

TASTE: A Context-Aware Explainable Recommendation Model

Anonymous ACL submission

Abstract

Travel recommendation involves complex, context-dependent decision-making that goes beyond static user preferences. While recent large language model (LLM)-based conversational recommender systems enable natural language interaction and explanation generation, most existing approaches treat user preferences as fixed and utilize user queries only for surface-level explanation, without adapting the underlying recommendation logic. This results in a mismatch between conversational intent and actual recommendation criteria.

In this paper, we propose TASTE (Query-Aware Travel Recommendation via Aspect-based Sentiment Profiling), a novel recommendation framework that dynamically adapts aspect-level preference weights according to natural language user queries while preserving long-term user preference profiles. TASTE constructs structured user preference representations using Aspect-Based Sentiment Analysis (ABSA) applied to review texts and combines learned aspect attention with query-derived aspect importance through a controllable weighting mechanism. Large language models are employed exclusively for query interpretation and explanation generation, ensuring that recommendation decisions remain deterministic, interpretable, and faithful to model-internal reasoning.

We conduct extensive experiments on three real-world datasets spanning restaurant, travel, and hotel domains (Yelp Restaurant, Yelp Travel-related, and TripAdvisor Hotel). Quantitative results demonstrate that TASTE achieves competitive rating prediction performance compared to strong rating- and review-based baselines, without sacrificing accuracy for interpretability. Qualitative and quantitative analyses further show that TASTE effectively captures query-aware preference shifts, producing meaningful changes in recommendation rankings for the same user under different queries. Finally, explanation quality evalua-

tion using an LLM-as-a-Judge protocol confirms that TASTE generates coherent, relevant, and aspect-aligned explanations across domains and languages.

Overall, TASTE provides a unified framework for accurate, query-aware, and explainable travel recommendation, addressing key limitations of existing conversational recommender systems.

1 Introduction

Travel recommendation systems have been widely adopted across online platforms and the travel industry, and recent advances in large language models (LLMs) have accelerated the emergence of conversational, chatbot-based travel recommendation services. In practice, many platforms now employ natural language interfaces that accept free-form user queries and provide personalized recommendations, as exemplified by systems such as OpenAI (OpenAI, 2025), Expedia (Expedia, 2025), HAIJEKO (Jeju Air, 2025), and the Malaysia Airlines chatbot. These trends indicate a shift from conventional item recommendation toward interactive, dialogue-oriented recommendation systems.

This transition has been driven by three key technological advances. First, progress in natural language processing enables direct understanding of unstructured textual data such as reviews and conversational utterances, allowing recommendation models to move beyond rating prediction toward reasoning-based and explainable recommendations. Second, user modeling has evolved to extract personalized preferences from review text rather than relying solely on structured interaction data. Third, these advances have facilitated the rapid adoption of LLM-based conversational recommendation systems across domains including travel, e-commerce, and entertainment. Consequently, natural language-based preference modeling and explainable recommendation have emerged

086	as central research challenges in modern recom-	138
087	mender systems.	139
088	Among recommendation domains, travel recom-	140
089	mendation is characterized by particularly	141
090	complex and context-dependent user preferences.	142
091	Travel decision-making involves a sequence of	143
092	interrelated choices—such as destination selec-	144
093	tion, accommodation, dining, and activity plan-	145
094	ning—rather than a single consumption decision	146
095	(Sarkar et al., 2023). Throughout this process, users	147
096	simultaneously consider functional utility and expe-	148
097	riential factors, including travel purpose, compan-	149
098	ions, budget constraints, and preferred atmosphere.	150
099	Prior studies have shown that even the same user	151
100	may apply substantially different preference cri-	152
101	teria depending on situational context, reflecting	153
102	differences between solo and group travel, leisure-	154
103	oriented versus sightseeing-focused trips, and other	155
104	contextual conditions (Renjith et al., 2020). As	156
105	a result, travel recommendation has long been re-	157
106	garded as a domain in which personal context and	158
107	subjective judgment play a more prominent role	159
108	than in conventional content recommendation tasks	160
109	such as music or movie recommendation (Lin and	161
110	Liu, 2024). Moreover, fine-grained preferences	
111	over factors such as price, location, service quality,	162
112	and atmosphere have been reported to exert a direct	
113	and significant impact on overall travel satisfaction	163
114	(Payandenick and Othman, 2026). These findings	
115	highlight the limitations of static user profiles and	164
116	fixed preference weights in modeling real-world	165
117	travel decision-making.	166
118	Meanwhile, prior research on review-based and	167
119	conversational recommendation has primarily fo-	168
120	cused on modeling long-term user preferences or	169
121	generating explanations for recommended items.	170
122	For example, studies have proposed constructing	171
123	interpretable user profiles from review text to sup-	172
124	port explainable recommendations, while TEARS	173
125	introduced a controllable recommendation frame-	174
126	work that allows users to influence recommenda-	175
127	tion outcomes (Ramos et al., 2024). In addition,	176
128	recent studies on ChatGPT-based travel recommen-	177
129	dation systems have reported that LLM-generated	178
130	explanations positively affect perceived relevance,	179
131	credibility, and usefulness.	180
132	However, most existing approaches still assume	181
133	that user preferences are static. Even when natural	182
134	language queries are supported in conversational	183
135	recommendation settings, queries are often used	184
136	only for explanation generation rather than for ac-	185
137	tively adjusting recommendation criteria. Conse-	186
	quently, when a single user issues queries reflecting	187
	different travel purposes or contextual needs, the	
	underlying preference weights remain unchanged,	
	leading to a mismatch between conversational input	
	and recommendation logic.	
	To address these limitations, we propose TASTE	
	(Query-Aware Travel Recommendation via	
	Aspect-based Sentiment Profiling) , a novel frame-	
	work for conversational travel recommendation.	
	TASTE constructs user preference profiles using	
	Aspect-Based Sentiment Analysis (ABSA) over re-	
	view texts and dynamically adjusts aspect-level im-	
	portance weights according to the user’s natural	
	language query. By preserving long-term preferences	
	while reweighting aspects based on query intent,	
	TASTE enables context-adaptive recommendation	
	behavior. Furthermore, aspect-level contribution	
	scores are explicitly exposed, and large language	
	models are used solely to generate natural language	
	explanations grounded in these quantitative signals,	
	thereby enhancing transparency and trustworthi-	
	ness.	
	Based on this motivation, this study addresses	
	the following research questions:	
	2 Related Works	
	2.1 LLM-based Recommendation	
	The rapid advancement of large language models	
	(LLMs) has introduced a new paradigm in recom-	
	mender systems research (Liu et al., 2024; Wu et al.,	
	2024; Wang et al., 2024a). Early studies primar-	
	ily explored the use of LLMs as text encoders or	
	auxiliary knowledge sources to enhance existing	
	neural recommendation models (Zhao et al., 2023).	
	More recent work has moved toward deeper inte-	
	gration of LLMs into the recommendation pipeline,	
	leveraging their capabilities in reasoning, summa-	
	rization, and content generation.	
	Prior research broadly categorizes LLM-based	
	recommendation approaches into LLM-enhanced	
	recommender systems, generative recommenda-	
	tion, and hybrid frameworks. LLM-enhanced rec-	
	ommender systems include knowledge enhance-	
	ment, interaction enhancement, and model en-	
	hancement strategies. Knowledge enhancement	
	approaches enrich user and item representations	
	through user history summarization, item attribute	
	generation, and knowledge graph expansion (De-	
	vlin et al., 2019; Wang et al., 2024b; Tian et al.,	
	2024; Cao et al., 2024; Ren et al., 2024; Sun et al.,	
	2025). Interaction enhancement methods mitigate	

