# ImpReSS: Implicit Recommender System for Support Conversations

Omri Haller
Ben Gurion University of The Negev
Be'er Sheva, Israel
haller@post.bgu.ac.il

Yair Meidan
Ben Gurion University of The Negev
Be'er Sheva, Israel
yair.meidan@gmail.com

Dudu Mimran
Ben Gurion University of The Negev
Be'er Sheva, Israel
dudu@dudumimran.com

Yuval Elovici
Ben Gurion University of The Negev
Be'er Sheva, Israel
elovici@bgu.ac.il

Asaf Shabtai
Ben Gurion University of The Negev
Be'er Sheva, Israel
shabtaia@bgu.ac.il

## Abstract

Following recent advancements in large language models (LLMs), LLM-based chatbots have transformed customer support by automating interactions and providing consistent, scalable service. While LLM-based conversational recommender systems (CRSs) have attracted attention for their ability to enhance the quality of recommendations, limited research has addressed the *implicit* integration of recommendations within customer support interactions. In this work, we introduce ImpReSS, an implicit recommender system designed for customer support conversations. ImpReSS operates alongside existing support chatbots, where users report issues and chatbots provide solutions. Based on a customer support conversation, ImpReSS identifies opportunities to recommend relevant solution product categories (SPCs) that help resolve the issue or prevent its recurrence – thereby also supporting business growth. Unlike traditional CRSs, ImpReSS functions entirely implicitly and does not rely on any assumption of a user's purchasing intent. Our empirical evaluation of ImpReSS's ability to recommend relevant SPCs that can help address issues raised in support conversations shows promising results, including an MRR@1 (and recall@3) of 0.72 (0.89) for general problem solving, 0.82 (0.83) for information security support, and 0.85 (0.67) for cybersecurity troubleshooting. To support future research, our data and code will be shared upon request.

## 1 Introduction

Large language models (LLMs) have become transformative tools in numerous domains, fundamentally changing the way we approach complex tasks in almost all areas of life [37]. Among their countless applications, the integration of LLMs in customer support has been particularly impactful [20], allowing businesses to provide 24/7 assistance and enhance customer satisfaction by improving the speed, consistency, and accuracy of customer interactions [6].

Previous research on LLM-based customer support chatbots has concentrated on improving service quality and reducing human agents' workload with self-service tools and automated responses [20]. However, the opportunity to leverage customer interactions for cross-selling and product recommendations [6] has largely been overlooked. As a result, businesses have missed out on value creation opportunities that could benefit them and their clients. Research on conversational recommender systems (CRSs) has led to advancements in the use of dialogue to understand user

preferences and suggest appropriate items [12]. However, these systems face challenges in customer support contexts, where interactions are primarily problem-focused rather than sales-oriented. In such cases, user preferences are often unavailable, and explicitly requesting them may be inappropriate or irrelevant to the situation, and potentially harm the user experience.

To address this gap, we developed **ImpReSS**: an **imp**licit **re**commender **s**ystem for **s**upport conversations. ImpReSS employs a novel approach that augments support-oriented conversations with product (or service) recommendations by implicitly identifying users' *needs* during the problem-solving stage, rather than users' preferences. In addition, unlike previous CRS research, ImpReSS does not assume any user purchasing intent. As shown in Figure 1, ImpReSS consists of three main steps: (1) *Query Generation*, in which an LLM generates a brief summary (with a diagnosis) of the conversation, as well as a set of preliminary solution product categories (SPCs) derived from the support conversation's summary and diagnosis; (2) *Candidate Retrieval*, in which the query is used to search designated catalogs for the most relevant SPCs; and (3) *Candidate Ranking*, in which the retrieved SPCs are prioritized.

To integrate ImpReSS into real-world business workflows, organizations can map general SPCs to their specific company products (brands), and implement various presentation strategies for the recommendations. In the (*"In-Chat"*) presentation strategy, the chatbot suggests the top-ranked candidate as a natural continuation of the support conversation, following the resolution of the user's issue. Alternatively, multiple top-ranked candidates can be displayed below the conversation interface, similar to e-commerce platforms, under a heading such as *"Users who encountered this problem found these products useful"*. This *"Related Items"* strategy is applicable to both chatbots and online forums operated by commercial entities.

We evaluate ImpReSS using three conversational support datasets from various domains. Our empirical results highlight ImpReSS's promising SPC recommendation capabilities and ability to make accurate recommendations without relying on user preference data. ImpReSS achieved MRR@1 values (and recall@3) of 0.72 (0.89) for general problem solving, 0.82 (0.83) for information security support, and 0.85 (0.67) for cybersecurity troubleshooting.

The contributions of this research can be summarized as follows:
(1) We introduce ImpReSS, a novel LLM-based method for implicit SPC recommendations in conversational support settings,

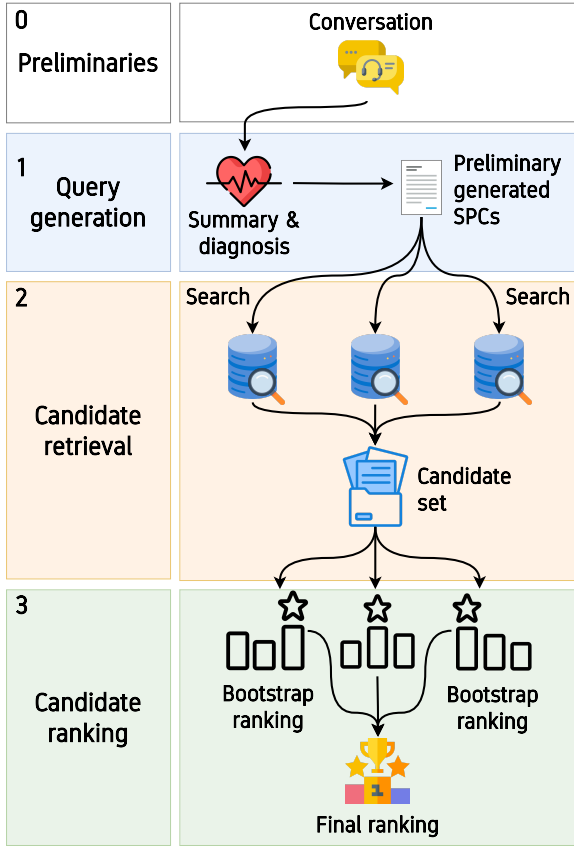Figure 1: The proposed method's pipeline[1].

which assumes no user purchasing intent and requires no user preferences or background information as input.

