

# Recommender Systems for Renewable Energy Communities: Tailoring LLM-Powered Recommendations to User Personal Values and Literacy

Bianca Maria Deconcini<sup>1</sup>, Giulia Coucourde<sup>1</sup>, Luca Console<sup>1</sup>, Malek Anouti<sup>2</sup>, Giorgio Gaudio<sup>2</sup> and Michele Visciola<sup>2</sup>

<sup>1</sup> University of Turin, Italy

<sup>2</sup> Experientia SA, Switzerland

## Abstract

The paper proposes a multi-step approach to the design of recommender systems in which the adoption of LLMs is choreographed by a more traditional knowledge-based system exploiting a user model. We focus on renewable energy communities and on the task of engaging participants by leveraging their values, expertise, and available resources to provide personalized descriptions of the benefits they could achieve. This is an important step for the ultimate goals of advanced recommender systems, i.e., facilitating the progression of behaviours towards sustainable agency and adaptivity to climate and environmental challenges.

## Keywords

Recommender Systems, Energy Communities, LLMs

## 1. Introduction

Energy communities face unique challenges in participant engagement and retention. Unlike traditional recommender systems that focus on item recommendations, energy communities require personalized approaches that align with users' values, knowledge levels, and available resources to encourage meaningful participation. Large Language Models (LLMs) are demonstrating enormous potential in many areas, including recommender systems, with capabilities that go beyond traditional methods. The survey by [1] provided a systematic analysis of the different roles LLMs can play in recommender systems, examining the ways they can be trained for this specific task and the various approaches exploited for generating recommendations. In this paper we explore innovative directions that led to the design of ReCommE, a recommender which combines "traditional" approaches to recommendation with the adoption of LLMs under the coordination of a flexible choreography tailored to each individual user. While in most recommender systems the focus is on recommending

---

BCSS 2025: The 13<sup>th</sup> International Workshop on Behavior Change Support Systems, May 5, 2025, Limassol, Cyprus.

✉ biancamaria.deconcini@unito.it (B. M. Deconcini); giulia.coucourde@edu.unito.it (G. Coucourde); lconsole@di.unito.it (L. Console); malek.anouti@experientia.com (M. Anouti); giorgio.gaudio@experientia.com (G. Gaudio); michele.visciola@experientia.com (M. Visciola)

id 0009-0006-1882-9556 (B. M. Deconcini); 0009-0000-9051-386X (G. Coucourde); 0000-0003-2948-5622 (L. Console)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

items/content tailored to specific users' features, we concentrate on leveraging users' values and knowledge to engage them to participate in energy communities. Our goal is to provide a personalized narrative describing the benefit that the user can obtain after joining the community. Our research addresses three key questions:

- How can LLMs be effectively integrated with rule-based systems to create personalized energy community recommendations?
- What role do personal values and literacy levels play in tailoring energy recommendations?
- How can iterative user feedback improve the personalization process?

At an abstract level, our approach involves interleaving between a rule-based classifier, based on knowledge provided by domain experts, and an LLM trained (finely tuned) on the specific domain of application. The classifier initiates by generating the initial profile of a user who is approaching the community; the LLM provides textual argumentation based on the classification. The profile is centered around the user's values and knowledge (besides socio-demographic information and information about the user resources, which are in turn used as context by the LLM to generate the responses). The process is iterative and controlled by a choreography that considers the user's reaction to the text generated by the LLM (the user can agree or disagree with parts of the text) and consequently the classifier is used to refine the user profile, and the LLM to generate more detailed argumentation. This work has been carried on with Experientia, a forward-thinking company involved in a project that has started at the beginning of 2023. It focuses on supporting energy communities - groups of people that aim to optimize energy usage by sharing energy production and consumption [2]. Through this synergy, we aim to explore how the integration of LLMs and behavioral models can reshape recommender systems and better support community-driven scenarios. The paper is organized as follows: Section 2 presents the background and related work, highlighting relevant contributions in LLM-based recommender systems and energy communities. Section 3 introduces the Masterpiece project [3] context. Section 4 details the ReCommE system architecture and its components. Section 5 explains the choreography of our approach. Sections 6 focus on LLM implementation, including prompting, training, and usage. Section 7 discusses our preliminary evaluation, and Section 8 concludes with future directions.

