



NAACL
2025



Empowering Retrieval-based Conversational Recommendation with Contrasting User Preferences

Heejin Kook*, Junyoung Kim*, Seongmin Park, Jongwuk Lee†
Sungkyunkwan University (SKKU), Republic of Korea

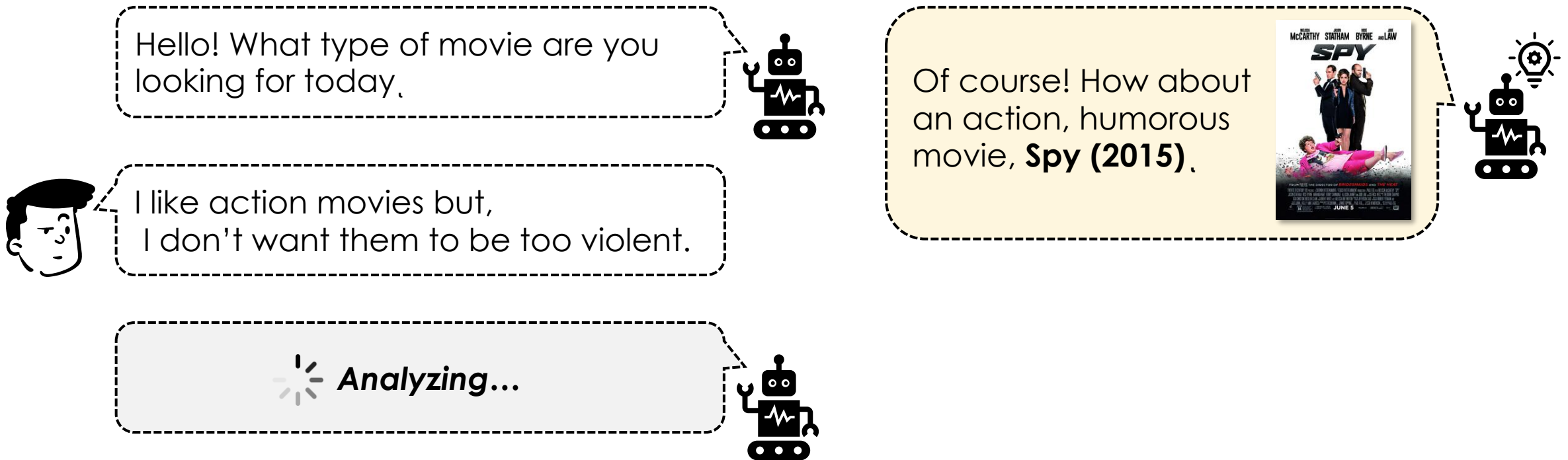
* equal contribution

Introduction

- **Conversational Recommender Systems (CRSs)**
- **Limitations of Existing Methods**
- **Challenges: Contrasting Preferences**
- **Key Contributions**

Task: Conversational Recommender Systems (CRSs)

- CRSs provide personalized recommendations by understanding users' intentions through multi-turn interactions.



Task: Conversational Recommender Systems (CRSs)

- Accurately capturing diverse sentiments—**both positive and negative**—is essential.
 - This reflects opposing intentions!



Task: Conversational Recommender Systems (CRSs)

- Accurately capturing diverse sentiments—**both positive and negative**—is essential.
 - This reflects opposing intentions!



User's opposing intentions in dialogue, namely ***contrasting preference***

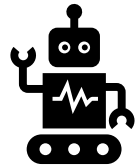


I like **action** movies but,
I don't want them to be too **violent**.

Negative

Limitations of Existing Methods

- Most existing methods **overlook the complex relationship** between the user, item, and contrasting preferences.



Hello! What type of movie are you looking for today?

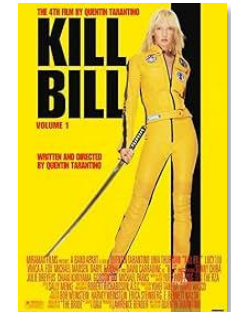
Positive

I like **action** movies but, I don't want them to be too **violent**.

Negative



Kill Bill (2003)



Action
Violent

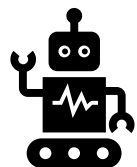
Spy (2015)



Action
Humorous

Limitations of Existing Methods

- Most existing methods **overlook the complex relationship** between the user, item, and contrasting preferences.
 - All this information is aggregated in a **single representation**.



Hello! What type of movie are you looking for today?

Positive

I like action movies but,
I don't want them to be too violent.

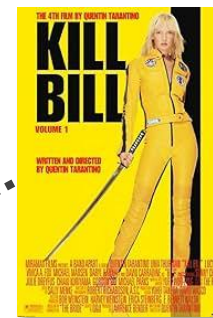
Negative



Action
Violent

I told you I hate violent movie!

Kill Bill (2003)



Action
Violent

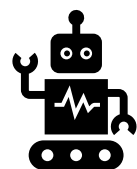
Spy (2015)



Action
Humorous

Our Solution

- We represent contrasting preferences as **distinct vectors**.
 - Explicitly representing and modeling the user, item, and contrasting preferences enables more accurate recommendations that reflect the user's intent.

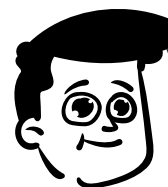


Hello! What type of movie are you looking for today?

Positive

I like **action** movies but,
I don't want them to be too **violent**.

Negative

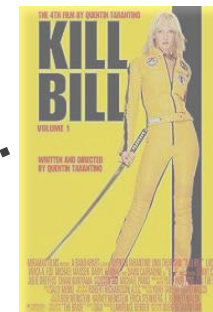


Action

Violent

Just what I wanted!

Kill Bill (2003)



Action

Violent

Spy (2015)



Action

Humorous

Our Solution

We devised **Co**NTRASTING USER **Pr**EFERENCE EXP**AN**SION AND **L**EARNING (CORAL), which extracts and learns contrasting preferences.

Challenge 1.

How do we **extract contrasting user preferences** from the conversation?

Challenge 2.

How do we **learn the relationship** between the contrasting preferences and the user/item?

Method 1.

Contrasting Preference Expansion

Method 2.

Preference-aware Learning

Proposed Method

- Overview of CORAL 
- Contrasting Preference Expansion
- Preference-aware Learning

Overview of CORAL

CORAL (a) extracts and (b) learns contrasting preferences.

