

# Parameter-Efficient Conversational Recommender System as a Language Processing Task

**Mathieu Ravaut**<sup>1</sup>, Hao Zhang<sup>1</sup>, Lu Xu<sup>2</sup>, Aixin Sun<sup>1</sup>, Yong Liu<sup>1</sup>

<sup>1</sup>Nanyang Technological University, Singapore

<sup>2</sup>Singapore University of Technology and Design, Singapore

# Conversational Recommender Systems (CRS)

Conversational Recommender Systems (CRS) jointly generate a natural language response to the user (**conversation task**) and recommend a list of items (**recommendation task**).

CRS approaches can be roughly divided into two categories :

- *Attribute-based* CRS : collect user preference on items attributes.
- *Generation-based* CRS : acquire feedback from users through language and generate natural responses.

We are focusing on **generation-based** CRS in this work.

# Challenges in CRS

CRS are challenging to build because item recommendation and language generation are two tasks of very different nature.

A long line of work relies on knowledge graphs to learn items representation [1, 24, 23]. Unfortunately, there are a few issues :

- Because they are learned separately, word representations and items representations are semantically misaligned.
- KG consist in an external source of knowledge, which may not be readily available in certain inference setups.
- This approach neglects rich text information available for items.

# Unifying CRS through language models

The recent MESE [21] approach uses pre-trained language models to learn items representations, and integrates them within the language response, bypassing the need for knowledge graph. However, it still relies on several models (two DistilBERT [15] and a GPT-2 [14]).

Overall, there does not exist yet a *truly unified* CRS model :

- UniCRS [20] uses a language model and a knowledge graph, and requires three training stages.
- BARCOR [19] and RecInDial [18] train in a single stage, but still need both a language model and a knowledge graph.
- MESE [21] discards the knowledge graph and still trains in a single stage, but relies on several pre-trained language models.

# Proposal

In this work, we push simplicity and unification to its finest and solve the CRS task with a **single pre-trained language model (LM) fine-tuned in a single stage, without using a knowledge graph (KG)**.

Besides, through parameter-efficient fine-tuning, we only update a small fraction of parameters.

# Dialogue Modeling

Let  $\mathcal{I} = \{I_1, I_2, \dots, I_{N_{\text{item}}}\}$  be the item database with  $N_{\text{item}}$  items.

Let  $\mathcal{D} = \{D_1, D_2, \dots, D_{N_{\text{dial}}}\}$  be the dataset with  $N_{\text{dial}}$  dialogues.

Let  $D = \{u_t\}_{t=1}^{n_{\text{utt}}}$  be a dialogue with  $n_{\text{utt}}$  utterances.

Conditioning on the dialogue history  $D_t = \{u_i\}_{i=1}^{t-1}$ , CRS predicts :

- The current utterance  $u_t = \{w_j\}_{j=1}^n$ , with  $n$  tokens.
- The set of recommended items  $\mathcal{I}_t$ , which may be empty.

Utterances are produced by the *seeker* or the *recommender*.

CRS only predicts the *recommender* utterances.

We use a decoder-only Transformer LM enhanced with special tokens :  
 “[ITEM]”, “[SEP]”, “[REC]” and “[REC\_END]”.

# Items Representation

We use the LM for both dialogue response and items representation.

Each movie item is described with a text in the template “*Movie title [SEP] Actors [SEP] Director(s) [SEP] Genre(s) [SEP] Plot*”.

We add an item head  $h_{\text{item}}$  and learnable pooling weight  $w$  to the LM. The  $j$ -th item representation is :

$$\mathbf{v}_j = h_{\text{item}}(w^T \cdot \mathbf{I}_j). \quad (1)$$

where  $\mathbf{I}_j$  is the LM contextual representation of the description.

# Context Representation

For each utterance of the context  $D_t = \{u_i\}_{i=1}^{t-1}$  we obtain contextual representation with the LM :  $\mathbf{u}_i = [\mathbf{c}_{i,1}, \dots, \mathbf{c}_{i,n}]$ .

Item names are replaced by the “[ITEM]” special token.

If the utterance is from the speaker, it becomes

$$\bar{\mathbf{u}}_i = \tilde{\mathbf{u}}_i = [\mathbf{v}_{\text{sep}}, \mathbf{v}_j, \mathbf{v}_{\text{sep}}, \mathbf{u}_i].$$

If it is from the recommender, it becomes

$$\bar{\mathbf{u}}_i = \tilde{\mathbf{u}}_i = [\mathbf{v}_{\text{rec}}, \mathbf{v}_j, \mathbf{v}_{\text{rec\_end}}, \mathbf{u}_i].$$

If there is no recommended item, it remains unchanged  $\bar{\mathbf{u}}_i = \mathbf{u}_i$ .

