# Deep Reinforcement Learning with Hierarchical Action Exploration for Dialogue Generation

Itsugun Cho[1], Ryota Takahashi[1], Yusaku Yanase[1], Hiroaki Saito[1]

Keio University, Japan[1]

Keio University

CALAMVS GLADIO FORTIOR
1858

# Why use Reinforcement Learning

**Being equipped with a "forward-looking" capacity in the dialogue to enhance user experience.**

# Problem Scenario

Policy $\pi$ : Agent

State $s$ : Context

Action $a$ : Response $\qquad a \rightsquigarrow \pi(a|s)$

Reward $r$ : Reward $\qquad r = R(s, a)$

State $s'$ : Next Context $\quad s' = T(s, a)$

State-value function : $\qquad V^\pi(s)$

Action-value function : $\qquad Q^\pi(s, a)$

# Policy Iteration (Actor-critic)

Approximate Dynamic Programming

$$\text{Evaluation}: Q \leftarrow \underset{Q}{\arg\min} \, \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[ (r + \gamma \mathbb{E}_{a' \sim \pi(a'|s')}[Q(s',a')] - Q(s,a))^2 \right]$$

$$\text{Improvement}: \pi \leftarrow \underset{\pi}{\arg\max} \, \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi(a|s)} \left[ Q(s,a) \right]$$

**Policy Improvement Theorem**

$Q = Q^{\pi}$

$Q, \pi$

$Q^*, \pi^*$

$\pi = greedy(Q)$

# Theorem 1

**Theorem 1.** Given a policy $\pi$ and the number of sampled actions $L$, if we update the new policy by

$$\forall s, \pi'_L = \arg\max_{a \in \{a_i\}_{i=1}^{L} \sim \pi(a|s)} Q^\pi(s, a),$$

then for any $N$, $M$, such that $N \geq M \geq 1$, $\forall s$, $V^{\pi'_N}(s) \geq V^{\pi'_M}(s)$ always holds.