(2) We empirically evaluate ImReSS using multiple datasets, and show promising performance, achieved early in a conversation.

(3) To promote future research, our data and code will be shared upon request.

## 2 Proposed Method

### 2.1 Approach

Unlike traditional recommendation scenarios [10], which center on product seeking, customer support conversations focus on problem solving. ImReSS addresses this distinction with an implicit recommendation approach that analyzes conversation content instead of relying on user preferences. This enables ImReSS to identify and recommend the most relevant SPCs (or specific products).

### 2.2 Assumptions

ImReSS leverages an existing online support platform (e.g., chatbot or forum) where users describe problems and request help. It can be implemented as an add-on to these platforms, analyzing each support conversation and returning a prioritized list of SPC recommendations.

---

[1]Figure icons obtained from flaticon.com

## 2.3 Key Steps

ImReSS's input is a support conversation, i.e., a set of utterances between a user and a knowledgeable entity, either human or virtual. As illustrated in Figure 1, ImReSS follows a three-step process in order to output a prioritized list of recommendations:

*2.3.1 Query Generation (Step 1).* In this step (illustrated in the example in Figure 2), given a conversation with multiple interactions, an LLM first generates a *conversation summary and diagnosis* object, which concisely outlines the issue raised by the user, a root cause diagnosis, and plausible measures to take. Then, based on the diagnosis, the LLM generates a *query* object, which is a preliminary list of relevant SPCs that can help resolve the issue or prevent its reoccurrence. Each SPC is also briefly explained, facilitating similarity search in the candidate retrieval step that follows.



Figure 2: Query generation example from $DS^{CT}$.

*2.3.2 Candidate Retrieval (Step 2).* In this step (illustrated in the example in Figure 3), ImReSS searches multiple designated databases (DBs) using the query generated in the previous step. These DBs, equipped with an L2 index for efficient similarity search using an embedding model, differ from one another in terms of the aspect of the SPCs they emphasize or their source, thus enabling more diverse candidate retrieval than a single-index approach. The final candidate set is formed by uniting the results from each index.

Note: Bold entries indicate true-labeled SPCs.

**Wi-Fi range extenders**, System repair tools, Network monitoring software, **Driver update software**, System optimization software

Figure 3: Candidate retrieval example from $DS^{CT}$.

*2.3.3 Candidate Ranking (Step 3).* Given a set of retrieved candidates, the LLM ranks the SPCs by their ability to help resolve the diagnosis generated in the query generation step (see Figure 4). To mitigate a possible position bias in ranking, ImReSS employs a bootstrap approach [8], concurrently repeating the ranking process three times with randomly shuffled orders of candidates.

**Table 1: Datasets used to empirically evaluate ImpReSS.**

| Dataset | Domain | Source | Setting | #Records | Example Question | Example SPC |
|---|---|---|---|---|---|---|
| $DS^{CT}$ | Cybersecurity Troubleshooting | Real conversations | Chatbot | 20 | "im in airport and i need to access my banking app do you think its safe to connect to available wifi" | VPN service |
| $DS^{IS}$ | Information Security | Stack Exchange network, IS community [9] | Online forum | 70 | "Precautions to secure (company) laptop & mobile while traveling" | Mobile antivirus |
| $DS^{GE}$ | General Problem Solving | Synthetic | Chatbot | 224 | "... how to cope with my partner's snoring – it's been super tough lately. Any tips on how to handle this without causing much tension?" | Sleep headphones |

Note: Bold entries indicate true-labeled SPCs.

**Wi-Fi range extenders**, **Driver update software**, Network monitoring software, System optimization software, System repair tools

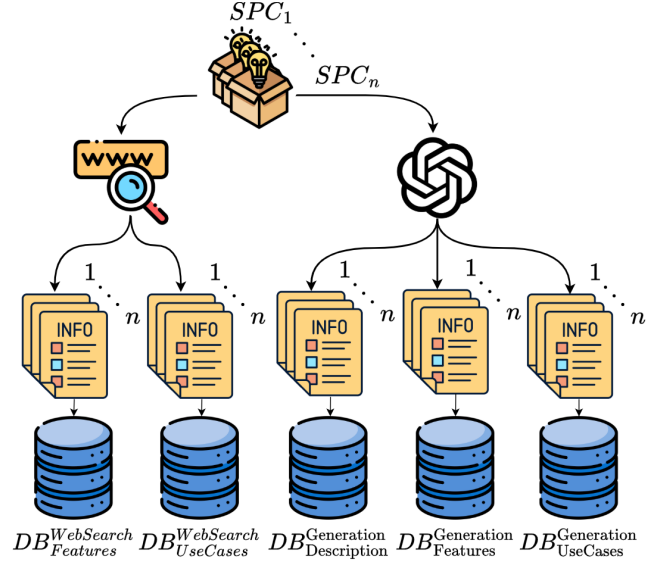**Figure 4: Candidate ranking example from $DS^{CT}$.**

## 3 Evaluation Method

### 3.1 Creation of Datasets

To evaluate ImpReSS, we constructed three datasets, which are summarized in Table 1:

*3.1.1 Cybersecurity Troubleshooting Dataset ($DS^{CT}$).* PC users may occasionally encounter various technical issues, some of which stem from cybersecurity threats. For example, PC slowness may be caused by software conflicts, outdated drivers, aging hardware, or cybersecurity attacks such as cryptojacking. To address this, we developed a cybersecurity-specialized chatbot in our lab and tasked students with troubleshooting a set of predefined complaints using only this chatbot. $DS^{CT}$ consists of these step-by-step cybersecurity troubleshooting conversations.