188 data sparsity by generating pseudo interaction logs
189 via prompt-based reasoning (Yue et al., 2023; Yang
190 et al., 2025b). Model enhancement approaches
191 distill semantic or reasoning signals derived from
192 LLMs into lightweight recommendation models
193 for efficient inference (Gheewala et al., 2024).

194 Generative recommendation adopts an alterna-
195 tive paradigm in which user histories and items are
196 represented as natural language token sequences
197 and recommendations are generated autoregres-
198 sively. To support this paradigm, techniques such
199 as ID tokenization, clustering-based indexing, and
200 VAE-based encoding have been proposed. In ad-
201 dition, several studies employ LLMs in re-ranking
202 stages or adopt retrieval-augmented generation
203 (RAG)-based recommendation frameworks (Liu
204 et al., 2025a; Zhou et al., 2025). Despite their
205 strong reasoning capability, these approaches gen-
206 erally overlook the problem of dynamically adjust-
207 ing internal preference structures according to user
208 query context.

209 2.2 Review-based and Aspect-based 210 Recommendation

211 User reviews provide a rich source of information
212 for modeling user preferences, experiences, and
213 emotions, and have been extensively studied in rec-
214 ommender systems. Early approaches enhanced
215 user and item representations using review embed-
216 dings, while later studies significantly improved
217 text understanding by incorporating pre-trained lan-
218 guage models such as BERT (Liu et al., 2019).
219 Models such as DAML jointly learn from ratings
220 and review text, achieving complementary effects
221 between numerical and textual signals (Liu et al.,
222 2025b).

223 Graph-based approaches have further ex-
224 plored structural relationships within review data.
225 APAHN models aspects as hyperedges to capture
226 interactions among multiple aspects (Zhang et al.,
227 2023), while the Multi-Aspect Enhanced GNN
228 constructs a user-item-review tripartite graph to
229 learn multi-relational representations (Hasan et al.,
230 2025). The Multi-Factor Collaborative Prediction
231 Model similarly improves recommendation perfor-
232 mance by extracting multiple latent semantic fac-
233 tors from review text (Zhang et al., 2022). A com-
234 prehensive survey on review-based recommenda-
235 tion identifies text processing, sentiment analysis,
236 aspect extraction, and multi-preference modeling
237 as core components (Shen et al., 2025). However,
238 many existing methods struggle to disentangle mul-

239 tiple coexisting aspects within a single review and
240 to explicitly model both positive and negative pref-
241 erences, motivating the adoption of ABSA-based
242 approaches.

243 2.3 Aspect-based Sentiment Analysis

244 Aspect-Based Sentiment Analysis (ABSA) aims to
245 identify aspects mentioned in text and determine
246 the sentiment polarity associated with each aspect.
247 ABSA has been widely adopted for fine-grained
248 user preference modeling in recommendation sys-
249 tems. While early ABSA models relied on LSTM-
250 and CNN-based architectures, recent approaches
251 have achieved substantial performance gains using
252 pre-trained language models.

253 Recent studies focus on more precise model-
254 ing of relationships between contextual semantics
255 and aspects. DMAN dynamically computes to-
256 ken importance to reduce contextual noise (Xu
257 et al., 2025), while DR-BERT enhances aspect-
258 specific representations through context-aware to-
259 ken reweighting (Zhao et al., 2025). In addition, Ze-
260 roABSA and DS2-ABSA enable effective ABSA
261 in low-resource and cross-domain settings through
262 LLM-based data synthesis and domain generaliza-
263 tion (Cheng et al., 2023; Yang et al., 2025a). These
264 advances provide a robust foundation for extracting
265 reliable aspect-level positive and negative prefer-
266 ence signals from review text.

267 2.4 Explainable Recommendation and Rating 268 Prediction Models

269 Explainability has become increasingly important
270 for improving user trust and system transparency
271 in recommender systems. Early studies proposed
272 expressing aspect importance based on model-
273 internal structures or learning interpretable user
274 profiles from review text (Lei et al., 2024; Shimizu
275 et al., 2025; Peng et al., 2024). With the emer-
276 gence of LLMs, natural language-based explana-
277 tion generation has been actively explored. Mod-
278 els such as MAPLE, RecExplainer, and Disentan-
279 gling Likes and Dislikes align user preferences and
280 model signals with natural language explanations
281 (Bang and Song, 2025; Penaloza et al., 2024; Koren
282 et al., 2009). More recent work further investigates
283 uncertainty-aware explanations and LLM-friendly
284 user profile representations (Mnih and Salakhutdi-
285 nov, 2007; Ramos et al., 2024; Huang et al., 2013).