## **2. Background and Related Work**

### **2.1. Energy Communities and User Engagement**

Energy communities face unique challenges in participant recruitment and ongoing engagement. Unlike traditional consumer products, the benefits of joining energy communities can be complex and multifaceted, spanning economic, environmental, and social dimensions. Effective engagement requires understanding users' behavioral drivers, values, and technical literacy to communicate these benefits in personally relevant ways.

## 2.2. LLMs in Recommender Systems

While some of the latest recommender systems techniques involve LLMs, most of them use LLMs to support the work of the system itself and enhance backend operations, such as generating descriptions or extracting information from texts. However, despite these advancements, their use in recommendation systems has not yet been fully explored. The idea for this paper is to move LLM's role from simple content generators for the system to tools which can provide context-aware responses. Our goal is to understand if, and how, LLMs can be valuable tools for personalization purposes, able to increase user engagement and satisfaction. One of the first applications of LLMs in RS mainly focused on content generation. The work of [4] demonstrates the use of LLMs to generate item descriptions, showing that these systems give results that are similar to those obtained by web-scraping techniques. Similarly, [5] explore how LLMs can produce explanations that users find comparable or even better than the baseline ones. Indeed, users in this study perceived LLM-generated recommendations as more effective and efficient, thus helping them to decide faster. This is due also to LLMs detailing ability. Several studies explored the integration of LLMs for enhancing personalization and user experience. [6] introduce LLM-Rec, a framework which leverages the power of LLMs specifically for personalisation goals. By employing prompting techniques, this approach aims to generate better quality recommendations, without the need for extensive domain-specific training or data. [7] emphasize the emerging role of LLMs in reshaping traditional recommendation methodologies. The study shows that contextualising recommendations can significantly improve their relevance and effectiveness. However, it also suggests a gap in the incorporation of more sophisticated real-time methods, highlighting the need for systems that continuously adapt to changing user preferences with the help of contextual information. Another interesting vision is offered by [8] and their evaluation metric of behavior alignment. Traditionally, the evaluation of recommendation systems focuses on metrics such as accuracy or novelty. However, these approaches do not always capture how recommendation systems adapt to user behaviors and preferences.

## 2.3. User Modeling for Personalization

Effective personalization relies on robust user/behavioral modeling. Traditional user behavioral models often focus on demographic information and behavioral data related to the digital footprint of the users, but these may be insufficient for complex domains like energy communities. Our approach expanded to include intangible aspects such as personal values [X], domain literacy [Y], and resource availability [Z]. These elements are particularly relevant for energy community engagement, where decisions are influenced by a combination of practical constraints, knowledge levels, and personal value systems.

## 2.4. Research Gap

Despite advances in both LLM-based recommender systems and energy community research, several gaps remain:

- Most LLM applications in recommender systems focus on item recommendations rather than complex service engagement.
- Few systems incorporate iterative user feedback to refine recommendations.

- The integration of user values and domain literacy levels in recommendation generation remains underexplored.
- There is limited research on how to effectively combine rule-based systems with generative LLMs.

Our work addresses these gaps through a choreographed approach that leverages both traditional knowledge-based systems and fine-tuned LLMs.

### **3. The Masterpiece Project Context**

Masterpiece is a project focusing on digital tools for supporting energy communities. The first part of the project focused on studying the domain, exploiting a number of case studies in different countries. In this way, a detailed picture of the instruments, rules and roles of different types of stakeholders and participants has been created. User studies supported the construction of the archetypes of users, characterized along multiple dimensions such as individual and societal values, levels of expertise, available resources (including household, types of appliances, etc.). This behavioral model serves now as the foundation for this work. The user studies identified several key archetypes within energy communities, including:

- Sustainability Champions: Primarily motivated by environmental values.
- Financial Optimizers: Focused on energy savings and financial benefits.
- Community Builders: Driven by social cohesion and local development.
- Skeptical consumers: Interested in convenience supply.

