

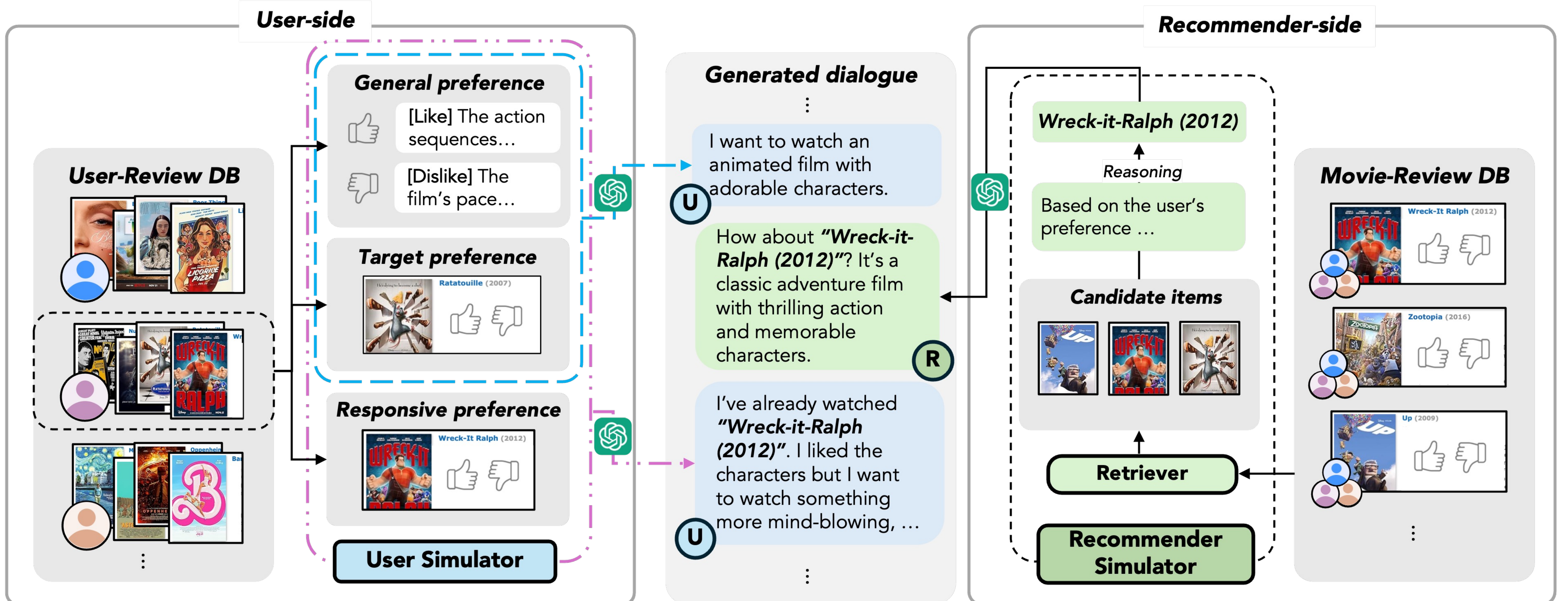
1 ABSTRACT

- We present a novel conversational recommendation dataset named PEARL, synthesized with persona- and knowledge-augmented LLM simulators.
- We obtain detailed persona and knowledge from real-world reviews and construct a dataset with 57k dialogues.
- We show the quality and utility of PEARL through human and automatic evaluations.

3 METHOD

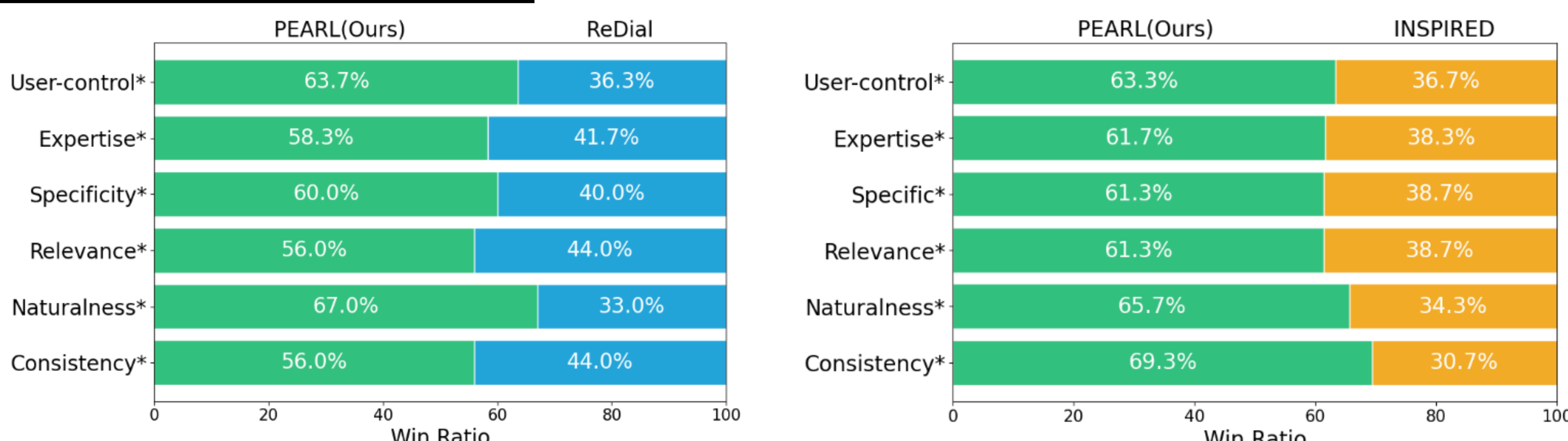
- Construct user-review and movie-review databases with metadata and reviews from IMDb
- Run persona-augmented user simulator
- Run knowledge-augmented recommender simulator
- Complete dialogues via turn-by-turn generation
- Filter out disqualified or low-quality dialogues

