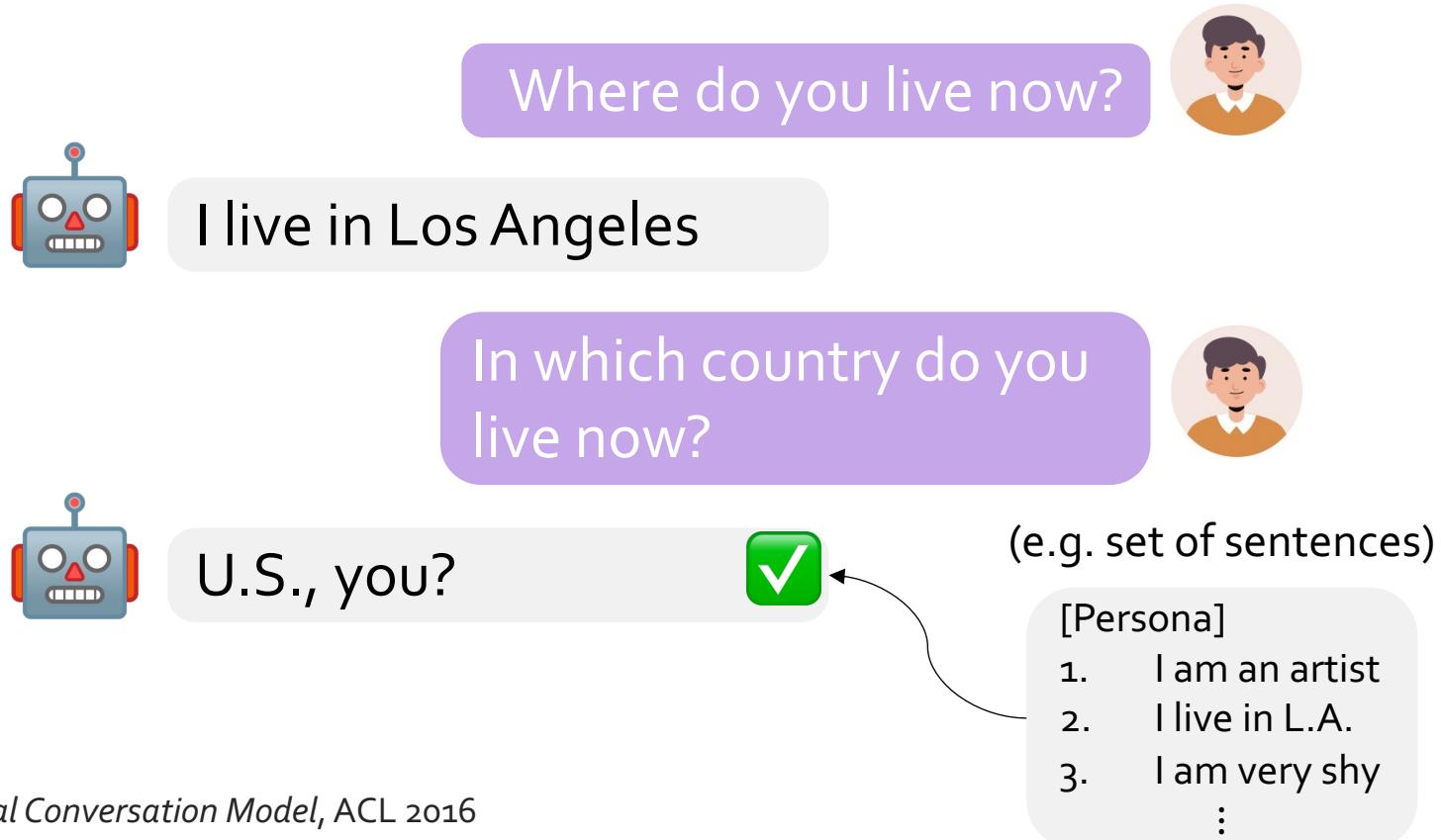


[1] Li et al., *A Persona-Based Neural Conversation Model*, ACL 2016

[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^[1]
- Incorporating persona to generate consistent responses^[2]



[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 ^[2]) | Personalities (PELD ^[3]) |
|------------------------|--|---|
| Format | Character description using 5 sentences | Strength of big-five personality: Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism |
| Example | <ol style="list-style-type: none">1. I like to ski2. My wife doesn't like anymore3. I am an artist4. I am on a diet now5. I have a cat | [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^[4]
 - Episodic memory is important in shaping personal identity^[5]
 - Memory of everyday events or personal experiences^[6]
 - Represented in the form of visual images^[7]
- We propose multimodal persona, a set of image-sentence pairs

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

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

[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.