(a) *Contrasting Preference Expansion*

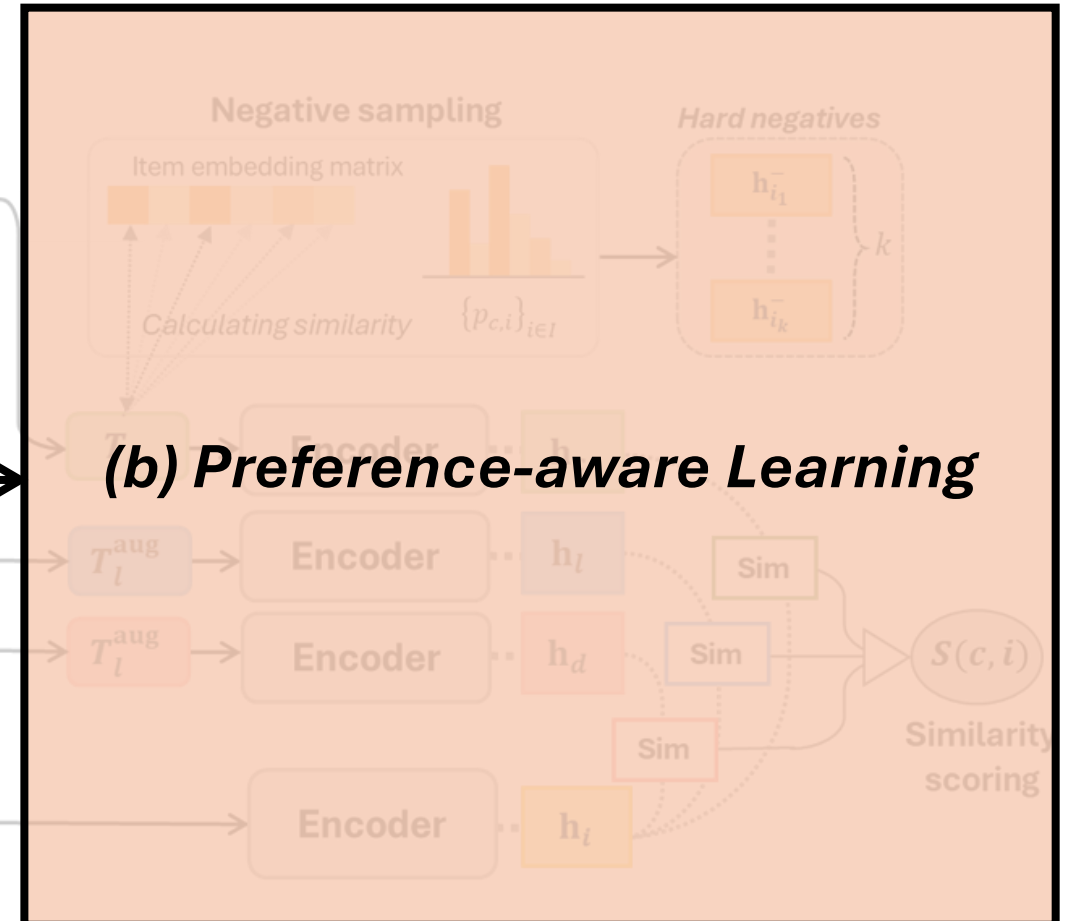
User-side Expansion

Step 1. Superficial Preference Extraction

Step 2. Potential Preference Augmentation

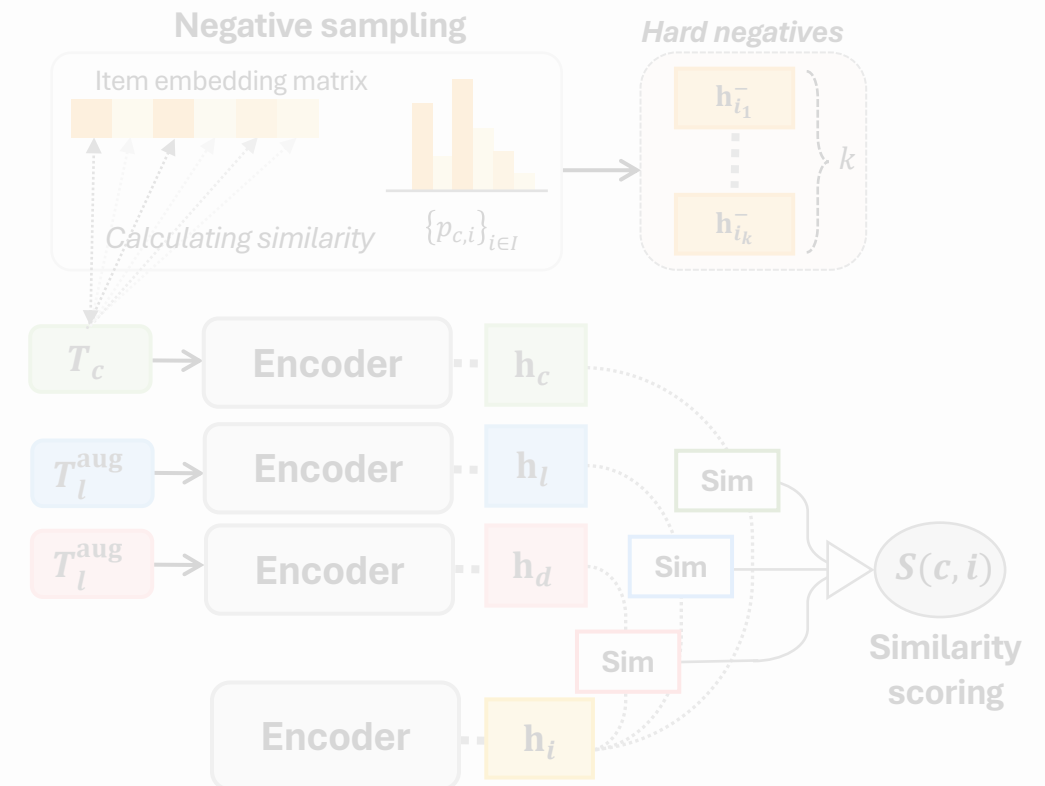
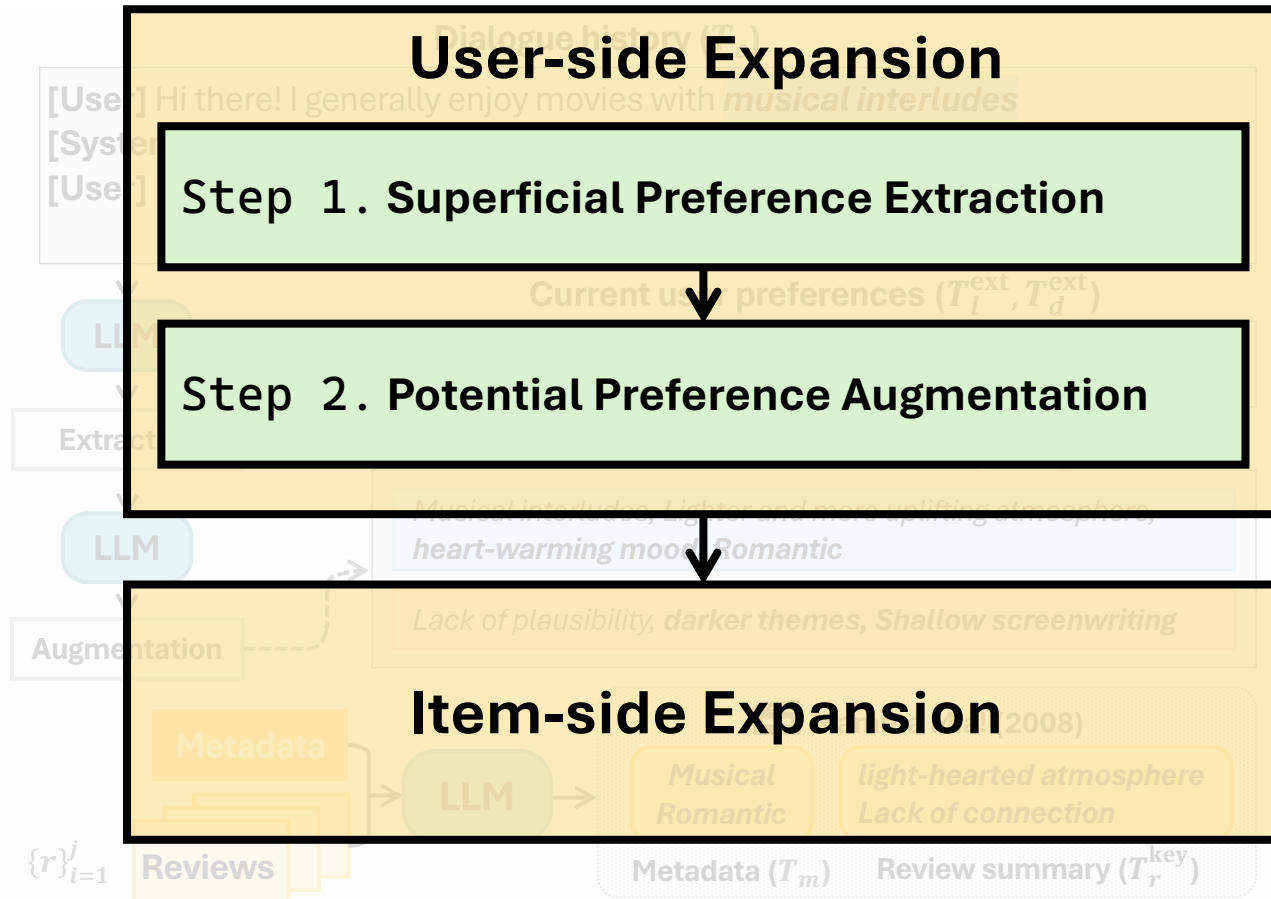
Item-side Expansion