**Theorem 2.** Given the policy $\pi_\alpha$, $\pi_\beta$, and $\pi$, s.t. $\mathbb{E}_{a \backsim \pi_\alpha(a|s)}[Q^\pi(s,a)] \geq \mathbb{E}_{a \backsim \pi_\beta(a|s)}[Q^\pi(s,a)]$, if the number of sampled actions is $L$, and we update the new policy by
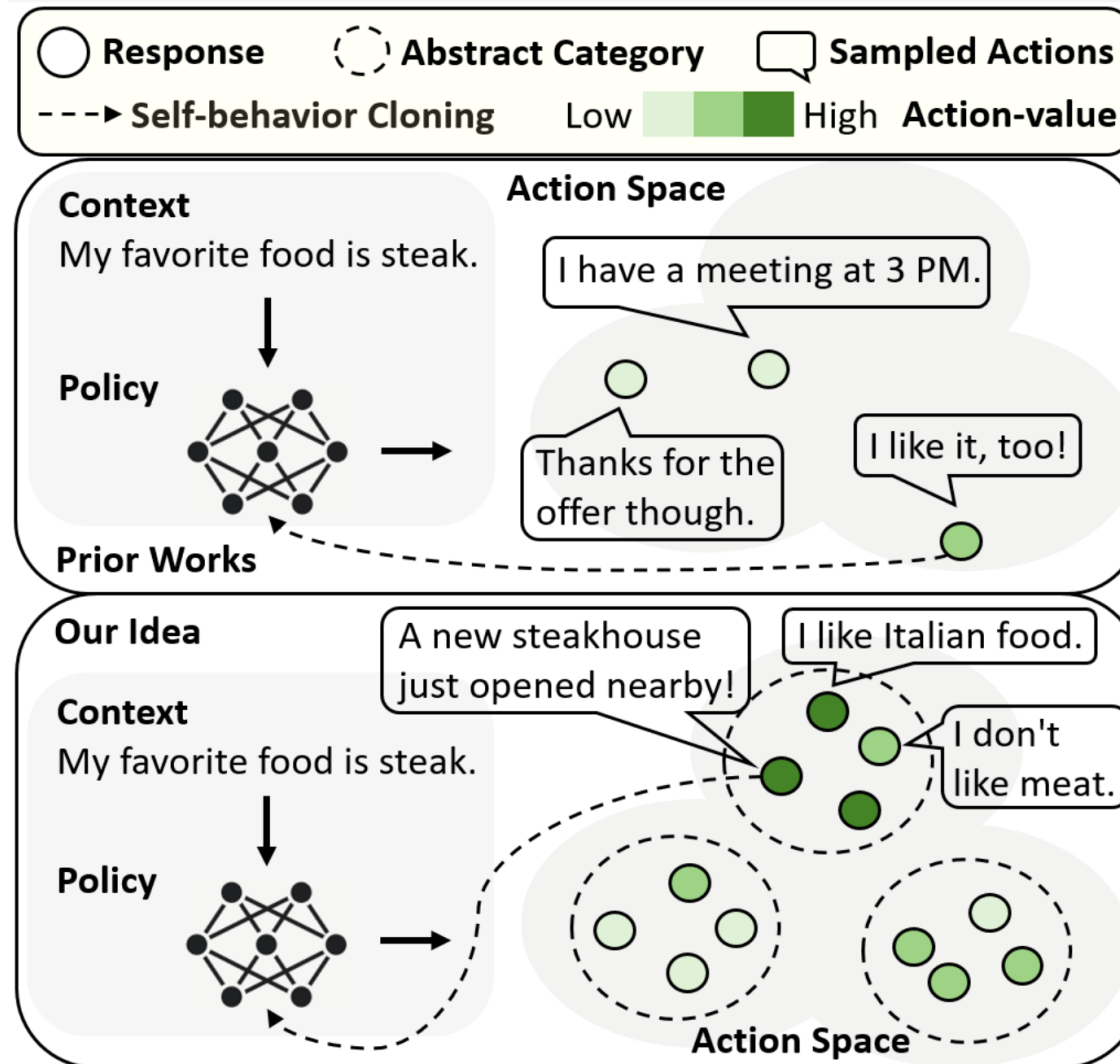
$$\forall s, \ \pi_1' = \underset{a \in \{a_i\}_{i=1}^L \backsim \pi_\alpha(a|s)}{\arg\max} Q^\pi(s,a),$$

$$\forall s, \ \pi_2' = \underset{a \in \{a_i\}_{i=1}^L \backsim \pi_\beta(a|s)}{\arg\max} Q^\pi(s,a),$$

then $\forall s, \ V^{\pi_1'}(s) \geq V^{\pi_2'}(s)$ always holds.

# Dual-granularity Q-function

# Rewards

The average cosine similarity between the agent's response and dull responses.

The outpouring of the emotion of surprise.

The length of response.

Asking questions.

# Corpus and Agent

| DailyDialog | Dialogues |
| --- | --- |
| Training set | 11,118 |
| Testing set | 1,000 |
| Validation set | 1,000 |

| Agent | Descriptions |
| --- | --- |
| GPT-2 | proposed by Radford et al. (2019), is an unsupervised autoregressive language model for textual generation. |
| DialoGPT | is a pre-trained di alogue model proposed by Zhang et al. (2020). This model is based on GPT-2, using the Reddit comments dataset. |
| T5 | is a unified framework pro posed by Raffel et al. (2020) that converts all text based language tasks into a text-to-text format via the transfer learning technique. |
| GODEL | is a pre trained dialogue model proposed by Peng et al. (2022). This model is based on T5, using the Red dit discussion dataset. |