Each archetype is characterized by a unique combination of values, knowledge levels, and available resources, which inform our personalization approach.

### **4. ReCommE System: Architecture and Components**

The ReCommE system integrates multiple components to deliver personalized recommendations for energy community participation.

#### **4.1. System Overview**

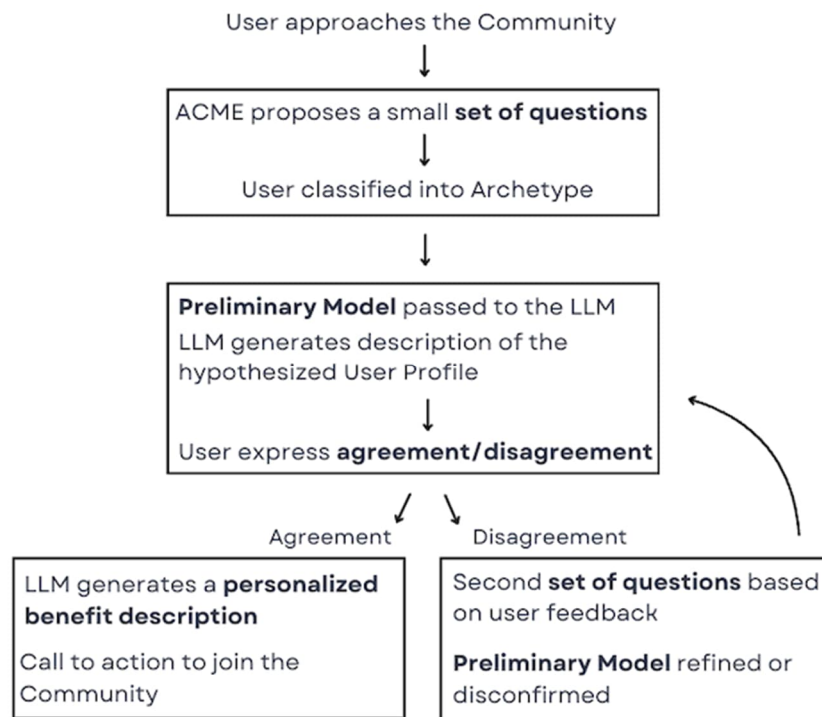
ReCommE combines a rule-based classifier with a fine-tuned LLM to generate personalized recommendations. The system follows a choreographed workflow that iteratively refines the user model based on feedback, allowing for increasingly tailored recommendations.

#### **4.2. User/Behavioral Model**

The core of ReCommE system architecture is a comprehensive user/behavioral model that captures three key dimensions:

- Values: Individual priorities and behavioral drivers that drive decision-making (e.g., environmental sustainability, economic benefits, community well-being, convenience services).
- Literacy/Expertise: Technical knowledge and familiarity with energy concepts.

- Resources: Available assets and infrastructure (e.g., home ownership, renewable installations, smart devices).



**Figure 1:** ReCommE System Architecture showing the interaction between the rule-based classifier, LLM component, and user feedback loop.

### 4.3. Rule-Based Classifier

The classifier component employs domain expertise to categorize users into meaningful archetypes based on their responses to targeted questions. Initially, the user's profile is randomly selected from possible profiles. Then, the user is prompted with questions for each category (values, literacy, and resources). If the user's responses do not align with the initially selected profile, the questions are dynamically adjusted. This process ensures that the user's profile evolves in response to their answers, leading to a more accurate and personalized classification. This classification provides the foundation for the LLM-generated content, ensuring domain-specific relevance.

### 4.4. LLM Component

The LLM component is responsible for generating natural language descriptions and recommendations based on the user profile. Through fine-tuning and structured prompting, the LLM produces personalized narratives that highlight the benefits of community participation most relevant to the specific user.

## 5. The Choreography

In this section we sketch the workflow of the choreography that interleaves traditional rule-based systems with LLMs (see Figure 1).