286 In parallel, traditional rating prediction models
287 include MF, PMF, SVD++, and NMF (He et al.,
288 2017; Li et al., 2021; Zheng et al., 2017), while

289	neural models such as NeuMF, DAML, and BERT-	dataset contains <i>user_id</i> , <i>hotel_id</i> , star ratings, and	338
290	MEF integrate textual and semantic information	and review text, and includes fine-grained evaluations	339
291	to improve performance (Tay et al., 2018; Chen	of accommodation attributes such as room quality,	340
292	et al., 2018; Cheng et al., 2018). However, existing	cleanliness, facilities, location, and staff service.	341
293	approaches remain limited in jointly leveraging	Although the datasets differ in review length dis-	342
294	aspect-level positive and negative preferences for	tributions, sparsity levels, and rating distributions,	343
295	both recommendation and explanation.	all provide sufficiently rich textual information re-	344
296	By integrating these research streams, this work	quired for ABSA-based multi-dimensional prefer-	345
297	distinguishes itself by constructing ABSA-based	ence profiling.	346
298	multi-aspect preference vectors from review text		
299	and dynamically adjusting aspect importance ac-	3.2.2 Filtering and Cleaning	347
300	cording to natural language query context. We	To ensure data quality and reliable preference esti-	348
301	empirically validate the effectiveness of the pro-	mation, we apply consistent preprocessing criteria	349
302	posed TASTE framework on large-scale real-world	across all datasets. We retain users with at least 10	350
303	datasets from the travel and dining domains.	reviews and items with at least 5 reviews, remove	351
304		duplicate entries, and inspect rating distributions to	352
	3 Datasets	exclude extreme outliers.	353
305			
	3.1 Data Sources and Domain Coverage	3.3 Aspect Definition and Data Splitting	354
306	To analyze the complex and context-dependent	Rather than relying on fully unsupervised as-	355
307	decision-making process in travel recommendation,	pect extraction, we predefine domain-specific	356
308	we employ three review-based datasets covering	aspect sets to enhance interpretability and enable	357
309	complementary travel-related domains. Travel ex-	consistent cross-domain comparison. For the	358
310	periences involve multiple interrelated decisions,	restaurant domain, we define the aspect set	359
311	including dining, sightseeing and activities, and ac-	$\mathcal{A} = \{food, service, ambience, price, location\}$.	360
312	commodation, each governed by distinct evaluation	For the travel and hotel domains, we define $\mathcal{A} =$	361
313	criteria and preference structures. Accordingly, we	$\{room_quality, service, facilities, location, value\}$.	362
314	select datasets that represent these major compo-	These aspects correspond to core evaluation	363
315	nents of travel decision-making.	factors that have been shown to strongly influence	364
316	Specifically, we use the Yelp Restaurant dataset	user satisfaction. Details of the ABSA model	365
317	to model dining experiences, a Yelp Travel-related	architecture and training procedure are described	366
318	dataset to capture tourist attractions and activity-	in Section 4.	367
319	related evaluations, and the TripAdvisor Hotel	All datasets are split into training, validation,	368
320	(Hong Kong) dataset to analyze accommodation	and test sets using an 80/10/10 ratio with a fixed	369
321	preferences. All datasets contain numerical ratings	random seed to ensure reproducibility. The split is	370
322	accompanied by relatively long and descriptive re-	performed at the user–item interaction level, with	371
323	view texts, making them well-suited for aspect-	the constraint that each user and item appears at	372
324	level sentiment analysis and personalized prefer-	least once in the training set to avoid cold-start	373
325	ence modeling.	issues during evaluation.	374
326			
	3.2 Data Description and Preprocessing	3.4 Summary Statistics	375
327		Table 1 summarizes the statistics of the final	376
	3.2.1 Dataset Characteristics	datasets. Due to computational constraints, we	377
328	The Yelp Restaurant dataset consists of <i>user_id</i> ,	randomly sample approximately 5,000 users from	378
329	<i>business_id</i> , star ratings, and review text, and pro-	the Yelp Restaurant and TripAdvisor Hotel datasets	379
330	vides rich evaluations of core dining aspects such	for experimental evaluation.	380
331	as food quality, service, ambiance, price, and loca-		
332	tion. The Yelp Travel-related dataset is constructed		
333	by filtering travel-, attraction-, and activity-related		
334	categories from the full Yelp corpus. It captures ex-		
335	periential and context-dependent factors, including		
336	sightseeing experiences, surrounding environments,		
337	and activity satisfaction. The TripAdvisor Hotel		

Dataset	#Users	#Items	#Reviews
Yelp Restaurant	5,000	31,321	129,489
Yelp Travel	948	2,942	14,658
TripAdvisor Hotel	5,000	6,280	177,726

Table 1: Summary statistics of the datasets used in this study.

4 Method

4.1 Aspect-Based Sentiment Profiling

We propose **TASTE** (Query-Aware Travel Recommendation via Aspect-based Sentiment Profiling), a recommendation framework that dynamically adapts aspect importance based on natural language queries while preserving long-term personalized preferences.

4.2 Problem Formulation

Let $\mathcal{U} = \{u_1, \dots, u_M\}$ denote a set of users and $\mathcal{I} = \{i_1, \dots, i_N\}$ denote a set of items (restaurants or travel entities). Each user-item interaction (u, i) is associated with a rating r_{ui} and a review text t_{ui} . We define a set of aspects $\mathcal{A} = \{a_1, \dots, a_K\}$ representing key preference dimensions (e.g., food, service, ambiance, price, location).