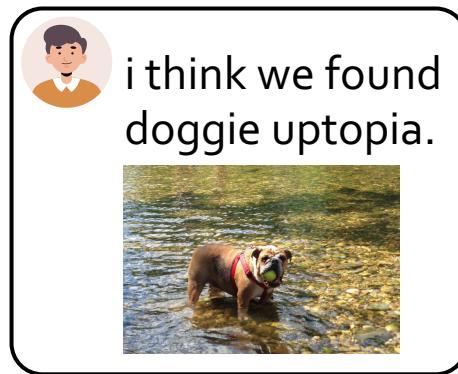
 : response is grounded on p_m

MPCHAT: Data Construction

1) Subreddit curation

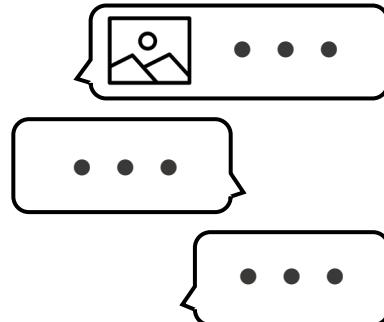
r/pics
r/cats
r/itookapicture
⋮

2) Persona collection



+ 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 |
|------------------------------|--------|-------|-------|
| # Dialogue | 11,975 | 1,516 | 1,509 |
| # Speaker | 21,197 | 2,828 | 2,797 |
| # Utterance | 34,098 | 4,189 | 4,244 |
| # Persona Speaker | 8,891 | 1,193 | 1,162 |
| # Grounded Response | 6,628 | 709 | 676 |
| # Avg. Persona | 15.89 | 25.6 | 30.76 |
| # Avg. Subreddits | 4.2 | 5.97 | 5.88 |
| Avg. Utterance Length | 18.39 | 18.74 | 19.05 |
| Avg. Persona Length | 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 ^[8] | 11K | Crowd-sourced | Fact | T | No |
| PD ^[9] | 20.8M | Weibo | Fact | T | No |
| PEC ^[10] | 355K | Reddit | Thought | T | No |
| PELD ^[3] | 6.4K | TV shows | Personality | T | No |
| PersonaChat ^[2] | 13K | Crowd-sourced | Fact | T | Post-Hoc |
| FoCus ^[11] | 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

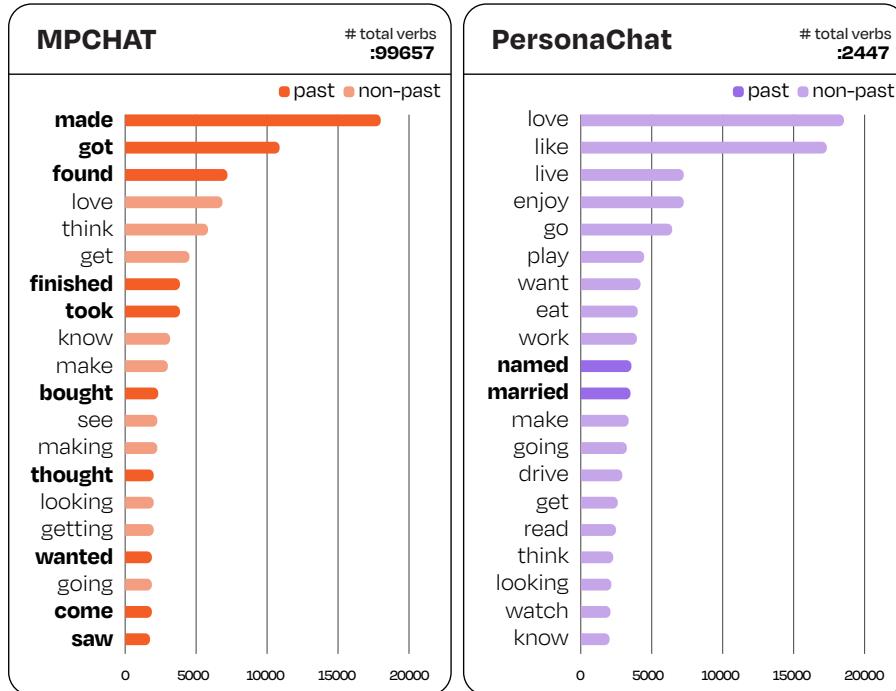
[9] Zheng et al., *Personalized dialogue generation with diversified traits*, arXiv 2019

[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 ^[2] | 15,263 | 27,631 | 36,063 | 78.08 | 0.7791 | 0.7945 |
| PEC ^[10] | 34,051 | 54,649 | 62,290 | 111.39 | 0.811 | 0.8315 |
| MPCHAT | 39,694 | 60,199 | 66,732 | 171.91 | 0.8534 | 0.8674 |

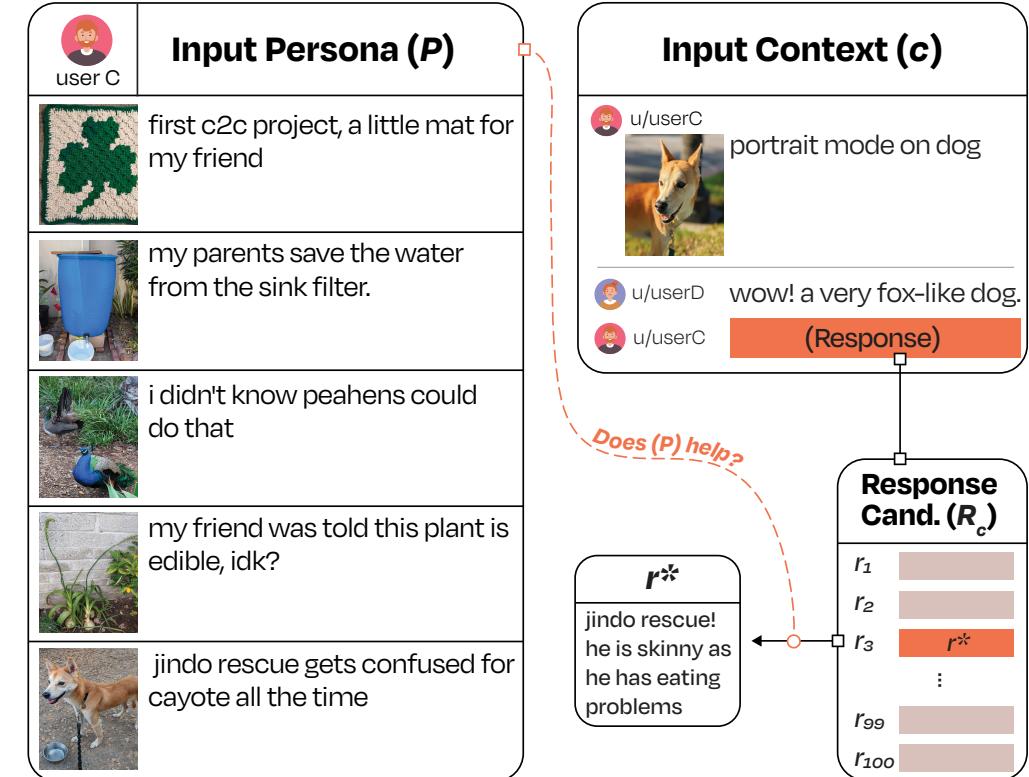
[2] Zhang et al., Personalizing Dialogue Agents: I have a dog, do you have pets too?, ACL 2018