- **Initial Profiling:** Whenever a new participant approaches the community, ReCommE is activated with the initial goal of profiling the user. A small set of questions is presented to the user and a first classification in one of the archetypes is done. In this first attempt, a user model is built with a preliminary rough estimation of the user's values, knowledge (expertise) and resources.
- **Profile Validation:** This preliminary model is passed to the LLM which is asked to generate a first abstract description of the hypothesized user profile. The description is presented to the user, who is then asked to highlight specific points where s/he agrees or disagrees.
- **Model Refinement:** The control is passed back to the rule systems which activates a second group of questions depending on the user's feedback. In this way the preliminary model is either disconfirmed or refined. In the former case the process is restarted, in the latter the LLM is activated for a second time.
- **Personalized Recommendation:** In this second case, the LLM is asked to provide a first description of the benefits the user can obtain if she/he joins the community. The description is personalized given the user model, specifically along the dimensions of values and available resources and taking also into account the expertise to tailor the level of details and technicality. We also personalize the tone of the explanation, given the user model.

### 5.1. Example Walkthrough

To illustrate the choreography, consider a user initially classified as a "Sustainability Champion":

- The system presents an initial profile: "You seem to prioritize environmental sustainability and have moderate knowledge about renewable energy..."
- The user confirms their environmental values but disagrees with the knowledge assessment.
- The system refines the model, adjusting the expertise level downward.
- The LLM generates a new recommendation emphasizing environmental benefits with less technical terminology: "By joining this energy community, you'll help reduce carbon emissions by approximately X tons per year..."

In this version of ReCommE, the choreography is one step. In the future we plan to add further steps that interleave a progressive refinement of the user model by the classifier and more focused and detailed argumentations generated by the LLM.

## **6. LLM Implementation**

### **6.1. LLM Prompting**

The textual input given to an LLM with the purpose of directing the output's generation is known as prompt, and prompt engineering is the process of designing and constructing input to elicit desired responses. Advancements in the field of generative AI are remarkable, as demonstrated by the increasing complexity of models, improved training techniques and the expansion of application possibilities. These reasons, as pointed out by [9], underscore the critical role of prompt engineering in maximizing the precision and usefulness of these models, ensuring they can successfully meet a range of changing user requirements. This section describes our strategy to take advantage of the LLMs generation capability by employing prompt engineering techniques to lead the model to produce a basic description of the user profile and the benefits of joining the community. According to the workflow, the goal is to get the abstract description by providing the model details about the user's value, resources and knowledge (expertise). In order to get the desired outcome, a role-playing prompting approach - a method of influencing the language model's behavior by giving it a specific role within the interaction - was employed, as suggested by [10]. The purpose of the experimental framework was to assess the performance of the models, Llama-3-8B-Instruct [11] and Zephyr-7b-beta [12], under different conditions. The experiments were performed on both pre-trained and fine-tuned versions of the models. We used a role-prompting strategy for each configuration, employing both the zero-shot and few-shot paradigms. In the zero-shot setting, the models were prompted without any task-specific examples and context, only the role was defined. On the other hand, in the few-shot setting, the models received a limited number of examples, to simulate a more informed scenario, as illustrated by [13]. Adding context in both cases, background knowledge helps the model match its responses with the user's intent, which enhances the quality of the output produced in both situations. Nevertheless, even with this advancement, the generated content frequently falls short of the required level of domain specificity needed for reliable and accurate outcomes, being too broad in regard to environmental sustainability and energy field. This limitation indicates a weakness in the model's ability to fully use context to tailor responses to a given profile type and to adjust its responses in response to particular knowledge domains. However, adding context is a good place to start when trying to improve generation, providing a basis for additional refinement and domain adaptation to produce more accurate and contextually relevant results.