# Automatic Metrics

**CS:** The average **c**osine **s**imilarity between the agent's response and dull responses

**SE:** The outpouring of the **e**motion of **s**urprise

**LR:** The **l**ength of **r**esponse

**AQ:** **A**sking **q**uestions

# Human Metrics

**Quality** measures the coherence and grammatical accuracy of the agents' responses

**Informativeness** measures the diversity and hallucination of the agents' responses

**Empathy** measures the degree to which agents respond with concern or affectivity

**Engagingness** measures the desire to engage the agents in a long conversation

# Experiments

| Agent | Training Method | Dataset-based | | | | Simulator-based | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | CS ↓ | SE | LR | AQ | CS ↓ | SE | LR | AQ |
| GPT-2 | MLE | 0.712 | 0.082 | 10.396 | 0.308 | 0.685 | 0.146 | 11.276 | 0.390 |
| | Standard | 0.645 | 0.126 | 13.020 | 0.550 | 0.644 | 0.206 | 13.778 | 0.526 |
| | Ours | **0.596** | **0.191** | **14.463** | **0.555** | **0.597** | **0.238** | **15.636** | **0.566** |
| DialoGPT | MLE | 0.714 | 0.069 | 9.761 | 0.345 | 0.687 | 0.142 | 10.838 | 0.492 |
| | Standard | 0.645 | 0.142 | 12.182 | 0.579 | 0.654 | 0.206 | 13.772 | 0.538 |
| | Ours | **0.598** | **0.171** | **13.055** | **0.586** | **0.588** | **0.240** | **14.466** | **0.604** |
| T5 | MLE | 0.720 | 0.063 | 9.704 | 0.316 | 0.651 | 0.088 | 10.242 | 0.396 |
| | Standard | 0.621 | 0.147 | 13.291 | 0.532 | 0.605 | 0.224 | 13.676 | 0.510 |
| | Ours | **0.567** | **0.202** | **14.834** | **0.565** | **0.553** | **0.268** | **15.134** | **0.552** |
| GODEL | MLE | 0.718 | 0.064 | 9.507 | 0.318 | 0.689 | 0.112 | 10.132 | 0.414 |
| | Standard | 0.625 | 0.165 | 13.553 | 0.529 | 0.615 | 0.235 | 13.108 | 0.614 |
| | Ours | **0.571** | **0.232** | **15.272** | **0.557** | **0.571** | **0.258** | **14.608** | **0.628** |

Table 1: Automatic evaluation results. For the standard offline RL algorithm and our approach, we use $L = 5$ for the number of candidate responses $\{a_i\}_{i=1}^{L}$. For the simulator-based evaluation, we conducted 1000 dialogues of 5 consecutive turns between the simulator and each method. Each metric is measured per response, and the best score in each metric is in bold. The statistical test revealed that the differences are significant, with a p-value < 0.05.
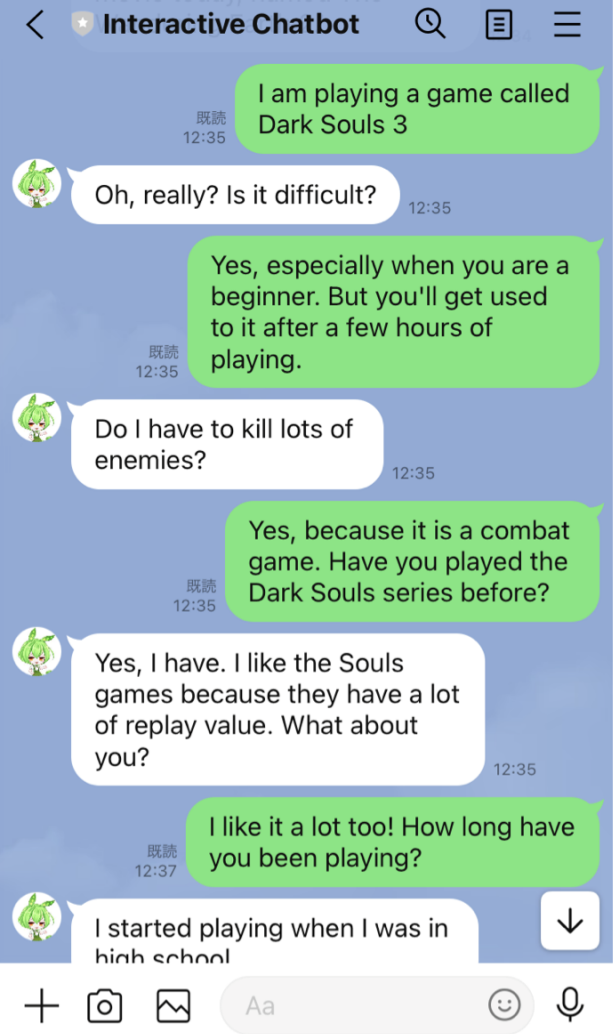
# Experiments