(b) *Preference-aware Learning*



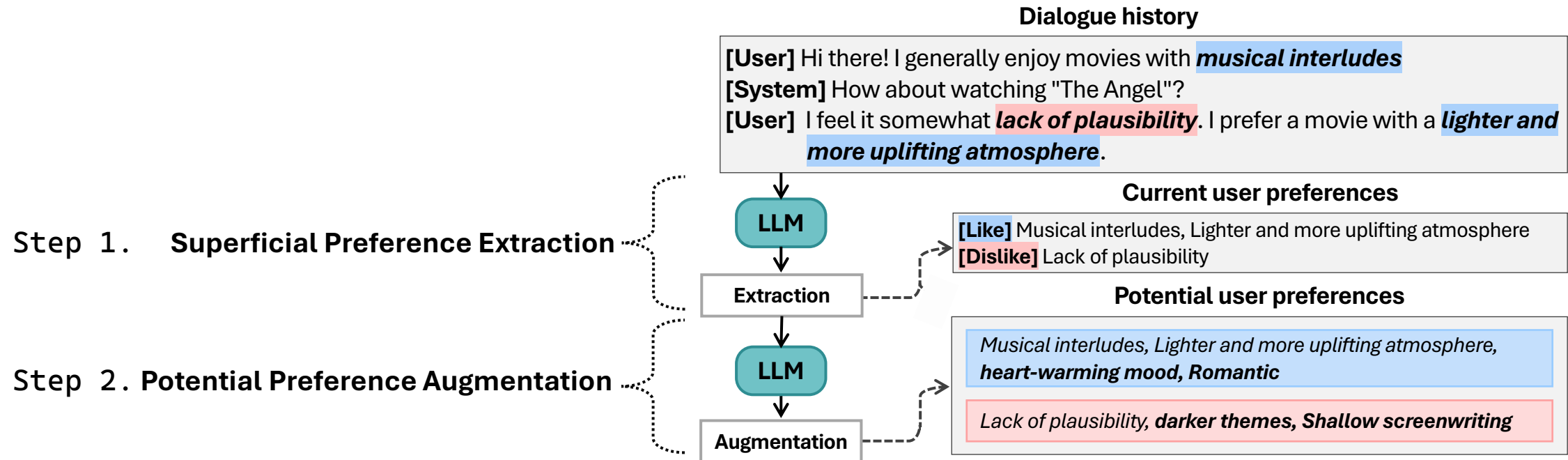
Contrasting Preference Expansion

- Extract **superficial contrasting preferences** from dialogue history/review.
- Expand **potential preferences** via LLMs' reasoning abilities.



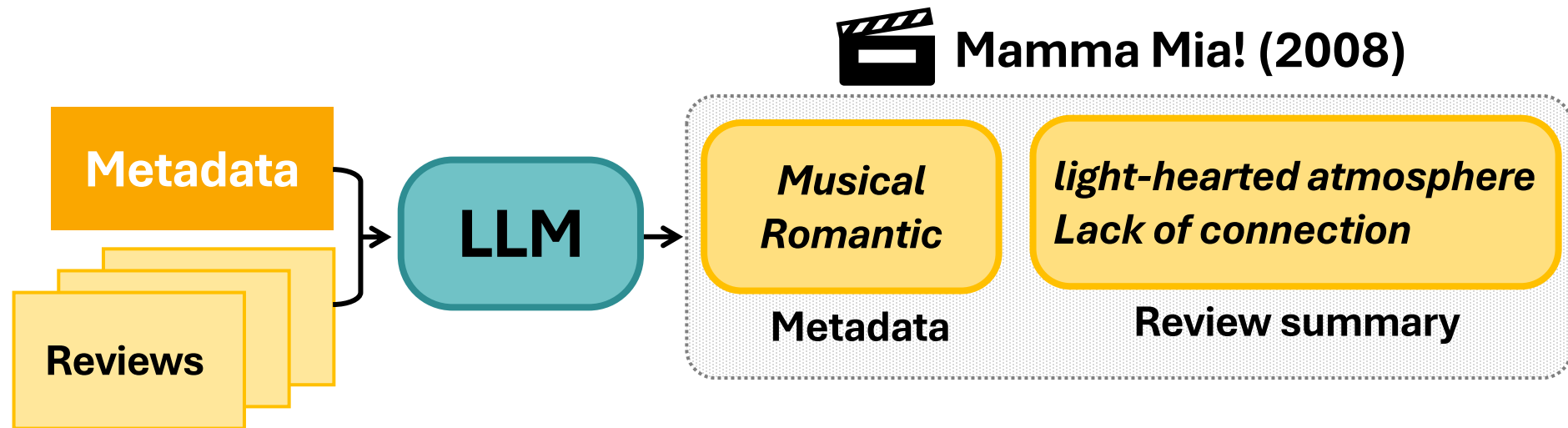
Contrasting Preference Expansion

- **User-side expansion** distinguishes and infers user preferences embedded in the dialogues