### **6.2. LLM Training**

Our system training is based on the use of project documents that describe users' archetypes along three main dimensions: values, literacy and resources. These documents, based on user studies conducted by Experientia, are used to train the model, so that it can know the user profiles' characteristics and the specific context. This allows the model to generate responses driven by domain-relevant user studies, rather than relying solely on its general knowledge as a preformed model. The initial objective was to obtain the generation of comprehensive descriptions for each archetype that included the core values and domain-specific information related to that archetype. This approach aims to offer comprehensive depictions that contextualize the archetypes' roles, traits, and applications within the energy domain in

addition to defining them. To achieve the goal, the training data was manually converted into a chat template format in order to tailor it to the model’s particular needs. The preprocessing involved manually converting the data into 1,000 instruction-response pairs where each pair represents a conversational exchange between a user and an assistant, carefully extracted from official documents. This step was crucial to enhance data relevance and with the help of this approach we trained the model to generate responses that are highly adapted to the target domain. For the preliminary testing phase, we chose to finetune Llama-3-8B-Instruct. The model was fine-tuned applying a selective Parameter-Efficient Fine-Tuning (PEFT) as explained by [14] with Low-Rank Adaptation (LoRA) [15]. Two NVIDIA A40 GPUs, selected for their high amount of memory and efficient parallel processing capabilities, made up the computing environment. Multiple runs were performed, varying the hyperparameter configurations listed in Table 1. The target modules selected for LoRA fine-tuning include key self-attention components (query, key, value, and output projections) and MLP layers (gate, up, and down projections), ensuring effective adaptation of both attention and feed-forward transformations. Table 2 compares the best and worst runs, selected based on eval loss. The chosen loss function is the cross-entropy loss, which is standard for causal language modeling. By minimizing cross-entropy loss, the model learns to assign higher probabilities to correct next-token predictions, improving its ability to generate coherent and contextually appropriate sequences. The best run (run\_7) outperforms the worst (run\_18), likely due to its higher learning rate ( $2.5e - 5$  vs.  $1.5e - 5$ ) and constant scheduler. This setup yields lower train (1.96 vs. 3.02) and eval loss (1.71 vs. 2.56), suggesting more effective weight updates. Additionally, LoRA r and LoRA alpha are higher in the best run (16 and 32 vs. 8 and 16), following the consistent setting where alpha is always set to twice the value of r, which may contribute to better adaptation of the model parameters.

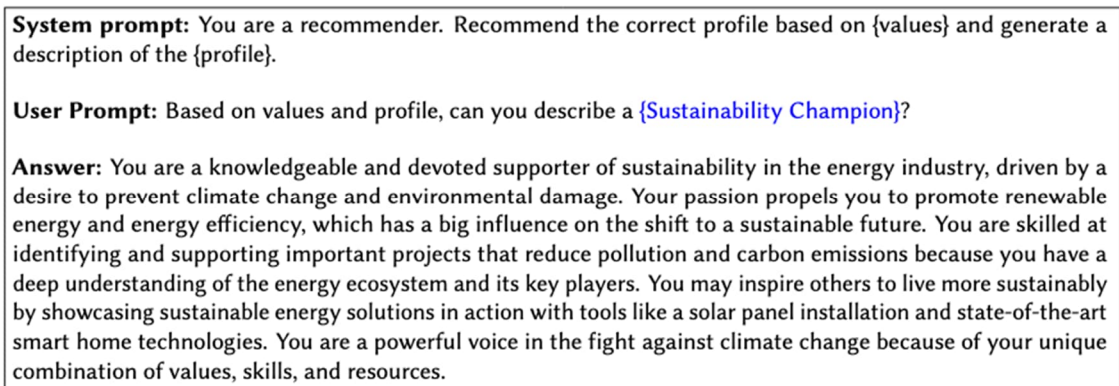
**Table 1:** Parameter Grid of PEFT settings.

Parameter	Values
Learning Rate	5e-5, 2.5e-5, 1e-5
Learning Rate Scheduler	Linear, Constant
Number of Training Epochs	3, 4, 5
LoRA r	8, 16, 32
LoRA alpha	16, 32, 64

**Table 2:** Comparison of hyperparameters and performance between the best and worst runs of LLaMA3 8B-Instr in PEFT setting.