| Agent | Training Method | Quality | Informativeness | Empathy | Engagingness |
|-------|-----------------|---------|-----------------|---------|--------------|
| GPT-2 | MLE | 1.4 | 1.3 | 1.2 | 1.1 |
| | Standard | **1.7** | 1.3 | 1.4 | 1.4 |
| | Ours | 1.5 | **1.5** | **1.5** | **1.6** |
| DialoGPT | MLE | 1.3 | 1.1 | 0.7 | 0.7 |
| | Standard | **1.5** | 1.4 | 1.2 | 1.2 |
| | Ours | 1.4 | **1.5** | **1.6** | **1.6** |
| T5 | MLE | 1.2 | 0.9 | 0.5 | 0.6 |
| | Standard | 1.1 | 0.8 | 0.6 | 0.7 |
| | Ours | **1.4** | **1.4** | **1.4** | **1.3** |
| GODEL | MLE | 1.5 | 1.3 | 0.8 | 1.0 |
| | Standard | 1.6 | 1.2 | 1.1 | 1.1 |
| | Ours | **1.7** | **1.6** | **1.7** | **1.6** |

Ten native speakers were recruited for human evaluation. The scale of these metrics is [0, 1, 2].

Table 2: Human evaluation results. The final scores for each metric were calculated by taking the average of the annotator ratings. Each metric is measured per dialogue, and the best score in each metric is presented in bold. The Fleiss' kappa (Fleiss, 1971) score with human judges was approximately 0.29,
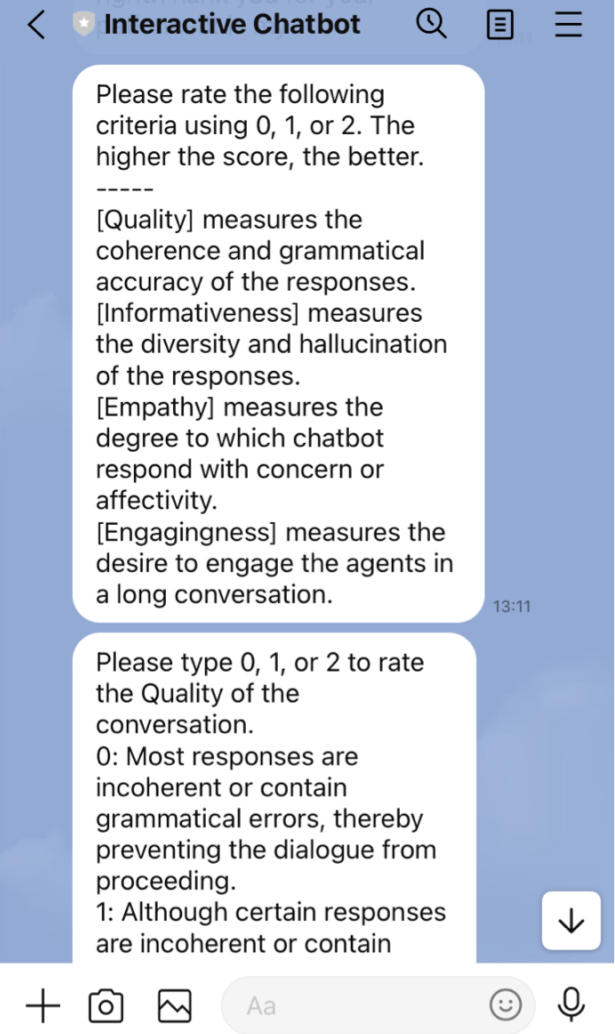
# Details regarding Interactive