Given a user u , a natural language query q , and the user’s historical reviews, our objectives are to: (1) generate a ranked list of top- K item recommendations, (2) provide natural language explanations grounded in aspect-level contributions, and (3) dynamically adjust recommendation criteria according to query context.

4.2.1 Aspect Sentiment Extraction

For each review text t_{ui} , we extract aspect-specific sentiment scores using a pre-trained DeBERTa-v3 model fine-tuned for Aspect-Based Sentiment Analysis (ABSA). The ABSA model outputs a probability distribution over sentiment polarities (negative, neutral, positive) for each aspect. We compute the sentiment score for aspect a as:

$$s_{ui}^a = p_{ui}^a(pos) - p_{ui}^a(neg), \quad (1)$$

where $p_{ui}^a(\cdot)$ denotes the predicted probability for aspect a . For each user-item pair, we obtain an aspect sentiment vector $\mathbf{s}_{ui} \in R^K$.

4.2.2 User Preference Normalization

To account for individual rating biases, we normalize ratings using user-specific statistics:

$$\tilde{r}_{ui} = \frac{r_{ui} - \mu_u}{\sigma_u}, \quad (2)$$

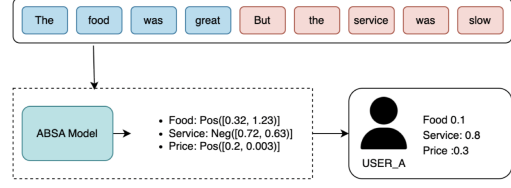


Figure 1: Illustration of the ABSA extraction process.

where μ_u and σ_u are the mean and standard deviation of user u ’s historical ratings. The normalized rating \tilde{r}_{ui} is used as the training target.

4.3 TASTE Model Architecture

The TASTE model consists of four components: (1) embedding layers, (2) aspect attention mechanism, (3) query-aware dynamic weighting, and (4) rating prediction.

4.3.1 Embedding Layer

We learn latent representations for users and items, $\mathbf{e}_u, \mathbf{e}_i \in R^d$, and user-specific aspect bias vectors $\mathbf{b}_u \in R^K$. Aspect biases capture each user’s inherent preference tendencies across aspects.

4.3.2 Aspect Attention Mechanism

Aspect sentiment scores are first adjusted by user biases:

$$\hat{\mathbf{s}}_{ui} = \mathbf{s}_{ui} + \mathbf{b}_u. \quad (3)$$

We compute aspect attention weights using a multi-layer perceptron (MLP) with softmax normalization:

$$\alpha_{ui} = \text{softmax}(MLP([\mathbf{e}_u; \mathbf{e}_i; \hat{\mathbf{s}}_{ui}])), \quad (4)$$

where $[\cdot; \cdot]$ denotes vector concatenation. The attention vector α_{ui} reflects the relative importance of aspects learned from historical interactions.

4.3.3 Query-Aware Dynamic Weighting

To incorporate situational context, we extract query-specific aspect weights β_q from the user query q using an LLM, subject to $\sum_a \beta_q^a = 1$. We combine learned attention and query weights via linear interpolation:

$$\gamma_{ui} = \lambda \alpha_{ui} + (1 - \lambda) \beta_q, \quad (5)$$

where $\lambda \in [0, 1]$ controls the balance between long-term preferences and query intent. We set $\lambda = 0.5$ in all experiments.

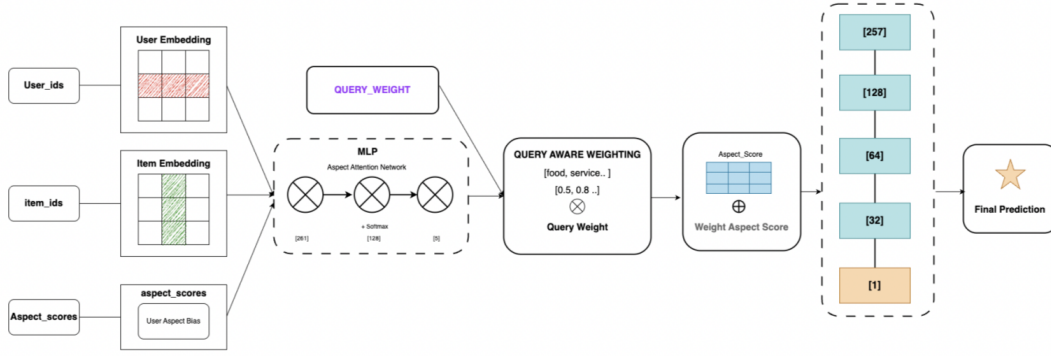


Figure 2: Overview of the TASTE model architecture.

4.3.4 Rating Prediction

We compute the query-aware weighted aspect score:

$$z_{ui} = \gamma_{ui}^T \hat{\mathbf{s}}_{ui}. \quad (6)$$

The final rating prediction is given by:

$$\hat{r}_{ui} = f(z_{ui}, \mathbf{e}_u, \mathbf{e}_i) + b_i + b, \quad (7)$$

where $f(\cdot)$ is an MLP with ReLU activations and dropout, b_i is an item bias, and b is a global bias. Predictions are denormalized to obtain ratings on the original scale.

4.4 Training Objective

We train the model using Smooth L1 loss on normalized ratings:

$$\mathcal{L} = \sum_{(u,i) \in \mathcal{D}_{train}} \text{SmoothL1}(\tilde{r}_{ui}, \hat{r}_{ui}) + \eta \Theta_2^2, \quad (8)$$

where Θ denotes all trainable parameters and η is the L2 regularization coefficient.

4.5 Explainable Recommendation with LLM

For each recommended item, we generate natural language explanations grounded in the query-aware aspect weights γ_{ui} . The LLM takes as input the user query, predicted rating, and aspect weights, and produces a concise explanation focusing on the most influential aspects. Crucially, the LLM does not make recommendation decisions; it only translates quantitative model outputs into human-readable explanations, ensuring faithfulness and controllability.

