

1 We thank all the reviewers for their positive comments, and address their major questions and comments below.
2 Clarifications will be added in the revision and we will keep improving our draft.

3 **Reviewer #1** We thank the reviewer for the positive reviews. The remarks raised are addressed below.

4 **Q: More details about model training**

5 **A:** λ_{\max} is selected via parameter search for $\lambda_{\max} \in \{0.1, 0.5, 1.0, 2.0, 5.0\}$, ending up with $\lambda_{\max} = 1.0$. For $\Gamma_{\theta}(\cdot)$,
6 we use PPO in our implementation to keep the stability as the parameters are updated within a ‘trust region’.

7 We are happy to release our code for better reproducibility.

8 **Reviewer #2** We appreciate reviewer’s acknowledgement of our novelty and suggestions provided.

9 Thanks for the suggestion on the paper improvements. In our updated version, we will add more results to show how
10 our approach can handle the more useful dialogue generation problem.

11 **Reviewer #3** We thank the reviewer for the positive reviews and appreciate the reviewer’s suggestions.

12 **Q: How to ensure that the sentences are short and meaningful.**

13 **A:** Our paper proposed a general constraint-augmented reinforcement learning framework. In the recommendation task,
14 natural language is only an example way for the user to express the preferences (*i.e.*, constraints), and we do not focus
15 too much on how to handle more complicated or even free-form language. Therefore, similar to [20], our sentences in
16 experiments are in the format of simple sentences with prefix. In this setting, the length of the sentences is implicitly
17 determined by the training data. That is, we train the GRU model (the user simulator) on short sentences collected by
18 human, and thus the trained model usually generates short sentences.

19 **Q: Whether historical behaviors are considered.**

20 **A:** Our approach considers the user historical behaviours within the current user session, although it does not model
21 the user historical behaviours in previous sessions. The GRU tracks the user behaviours and the discriminator considers
22 all previous user preferences in the current user session. However, we do agree that user historical behaviours from
23 previous sessions can be employed to further enhance the performance.

24 **Q: Comparison with traditional methods considering historical behaviors.**

25 **A:** Our approach is not directly comparable to the traditional methods considering historical behaviors. Traditional
26 recommendation models are usually trained in an offline manner, *i.e.*, the model is trained on a pre-collected dataset.
27 By contrast, our method is proposed in a different setting, where the pre-collected dataset is not required and our
28 recommender is interactively learned when the user interacts with the system. Moreover, it is not clear how to handle
29 the interactive natural language feedback to provide interactive recommendations by traditional models.

30 **Q: Incorporating traditional methods into our framework.**

31 **A:** A simple approach can be developed to incorporate the traditional recommendation methods into our framework,
32 to leverage historical behaviours from previous sessions. Assume we have users’ historical behavior data and train a
33 traditional model. In each user session, we make initial recommendations by the traditional model, collect natural-
34 language feedback from users, and make further interactive recommendations by our framework. After a number of
35 user sessions, we can update the traditional model based on the recently collected users’ feedback to the items.

36 Thank you very much for providing us with the related work [49, 50, 51, 52]. We will definitely discuss these references
37 in our related work section, and fix all the typos in our minor revision.

38 [49] Konstantina Christakopoulou, Alex Beutel, Rui Li, Sagar Jain, and Ed H Chi. Q&R: A two-stage approach toward
39 interactive recommendation. In KDD, 2018.

40 [50] Yu Zhu, Yu Gong, Qingwen Liu, Yingcai Ma, Wenwu Ou, Junxiong Zhu, Beidou Wang, Ziyu Guan, and Deng Cai.
41 Query-based interactive recommendation by meta-path and adapted attention-gru.arXiv:1907.01639, 2019.

42 [51] Yu Zhu, Hao Li, Yikang Liao, Beidou Wang, Ziyu Guan, Haifeng Liu, and Deng Cai. What to do next: modeling
43 user behaviors by time-lstm. In IJCAI, 2017.

44 [52] Yu Zhu, Junxiong Zhu, Jie Hou, Yongliang Li, Beidou Wang, Ziyu Guan, and Deng Cai. A brand-level ranking
45 system with the customized attention-gru model. In IJCAI, 2018.