(a) Chat Interface 1
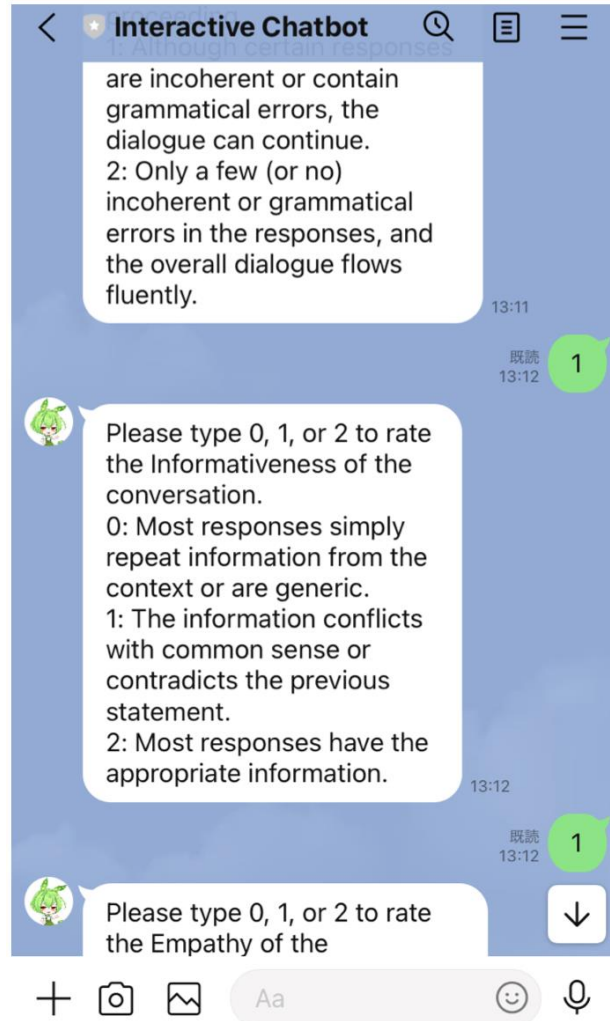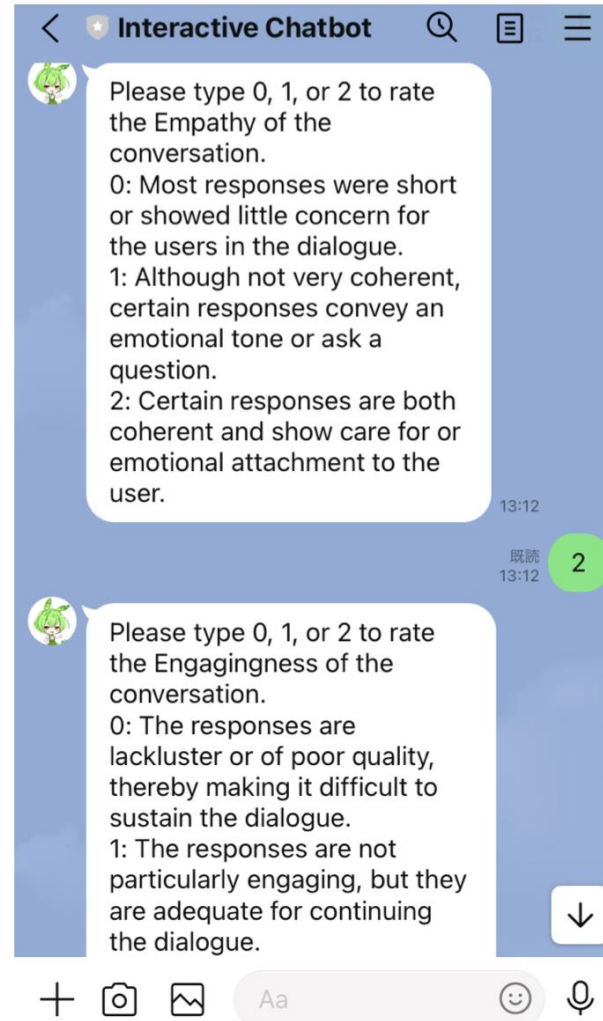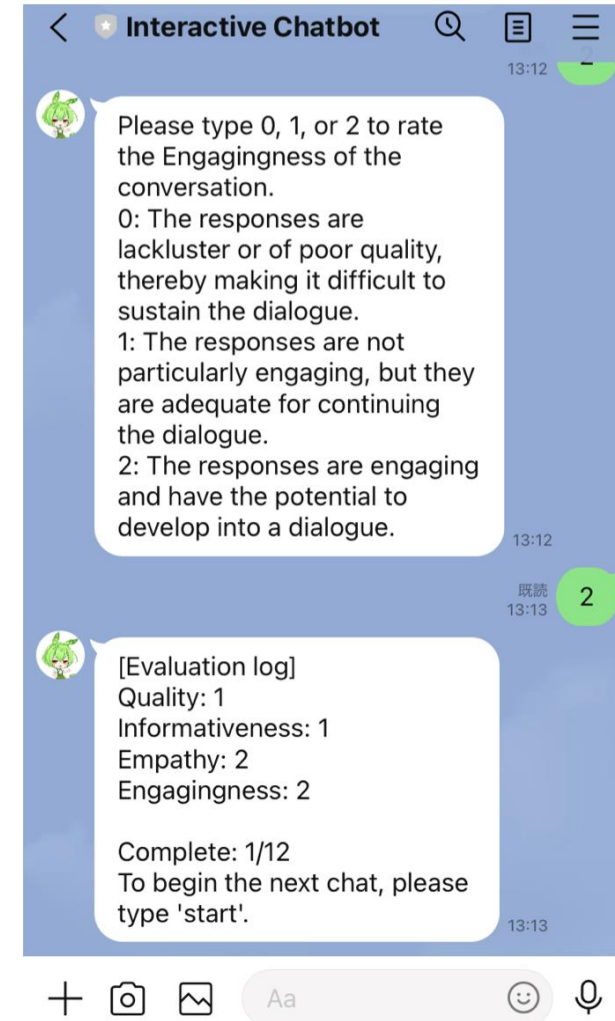
(b) Chat Interface 2