*3.1.2 Information Security Dataset ($DS^{IS}$).* Stack Exchange [22] is a network of nearly 200 question-and-answer (Q&A) communities where millions collaborate monthly to ask questions, share knowledge, and solve problems across various technical and professional domains. Its most prominent community is the Stack Overflow community [23], which focuses on programming. The Information Security (IS) community [9] on Stack Exchange, features discussions on cryptography, encryption, network security, and related topics. From this community, we extracted multiple Q&A pairs in which the answer was accepted by the question's author. Out of the 239 Q&A pairs that we managed to annotate, we identified 70 pairs containing at least one product recommendation – and these formed our final dataset. In order to (1) prevent data leakage during testing and (2) ensure that the Candidate Generation step (Sec. 2.3.2) relies solely on ImpReSS's conversation summary and diagnosis, we made sure to remove any mentions of specific product recommendations or specific solutions from the answers.

*3.1.3 General Problem-Solving Dataset ($DS^{GE}$).* Following the LLM-based student-teacher interaction simulation framework proposed by Abbasiantaeb et al. [1], we simulated support conversations with a chatbot specialized in a broad range of consultation topics, e.g., kitchen, pet care, and camping. To simulate these conversations, we (1) generated synthetic users based on persona distributions, characterized by attributes such as age, gender, and occupation; and (2) instructed them to interact with an AI assistant (chatbot).



**Figure 5: Catalog DB creation process[1].**

Each conversation begins with the user describing a general problem encountered, followed by up to four Q&A exchanges with the chatbot, enabling it to gather sufficient information. The ground-truth SPC label is derived from the conversation generation prompt, which encapsulates both the essence of the problem and its root cause, to be communicated to the chatbot by the (synthetic) user. Although the user is aware of the cause, they are instructed not to disclose it explicitly. This setup helps maintain coherence and prevents divergence in the conversation.

To ensure conversation quality, we asked three experts to manually rate a random sample containing 10% of $DS^{GE}$. The rating was performed using the USR metrics for dialog generation [13] and achieved very high scores across all assessment dimensions. The generated conversations received perfect scores for Understandable (range [0, 1], mean±standard deviation 1.00±0.00) and Maintains Context ([1, 3], 3.00±0.00), and high rating for Natural ([1, 3], 2.33±0.50) and Overall Quality ([1, 5], 4.01±0.66). Gwet's AC2 coefficient ranged [0.84, 1.00], indicating near-perfect inter-rater agreement on these conversation quality metrics.

### 3.2 Creation of Catalog DBs

As illustrated in Figure 5, we constructed five catalog DBs using both web search-based and LLM-based text generation methods. These approaches are complementary: search yields web pages deemed

relevant by search engines, which likely contain useful information but may (1) include irrelevant content, or (2) lack comprehensive coverage by presenting only a few specific results. In contrast, LLM-based text generation can synthesize information across multiple sources, offering broader summaries, which may come at the cost of factual inaccuracies. In Sec. 4.4, we present an ablation study to empirically assess the contribution of each catalog DB.

*3.2.1 Web Search-Based Catalog DBs Creation.* Using Tavily [25], a web search engine tailored for LLM retrieval-augmented generation (RAG) use cases, we constructed two catalog DBs: $DB_{Features}^{WebSearch}$ and $DB_{UseCases}^{WebSearch}$. For each SPC, Tavily was queried separately for each SPC's features and use cases, and each query returned five results. The retrieved results for each SPC were then concatenated into two separate documents: one document containing all feature-related results, and the other containing all use case-related results. The documents for each of the SPCs formed the respective DB.

*3.2.2 Generation-Based Catalog DBs Creation.* Using GPT-4o, we generated additional three catalog DBs. For each SPC, GPT-4o produced a brief description, a list of key features, and three example use cases. These outputs were stored in $DB_{Descriptions}^{Generation}$, $DB_{Features}^{Generation}$, and $DB_{UseCases}^{Generation}$, respectively.

## 3.3 Experimental Setup

We used GPT-4o mini to create $DS^{GE}$ – both to generate users and the support conversations, where each user query and assistant response were generated independently. We chose to use a high temperature (1.0) for these tasks to elicit diverse responses that better represent conversational scenarios. In contrast, for catalog DB creation and candidate search we used GPT-4o (and text-embedding-3-small [14], with a low temperature (0.3), to keep the outputs grounded and reduce variability in the content generated. We implemented ImpReSS's pipeline on a laptop with an Intel i7-1365U 13th generation processor and 32 GB of RAM, using LangChain and LangSmith. Access to GPT and Llama models was provided via Azure OpenAI Service and Amazon Bedrock, respectively.

## 3.4 Performance Metrics

We evaluated ImpReSS's performance using two metrics commonly used in recommender systems research [36, 38]: MRR@k, which measures how high a relevant SPC ranks among the first k suggestions; and Recall@k (abbreviated as R@k), which measures how many of the total relevant SPCs are successfully retrieved. To align with both presentation strategies, we focused mainly on MRR@1 and R@3. We chose k=1 as most relevant for the "In-Chat" strategy for two reasons: (1) it may not be appropriate for a chatbot to recommend more than one SPC in a support conversation since recommendations are displayed directly within the conversation interface, and (2) most conversations contain one relevant SPC. Thus, MRR@1 effectively reflects how well ImpReSS retrieves and ranks this key recommendation. For "Related Items" representation strategy, which is applicable for multiple SPCs, and for a more comprehensive assessment, we also compute the MRR and recall at larger k values.