Dialogue representation is  $\mathbf{D}_t = [\bar{\mathbf{u}}_1, \dots, \bar{\mathbf{u}}_{t-1}, \mathbf{v}_{\text{rec}}]$ , and we use the output of the last “[REC]” token, noted  $\mathbf{d}_t$ .



# Recommendation task (1/2)

We use a **contrastive learning** approach [4, 13, 12] to bring closer the query  $\mathbf{d}_t$  and the positive item  $\mathbf{v}_p$ ; while pushing apart  $\mathbf{d}_t$  and  $M$  sampled negative items  $\{\mathbf{v}'_j\}_{j=1}^M$ .

$$\mathcal{L}_{\text{recall}} = -\frac{1}{|\mathcal{D}|} \sum_{D_t \in \mathcal{D}} \log(\mathcal{E}_{D_t}). \quad (2)$$

where :

$$\mathcal{E}_{D_t} = \frac{e^{f(\mathbf{d}_t)^\top \odot \mathbf{v}_p}}{e^{f(\mathbf{d}_t)^\top \odot \mathbf{v}_p} + \sum_{(\mathbf{d}_t, \mathbf{v}'_j) \sim \mathcal{N}} e^{f(\mathbf{d}_t)^\top \odot \mathbf{v}'_j}}, \quad (3)$$

where where  $f$  is a projection head MLP.

We **stop the gradients** of LM and only optimize the pooling ( $w$ ) and MLP layers ( $h_{\text{item}}, f$ ).

# Recommendation task (2/2)

To refine item selection, we use a **re-ranking** approach.

We concatenate the context and all items,  $[D_t, v_p, v'_1, \dots, v'_M]$ .

This input is fed into LM then MLP  $f$ , with attention mask blocking attention between items, yielding representations  $[q_p, q_1, \dots, q_M]$ .

Another MLP layer  $g$  is applied to compute the final item scores as  $\mathbf{r} = [g(q_p), g(q_1), \dots, g(q_M)] = [r_p, r_1, \dots, r_M]$ .

Items are re-ranked through a cross-entropy loss :

$$\mathcal{L}_{\text{rerank}} = \frac{1}{|\mathcal{D}|} \sum_{D_t \in \mathcal{D}} f_{\text{XE}}(\mathbf{r}, \mathbf{Y}), \quad (4)$$

where  $\mathbf{Y} = [1, 0, \dots, 0]$  and  $f_{\text{XE}}$  denotes cross-entropy loss.

# Response generation task

If  $u_t$  contains an item to be recommended, it is appended to the context :

$$\tilde{\mathbf{D}}_t = [\bar{\mathbf{u}}_1, \dots, \bar{\mathbf{u}}_{t-1}, \mathbf{v}_{\text{rec}}, \mathbf{v}_p, \mathbf{v}_{\text{rec\_end}}]. \quad (5)$$

otherwise,  $\tilde{\mathbf{D}}_t = [\bar{\mathbf{u}}_1, \dots, \bar{\mathbf{u}}_{t-1}]$ .

Response generation is optimized by the standard next-token prediction objective :

$$\mathcal{L}_{\text{gen}} = -\frac{1}{|\mathcal{D}|} \sum_{D_t \in \mathcal{D}} \frac{1}{n} \sum_{j=1}^n \log(p_{\theta}(w_j | w_{1:(j-1)}, \tilde{\mathbf{D}}_t)). \quad (6)$$

# Parameter-efficiency

We keep the backbone LM frozen, and instead add LoRA [7] layers to be updated.

This prevents catastrophic forgetting of the LM's text generation capability, while adapting it to the CRS task [8, 22].

The only learnable weights are : task-specific MLP layers  $f$ ,  $g$ ,  $h_{\text{item}}$ , pooling weights  $w$ , and the special tokens embeddings.

Our model is dubbed Parameter-Efficient Conversational Recommender System (**PECRS**).

# Training

We train in a *single-stage* end-to-end manner by minimizing the following loss :

$$\mathcal{L} = \alpha \times \mathcal{L}_{\text{recall}} + \beta \times \mathcal{L}_{\text{rerank}} + \gamma \times \mathcal{L}_{\text{gen}}, \quad (7)$$

During training :

- Sample  $M_{\text{train}}$  negative items, and **share them across losses**  $\mathcal{L}_{\text{recall}}$  and  $\mathcal{L}_{\text{rerank}}$  and **across batch items**.
- Append the ground-truth item to the dialogue context.