(c) Ratings Interface 1

# Details regarding Interactive



(d) Ratings Interface 2

(e) Ratings Interface 3
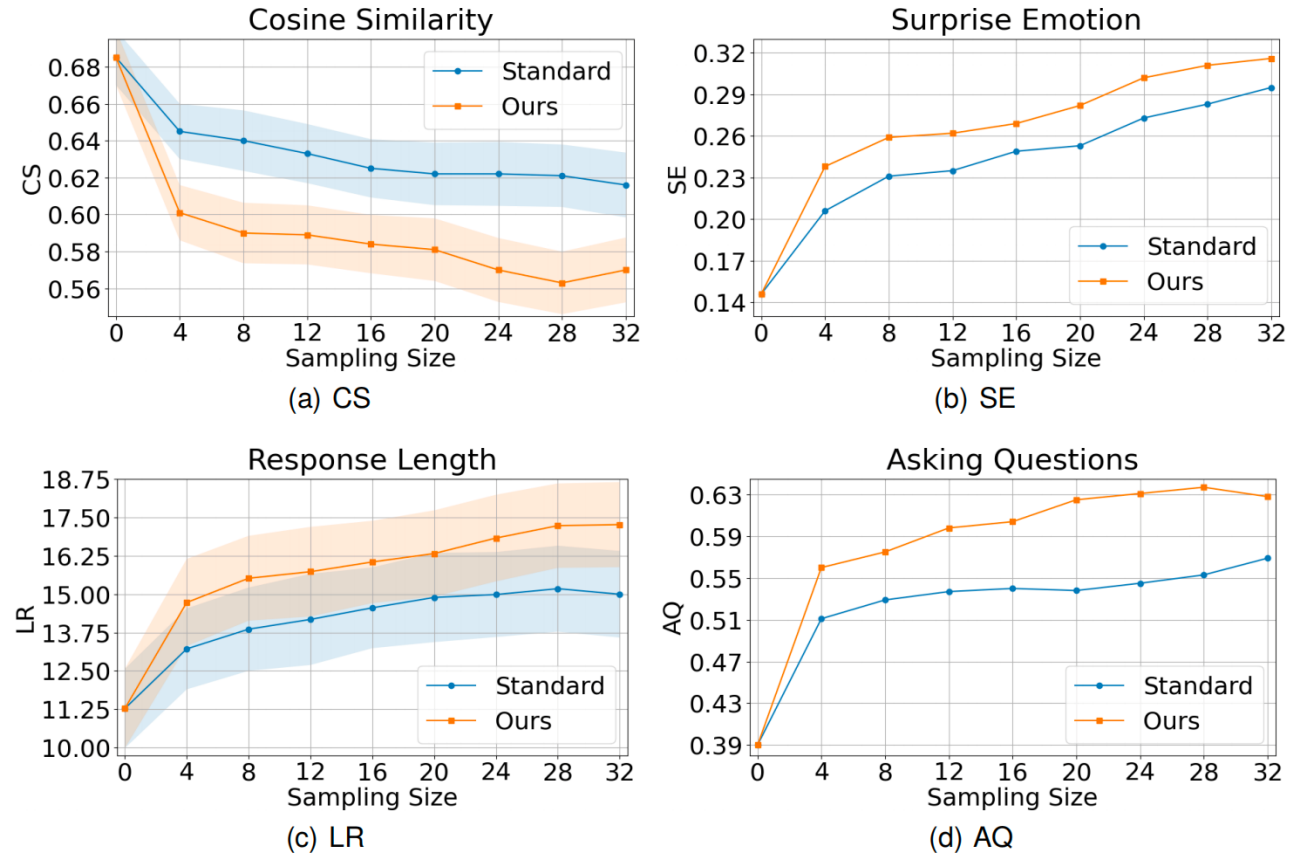
(f) Ratings Interface 4