4.6 Query-Aware Recommendation Pipeline

Given a user u and query q , we extract query-specific aspect weights, compute learned attention

for candidate items, and fuse both signals to obtain query-aware weights. Items are ranked by predicted ratings, and explanations are generated for the top- K recommendations. Explanation quality is evaluated using an LLM-as-a-Judge protocol.

5 Experiments and Evaluation

We evaluate the proposed TASTE framework through a combination of quantitative and qualitative experiments. The evaluation focuses on three aspects: (i) rating prediction accuracy, (ii) query-aware adaptability, and (iii) explanation quality. Specifically, we address the following research questions.

RQ1 (Performance): Does TASTE achieve competitive rating prediction performance compared to existing baselines?

RQ2 (Query-awareness): Does TASTE dynamically adjust aspect importance and recommendation outcomes in response to different user queries?

RQ3 (Explanation Quality): Do aspect-grounded, LLM-based explanations improve explanation quality and faithfulness?

All experiments are conducted on three real-world datasets: Yelp Restaurant, Yelp Travel-related, and TripAdvisor Hotel. Dataset details and preprocessing procedures are described in Section 3.

5.1 Experimental Setup

5.1.1 Baselines

We compare TASTE with representative models from three categories.

Traditional Collaborative Filtering: MF (Zheng et al., 2017), NeuMF (Tay et al., 2018).

Review-based Models: DeepCoNN (?), MPCN (?), DAML (?), and NARRE (?).

Parameter	Value	Description
Embedding Dimension	128	Latent vector size
Number of Aspects	5	Domain-specific aspects
Batch Size	256	Mini-batch size
Learning Rate	1e-3	AdamW optimizer
Weight Decay	1e-3	L2 regularization
Max Epochs	30	Training iterations
Early Stopping	5	Validation patience
Query Weight λ	0.5	Attention-query balance
LLM Model	GPT-4	Query parsing and explanation

Table 2: Hyperparameters for TASTE.

Aspect-aware Models: ANCF (?), which incorporates aspect attention but does not support query-aware adaptation.

All models are trained on identical data splits with the same rating normalization strategy.

5.1.2 Training Configuration

Table 2 summarizes the hyperparameters used for TASTE.

5.2 Rating Prediction Performance (RQ1)

5.2.1 Evaluation Metrics

We evaluate rating prediction using Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE):

$$MAE = \frac{1}{|\mathcal{D}_{test}|} \sum_{(u,i) \in \mathcal{D}_{test}} |r_{ui} - \hat{r}_{ui}|, \quad (9)$$

$$RMSE = \sqrt{\frac{1}{|\mathcal{D}_{test}|} \sum_{(u,i) \in \mathcal{D}_{test}} (r_{ui} - \hat{r}_{ui})^2}. \quad (10)$$

All metrics are computed on denormalized ratings (1–5 scale).

5.2.2 Results

Table 3 reports rating prediction performance across Yelp Restaurant, Yelp Travel-related, and TripAdvisor Hotel datasets. TASTE achieves the best or highly competitive MAE and RMSE on all three datasets, demonstrating that query-aware aspect modeling improves robustness without sacrificing predictive accuracy. TASTE achieves consistently competitive performance across all datasets. While it does not always achieve the absolute best score, it outperforms most review-based models and maintains stable performance across domains. This suggests that explicit aspect modeling improves robustness without sacrificing accuracy.

5.3 Query-aware Recommendation Effects (RQ2)

We qualitatively evaluate query-aware behavior by issuing different queries for the same user and ob-

Model	Yelp Restaurant		Yelp Travel		TripAdvisor Hotel	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
MF	1.0925	0.8398	0.9522	0.7318	0.8260	0.6286
NeuMF	1.0819	0.8261	0.9465	0.7110	0.8451	0.6373
DeepCoNN	0.9865	0.7486	0.8953	0.6803	0.5839	0.7843
MPCN	0.9936	0.7566	0.9013	0.6832	0.5856	0.7784
DAML	0.9865	0.7486	0.9040	0.6832	0.5866	0.7795
NARRE	0.7600	0.5789	0.7383	0.5590	0.5029	0.6703
ANCF	0.7644	0.5920	0.7714	0.5709	0.5161	0.6903
TASTE	0.7577	0.5752	0.7524	0.5709	0.5157	0.6940

Table 3: Rating prediction performance (MAE and RMSE) across three datasets. Lower values indicate better performance. Best results for each metric are highlighted in bold.

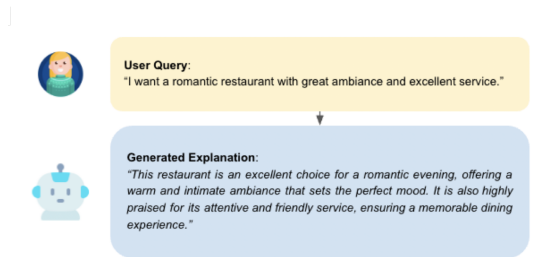


Figure 3: High-scoring example.

serving changes in aspect weights and top- K rankings. Table 4 illustrates representative query scenarios and resulting aspect weight shifts on Yelp Restaurant.

Aspect importance shifts align well with query intent and lead to meaningful ranking changes. Additional multilingual case studies are provided in Appendix A.

5.4 Explanation Quality (RQ3)

5.4.1 LLM-as-a-Judge Evaluation

We adopt the LLM-as-a-Judge protocol (?) to evaluate explanation quality. Explanations are scored along five dimensions: relevance, coherence, informativeness, aspect alignment, and naturalness.

5.4.2 Results

Table 5 reports average scores across datasets. Results indicate consistently high explanation quality, with strong coherence and relevance. Aspect alignment scores validate that explanations faithfully reflect model-internal reasoning. Representative examples are shown in Figure 3, with additional cases in the appendix.