**Table 2: ImpReSS's performance across datasets.**

| Dataset | MRR@1 | MRR@3 | MRR@5 | R@1 | R@3 | R@5 |
|---------|-------|-------|-------|-----|-----|-----|
| $DS^{IS}$ | 0.82 | 0.87 | 0.87 | 0.68 | 0.83 | 0.87 |
| $DS^{CT}$ | 0.85 | 0.87 | 0.88 | 0.56 | 0.67 | 0.74 |
| $DS^{GE}$ | 0.72 | 0.78 | 0.81 | 0.72 | 0.89 | 0.93 |

**Table 3: Performance across LLMs and embedding models.**

| LLM | Embed. model | $DS^{IS}$ | | $DS^{CT}$ | | $DS^{GE}$ | |
|-----|-------|-------|-----|-------|-----|-------|-----|
| | | MRR@1 | R@3 | MRR@1 | R@3 | MRR@1 | R@3 |
| 4o | multi-e5 | 0.74 | **0.83** | 0.60 | 0.56 | 0.70 | **0.90** |
| | embed-3 | **0.82** | **0.83** | **0.85** | 0.67 | **0.72** | 0.89 |
| 4o mini | multi-e5 | 0.51 | 0.60 | 0.55 | 0.55 | 0.63 | 0.85 |
| | embed-3 | 0.55 | 0.67 | 0.65 | 0.64 | 0.61 | 0.84 |
| Llama | multi-e5 | 0.62 | 0.66 | 0.55 | 0.56 | 0.66 | 0.85 |
| | embed-3 | 0.62 | 0.73 | 0.80 | **0.69** | 0.67 | 0.84 |

Note: Bold values indicate the best configuration for each metric.

## 4 Results

## 4.1 Performance Across Datasets

As can be seen in Table 2, for each evaluated dataset both the MRR@k and R@k increase with k, although with diminishing effect. The improvement in performance makes sense, since the likelihood of missing relevant SPCs decreases as more SPCs are retrieved. The diminishing effect is a good indication that ImpReSS has already identified and ranked the most important recommendations, so as k increases, we see smaller improvements, if at all.

## 4.2 Sensitivity to LLM and Embedding Models

The results presented in Sec. 4.1 were obtained using a proprietary, state-of-the-art LLM and embedding model, respectively OpenAI's GPT-4o and text-embedding-3-small. To evaluate ImpReSS's sensitivity to the underlying LLM and embedding model, we repeated the previous experiment with additional configurations. That is, for the steps presented in Sec. 2 we used two additional LLMs, GPT-4o mini (less costly) and Llama-3.3-70B-Instruct (open sourced), along with Multilingual-E5-Large-Instruct [28] as an additional embedding model (open sourced, 560M parameters). As can be seen in Table 3, GPT-4o outperformed GPT-4o mini and Llama-3.3-70B-Instruct in almost all cases (frequently with a high margin), as did text-embedding-3-small compared to Multilingual-E5-Large-Instruct (though with a smaller margin). Moreover, the combination of GPT-4o with text-embedding-3-small yielded the highest MRR@1 for 3 out of 3 datasets and the highest R@3 for 2 out of 3 datasets. The superiority of GPT-4o in this experiment is consistent with prior research [17], including recommender systems research [26].

## 4.3 Sensitivity to Conversation Length

The experimental results discussed thus far were achieved using all available utterances in each conversation. To assess ImpReSS's sensitivity to conversation length, Figure 6 shows the MRR@1 and R@3 scores as functions of the number of utterances in the conversation. Since $DS^{IS}$ contains only one user message followed
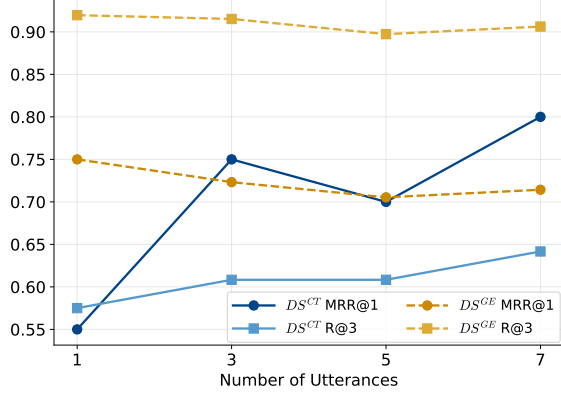
Figure 6: Performance as a function of conversation length.



Figure 8: Performance as a function of bootstrap iterations.

by an assistant response, without further utterances, it is not included in this analysis, which only pertains to $DS^{CT}$ and $DS^{GE}$. The experimental results show that on $DS^{CT}$, both the MRR@1 and R@3 increase with the number of utterances. This makes sense, as cybersecurity troubleshooting often requires iteratively ruling out possible root causes. In other words, the less irrelevant root causes considered by the chatbot, the fewer irrelevant SPCs recommended. In comparison, the user complaints and questions in $DS^{GE}$ were simpler (e.g., planning holiday meals), with most of the information provided in the first question, so high MRR@1 and R@3 values were quickly achieved.