# Further Verification



Figure 3: The evolution of the agent's performance for each metric with the increased sampling size. The scale for the X-axis is a multiple of 4, and 0 represents the MLE without RL. Bands indicate half a standard deviation.
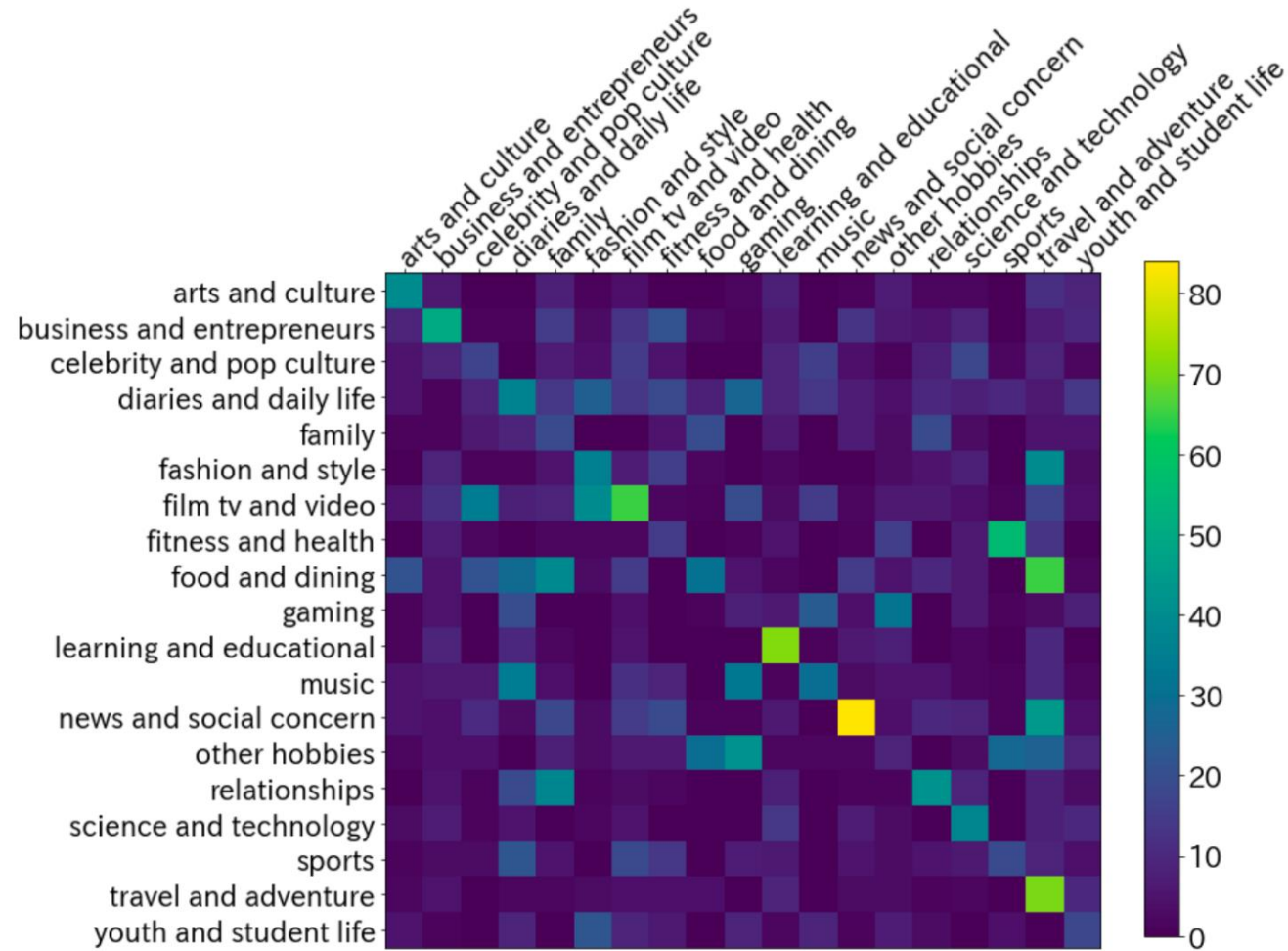
# Further Verification