# Inference

During inference :

- Encode every single item.
- Retrieve the closest  $M_{\text{inference}}$  items to the dialogue query via  $f(\mathbf{d}_t)^\top \odot \mathbf{v}_j$ .
- Re-rank them and output the highest score one as prediction.
- Append the predicted item to the context.
- The presence of “[ITEM]” in the generated response assesses recommendation.

# Experimental Setup

We apply PECRS to **movie** recommendation on ReDial [9] and INSPIRED [5] datasets.

For the backbone LM, we use GPT-2 (**PECRS-small**) and GPT-2-medium (**PECRS-medium**).

We train with AdamW and LR as  $3e - 5$ , warming up one epoch.

We set  $M_{\text{train}} = 150$  for training and  $M_{\text{infer}} = 700$  for inference.

We balance losses with  $\alpha = 0.15$ ,  $\beta = 0.85$ , and  $\gamma = 1.0$ .

# Evaluation Setup

We measure **recommendation** performance with *Recall@K* ( $R@K$ ) metric, taking  $K \in \{1, 10, 50\}$  and *Unique*, the number of unique recommended items throughout the test set.

We measure **conversation** with *Perplexity* ( $PPL$ ) (fluency), *Distinct@K* ( $Dist@K$ ) with  $K \in \{2, 3, 4\}$  (diversity),  $F-1$  score of the presence of "[ITEM]" (recommendation decision) and *ROUGE-K* ( $RG-K$ ), taking  $K \in \{1, 2\}$  (closeness to the ground truth).



# Recommendation Results

Model	Metadata			Model Properties			ReDial				INSPIRED				
	KG	Reviews	Description	Extra	Model	PEFT	Rounds	R@1	R@10	R@50	Unique	R@1	R@10	R@50	Unique
ReDial (Li et al., 2018)	✗	✗	✗	✓	✗	✗	3	2.4	14.0	32.0	–	–	–	–	–
KBRD (Chen et al., 2019)	✓	✗	✗	✓	✗	✗	2	3.0	16.3	33.8	–	–	–	–	–
KGSF (Zhou et al., 2020a)	✓	✗	✗	✓	✗	✗	3	3.9	18.3	37.8	–	–	–	–	–
KECRS (Zhang et al., 2022)	✓	✗	✗	✓	✗	✗	2	2.3	15.7	36.6	–	–	–	–	–
BARCOR (Wang et al., 2022b)	✓	✗	✗	✓	✗	✗	1	2.5	16.2	35.0	–	–	–	–	–
UniCRS (Wang et al., 2022c)	✓	✗	✗	✓	✓	✗	3	5.1	22.4	42.8	–	<b>9.4</b>	<b>25.0</b>	<b>41.0</b>	–
RecInDial (Wang et al., 2022a)	✓	✗	✗	✓	✗	✗	1	3.1	14.0	27.0	–	–	–	–	–
VRICR (Zhang et al., 2023b)	✓	✗	✗	✓	✗	✗	3	5.7	25.1	41.6	–	–	–	–	–
RevCore (Lu et al., 2021)	✓	✓	✗	✓	✗	✗	2	<b>6.1</b>	23.6	<b>45.4</b>	–	–	–	–	–
C <sup>2</sup> -CRS (Zhou et al., 2022)	✓	✓	✗	✓	✗	✗	2	5.3	23.3	40.7	–	–	–	–	–
MESE (Yang et al., 2022)	✗	✗	✓	✓	✗	✗	1	5.6	<b>25.6</b>	<b>45.5</b>	–	4.8	13.5	30.1	–
PECRS-small	✗	✗	✓	✗	✓	✗	1	4.7	20.8	40.5	<b>463</b>	5.4	16.1	33.3	<b>34</b>
PECRS-medium	✗	✗	✓	✗	✓	✗	1	<b>5.8</b>	22.5	41.6	<b>634</b>	<b>5.7</b>	<b>17.9</b>	<b>33.7</b>	<b>72</b>

PECRS-medium is on par with previous leading approaches (RevCore [11], MESE [21]) for Recall@1 on ReDial.

Scaling up LM size improves recall and items diversity (Unique).