## 4.4 Ablation Study

Our literature review did not reveal any prior research (with or without an experimental dataset made public) on conversational recommendations based on implicit user needs; in fact, we found no research whatsoever on conversational recommendations that are not based on user preferences. Hence, since we were unable to directly compare ImpReSS to a public benchmark, and in order to better explore ImpReSS's potential while assessing the relative importance of its components, we performed an ablation study, which is described below.
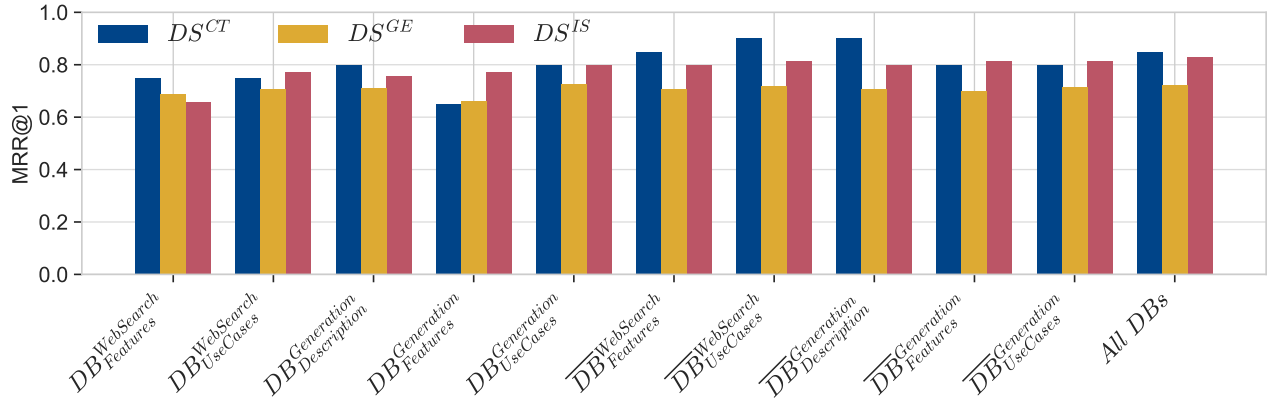
*4.4.1 Importance of the Various Catalog DBs.* As described in Sec. 3.2, the catalog DBs contain SPC features, descriptions, and use cases, obtained using web search and LLM-based text generation. In this experiment, we compared the performance across all three evaluated datasets when using all five DBs together (denoted as *All DBs*) against using each DB individually, as well as against using every combination of four DBs. For example, $DB_{Features}^{WebSearch}$ denotes including *only* this DB, while $\overline{DB}_{Features}^{WebSearch}$ denotes including all other (four) DBs *except for* this DB. As can be seen in Figure 7, on $DS^{IS}$ the *All DBs* configuration outperforms all other configurations, similarly to $DS^{GE}$, and on $DS^{CT}$, the *All DBs* configuration works best in 9 out of 11 of tested configurations; based on this, we can conclude that each of the five DBs contributes to ImpReSS's performance. In future implementations, if overhead limitations, such as token consumption costs, are more restrictive, the best option would be to use any of the Use Cases DBs (based on web search or LLM generation), as retrieving SPC candidates based on their typical use cases is more effective than retrieval based on the SPCs' features or descriptions.

*4.4.2 Importance of Candidates' Bootstrap Ranking.* As discussed in Sec. 4.5, the third step of ImpReSS– candidate ranking – introduces considerable overhead in both cost and latency. To assess



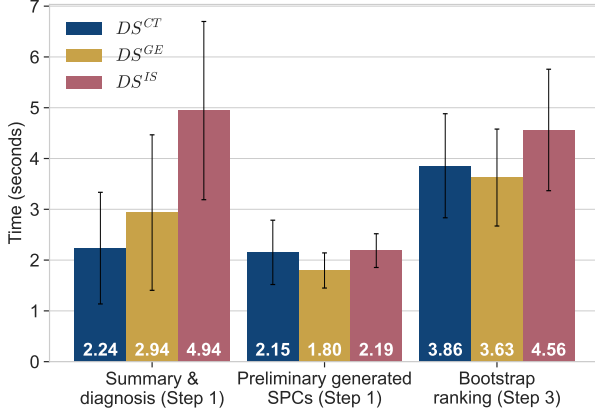Figure 7: Ablation study results for various combinations of catalog DBs.

Figure 9: Time overhead analysis.

the necessity of this step, we evaluated ImpReSS with zero to three bootstrap ranking iterations. The results, shown in Figure 8, reveal a marked performance drop across all datasets when the candidate ranking is disabled, i.e., when ImpReSS relies solely on the candidate SPCs retrieved from the catalog DBs, ranked by their similarity to the generated query (Sec. 2.3.1). Unlike Hou et al. [8], who advocate for at least three ranking iterations, our findings indicate that even one or two iterations substantially improve performance, though still falling short of the optimal results achieved with three.

## 4.5 Overhead

Figures 9 and 10 present ImpReSS's overhead in terms of time (seconds) and token consumption, respectively indicating computation and monetary expenses. Although the overhead was typically higher for $DS^{IS}$ than it was for $DS^{CT}$ and $DS^{GE}$, in most cases the differences were relatively small.

While bootstrap ranking the retrieved SPC candidates is the most resource-consuming step, this step also greatly affects ImpReSS's performance (as discussed in Sec. 4.4.2). The overhead of that step can be reduced in several ways, including the use of a local LLM, which may result in shorter and more consistent waiting times.

## 5 Related Work

### 5.1 Chatbots in Customer Service and Support

Live chat interfaces have become a widely adopted channel for delivering real-time customer service [19]. Customers increasingly rely on these platforms to access information or receive assistance, and the immediacy of live chat strongly influences customer trust and satisfaction [2]. Hardalov et al. [7] examined the automation of customer support using conversational agents, comparing methods such as information retrieval, sequence-to-sequence (Seq2Seq) models, attention mechanisms, and the transformer architecture. Other studies have investigated domain-specific applications in sectors such as banking [32] and healthcare [3], employing various deep learning and natural language processing techniques. More recently, LLMs have been applied to customer support tasks, including fine-tuning OpenAI's GPT-4 for e-commerce [18] and deploying Google's Flan-T5 XXL model for real-time assistance [15]. Additional research has explored LLM-driven approaches to create context-aware and personalized support chatbots [18]. While these methods aim to enhance the quality of customer service, they typically do not utilize conversation data to recommend relevant products or services – an approach that could further support users and drive business growth at the same time.