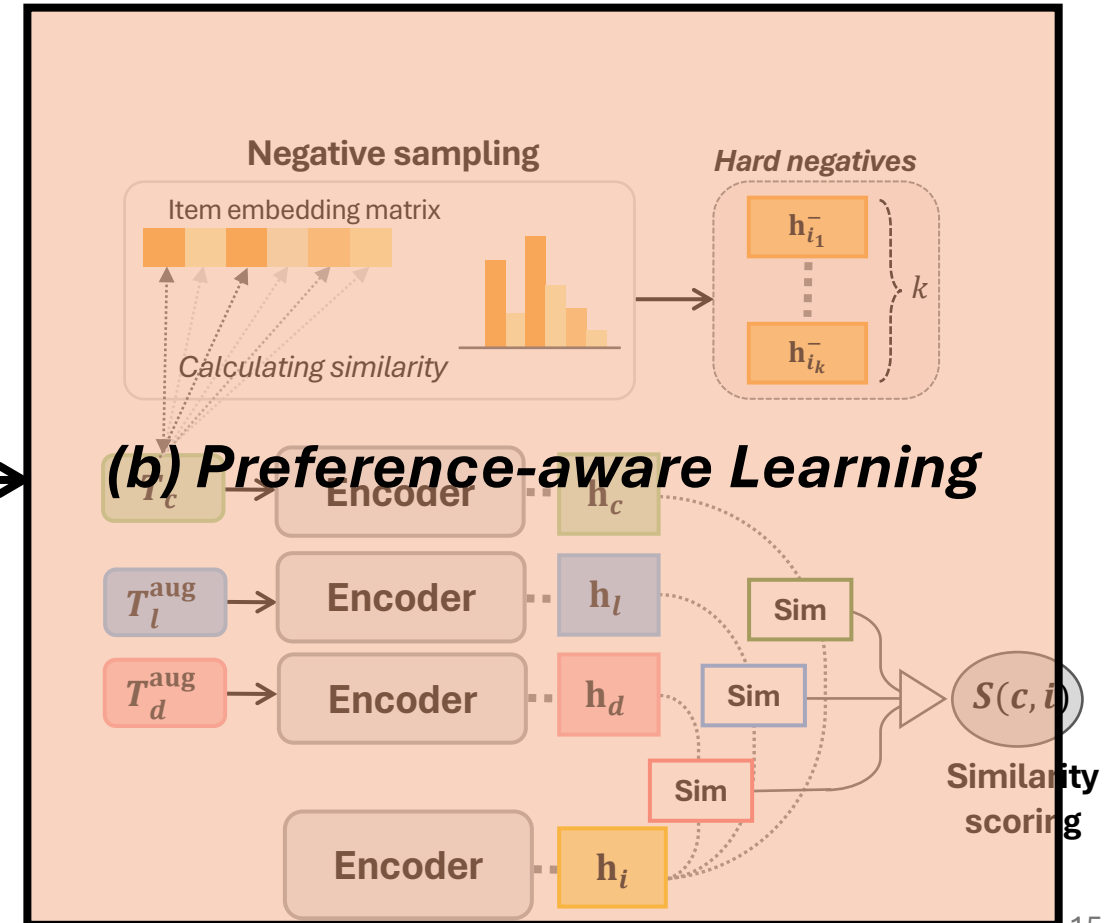
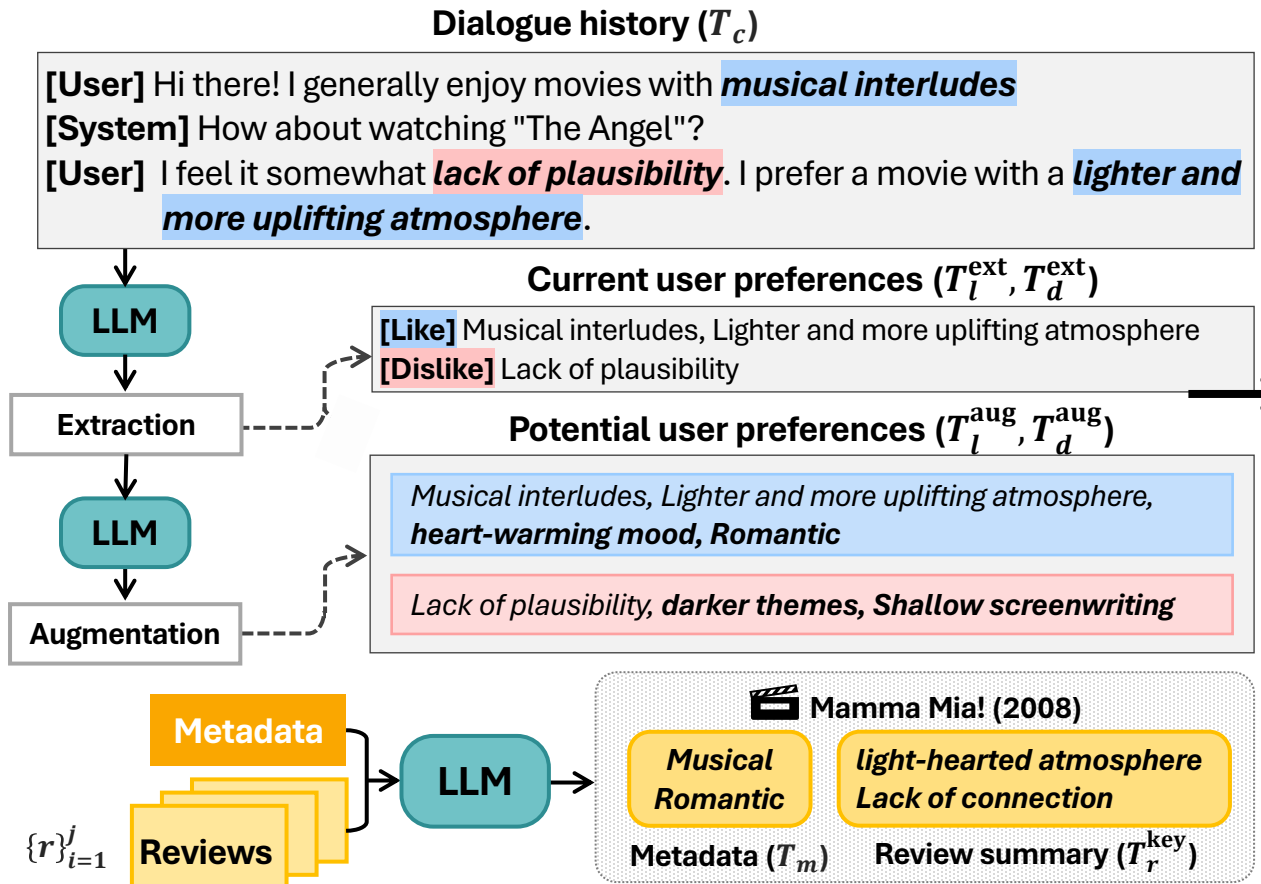
Contrasting Preference Expansion