# Conversation Results

Model	Reference-based			Reference-free			
	RG-1	RG-2	F-1	PPL	Dist@2	Dist@3	Dist@4
C <sup>2</sup> -CRS	-	-	-	-	0.163	0.291	0.417
UniCRS	-	-	-	-	0.492	0.648	0.832
RecInDial	-	-	-	-	0.518	0.624	0.598
MESE	-	-	-	12.9	<b>0.822</b>	1.152	1.313
<b>PECRS-small</b>	<u>36.28</u>	<u>14.77</u>	<u>86.04</u>	<u>9.89</u>	0.745	<u>1.462</u>	<u>2.132</u>
<b>PECRS-medium</b>	<b>36.86</b>	<b>15.27</b>	<b>86.36</b>	<b>8.98</b>	<u>0.820</u>	<b>1.552</b>	<b>2.154</b>

PECRS-medium reaches SOTA generation capability on ReDial.

Aspect	MESE	PECRS-small	Tie
Fluency	28.00 (1.63)	<b>46.67</b> (5.91)	25.33 (6.24)
Relevancy	26.33 (2.62)	<b>46.00</b> (0.82)	27.67 (2.87)

Which is confirmed by a human evaluation for fluency and relevancy.

# Comparison with LLMs

We compare against popular instruction-tuned LLMs used in zero-shot [16, 2] on INSPIRED :

Model	Rec.			Conv.		
	R@1	R@10	R@50	Unique	RG-1	RG-2
PECRS-small	5.4	<b>16.1</b>	<b>33.3</b>	<b>34</b>	<b>29.72</b>	<b>8.26</b>
Llama-2-7B-chat	<b>9.3</b>	9.3	9.3	26	19.88	2.88
Vicuna-1.5-7B	8.2	8.2	8.2	23	21.18	3.50

LLMs used in this fashion tend to always recommend among the same small subset of items.

It is not straightforward how to score multiple items with LLMs in zero-shot.

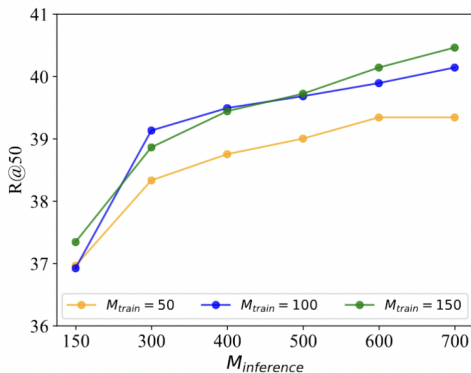
# Analysis

Decoding Strategy	Reference-based		Reference-free		
	RG-1	RG-2	Dist@2	Dist@3	Dist@4
Greedy decoding	38.54	16.25	0.208	0.311	0.390
Beam search	38.23	16.83	0.235	0.353	0.444
Diverse beam search (diversity=0.5)	39.94	<u>17.30</u>	0.190	0.287	0.361
Diverse beam search (diversity=1.0)	<b>40.29</b>	<b>17.40</b>	0.179	0.264	0.320
Diverse beam search (diversity=1.5)	40.07	17.23	0.172	0.246	0.290
Top-k sampling (k=25)	33.54	14.40	0.593	1.177	1.806
Top-k sampling (k=50)	33.37	14.17	<b>0.647</b>	<u>1.300</u>	<u>1.989</u>
Top-k sampling (k=75)	33.48	14.15	<u>0.644</u>	<b>1.303</b>	<b>1.992</b>
Nucleus sampling (p=0.90)	36.35	16.04	0.329	0.555	0.760
Nucleus sampling (p=0.95)	36.44	16.02	0.351	0.594	0.804
Nucleus sampling (p=0.99)	36.60	16.07	0.352	0.593	0.809

Different decoding methods [17, 3, 6] yield very inconsistent Dist@K results.

We advocate for using reference-based methods like ROUGE [10], which are much more stable.

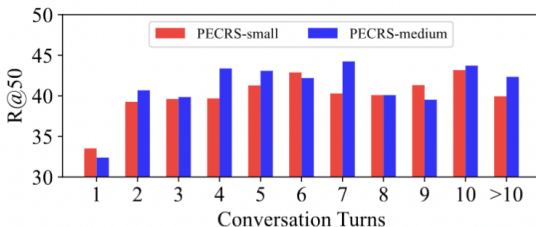
# Analysis



The  $M$  parameter controlling the number of negatives is crucial.

A higher  $M$  is better, albeit at greater computational cost.

# Analysis



Beyond the first turn, recall is relatively stable w.r.t the number of turns in the context.

# Conclusion

In brief, we have introduced PECRS, a simple model fine-tuning a pre-trained LM in a single stage for the CRS task.