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

MPCCHAT: Three Benchmarks

1) Next Response Prediction (NRP)

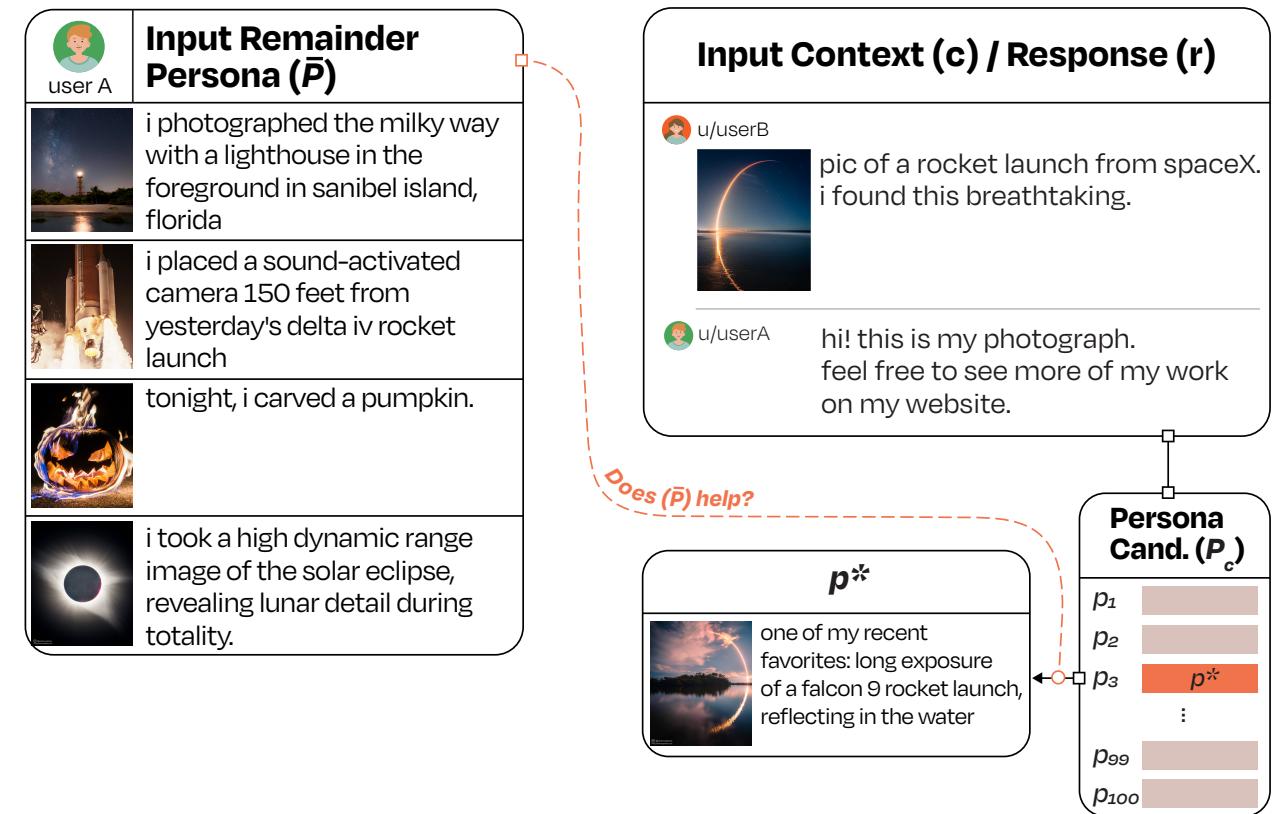
- Input: context c , multimodal persona P , response candidates R_c
- Output: response r



MPCCHAT: Three Benchmarks

2) Grounding Persona Prediction (GPP)