Aspect	No Query	Romantic	Value	Convenience
Food	0.20	0.15	0.35	0.18
Service	0.20	0.30	0.15	0.20
Ambiance	0.20	0.40	0.10	0.12
Price	0.20	0.05	0.35	0.32
Location	0.20	0.10	0.05	0.18

Table 4: Aspect weight changes under different queries.

Dataset	Overall	Rel.	Coh.	Align.	Nat.
Yelp Restaurant	4.85	5.00	5.00	4.80	5.00
Yelp Travel	4.87	5.00	5.00	4.67	5.00
TripAdvisor	4.89	4.98	5.00	4.78	5.00

Table 5: LLM-as-a-Judge evaluation results.

6 Conclusion and Limitations

6.1 Conclusion

This paper proposed TASTE, a query-aware recommendation framework that integrates aspect-based sentiment profiling with dynamic weighting derived from natural language queries. Across three real-world datasets, TASTE achieved competitive rating prediction performance (MAE/RMSE) compared to strong rating- and review-based baselines, while additionally providing aspect-level interpretability and query-aware adaptation. Qualitative analyses further showed that different queries from the same user yield meaningful shifts in aspect importance and corresponding changes in top- K recommendations, indicating that TASTE can capture situational preference variations beyond static user profiles. Moreover, TASTE generates faithful and interpretable explanations grounded in quantitative aspect contributions. By restricting the role of LLMs to query interpretation and explanation generation (rather than decision-making), the framework maintains alignment between model reasoning and natural language explanations, resulting in consistently high explanation quality.

6.2 Limitations and Future Work

Despite its effectiveness, this work has several limitations. First, aspect granularity is constrained by a fixed, predefined aspect set (e.g., five aspects per domain), which may miss finer-grained or latent factors (e.g., quietness, hygiene, family-friendliness). Future work could explore automatic aspect discovery or hierarchical aspect representations.

Second, the fusion coefficient λ is fixed (we use $\lambda = 0.5$), which may not be optimal across users and query types. A natural extension is to learn

λ dynamically from user characteristics, query semantics, or uncertainty estimates.

Third, TASTE prioritizes interpretability and query-aware controllability, and thus does not consistently achieve the lowest prediction error across all settings. Future work could investigate hybrid architectures that preserve aspect interpretability while improving predictive accuracy.

Fourth, explanation evaluation is limited by the lack of strong, directly comparable baselines for query-aware, aspect-grounded explanations. Future studies could develop controlled explanation baselines or compare against emerging LLM-based explanation frameworks.

Finally, we rely on LLM-as-a-Judge for scalable evaluation; however, human-centered studies are necessary to assess perceived usefulness, trust, and satisfaction. Future work should incorporate user studies (e.g., surveys or A/B tests) to validate the practical impact of query-aware explanations.

Overall, TASTE demonstrates that a lightweight, weight-based aspect modeling approach can balance accuracy, contextual adaptability, and explanation faithfulness, providing a foundation for human-centered conversational recommendation.

References

- S. Bang and H. Song. 2025. Llm-based user profile management for recommender systems. *arXiv preprint arXiv:2502.14541*.
- Y. Cao, N. Mehta, X. Yi, R. H. Keshavan, L. Heldt, L. Hong, and M. Sathiamoorthy. 2024. Aligning large language models with recommendation knowledge. In *Findings of NAACL*, pages 1051–1066.
- C. Chen, M. Zhang, Y. Liu, and S. Ma. 2018. Neural attentional rating regression with review-level explanations. In *WWW*, pages 1583–1592.
- H. Cheng, S. Wang, W. Lu, W. Zhang, M. Zhou, K. Lu, and H. Liao. 2023. Explainable recommendation with personalized review retrieval and aspect learning. *arXiv preprint arXiv:2306.12657*.
- Z. Cheng, Y. Ding, X. He, L. Zhu, X. Song, and M. S. Kankanhalli. 2018. a^3nfc : An adaptive aspect attention model for rating prediction. In *IJCAI*, pages 3748–3754.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186.
- Expedia. 2025. Expedia Korea. <https://www.expedia.co.kr/>. Accessed: Dec. 23, 2025.