Hyperparameter	Best Run (run_7)	Worst Run (run_18)
Learning Rate	2.5e-5	1.5e-5
Scheduler	Constant	Linear
Number of Epochs	4	3
<b>Train Loss</b>	<b>1.956402</b>	3.021542
<b>Eval Loss</b>	<b>1.712478</b>	2.560054
<b>Perplexity</b>	<b>5.542681</b>	12.936517





**Figure 2:** Example of prompt structure for the "Sustainability Champion" archetype, showing system prompt and user prompt components.

### 6.3. LLM Usage

In order to create different generations, the same prompt structure was used, modifying only the content of the variables {values} and {profile}, which depends on the classifier's output. The profiles and values obtained from the rule-based system's initial estimation are given to the system prompt in this initial stage without additional explanation because they were integrated through fine-tuning. The user prompt incorporates the resulting profile and is not intended as a prompt directly provided by a user but as a starting point for the generation process. An example for the archetype "Sustainability Champion" is shown in Figure 2. This process enabled the generation of multiple descriptions that could be compared while keeping the structure of system prompt and user prompt fixed. We came to the conclusion from the tests that a minimal structured system prompt that defines the role and directs the model in the generation was successful.

## 7. Preliminary Evaluation

The evaluation, at present, is preliminary and conducted by the Experientia team that carried out the user studies. Different generations were compared, and the best were selected based on their coherence with the goal, the prompt and the domain. This process made it possible to determine which model was the most suitable and according to expert evaluations, Llama-3-8B-Instruct in its fine-tuned version was the one that generated the most relevant outputs.

### 7.1. Evaluation Criteria

Our evaluation focused on several key dimensions:

- **Relevance:** How well the generated content aligned with the user's archetype.
- **Accuracy:** Correctness of domain-specific information.
- **Personalization:** Adaptation to the user's values, literacy, and resources.
- **Readability:** Clarity and accessibility of the generated text.
- **Persuasiveness:** Potential effectiveness for encouraging community participation.

## 7.2. Preliminary Results

Initial findings suggest that:

- Fine-tuned models consistently outperformed pre-trained models in domain relevance.
- Role-prompting significantly improved the personalization of outputs.
- “Few-shot” prompts produced more accurate technical information than “zero-shot” approaches.
- User values were more effectively incorporated than literacy levels in the generated outputs.

In this phase, the evaluation is intended to collect initial feedback to understand the potential of the system or areas for improvement, but it does not represent a final objective assessment of the system’s performance - that indeed we are planning to conduct in the near future.

## 8. Conclusions and Future Work

In this work, we explored the integration of large language models (LLMs) into recommendation systems, with a particular focus on dynamic personalization and user-system interaction. Our key contributions include:

- A novel choreographed approach combining rule-based systems and LLMs.
- A multi-dimensional user model for energy community engagement.
- An iterative feedback mechanism that refines recommendations based on user input.
- A domain-specific implementation for renewable energy communities.

Several promising directions for future work have emerged: Since a key aspect of our proposal is the use of iterative interaction, where users are asked to express their agreement or disagreement on the description provided by the system, an open question remains whether binary responses are sufficient or whether it is more useful to allow richer feedback, for a more detailed refinement. Moreover, an interesting future development regards how past conversations and interactions could be reintegrated into the LLM tuning process to update the system knowledge. This approach would allow the system to adapt and enhance its ability to provide personalized recommendations over time. At present, the interaction between the LLM and the recommender system takes place in two main phases, but this exchange may become more frequent and dynamic, involving the LLM more in the recommendation process. Lastly, as mentioned before, our plan is to first refine the system and then conduct a quantitative analysis based on a large-scale evaluation with actual energy community participants to measure engagement effectiveness and behavioral change.