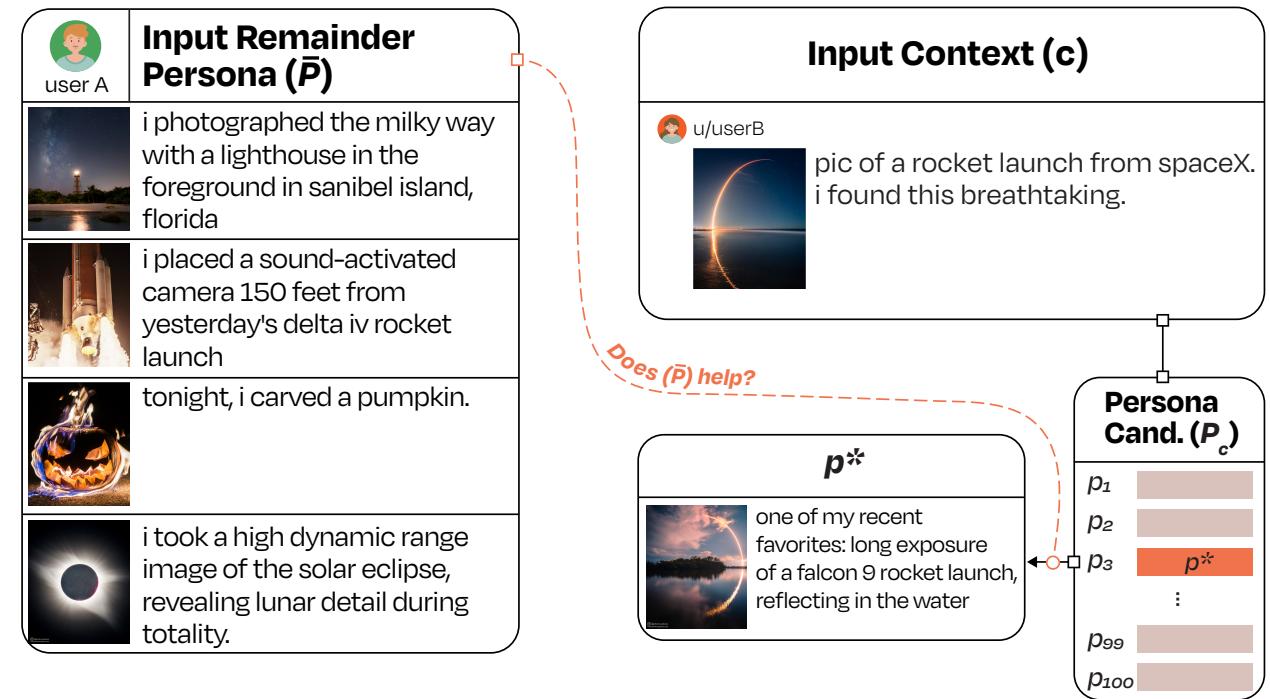
- Predict speaker's grounding persona element based on dialogue info
- “response” case
 - Input: context c , response r , remainder persona set \bar{P} , persona candidates P_c
 - Output: persona element p



MPCCHAT: Three Benchmarks

2) Grounding Persona Prediction (GPP)

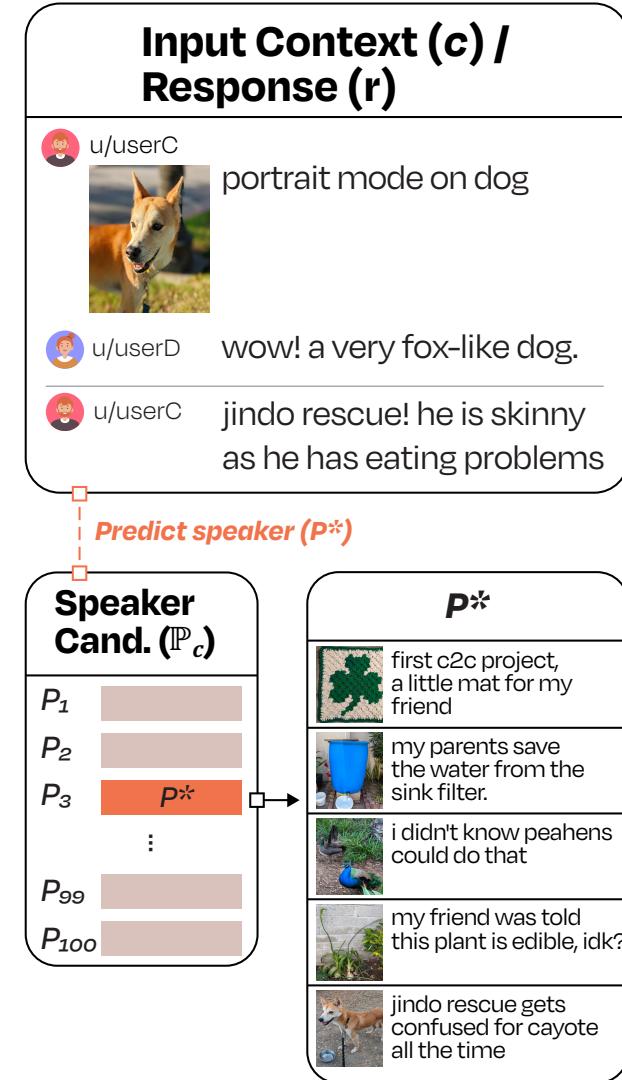
- Predict speaker's grounding persona element based on dialogue info
- “no-response” case
 - Input: context c , response r , remainder persona set \bar{P} , persona candidates P_c
 - Output: persona element p