663	S. Gheewala, S. Xu, S. Yeom, and S. Maqsood. 2024.	M. Payandenick and M. K. Othman. 2026. A multi-	716
664	Exploiting deep transformer models in textual review	criteria recommendation system for personalised	717
665	based recommender systems. <i>Expert Systems with</i>	tourism experiences with user query analysis. <i>In-</i>	718
666	<i>Applications</i> , 235:121120.	<i>formation Technology & Tourism</i> , 28(1):3.	719
667	E. Hasan, M. Rahman, C. Ding, J. X. Huang, and	E. Penaloza, O. Gouvert, H. Wu, and L. Charlin. 2024.	720
668	S. Raza. 2025. Review-based recommender systems:	Tears: Textual representations for scrutable recom-	721
669	A survey of approaches, challenges and future per-	mendations. <i>arXiv preprint arXiv:2410.19302</i> .	722
670	spectives. <i>ACM Computing Surveys</i> , 58(1).	Y. Peng, H. Chen, C.-S. Lin, G. Huang, J. Hu, H. Guo,	723
671	X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua.	and X. Wang. 2024. Uncertainty-aware explainable	724
672	2017. Neural collaborative filtering. In <i>WWW</i> , pages	recommendation with large language models. In	725
673	173–182.	<i>IJCNN</i> , pages 1–8.	726
674	K. Huang, N. D. Sidiropoulos, and A. Swami. 2013.	J. Ramos, H. A. Rahmani, X. Wang, X. Fu, and A. Li-	727
675	Non-negative matrix factorization revisited. <i>IEEE</i>	pani. 2024. Transparent and scrutable recommenda-	728
676	<i>Transactions on Signal Processing</i> , 62(1):211–224.	tions using natural language user profiles. In <i>Pro-</i>	729
677	Jeju Air. 2025. Jeju Air AI Chatbot “Hi-	<i>ceedings of the Annual Meeting of the Association</i>	730
678	Jeko”. https://cschatbot-jwhiz.jejuair.net/	<i>for Computational Linguistics</i> .	731
679	chatbot?language=KO . Accessed: Dec. 23, 2025.	Y. Ren, Z. Chen, X. Yang, L. Li, C. Jiang, L. Cheng, and	732
680	Y. Koren, R. Bell, and C. Volinsky. 2009. Matrix factor-	J. Zhou. 2024. Enhancing sequential recommenders	733
681	ization techniques for recommender systems. <i>Com-</i>	with augmented knowledge from aligned large lan-	734
682	<i>puter</i> , 42(8):30–37.	guage models. In <i>SIGIR</i> , pages 345–354.	735
683	Y. Lei, J. Lian, J. Yao, X. Huang, D. Lian, and X. Xie.	S. Renjith, A. Sreekumar, and M. Jathavedan. 2020. An	736
684	2024. Recexplainer: Aligning large language models	extensive study on the evolution of context-aware per-	737
685	for explaining recommendation models. In <i>KDD</i> ,	sonalized travel recommender systems. <i>Information</i>	738
686	pages 1530–1541.	<i>Processing & Management</i> , 57(1):102078.	739
687	L. Li, Y. Zhang, and L. Chen. 2021. Personalized	J. L. Sarkar, A. Majumder, C. R. Panigrahi, S. Roy, and	740
688	transformer for explainable recommendation. <i>arXiv</i>	B. Pati. 2023. Tourism recommendation system: A	741
689	<i>preprint arXiv:2105.11601</i> .	survey and future research directions. <i>Multimedia</i>	742
690	J. Lin and X. Liu. 2024. Context-awareness-based in-	<i>Tools and Applications</i> , 82(6):8983–9027.	743
691	telligent recommendation method of tourism service	C. Shen, W. Wei, D. Wang, and Z. Wang. 2025. Zero-	744
692	resources. In <i>Proceedings of the 8th International</i>	shot cross-domain aspect-based sentiment analysis	745
693	<i>Conference on Electronic Information Technology</i>	via domain-contextualized chain-of-thought reason-	746
694	<i>and Computer Engineering</i> , pages 5–10.	ing. In <i>Findings of EMNLP</i> , pages 4558–4573.	747
695	C.-C. Liu, H.-R. Yao, D.-C. Chang, and O. Frieder.	R. Shimizu, T. Wada, Y. Wang, J. Kruse, S. O’Brien, and	748
696	2025a. Treatrag: A framework for personalized treat-	J. McAuley. 2025. Disentangling likes and dislikes in	749
697	ment recommendation. In <i>Proceedings of the ACM</i>	personalized generative explainable recommendation.	750
698	<i>Conference on Recommender Systems</i> , pages 690–	In <i>The Web Conference</i> , pages 4793–4809.	751
699	695.	J. Sun, S. Qian, Z. Han, W. Li, Z. Qian, D. Yang, and	752
700	D. Liu, J. Li, B. Du, J. Chang, and R. Gao. 2019. Daml:	G. Xue. 2025. Lkd-kgc: Domain-specific kg con-	753
701	Dual attention mutual learning between ratings and	struction via llm-driven knowledge dependency pars-	754
702	reviews for item recommendation. In <i>KDD</i> , pages	ing. <i>arXiv preprint arXiv:2505.24163</i> .	755
703	344–352.	Y. Tay, A. T. Luu, and S. C. Hui. 2018. Multi-pointer	756
704	J. Liu, T. Li, D. Wu, Z. Tang, Y. Fang, and Z. Yang.	co-attention networks for recommendation. In <i>KDD</i> ,	757
705	2025b. An aspect performance-aware hypergraph	pages 2309–2318.	758
706	neural network for review-based recommendation.	J. Tian, Z. Wang, J. Zhao, and Z. Ding. 2024. Mmrec:	759
707	In <i>WSDM</i> , pages 503–511.	Llm-based multi-modal recommender system. In	760
708	Q. Liu, X. Zhao, Y. Wang, Y. Wang, Z. Zhang, Y. Sun,	<i>SMAP</i> , pages 105–110.	761
709	and F. Tian. 2024. Large language model enhanced	Q. Wang, J. Li, S. Wang, Q. Xing, R. Niu, H. Kong, and	762
710	recommender systems: A survey. <i>arXiv preprint</i>	C. Zhang. 2024a. Towards next-generation llm-based	763
711	<i>arXiv:2412.13432</i> .	recommender systems: A survey and beyond. <i>arXiv</i>	764
712	A. Mnih and R. R. Salakhutdinov. 2007. Probabilistic	<i>preprint arXiv:2410.19744</i> .	765
713	matrix factorization. <i>NeurIPS</i> , 20.	Y. Wang, Y. Wang, Z. Fu, X. Li, W. Wang, Y. Ye,	766
714	OpenAI. 2025. ChatGPT Image Generation. https://	and R. Tang. 2024b. Llm4msr: An llm-enhanced	767
715	chatgpt.com/images/ . Accessed: Dec. 23, 2025.	paradigm for multi-scenario recommendation. In	768
		<i>CIKM</i> , pages 2472–2481.	769