### 5.2 Conversational Recommender Systems

CRSs collect user preferences and deliver personalized suggestions through interactive natural language dialogues [24]. CRSs are broadly categorized into attribute-based and open-ended systems. Attribute-based CRSs refine user preferences through explicit, attribute-driven Q&A interactions, aiming to identify the optimal item(s) in minimal rounds [11, 16, 33]. In contrast, open-ended CRSs such as RecInDial [27] and UniCRS [29] support free-form conversations, enabling more flexible user interactions. To enhance performance, UniCRS incorporates dialogue history as contextual input and knowledge graph entities as external information to jointly address recommendation and dialogue tasks. Deng et al. [5] proposed UniMIND, a Seq2Seq model that unifies multiple CRS goals – including chitchat, question answering, topic prediction, and recommendation—via prompt-based learning. However, UniMIND's multi-goal approach can detract from recommendation progress,
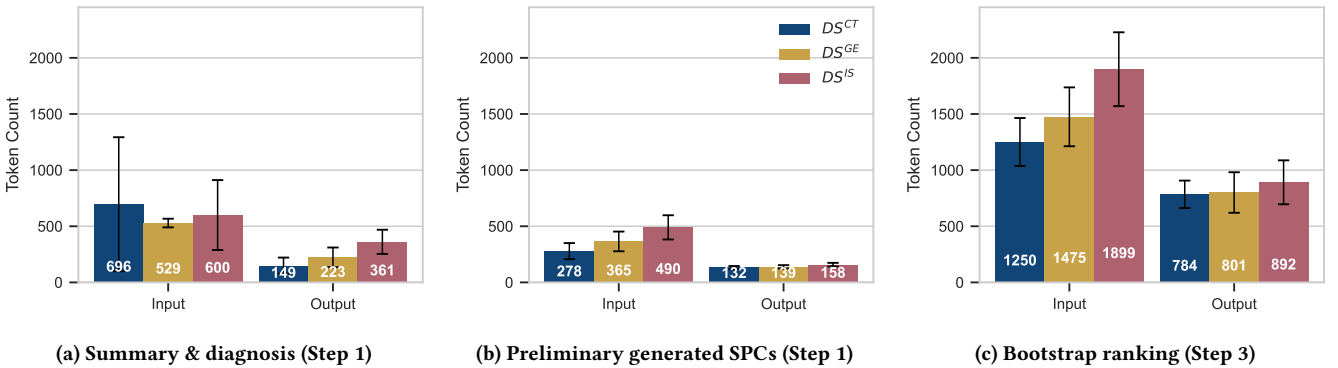


(a) Summary & diagnosis (Step 1)

(b) Preliminary generated SPCs (Step 1)

(c) Bootstrap ranking (Step 3)

Figure 10: Token consumption overhead analysis.

Table 4: Comparison of reviewed conversational and LLM-based recommender systems.

| Paper | Name | Recommender Type | Paradigm | Application Domains | Rec. Intent Assumption? | Conversation Control? |
|-------|------|------------------|----------|---------------------|-------------------------|----------------------|
| [29] | UniCRS | Conversational | Preference-based | Movies | Yes | Yes |
| [5] | MG-CRS | Conversational | Preference-based | Movies, music, restaurants | No | Yes |
| [27] | RecInDial | Conversational | Preference-based | Movies | Yes | Yes |
| [33] | FPAN | Conversational | Preference-based | Music, Yelp dataset | Yes | Yes |
| [30] | RecMind | LLM-powered agent | Preference-based | E-commerce, Yelp dataset | Yes | N.A. |
| [31] | MACRec | Multi-agent collaborative | Preference-based | E-commerce, movies, social media | Yes | N.A. |
| [21] | RAH! | Human-centered, LLM agents | Personality-based | Movies, books, games | Yes | N.A. |
| This | ImpReSS | Support conversation add-on | Need-based | Support conversations | No | No |

especially when chitchat dominates. Fundamentally, CRSs drive recommendations by steering conversations to elicit user preferences and background, assuming that users engage with the intent of receiving suggestions. In contrast, ImpReSS assumes no underlying purchasing intent and exerts no control over the conversation's direction. Functioning as an add-on, ImpReSS analyzes the support conversation, and recommends the most suitable SPCs – the ones capable of addressing the identified issue.

## 5.3 LLM-Based Recommender Systems

LLMs have been integrated in various components of recommender systems, including feature engineering, user and item embeddings, scoring, ranking, or even functioning as agents that guide the recommendation process itself [12]. For item embedding, TedRec [34] performs sequence-level semantic fusion of textual and ID features for sequential recommendation, while NoteLLM [35] combines note semantics with collaborative signals to produce note embeddings. In contrast, Chen et al. [4] proposed a hierarchical approach where an LLM extracts features from item descriptions and converts them into compact embeddings, reducing computational overhead. Other studies leveraged autonomous, LLM-powered agents. Unlike methods that merely prompt LLMs with user history, these approaches exploit advanced agentic capabilities such as planning and tool use. RecMind [30] performs multi-path reasoning through LLM planning, MACRec [31] features collaboration among specialized agents, and RAH [21] is a human-centered framework in which multiple LLM-based agents mediate between users and recommender systems to align suggestions with a user's personality and reduce cognitive load.

While most LLM-based systems focus on predicting what users might *like* based on past behavior or preferences (see Table 4), we propose a paradigm shift: understanding what the user *needs*. ImpReSS identifies recommendation opportunities by detecting needs – explicitly expressed, implicitly inferred, or derived from interactions between a user and a knowledgeable entity, either human or virtual.

## 6 Discussion

### 6.1 Key Insights

ImpReSS's results are promising, with a mean of 0.8 for both MRR@1 and R@3 achieved using a state-of-the-art LLM and embedding model. These performance levels are typically reached early in a conversation, highlighting ImpReSS's ability to produce accurate

recommendations quickly, after just a short interaction with a customer/user. Our ablation study confirms that each of ImpReSS's components contributes to its overall performance, although some overhead is incurred (Sec. 4.5). To mitigate time overhead, we recommend using more powerful hardware, and using parallel LLM calls for bootstrap ranking. To reduce token consumption, our findings indicate that even one or two iterations of bootstrap ranking substantially improve performance, though still falling short of the optimal results achieved with three.

### 6.2 Research Limitations and Future Work

While the results of this study are promising, their generalizability might be affected by factors such as the modest size of the datasets and the limited domain diversity. In future work we will evaluate ImpReSS on additional, larger conversational datasets drawn from a variety of support scenarios. We have also identified three other interesting research directions: (1) conducting an online experiment to examine user conversion rates after receiving recommendations from ImpReSS; (2) optimizing the *timing* of the recommendation within the support conversation; and (3) refining the chatbot's phrasing of recommendations within the "In-Chat" presentation strategy. In addition, given ImpReSS's sensitivity to the underlying LLM selection (Sec. 4.2), future research could explore the use of newly released LLMs, particularly open-source models, that can run locally to both reduce costs and enhance privacy.