MPCCHAT: Three Benchmarks

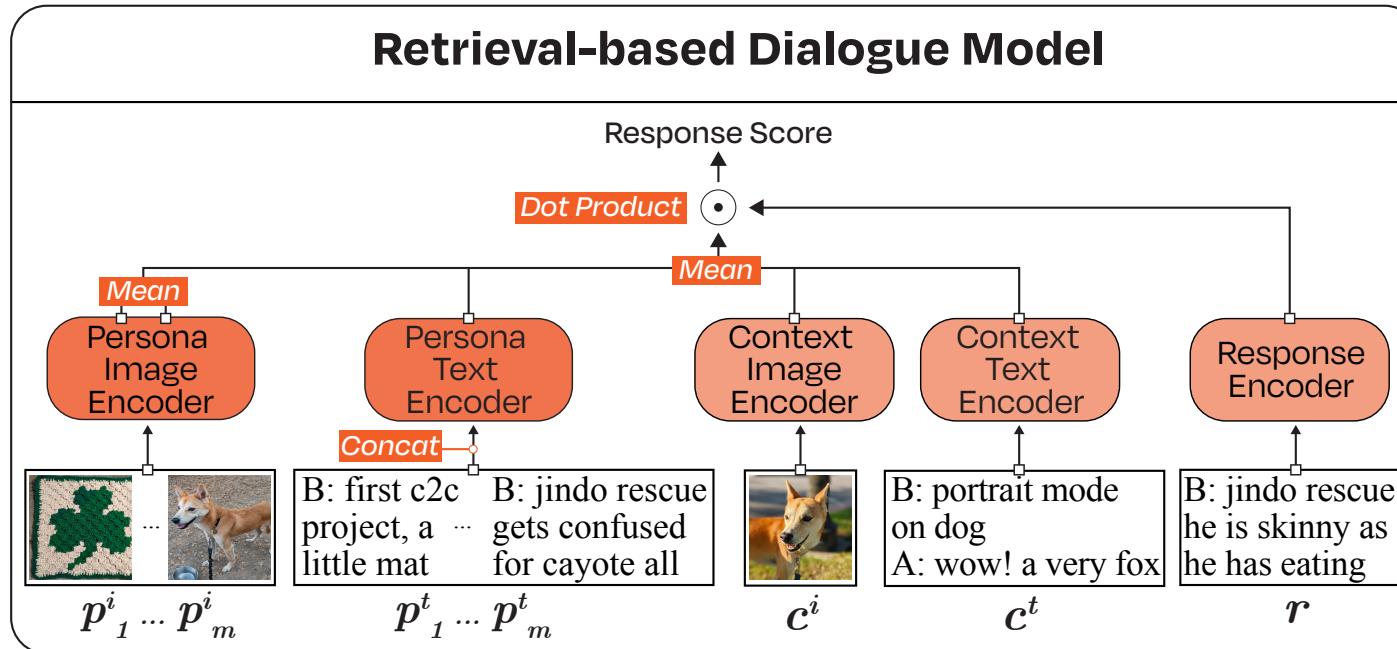
3) Speaker Identification (SI)

- 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^[12], CLIP-ViT-B/32^[13] vision model
 - Text encoder: SBERT^[14], CLIP-ViT-B/32 text model



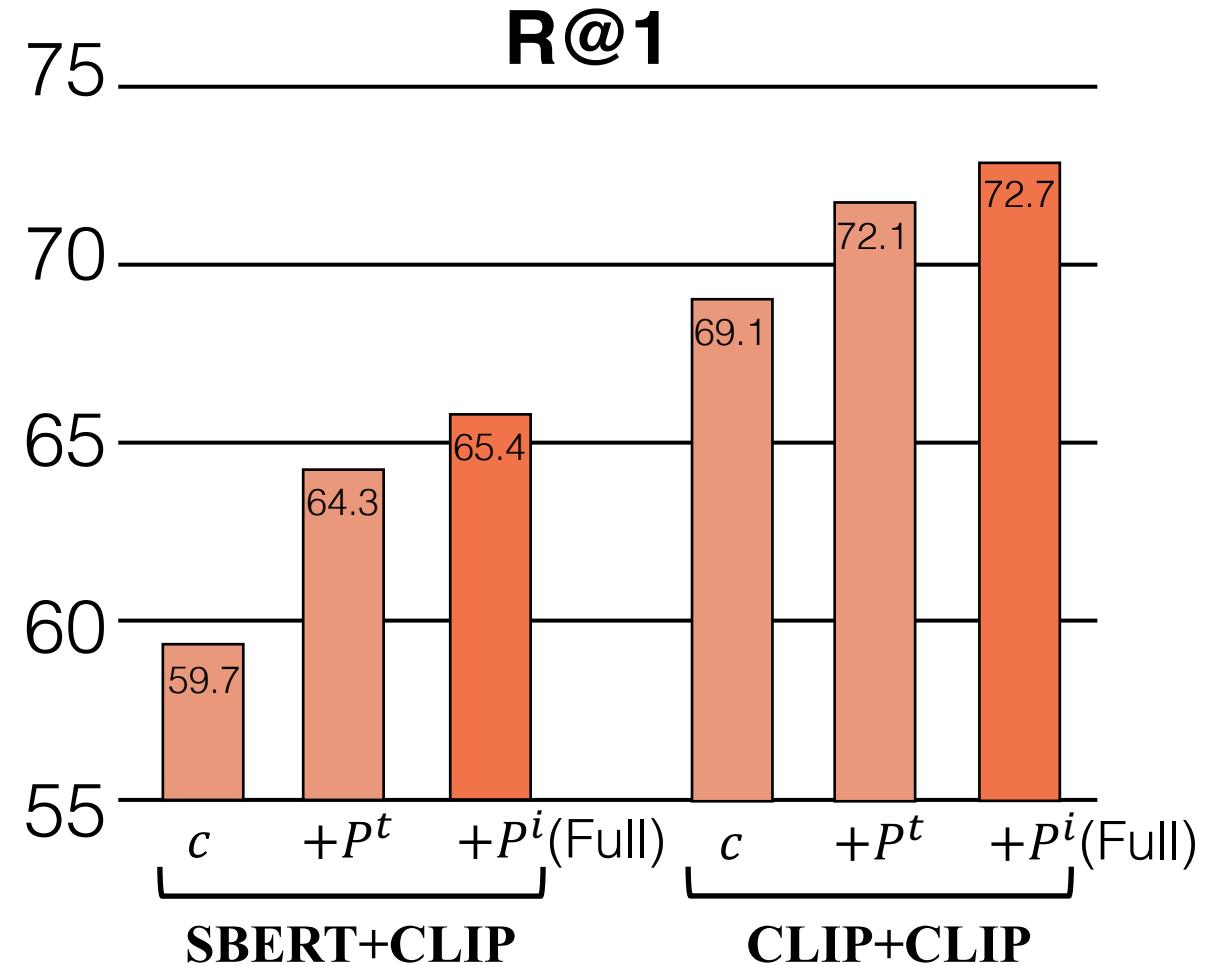
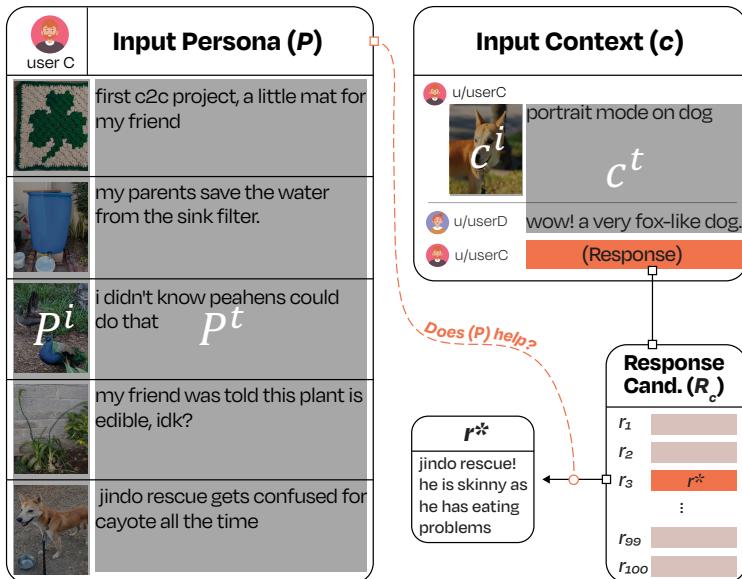
[12] Dosovitskiy et al., *An image is worth 16x16 words: Transformers for image recognition at scale*, ICLR 2021