PEARL Construction Overview



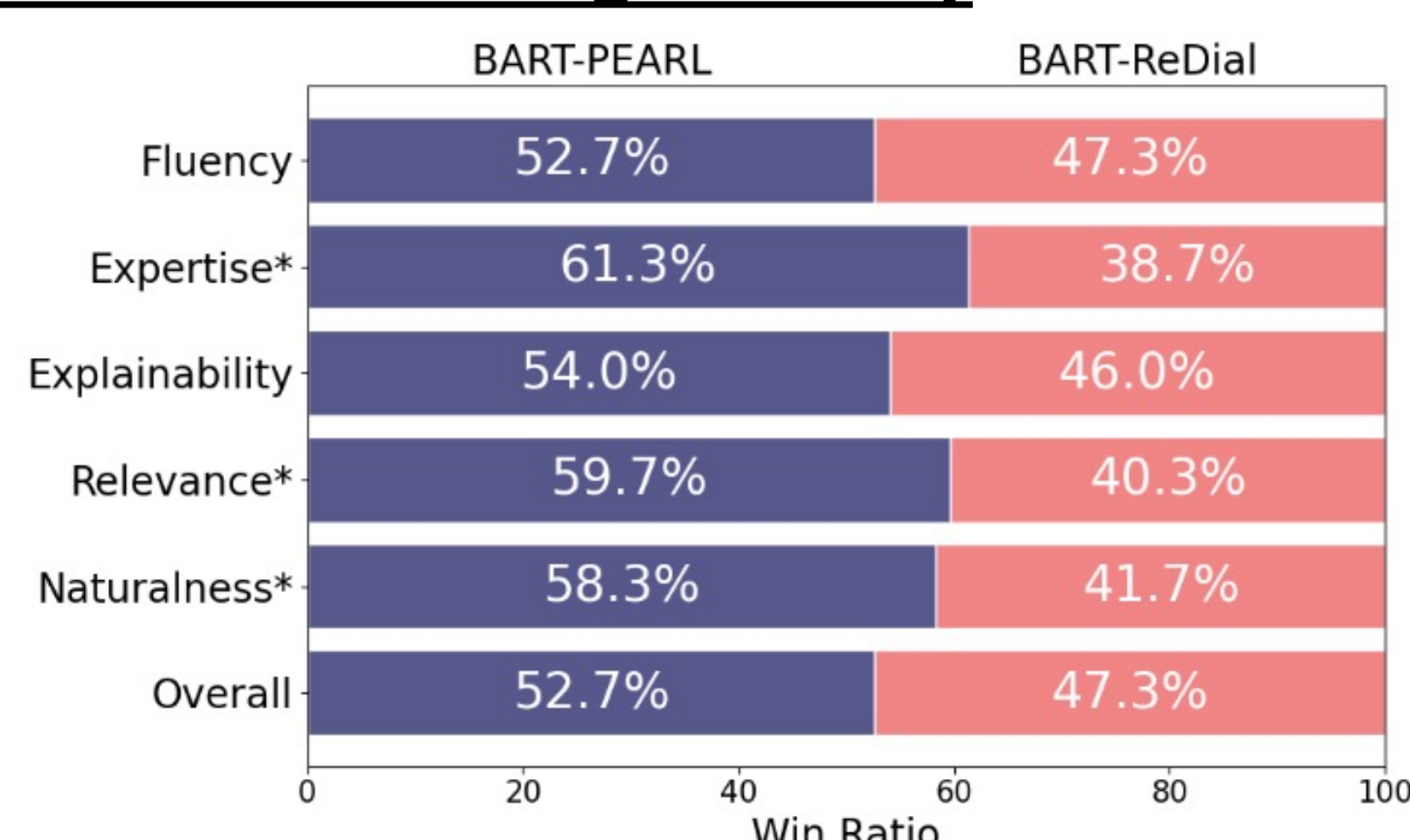
4 EXPERIMENTS

Evaluation on Dialogue Quality



- Dialogues of PEARL shows superiority in expertise, specificity, etc. over those of baseline datasets.

Evaluation on Dialogue Utility



- Models trained on PEARL generate better responses than those trained on a baseline dataset.

Model	Dist-3	Dist-4
BART-ReDial	0.6220	0.5057
BART-PEARL	0.9241	0.8861
UniCRS-ReDial	0.5413	0.3667
UniCRS-PEARL	0.9338	0.9007
PECRS-ReDial	0.6798	0.5906
PECRS-PEARL	0.9132	0.8947
GPT-3.5	0.9256	0.8910

- PEARL improves the diversity of responses generated by downstream models.

	ReDial	INSPIRED	PEARL
# of dialogues	10,006	1,001	57,277
# of utterances	182,150	35,811	548,061
2-gram specificity	65.44	119.56	141.79
3-gram specificity	65.97	123.01	149.75
4-gram specificity	65.37	122.81	153.00

- PEARL is superior to previous datasets in scalability and specificity in user preferences.

Model	R@1	R@10	R@50
BERT-PEARL	0.0018	0.0208	0.0736
UniCRS-PEARL	0.0310	0.0697	0.1202
PECRS-PEARL	0.0151	0.0339	0.0798
GPT-3.5	0.0071	0.0355	0.0709

- All models, including GPT-3.5, show low performances.
- Future research on utilizing PEARL, built based on real-world data rather than on parametric knowledge, is necessary.