## 7 Conclusion

In this paper, we introduced ImpReSS, a novel LLM-based recommender system that, unlike traditional preference-based approaches, implicitly infers user needs through conversation. It then suggests relevant SPCs (or specific products) potentially delivering value to both the customer and the business operating the support chatbot. Our comprehensive quantitative evaluation performed on datasets from three application domains demonstrated ImpReSS's effectiveness, analyzed its sensitivity to various parameters, and assessed the relative contributions of its components. In future work, we plan to broaden our experimental scope and collect online behavioral data, so that we can optimize the timing and refine phrasing of recommendations for the "In-Chat" presentation strategy to maximize conversion rates.

## References

[1] Zahra Abbasiantaeb, Yifei Yuan, Evangelos Kanoulas, and Mohammad Aliannejadi. 2024. Let the llms talk: Simulating human-to-human conversational qa via zero-shot llm-to-llm interactions. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*. 8–17.

[2] Martin Adam, Michael Wessel, and Alexander Benlian. 2021. AI-based chatbots in customer service and their effects on user compliance. *Electronic Markets* 31, 2 (2021), 427–445.

[3] Soufyane Ayanouz, Boudhir Anouar Abdelhakim, and Mohammed Benhmed. 2020. A Smart Chatbot Architecture based NLP and Machine Learning for Health Care Assistance. In *Proceedings of the 3rd International Conference on Networking, Information Systems & Security* (Marrakech, Morocco) *(NISS '20)*. Association for Computing Machinery, New York, NY, USA, Article 78, 6 pages. doi:10.1145/3386723.3387897

[4] Junyi Chen, Lu Chi, Bingyue Peng, and Zehuan Yuan. 2024. HLLM: Enhancing Sequential Recommendations via Hierarchical Large Language Models for Item and User Modeling. arXiv:2409.12740 [cs.IR] https://arxiv.org/abs/2409.12740

[5] Yang Deng, Wenxuan Zhang, Weiwen Xu, Wenqiang Lei, Tat-Seng Chua, and Wai Lam. 2023. A Unified Multi-task Learning Framework for Multi-goal Conversational Recommender Systems. *ACM Trans. Inf. Syst.* 41, 3, Article 77 (Feb. 2023), 25 pages. doi:10.1145/3570640

[6] Esha Manideep Dinne. 2024. How LLMs Are Transforming The Customer Support Industry. https://www.forbes.com/councils/forbestechcouncil/2024/09/20/how-llms-are-transforming-the-customer-support-industry/

[7] Momchil Hardalov, Ivan Koychev, and Preslav Nakov. 2018. Towards Automated Customer Support. In *Artificial Intelligence: Methodology, Systems, and Applications: 18th International Conference, AIMSA 2018, Varna, Bulgaria, September 12–14, 2018, Proceedings* (Varna, Bulgaria). Springer-Verlag, Berlin, Heidelberg, 48–59. doi:10.1007/978-3-319-99344-7_5

[8] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. 2024. Large Language Models are Zero-Shot Rankers for Recommender Systems. In *Advances in Information Retrieval: 46th European Conference on Information Retrieval, ECIR 2024, Glasgow, UK, March 24–28, 2024, Proceedings, Part II* (Glasgow, United Kingdom). Springer-Verlag, Berlin, Heidelberg, 364–381. doi:10.1007/978-3-031-56060-6_24

[9] Information Security Stack Exchange. 2025. Information Security Stack Exchange. https://security.stackexchange.com

[10] MuhammadZaid Katlariwala and Aakash Gupta. 2024. Product Recommendation System Using Large Language Model: Llama-2. In *2024 IEEE World AI IoT Congress (AIIoT)*. IEEE, 0491–0495.

[11] Wenqiang Lei, Gangyi Zhang, Xiangnan He, Yisong Miao, Xiang Wang, Liang Chen, and Tat-Seng Chua. 2020. Interactive Path Reasoning on Graph for Conversational Recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) *(KDD '20)*. Association for Computing Machinery, New York, NY, USA, 2073–2083. doi:10.1145/3394486.3403258

[12] Jianghao Lin, Xinyi Dai, Yunjia Xi, Weiwen Liu, Bo Chen, Hao Zhang, Yong Liu, Chuhan Wu, Xiangyang Li, Chenxu Zhu, Huifeng Guo, Yong Yu, Ruiming Tang, and Weinan Zhang. 2025. How Can Recommender Systems Benefit from Large Language Models: A Survey. *ACM Trans. Inf. Syst.* 43, 2, Article 28 (Jan. 2025), 47 pages. doi:10.1145/3678004

[13] Shikib Mehri and Maxine Eskenazi. 2020. USR: An Unsupervised and Reference Free Evaluation Metric for Dialog Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 681–707.

[14] OpenAI. 2023. text-embedding-3-small. https://openai.com/blog/new-embedding-models-and-api-updates. Accessed: May 4, 2025.

[15] Keivalya Pandya and Mehfuza Holia. 2023. Automating Customer Service using LangChain: Building custom open-source GPT Chatbot for organizations. *arXiv preprint arXiv:2310.05421* (2023).

[16] Xuhui Ren, Hongzhi Yin, Tong Chen, Hao Wang, Zi Huang, and Kai Zheng. 2021. Learning to Ask Appropriate Questions in Conversational Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 808–817. doi:10.1145/3404835.3462839

[17] Rohan Sanghera, Arun James Thirunavukarasu, Marc El Khoury, Jessica O'Logbon, Yuqing Chen, Archie Watt, Mustafa Mahmood, Hamid Butt, George Nishimura, and Andrew AS Soltan. 2025. High-performance automated abstract screening with large language model ensembles. *Journal of the American Medical Informatics Association* (2025), ocaf050.

[18] Kathari Santosh, Timur Kholmukhamedov, M Sandeep Kumar, Mohd Aarif, Iskandar Muda, and B Kiran Bala. 2024. Leveraging GPT-4 Capabilities for Developing Context-Aware, Personalized Chatbot Interfaces in E-commerce Customer Support Systems. In *2024 10th International Conference on Communication and Signal Processing (ICCSP)*. 1135–1140. doi:10.1109/ICCSP60870.2024.10544016

[19] Fu Shang, Fanyi Zhao, Mingxuan Zhang, Jun Sun, and Jiatu Shi. 2024. Personalized recommendation systems powered by large language models: Integrating semantic understanding and user preferences. *International Journal of Innovative Research in Engineering and Management* 11, 4 (2024), 39–49.