- Our model uses GPT-2 for both response generation and item encoding. This is rendered possible through :
  - Projection heads for items and items re-ranking.
  - Stop gradient operator on the backbone.
  - Parameter-efficiency LoRA.
- Optimization is streamlined through re-using the same negative samples across batch items and losses.
- For conversation evaluation, we advocate for not using the popular Dist@K metrics, and use reference-based metrics instead.

- [1] Qibin Chen, Junyang Lin, Yichang Zhang, Ming Ding, Yukuo Cen, Hongxia Yang, and Jie Tang. Towards knowledge-based recommender dialog system. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1803–1813, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [2] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna : An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023.
- [3] Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1 : Long Papers)*, pages 889–898. Association for Computational Linguistics, July 2018.
- [4] Michael U. Gutmann and Aapo Hyvärinen. Noise-contrastive estimation of unnormalized statistical models, with applications



to natural image statistics. *J. Mach. Learn. Res.*, 13 :307–361, feb 2012.

- [5] Shirley Anugrah Hayati, Dongyeop Kang, Qingxiaoyang Zhu, Weiyan Shi, and Zhou Yu. INSPIRED : Toward sociable recommendation dialog systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8142–8152, Online, November 2020. Association for Computational Linguistics.
- [6] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020.
- [7] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA : Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [8] Zhiqiang Hu, Yihuai Lan, Lei Wang, Wanyu Xu, Ee-Peng Lim, Roy Ka-Wei Lee, Lidong Bing, Xing Xu, and Soujanya Poria. Llm-adapters : An adapter family for parameter-efficient fine-tuning of large language models. *ArXiv*, abs/2304.01933, 2023.

- [9] Raymond Li, Samira Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. Towards deep conversational recommendations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, page 9748–9758. Curran Associates Inc., 2018.
- [10] Chin-Yew Lin. ROUGE : A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [11] Yu Lu, Junwei Bao, Yan Song, Zichen Ma, Shuguang Cui, Youzheng Wu, and Xiaodong He. RevCore : Review-augmented conversational recommendation. In *Findings of the Association for Computational Linguistics : ACL-IJCNLP 2021*, pages 1161–1173, Online, August 2021. Association for Computational Linguistics.
- [12] Andriy Mnih and Koray Kavukcuoglu. Learning word embeddings efficiently with noise-contrastive estimation. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.

- [13] Andriy Mnih and Yee Whye Teh. A fast and simple algorithm for training neural probabilistic language models. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, page 419–426, 2012.
- [14] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8) :9, 2019.
- [15] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert : smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108, 2019.
- [16] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2 : Open foundation and fine-tuned chat models. *arXiv preprint arXiv :2307.09288*, 2023.
- [17] Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R. Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. Diverse beam search : Decoding diverse solutions from neural sequence models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.

- [18] Lingzhi Wang, Huang Hu, Lei Sha, Can Xu, Daxin Jiang, and Kam-Fai Wong. RecInDial : A unified framework for conversational recommendation with pretrained language models. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1 : Long Papers)*, pages 489–500, Online only, November 2022. Association for Computational Linguistics.
- [19] Ting-Chun Wang, Shang-Yu Su, and Yun-Nung Chen. Barcor : Towards a unified framework for conversational recommendation systems. *ArXiv*, abs/2203.14257, 2022.
- [20] Xiaolei Wang, Kun Zhou, Ji-Rong Wen, and Wayne Xin Zhao. Towards unified conversational recommender systems via knowledge-enhanced prompt learning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, page 1929–1937. Association for Computing Machinery, 2022.
- [21] Bowen Yang, Cong Han, Yu Li, Lei Zuo, and Zhou Yu. Improving conversational recommendation systems’ quality with context-aware item meta-information. In *Findings of the*

*Association for Computational Linguistics : NAACL 2022*, pages 38–48, Seattle, United States, July 2022. Association for Computational Linguistics.

- [22] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter : Efficient fine-tuning of language models with zero-init attention. *ArXiv*, abs/2303.16199, 2023.
- [23] Tong Zhang, Yong Liu, Boyang Li, Peixiang Zhong, Chen Zhang, Hao Wang, and Chunyan Miao. Toward knowledge-enriched conversational recommendation systems. In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 212–217, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [24] Kun Zhou, Wayne Xin Zhao, Shuqing Bian, Yuanhang Zhou, Ji-Rong Wen, and Jingsong Yu. Improving conversational recommender systems via knowledge graph based semantic fusion. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 1006–1014. Association for Computing Machinery, 2020.