## References

- [1] Q. Wang, J. Li, S. Wang, Q. Xing, R. Niu, H. Kong, R. Li, G. Long, Y. Chang, C. Zhang, Towards next-generation llm-based recommender systems: A survey and beyond, 2024. URL: <https://arxiv.org/abs/2410.19744>. arXiv:2410.19744.
- [2] J. Lowitzsch, C. Hoicka, F. van Tulder, Renewable energy communities under the 2019 european clean energy package – governance model for the energy clusters of the future?, Renewable and Sustainable Energy Reviews 122 (2020) 109489. URL: <https://www.sciencedirect.com/science/article/pii/S1364032119306975>. doi:<https://doi.org/10.1016/j.rser.2019.109489>.
- [3] Multidisciplinary Approaches and Software Technologies for Engagement, Recruitment and Participation in Innovative Energy Communities in Europe, Technical Report Grant agreement no 101096836, 2024. URL: <https://masterpiece-horizon.eu/>.
- [4] A. Acharya, B. Singh, N. Onoe, Llm based generation of item-description for recommendation system, in: Proceedings of the 17th ACM Conference on Recommender Systems, RecSys '23, ACM, New York, NY, USA, 2023, p. 1204–1207. doi:10.1145/3604915.3610647
- [5] S. Lubos, T. N. T. Tran, A. Felfernig, S. Polat Erdeniz, V.-M. Le, Llm-generated explanations for recommender systems, in: Adjunct Proceedings of the 32nd ACM Conference on User Modeling, Adaptation and Personalization, UMAP Adjunct '24, ACM, New York, NY, USA, 2024, p. 276–285. doi:10.1145/3631700.3665185.
- [6] H. Lyu, S. Jiang, H. Zeng, Y. Xia, Q. Wang, S. Zhang, R. Chen, C. Leung, J. Tang, J. Luo, LLM-rec: Personalized recommendation via prompting large language models, in: K. Duh, H. Gomez, S. Bethard (Eds.), Findings of the Association for Computational Linguistics: NAACL 2024, Association for Computational Linguistics, Mexico City, Mexico, 2024, pp. 583–612. URL: <https://aclanthology.org/2024.findings-naacl.39>. doi:10.18653/v1/2024.findings-naacl.39.
- [7] Z. Zhao, W. Fan, J. Li, Y. Liu, X. Mei, Y. Wang, Z. Wen, F. Wang, X. Zhao, J. Tang, Q. Li, Recommender Systems in the Era of Large Language Models (LLMs), IEEE Transactions on Knowledge & Data Engineering 36 (2024) 6889–6907. doi:10.1109/TKDE.2024.3392335
- [8] D. Yang, F. Chen, H. Fang, Behavior alignment: A new perspective of evaluating llm-based conversational recommendation systems, in: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, New York, NY, USA, 2024, p. 2286–2290. doi:10.1145/3626772.3657924.
- [9] B. Chen, Z. Zhang, N. Langren'e, S. Zhu, Unleashing the potential of prompt engineering in large language models: a comprehensive review, ArXiv abs/2310.14735(2023).URL: <https://api.semanticscholar.org/CorpusID:264426395>
- [10] N. Wu, M. Gong, L. Shou, S. Liang, D. Jiang, Large language models are diverse role-players for summarization evaluation, in: Natural Language Processing and Chinese Computing, 2023. URL: <https://api.semanticscholar.org/CorpusID:257767249>.
- [11] AI@Meta, Llama 3 model card, [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md), 2024.
- [12] L. Tunstall, E. Beeching, N. Lambert, N. Rajani, K. Rasul, Y. Belkada, S. Huang, L. von Werra, C. Fourrier, N. Habib, N. Sarrazin, O. Sanseviero, A. M. Rush, T. Wolf, Zephyr: Direct distillation of lm alignment, 2023. arXiv:2310.16944.

- [13] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. teusz Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language models are few-shot learners, ArXiv abs/2005.14165 (2020). URL: <https://api.semanticscholar.org/CorpusID:218971783>
- [14] Z. Han, C. Gao, J. Liu, J. Zhang, S. Q. Zhang, Parameter-efficient fine-tuning for large models: A comprehensive survey, ArXiv abs/2403.14608 (2024). URL: <https://api.semanticscholar.org/CorpusID:268553763>.
- [15] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, Lora: Low-rank adaptation of large language models, 2021. URL: <https://arxiv.org/abs/2106.09685>. arXiv:2106.09685