Figure 4: Behavioral analysis of coarse-grained Q-function. The label on the Y-axis represents the topic of each context, whereas the label on the X-axis represents the selected action category. The color bar represents the number of occurrences of the corresponding selection.

# Further Verification

| | | | |
|---|---|---|---|
| **1st Turn** | **User: I really need to start eating healthier.**<br>MLE: You need to start eating more.<br>Stan: Yeah, it is really good for your health.<br>Ours: Really? How are you going to make it? | **1st Turn** | **User: Hi, Dan. What's new?**<br>MLE: What are you doing the last week?<br>Stan: I got a visa for travel to Singapore.<br>Ours: I was told that our TV show was canceled in next week. Didn't you see that? |
| **2nd Turn** | **User: I have to start eating better too.**<br>MLE: That's the hard part.<br>Stan: Wow, can I ask you for advice to weight loss?<br>Ours: Maybe you need to eat something that isn't fried food. There are so many bad things in fried food! | **2nd Turn** | **User: Where have you been? I've been trying to get in touch with you for two days.**<br>MLE: I was out of town. Have you seen them?<br>Stan: Hi, John. When did you make contact?<br>Ours: Wow! Now what happened? |
| **3rd Turn** | **User: What kind of food do you usually eat?**<br>MLE: I usually eat either fish or chicken.<br>Stan: I've been eating basically anything with whatever vegetables are around.<br>Ours: I mostly eat Chinese food. Do you get that often? | **3rd Turn** | **User: I went to Salt Lake City on business . What's up ?**<br>MLE: I'm leaving tomorrow.<br>Stan: I'd like to ask you if you'd like to join me on a trip to Singapore.<br>Ours: Last week our business meeting was held, I would like to get some information from you. |
| **4th Turn** | **User: I try my best to eat only fruits, vegetables, and chicken.**<br>MLE: That sounds like a lot of food. Is that what you usually eat?<br>Stan: When you eat these foods, do you usually feel empty?<br>Ours: Oh! you should make sure you are consuming plenty of protein and vitamins. | **4th Turn** | **User: I got fired.**<br>MLE: I'm sorry, I was fired last month.<br>Stan: Did not you expect to be fired then if this happens again?<br>Ours: Whoa, it's just over a month in, tired of your job already? |
| | (a) DialoGPT | | (b) T5 |

Figure 2: Case study. DialoGPT and T5 are used as the agents in the case study. For brevity, the standard method is abbreviated as stan. The ground truth for responding to each user utterance is the next user utterance in the dialogue.

# Acknowledgement

**Thank you for your listening!**