[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

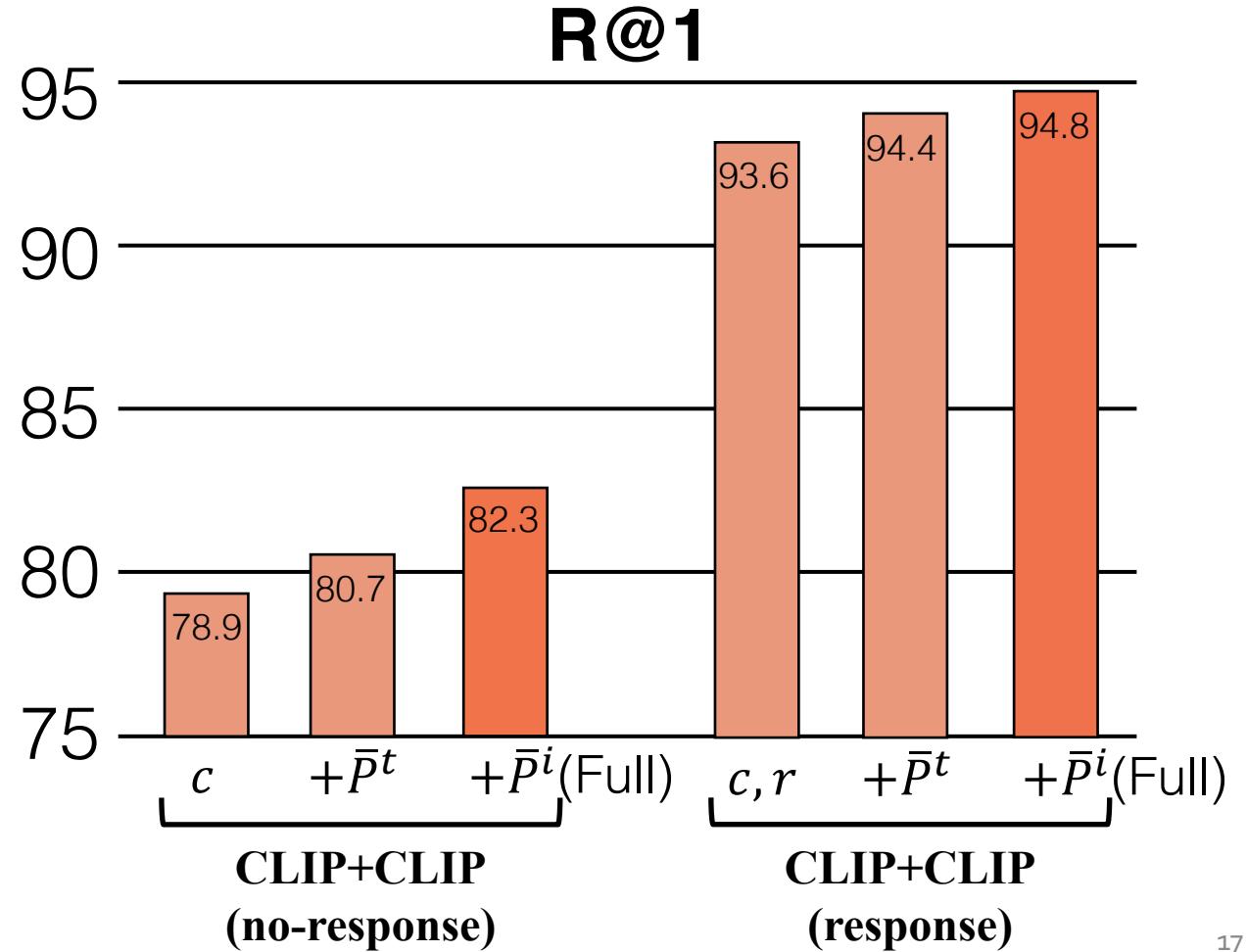
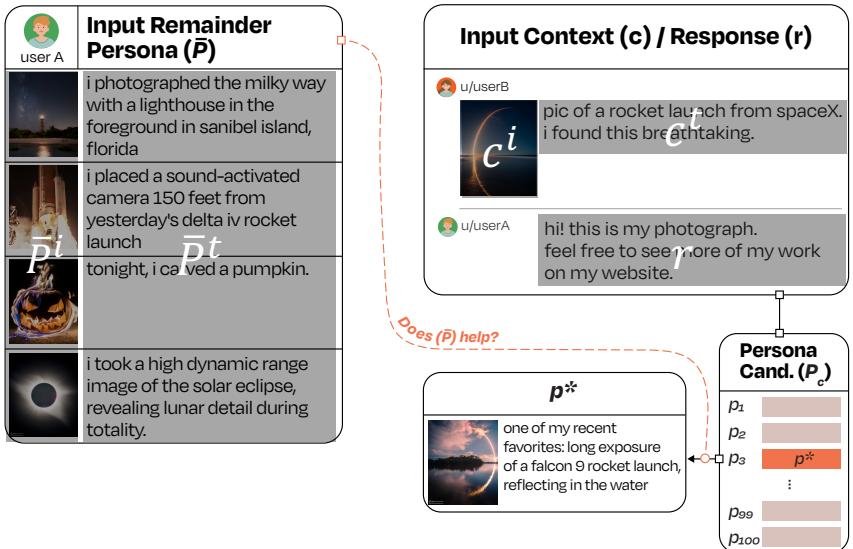
Quantitative Results on NRP