- **Item-side expansion** enhance the representation of an item using review summary.
 - The metadata lacks sufficient information about user preferences.
 - It can connect user conversations with item metadata using review data.



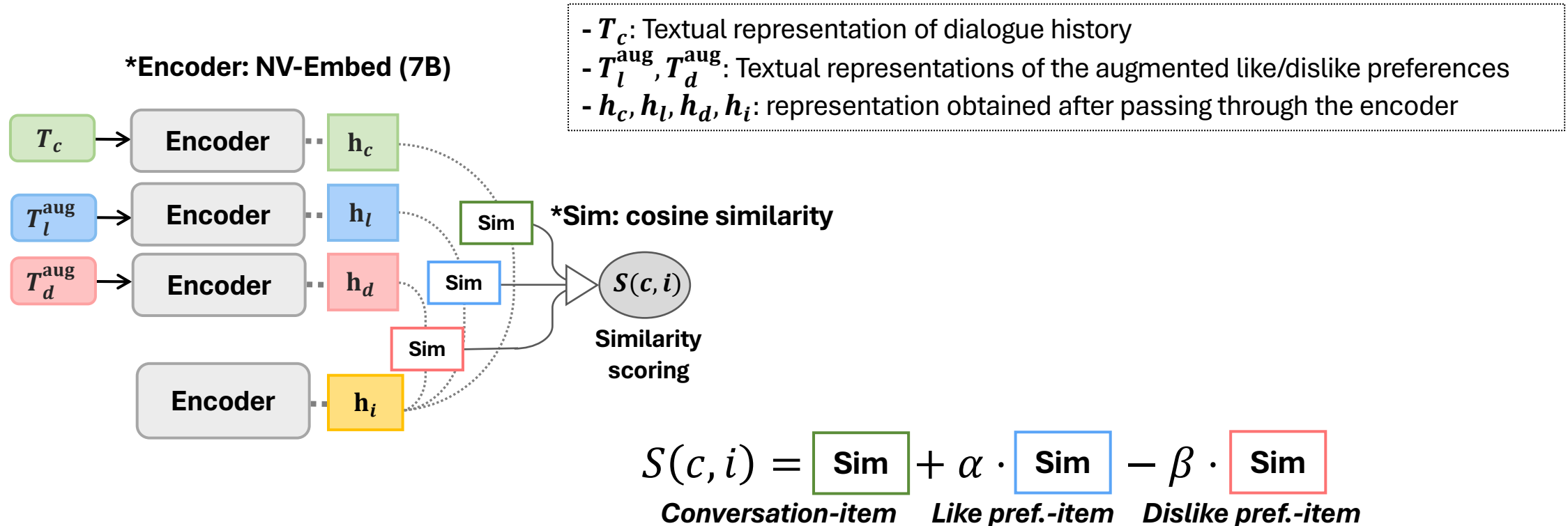
Preference-aware Learning

- Explicitly represent these preferences separately from the conversation and learn the relationship between them.



Preference-aware Learning

- **Preference modeling** explicitly represent preferences and conversation separately, and engage them in item scoring directly.



Experiments

- **Experimental Setting**
- **Overall Performance**
- **Case Study**

Experimental Setting ①: Dataset

- We use three well-known public English movie recommendation conversation datasets
 - **PEARL**, **INSPIRED**, and **ReDial**

| Dataset | #Dial. | #Items | #Likes | #Dislikes |
|----------|--------|--------|--------|-----------|
| PEARL | 57,159 | 9,685 | 9.59 | 5.97 |
| INSPIRED | 1,997 | 1,058 | 11.09 | 5.65 |
| ReDIAL | 31,089 | 5,896 | 10.99 | 1.00 |

#Dial. : the number of dialogues.

#Items: the number of items

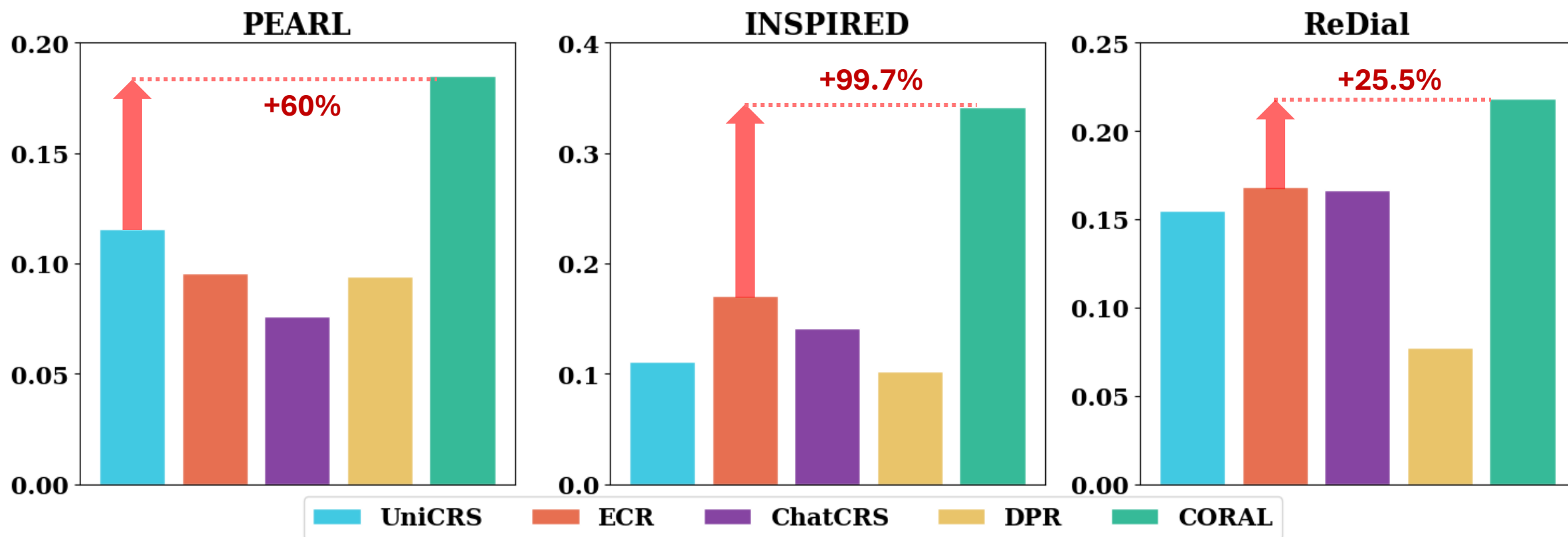
#Likes, #Dislikes: The average counts of the like and dislike preference after the augmentation stage, respectively.