[20] Farooq Shareef. 2024. Enhancing Conversational AI with LLMs for Customer Support Automation. In *2024 2nd International Conference on Self Sustainable Artificial Intelligence Systems (ICSSAS)*. 239–244. doi:10.1109/ICSSAS64001.2024.10760403

[21] Yubo Shu, Haonan Zhang, Hansu Gu, Peng Zhang, Tun Lu, Dongsheng Li, and Ning Gu. 2024. RAH! RecSys–Assistant–Human: A Human-Centered Recommendation Framework With LLM Agents. *IEEE Transactions on Computational Social Systems* 11, 5 (2024), 6759–6770. doi:10.1109/TCSS.2024.3404039

[22] Stack Exchange. 2025. Stack Exchange. https://www.stackexchange.com

[23] Stack Exchange. 2025. Stack Overflow. https://stackoverflow.com

[24] Yueming Sun and Yi Zhang. 2018. Conversational Recommender System. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval* (Ann Arbor, MI, USA) *(SIGIR '18)*. Association for Computing Machinery, New York, NY, USA, 235–244. doi:10.1145/3209978.3210002

[25] Tavily. 2025. Tavily Search API. https://tavily.com/

[26] Ajay Krishna Vajjala, Dipak Meher, Ziwei Zhu, and David S Rosenblum. 2024. Cross-Domain Recommendation Meets Large Language Models. *arXiv preprint arXiv:2411.19862* (2024).

[27] Lingzhi Wang, Huang Hu, Lei Sha, Can Xu, Kam-Fai Wong, and Daxin Jiang. 2022. RecInDial: A Unified Framework for Conversational Recommendation with Pretrained Language Models. arXiv:2110.07477 [cs.CL] https://arxiv.org/abs/2110.07477

[28] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Multilingual E5 Text Embeddings: A Technical Report. arXiv:2402.05672 [cs.CL] https://arxiv.org/abs/2402.05672

[29] Xiaolei Wang, Kun Zhou, Ji-Rong Wen, and Wayne Xin Zhao. 2022. Towards Unified Conversational Recommender Systems via Knowledge-Enhanced Prompt Learning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Washington DC, USA) *(KDD '22)*. Association for Computing Machinery, New York, NY, USA, 1929–1937. doi:10.1145/3534678.3539382

[30] Yancheng Wang, Ziyan Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Xiaojiang Huang, Yanbin Lu, and Yingzhen Yang. 2024. RecMind: Large Language Model Powered Agent For Recommendation. arXiv:2308.14296 [cs.IR] https://arxiv.org/abs/2308.14296

[31] Zhefan Wang, Yuanqing Yu, Wendi Zheng, Weizhi Ma, and Min Zhang. 2024. MACRec: A Multi-Agent Collaboration Framework for Recommendation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Washington DC, USA) *(SIGIR '24)*. Association for Computing Machinery, New York, NY, USA, 2760–2764. doi:10.1145/3626772.3657669

[32] Dinithi Weerabahu, Agra Gamage, Chathurya Dulakshi, Gamage Upeksha Ganegoda, and Thanuja Sandanayake. 2019. Digital Assistant for Supporting Bank Customer Service. In *Artificial Intelligence*, Jude Hemanth, Thushari Silva, and Asoka Karunananda (Eds.). Springer Singapore, Singapore, 177–186.

[33] Kerui Xu, Jingxuan Yang, Jun Xu, Sheng Gao, Jun Guo, and Ji-Rong Wen. 2021. Adapting User Preference to Online Feedback in Multi-round Conversational Recommendation. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining* (Virtual Event, Israel) *(WSDM '21)*. Association for Computing Machinery, New York, NY, USA, 364–372. doi:10.1145/3437963.3441791

[34] Lanling Xu, Zhen Tian, Bingqian Li, Junjie Zhang, Daoyuan Wang, Hongyu Wang, Jinpeng Wang, Sheng Chen, and Wayne Xin Zhao. 2024. Sequence-level Semantic Representation Fusion for Recommender Systems. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management* (Boise, ID, USA) *(CIKM '24)*. Association for Computing Machinery, New York, NY, USA, 5015–5022. doi:10.1145/3627673.3680037

[35] Chao Zhang, Shiwei Wu, Haoxin Zhang, Tong Xu, Yan Gao, Yao Hu, and Enhong Chen. 2024. NoteLLM: A Retrievable Large Language Model for Note Recommendation. In *Companion Proceedings of the ACM Web Conference 2024* (Singapore, Singapore) *(WWW '24)*. Association for Computing Machinery, New York, NY, USA, 170–179. doi:10.1145/3589335.3648314

[36] Wayne Xin Zhao, Shanlei Mu, Yupeng Hou, Zihan Lin, Yushuo Chen, Xingyu Pan, Kaiyuan Li, Yujie Lu, Hui Wang, Changxin Tian, Yingqian Min, Zhichao Feng, Xinyan Fan, Xu Chen, Pengfei Wang, Wendi Ji, Yaliang Li, Xiaoling Wang, and Ji-Rong Wen. 2021. RecBole: Towards a Unified, Comprehensive and Efficient Framework for Recommendation Algorithms. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management* (Virtual Event, Queensland, Australia) *(CIKM '21)*. Association for Computing Machinery, New York, NY, USA, 4653–4664. doi:zhao2021recbole

[37] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2025. A Survey of Large Language Models. arXiv:2303.18223 [cs.CL] https://arxiv.org/abs/2303.18223

[38] Yuhan Zhao, Rui Chen, Qilong Han, Hongtao Song, and Li Chen. 2024. Unlocking the Hidden Treasures: Enhancing Recommendations with Unlabeled Data. In *Proceedings of the 18th ACM Conference on Recommender Systems* (Bari, Italy) *(RecSys '24)*. Association for Computing Machinery, New York, NY, USA, 247–256. doi:10.1145/3640457.3688149