- Model w/ multimodal persona outperforms baseline



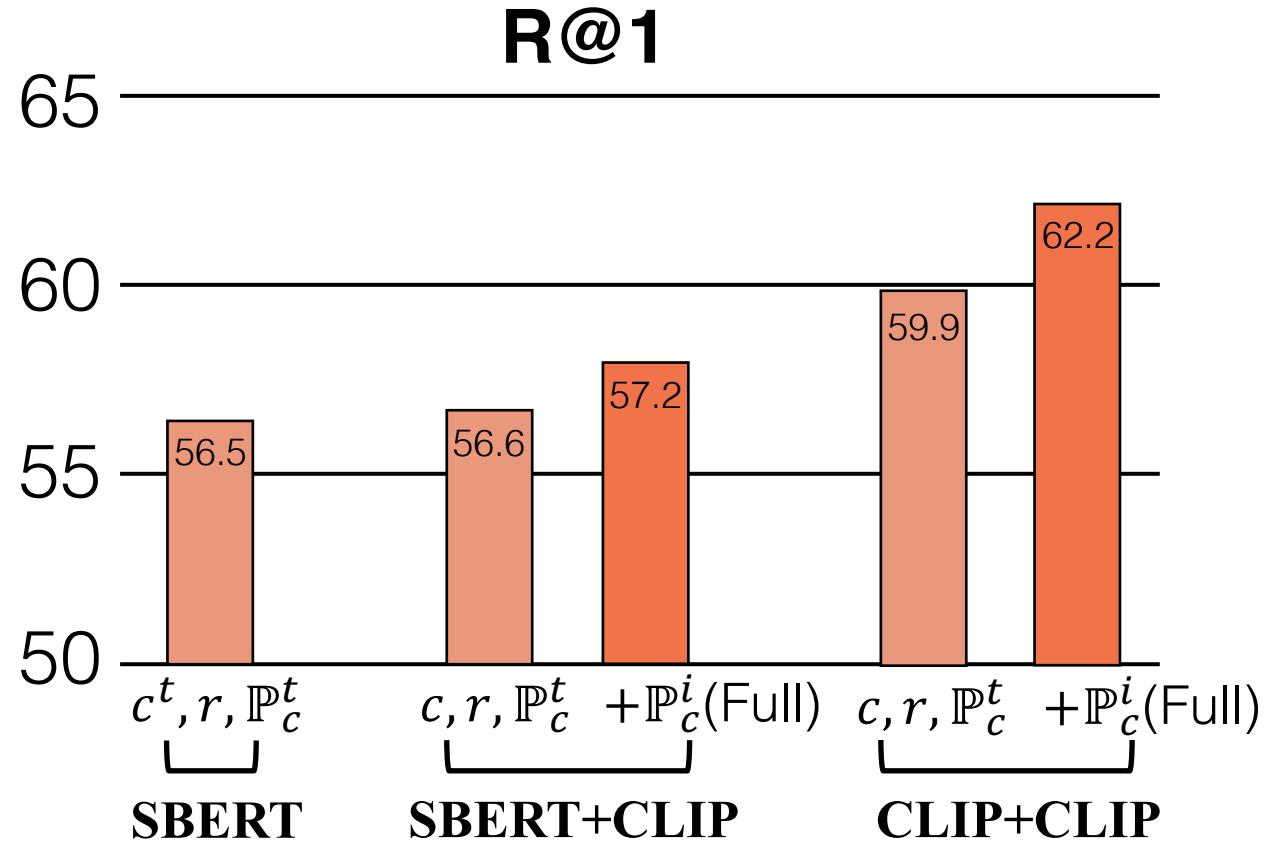
Quantitative Results on GPP

- Model w/ multimodal persona outperforms baseline



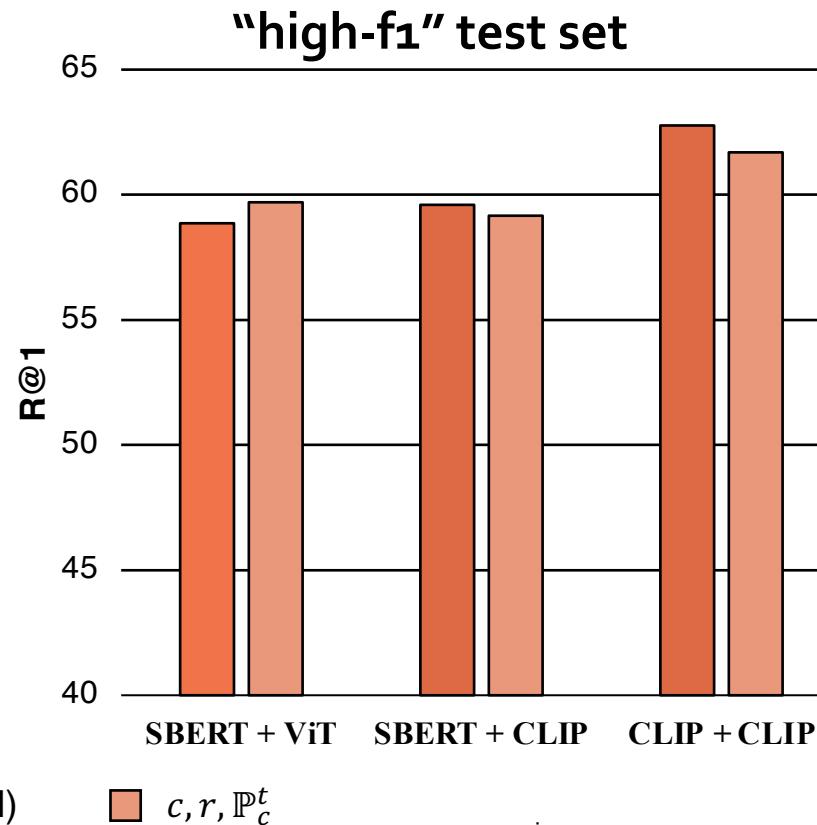
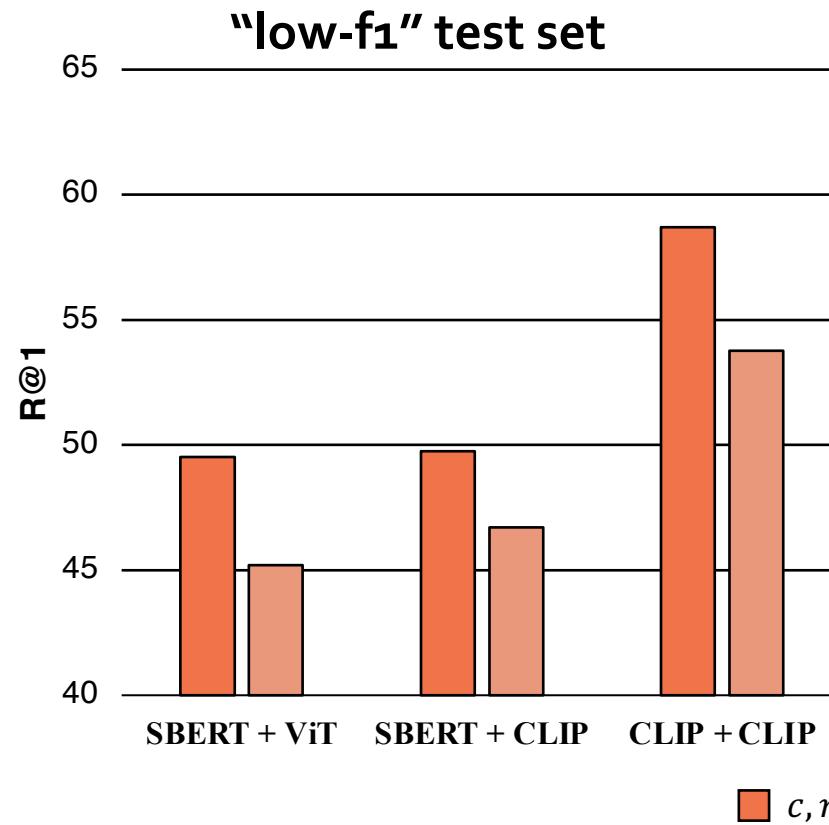
Quantitative Results on SI

- Model w/ multimodal persona outperforms baseline



When is multimodal persona helpful?

- SI: Larger gap in “low-f₁” 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

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

| GPP (no-response) | Context & persona | Persona-only |
|------------------------------|-------------------|--------------|
| Multimodal understanding | 15 (50%) | 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_HAT: Multimodal persona-grounded dialogue dataset & propose three benchmarks: NRP, GPP, SI
- Outperforms baselines on all tasks w/ multimodal persona
 - MPC_HAT 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