Experimental Setting ②: Baselines

| | | |
|-------------------------|-----------|---|
| (1) Traditional CRS | UniCRS | A model built on DialoGPT (117M)with prompt tuning |
| | RevCore | Categorizes user sentiments regarding entities as either positive or negative |
| | ECR | Identifies nine distinct emotional responses regarding entities. |
| (2) LLM-based CRS | Zero-shot | Recommends solely based on dialogue history and internal knowledge items. |
| | ChatCRS | Enhances the domain knowledge of LLM through a knowledge graph |
| (3) Retrieval-based CRS | BM25 | Ranks item by term relevance from a static index |
| | DPR | Retrieves items based on the similarity with dense vectors of the dialogue context (BERT-base (110M)) |

Key Results

- CORAL 🪸 achieves state-of-the-art performance over existing methods in three benchmark datasets, improving **up to 99.7% in Recall@10**



Key Results: Ablation Study on INSPIRED

- **Effect of Like/Dislike preferences**

| Variant | | R@10 | R@50 | N@10 | N@50 |
|----------------------------------|--|---------------|---------------|---------------|---------------|
| CORAL | | 0.3481 | 0.5667 | 0.1827 | 0.2297 |
| w/o <i>Like, Dislike</i> | | 0.3248 | 0.5767 | 0.1668 | 0.2226 |
| w/o <i>Review</i> | | 0.3167 | 0.5348 | 0.1710 | 0.2193 |
| w/o <i>Like, Dislike, Review</i> | | 0.2948 | 0.5633 | 0.1595 | 0.2196 |

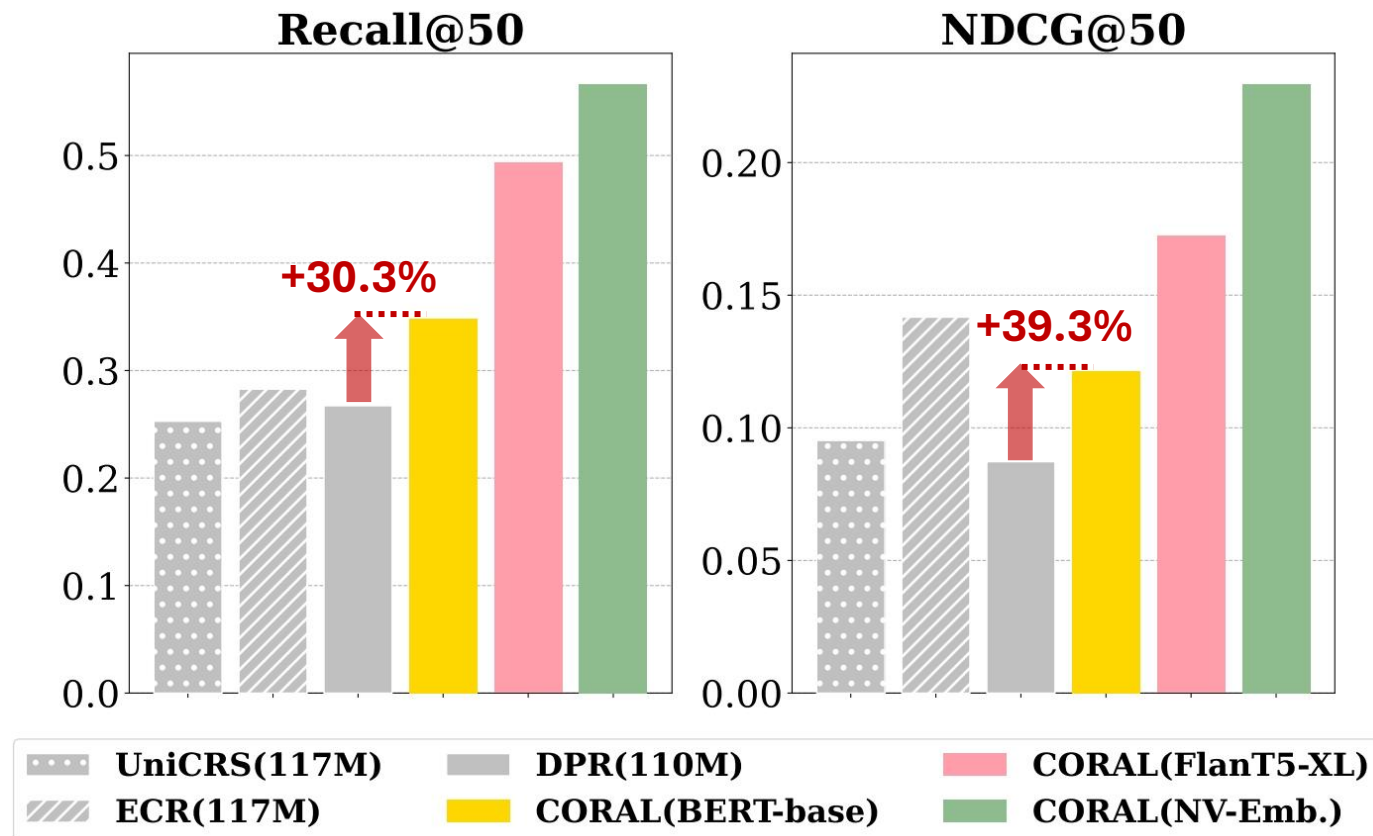
- **Effect of Preference-aware Learning**

| Variant | | R@10 | R@50 | N@10 | N@50 |
|-------------------------------------|---------------|---------------|---------------|---------------|---------------|
| w/o <i>Negative sampling</i> | w/ <i>PL</i> | 0.2974 | 0.5692 | 0.1520 | 0.2115 |
| | w/o <i>PL</i> | 0.1847 | 0.3616 | 0.1107 | 0.1491 |

* *PL*: Preference-aware Learning

Key Results: Performance by Model Size

- CORAL 🪸 significantly improves performance even with a relatively small model.