- 770 L. Wu, Z. Zheng, Z. Qiu, H. Wang, H. Gu, T. Shen, and
771 E. Chen. 2024. A survey on large language models
772 for recommendation. *World Wide Web*, 27(5).
- 773 H. Xu, Y. Zhang, Q. Wang, and R. Xu. 2025. Ds2-absa:
774 Dual-stream data synthesis with label refinement for
775 few-shot aspect-based sentiment analysis. In *ACL*,
776 pages 15460–15478.
- 777 C.-W. Yang, Z.-Q. Feng, Y.-J. Lin, C.-W. Chen, K.-D.
778 Wu, H. Xu, and H.-Y. Kao. 2025a. Maple: Enhanc-
779 ing review generation with multi-aspect prompt learn-
780 ing in explainable recommendation. In *ACL*, pages
781 31803–31821.
- 782 W. Yang, W. Zhang, Y. Liu, Y. Han, Y. Wang, J. Lee,
783 and P. S. Yu. 2025b. Cold-start recommendation with
784 knowledge-guided retrieval-augmented generation.
785 *arXiv preprint arXiv:2505.20773*.
- 786 Z. Yue, S. Rabhi, G. D. S. P. Moreira, D. Wang, and
787 E. Oldridge. 2023. Llamarec: Two-stage recom-
788 mendation using large language models for ranking.
789 *arXiv preprint arXiv:2311.02089*.
- 790 C. Zhang, S. Xue, J. Li, J. Wu, B. Du, D. Liu, and
791 J. Chang. 2023. Multi-aspect enhanced graph neural
792 networks for recommendation. *Neural Networks*,
793 157:90–102.
- 794 K. Zhang, M. Zhang, H. Zhao, Q. Liu, W. Wu, and
795 E. Chen. 2022. Incorporating dynamic semantics
796 into pre-trained language model for aspect-based sen-
797 timent analysis. *arXiv preprint arXiv:2203.16369*.
- 798 W. Zhao, P. Zhang, H. Gu, D. Li, T. Lu, and N. Gu.
799 2025. Aotree: Aspect order tree-based model for
800 explainable recommendation. In *AAAI*, pages 2273–
801 2286.
- 802 X. Zhao, M. Wang, X. Zhao, J. Li, S. Zhou, D. Yin, and
803 R. Guo. 2023. Embedding in recommender systems:
804 A survey. *arXiv preprint arXiv:2310.18608*.
- 805 L. Zheng, V. Noroozi, and P. S. Yu. 2017. Joint deep
806 modeling of users and items using reviews for recom-
807 mendation. In *WSDM*, pages 425–434.
- 808 H. Zhou, H. Gu, Z. Zhan, X. Liu, K. Zhou, Y. Xiao, and
809 T. Chen. 2025. The efficiency vs. accuracy trade-off:
810 Optimizing rag-enhanced llm recommender systems
811 using multi-head early exit. In *ACL*, pages 26443–
812 26458.

813	A Additional Qualitative Examples	• Strong relevance to the user query	855
814	A.1 Appendix A: LLM-as-a-Judge Evaluation Example	• Natural and human-like language	856
815		• Clear and coherent structure	857
816	Item ID: 7Cecref3HqLLpVLhMB1oTQ	• Informative and decision-supportive content	858
817	Evaluation Scores:		
818	• Overall Score: 4.80 / 5.0		
819	• Relevance: 5.0 / 5.0		
820	• Coherence: 5.0 / 5.0		
821	• Informativeness: 5.0 / 5.0		
822	• Aspect Alignment: 4.0 / 5.0		
823	• Naturalness: 5.0 / 5.0		
824	User Query: “I want a romantic restaurant with great ambiance and excellent service.”	A.2 Appendix B: Query-Aware Recommendation Case Study	859
825			860
826	Generated Explanation:	This appendix presents a detailed case study demonstrating the query-aware recommendation mechanism of TASTE using a non-English (Korean) user query. The example illustrates the complete pipeline from query interpretation to final recommendations and explanations.	861
827	If you’re looking for a romantic spot, this restaurant is an absolute gem. The ambiance is enchanting, with dim lighting and cozy seating that sets the perfect mood for a special evening. Paired with exceptional service, where the staff goes above and beyond to make guests feel welcomed and cared for, this place offers a truly memorable dining experience. The thoughtfully curated menu further enhances the overall experience, making it an ideal choice for a special night out.		862
828			863
829			864
830			865
831			866
832		A.2.1 B.1 Query Information	867
833		User Query (English Translation): “I am looking for an affordable restaurant that is suitable for dining with my family.”	868
834		User ID: _BcWyKQL16ndpBdggh2kNA	869
835			870
836		A.2.2 B.2 Query Parsing via LLM	871
837		The primary intent of this query is to find a family-friendly restaurant with good value for money. Accordingly, the highest weight is assigned to the <i>value</i> aspect. Since the query involves family dining, food quality and service are also emphasized. In addition, cleanliness and convenient location are considered important factors for family visits.	872
838		Derived Aspect Weights:	873
839	Aspect Weights:	• Value: 0.30	874
840	• Food: Query = 0.10, Attention = 0.18	• Food: 0.20	875
841	• Service: Query = 0.25, Attention = 0.24	• Service: 0.20	876
842	• Ambiance: Query = 0.50, Attention = 0.33	• Cleanliness: 0.15	877
843	• Price: Query = 0.05, Attention = 0.15	• Location: 0.15	878
844	• Location: Query = 0.10, Attention = 0.09		879
845	LLM-as-a-Judge Justification: The explanation directly addresses the user’s intent by emphasizing romantic ambiance and high-quality service, which are the most important aspects highlighted in the query. The content is well-structured, logically coherent, and provides rich descriptive details that support decision-making. While the explanation could have mentioned food quality more explicitly, it remains highly engaging and natural overall.	A.2.3 B.3 Top-5 Recommendation Example	880
846		Top Recommendation (Rank #1): Item ID: 2o4GOsWEmcdwjgnEw0w9A	881
847		Generated Explanation:	882
848		This restaurant is highly rated for value for money, which directly aligns with your preference for an affordable family dining option. It also receives strong ratings across key aspects such as food quality, service, cleanliness, and location, making it an excellent choice for a comfortable and budget-friendly family meal.	883
849			884
850			885
851			886
852			887
853			888
854			889

898

Reason for Ranking #1:

899

- Highest predicted rating among candidate items (1.01)

900

901

- Strong value aspect score (0.895), matching query emphasis

902

903

- High combined weight on the value aspect (0.282)

904

905

- Consistently high scores across all aspects (0.776–0.89)

906