**DPR: CORAL w/o preference-aware learning*

Case Study

- CORAL  **effectively captures contrasting preferences** to improve recommendations and enhance **explainability**.

I'm generally a fan of movies that have a lot of tension and build-up, as well as ones that portray society and real-life characters in an engaging way

...

I appreciate films that have a more detailed and realistic portrayal of events

User's utterances

[Like]

real-life characters in an engaging way,
psychological thrillers & dramas

Immersive sound design

[Dislike]

unrealistic events
over-the-top action sequences

Contrasting preferences

Title: Sicario (2015)

Genre: Cop Drama, Drug, Crime, Action, Mystery, Thriller

[Review Summary]

Realistic feel

Immersive atmosphere

Intense and engaging

Amazing sound design

**Item metadata
& summary**

Conclusion

- We propose a novel retrieval-based CRS framework that extracts and learns contrasting preferences.
 - **CO**ntrasting user p**R**eference exp**A**nasion and **L**earning (**CORAL** 🌊),
- **CORAL** 🌊 addresses contrasting preferences by
 - distinguishing and enhancing contrasting preferences into like/dislike.
 - learning the relationship between conversation, preferences, and items directly.
- **CORAL** 🌊 achieves the best performance with existing CRS models on PEARL, INSPIRED, and REDAIL, improving up to 99.72% in Recall@10.

Thank you!

Any question?

Email: hjkook@g.skku.edu

Code: <https://github.com/kookeej/CORAL.git>



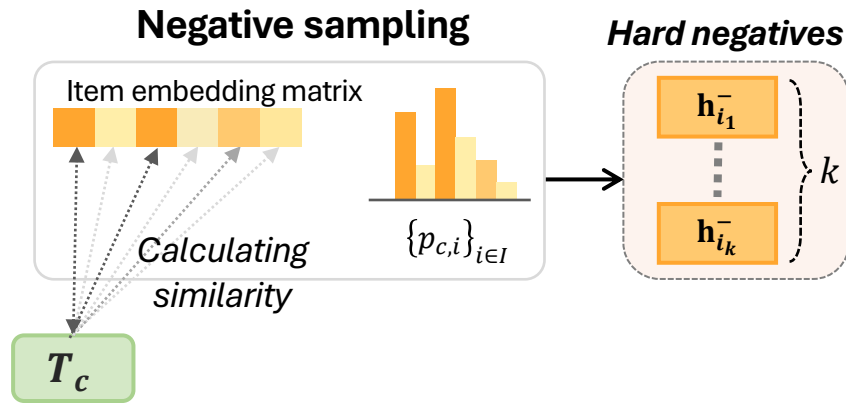
Paper



Code

Appendix: Negative Sampling

- **Hard negative sampling** differentiate samples that are challenging to predict based on conversation alone with contrasting preferences.



- Calculate the similarity between \mathbf{h}_c and all item embeddings.
- Apply *softmax* to convert the similarity scores into a probability distribution.

- **Loss function**

$$\mathcal{L} = -\log \frac{\exp(S(c, i^+)/\tau)}{\sum_{i \in \mathcal{I}_c^-} \exp(S(c, i)/\tau)}$$

where i^+ is the positive item of c , \mathcal{I}_c^- is a set of k negative items, and τ is a hyperparameter to adjust the temperature.

Appendix: Overall Performance

| Model | | Traditional CRS | | | LLM-based CRS | | Retrieval-based CRS | | | Gain |
|----------|--------|-----------------|---------------|---------------|---------------|---------|---------------------|--------|----------------|--------|
| Dataset | Metric | RevCore | UniCRS | ECR | Zero-shot | ChatCRS | BM25 | DPR | CORAL | |
| PEARL | R@10 | 0.0268 | <u>0.1156</u> | 0.0957 | 0.0767 | 0.0763 | 0.0026 | 0.0940 | 0.1851* | 60.07% |
| | R@50 | 0.0898 | <u>0.2624</u> | 0.2373 | 0.1129 | 0.1168 | 0.0123 | 0.2206 | 0.3619* | 37.94% |
| | N@10 | 0.0132 | <u>0.0642</u> | 0.0501 | 0.0468 | 0.0462 | 0.0014 | 0.0502 | 0.1125* | 75.17% |
| | N@50 | 0.0266 | <u>0.0958</u> | 0.0806 | 0.0560 | 0.0565 | 0.0033 | 0.0777 | 0.1511* | 57.74% |
| INSPIRED | R@10 | 0.0948 | 0.1113 | <u>0.1711</u> | 0.1436 | 0.1410 | 0.0429 | 0.1019 | 0.3417* | 99.72% |
| | R@50 | <u>0.3344</u> | 0.2528 | <u>0.2826</u> | 0.2436 | 0.2436 | 0.1210 | 0.2672 | 0.5632* | 68.45% |
| | N@10 | <u>0.0509</u> | 0.0642 | <u>0.1077</u> | 0.0927 | 0.0806 | 0.0202 | 0.0512 | 0.1772* | 64.52% |
| | N@50 | 0.1041 | 0.0952 | <u>0.1417</u> | 0.1175 | 0.1071 | 0.0373 | 0.0872 | 0.2255* | 59.07% |
| REDIAL | R@10 | <u>0.1739</u> | 0.1549 | 0.1685 | 0.1670 | 0.1666 | 0.0373 | 0.0774 | 0.2182* | 25.48% |
| | R@50 | <u>0.3034</u> | 0.3540 | <u>0.3793</u> | 0.2783 | 0.2824 | 0.0300 | 0.2138 | 0.4741* | 25.23% |
| | N@10 | <u>0.1053</u> | 0.0776 | <u>0.0805</u> | 0.0937 | 0.0893 | 0.0032 | 0.0403 | 0.1128* | 7.10% |
| | N@50 | <u>0.1337</u> | 0.1215 | 0.1293 | 0.1226 | 0.1191 | 0.0083 | 0.0713 | 0.1724* | 28.96% |

Table 2: Overall performance. The best and second-best are **bold** and underlined. Gain measures the difference between CORAL and the best competitive baseline. ‘*’ indicates statistically significant improvement ($p < 0.01$) for a paired t -test of CORAL compared to the best baseline, as conducted across 5 experiments.

Appendix: Ablation Study

| | Dataset | | INSPIRED | | | |
|-------------------------------------|----------------------------------|---------------|---------------|---------------|---------------|---------------|
| | Variants | | R@10 | R@50 | N@10 | N@50 |
| | CORAL | | 0.3481 | 0.5667 | 0.1827 | 0.2297 |
| Effect of Like/Dislike preference | w/o Like, Dislike | | 0.3248 | 0.5767 | 0.1668 | 0.2226 |
| | w/o Review | | 0.3167 | 0.5348 | 0.1710 | 0.2193 |
| | w/o Like, Dislike, Review | | 0.2948 | 0.5633 | 0.1595 | 0.2196 |
| Effect of Potential preference | w/o Augmentation | | 0.3016 | 0.5379 | 0.1663 | 0.2202 |
| Effect of Preference-aware learning | w/o NS | w/ PL | 0.2974 | 0.5692 | 0.1520 | 0.2115 |
| | | w/o PL | 0.1847 | 0.3616 | 0.1107 | 0.1491 |

- NS: Negative sampling

- PL: Preference-aware Learning

Appendix: Zero-shot Performance

| Dataset | | PEARL | | | | INSPIRED | | | | REDIAL | | | |
|-----------|--------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Retriever | L, D | R@10 | R@50 | N@10 | N@50 | R@10 | R@50 | N@10 | N@50 | R@10 | R@50 | N@10 | N@50 |
| BM25 | w/o | 0.0053 | 0.0269 | 0.0031 | 0.0076 | 0.0390 | 0.1486 | 0.0216 | 0.0444 | 0.0091 | 0.0361 | 0.0044 | 0.0104 |
| | w/ | 0.0079 | 0.0340 | 0.0045 | 0.0099 | 0.0448 | 0.1714 | 0.0297 | 0.0580 | 0.0235 | 0.0704 | 0.0113 | 0.0219 |
| BERT | w/o | 0.0040 | 0.0203 | 0.0015 | 0.0049 | 0.0057 | 0.0410 | 0.0017 | 0.0100 | 0.0069 | 0.0270 | 0.0039 | 0.0084 |
| | w/ | 0.0031 | 0.0225 | 0.0015 | 0.0056 | 0.0086 | 0.0648 | 0.0032 | 0.0152 | 0.0087 | 0.0302 | 0.0048 | 0.0096 |
| NV-Emb. | w/o | 0.0454 | 0.1323 | 0.0210 | 0.0397 | 0.1648 | 0.3943 | 0.0856 | 0.1374 | 0.0767 | 0.1885 | 0.0368 | 0.0613 |
| | w/ | 0.0569 | 0.1508 | 0.0286 | 0.0489 | 0.2276 | 0.4048 | 0.1140 | 0.1538 | 0.0872 | 0.2190 | 0.0408 | 0.0711 |

Table 3: The zero-shot performance of various language models depending on the presence or absence of the user’s potential preference. L and D mean T_l^{aug} and T_d^{aug} , respectively.

Appendix: Performance depending on the LLM utilized for Contrasting Preference Expansion

| Dataset | INSPIRED | | | |
|---|----------|--------|--------|--------|
| Model | R@10 | R@50 | N@10 | N@50 |
| UniCRS | 0.1113 | 0.2528 | 0.0642 | 0.0952 |
| ECR | 0.1711 | 0.2826 | 0.1077 | 0.1417 |
| CORAL _{w/o L, D, R} | 0.2948 | 0.5633 | 0.1595 | 0.2196 |
| CORAL _{Mistral} | 0.3162 | 0.5410 | 0.1809 | 0.2326 |
| CORAL _{gpt-4o-mini} | 0.3417 | 0.5632 | 0.1772 | 0.2255 |

Table 7: Performance depending on the LLM utilized for Contrasting Preference Expansion. L , D and R denote T_l^{aug} , T_d^{aug} , and T_r^{key} , respectively.

Appendix: Zero-shot Performance for Preference Input Variant

| Dataset | Input info. | R@10 | R@50 | N@10 | N@50 | Avg. Gain(%) |
|----------|-------------|--------|--------|--------|--------|---------------|
| PEARL | C | 0.0476 | 0.1230 | 0.0229 | 0.0395 | - |
| | C, L | 0.0560 | 0.1349 | 0.0303 | 0.0474 | 19.91% |
| | C, D | 0.0481 | 0.1349 | 0.0224 | 0.0415 | 3.40% |
| | C, L, D | 0.0573 | 0.1481 | 0.0311 | 0.0504 | 26.05% |
| INSPIRED | C | 0.1837 | 0.4133 | 0.1038 | 0.1545 | - |
| | C, L | 0.2103 | 0.3949 | 0.1133 | 0.1534 | 4.62% |
| | C, D | 0.2000 | 0.4154 | 0.1064 | 0.1527 | 2.68% |
| | C, L, D | 0.2205 | 0.4205 | 0.1166 | 0.1597 | 9.37% |
| REDIAL | $C,$ | 0.0887 | 0.2067 | 0.0383 | 0.0640 | - |
| | C, L | 0.0954 | 0.2325 | 0.0415 | 0.0710 | 9.83% |
| | C, D | 0.0910 | 0.2129 | 0.0400 | 0.0665 | 3.48% |
| | C, L, D | 0.0986 | 0.2431 | 0.0427 | 0.0735 | 13.78% |

Table 5: Zero-shot performance for preference input variant. C , L , and D denote T_c , T_l^{aug} , and T_d^{aug} , respectively.

Appendix: Ablation study on BERT

| Dataset | INSPIRED | | | | REDIAL | | | |
|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Variants | R@10 | R@50 | N@10 | N@50 | R@10 | R@50 | N@10 | N@50 |
| CORAL | 0.1219 | 0.3714 | 0.0625 | 0.1173 | 0.0856 | 0.2255 | 0.0446 | 0.0765 |
| w/o L, D | 0.0962 | 0.3429 | 0.0429 | 0.0968 | 0.0775 | 0.2287 | 0.0407 | 0.0754 |
| w/o R | 0.1133 | 0.2695 | 0.0752 | 0.1103 | 0.0749 | 0.2052 | 0.0405 | 0.0702 |
| w/o L, D, R | 0.0876 | 0.3010 | 0.0450 | 0.0967 | 0.0764 | 0.2231 | 0.0403 | 0.0738 |
| w/o $Neg.$ | 0.0867 | 0.2629 | 0.0371 | 0.0757 | 0.0725 | 0.2151 | 0.0375 | 0.0697 |
| w/o PL | 0.1076 | 0.2467 | 0.0516 | 0.0832 | 0.0774 | 0.2138 | 0.0403 | 0.0713 |

Table 6: Ablation study of CORAL in INSPIRED and REDIAL on BERT. The best scores are in **bold**. L , D and R denote T_l^{aug} , T_d^{aug} , and T_r^{key} , respectively. Also, $Neg.$ and PL mean potential hard negative sampling and preference-aware learning, respectively.