# Research on the Application of Large Language Models in Recommendation Systems

**Yan Yang***, Rong Li

*Computer School, Central China Normal University, Wuhan, China*

*\*Corresponding Author*

**Abstract: In the context of information explosion and diversified user demands, large language models (LLMs) provide new ideas for optimizing recommendation systems. This paper studies the application of LLMs in recommendation systems. Firstly, the paper outlines the basic framework of LLM-based recommendation systems, and then focuses on analyzing three advanced recommendation methods based on fine-tuning, prompt learning, and instruction tuning. The paper further elaborates on the advantages of LLMs compared to traditional recommendation systems in terms of accuracy, cold start, diversity, and other aspects, and points out their limitations in data processing, computational efficiency, and privacy security. Finally, the paper provides an outlook on future research directions.**

**Keywords: Recommendation System; Large Language Models; Instruction Tuning**

## 1. Introduction

In the era of big data, the rapid growth of diverse types, modalities, and sources of data has made it increasingly difficult for people to effectively and accurately locate the information they need. Consequently, the demand for systems that can quickly and accurately retrieve target information has become more urgent than ever. Against this backdrop, the roles of information retrieval and recommendation systems have become increasingly prominent. The former focuses on selecting the most relevant content from vast amounts of data, while the latter aims to provide personalized suggestions to help users find desired products or services [1]. Currently, recommendation systems have become a critical component across various domains—from product promotion and social media content filtering to precision recommendations in the entertainment industry. By deeply analyzing user behavior, preferences, and historical records, these systems present the most relevant choices, thereby enhancing information matching accuracy and user experience [2].

Traditional recommendation systems primarily encompass collaborative filtering-based methods, content-based methods, and knowledge graph-based methods, focusing on delivering personalized recommendations through user behavior and item attributes [3]. Although traditional recommendation systems have advanced scientific research across numerous fields, they still face a series of challenges, such as poor interactivity and insufficient explainability. These issues not only impact user experience but also relate to core concerns like system transparency and trustworthiness, thereby limiting the widespread deployment of these systems in practical applications [4].
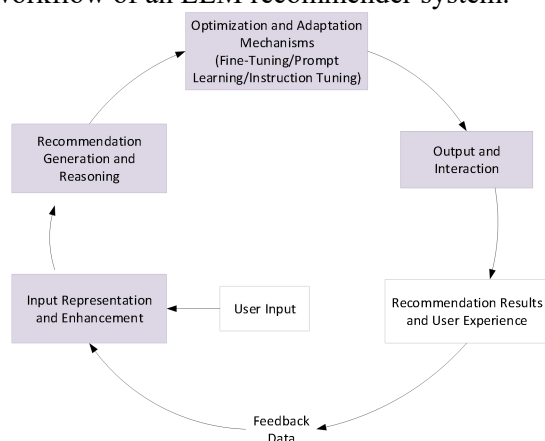
The 2022 launch of ChatGPT marked the rise of large language models (LLMs). This emergence represents a groundbreaking breakthrough in Natural Language Processing (NLP), showcasing remarkable advancements in neural network technology for NLP tasks [5]. In recent years, the exceptional performance of pre-trained language models like BERT and GPT in NLP has also provided new insights for the development of recommendation systems [6].

Against the backdrop of evolving NLP models, the recommendation system domain has undergone continuous transformation. Relevant NLP models, leveraging their exceptional language comprehension and generation capabilities, have gradually become the core of recommendation systems [7]. A key reason for the applicability of large language models in recommendation systems lies in their outstanding ability to understand complex user intent and preferences. They can analyze user search histories, purchase records, review

content, and other aspects to capture subtle shifts in user interests, thereby generating personalized recommendations. For instance, large language models can extract emotional tendencies, specific preferences, needs, and even implicit expectations from user reviews, aspects that traditional recommendation systems struggle to capture accurately [8]. Simultaneously, these models account for linguistic diversity and cultural differences during the recommendation process, delivering more precise recommendations to users across regions and language backgrounds. This significantly enhances the relevance and user satisfaction of recommendation systems, offering richer and more personalized experiences.

## 2. Framework of Large Language Model Recommendation System

LLM recommender systems integrate natural language understanding and generation capabilities to build a personalized recommendation framework centered on user semantic intent analysis, driven by collaborative data and model synergy. This framework typically comprises four core modules: input representation and enhancement, recommendation generation and inference, optimization and adaptation mechanisms, and output and interaction. Together, they form a complete closed-loop process from user demand comprehension to recommendation output generation. Figure 1 illustrates the general workflow of an LLM recommender system.



**Figure 1 A Framework of LLM Recommendation System**

Input Representation and Enhancement Module: This module standardizes and semantically fuses multi-source heterogeneous user information—such as query text, historical behavior sequences, item descriptions, and contextual data. Its core tasks involve text cleansing, key information extraction, and constructing structured task prompts. This transforms raw data into text sequences rich with user intent and recommendation instructions, providing high-quality, interpretable inputs for the large language model.

Recommendation Generation and Reasoning Module: As the system's core, this module leverages the deep semantic understanding and generative capabilities of large language models to parse and reason over input prompts. By identifying users' explicit needs and implicit preferences, combined with internal knowledge or external information, it performs candidate item generation, ranking, or comparative analysis, thereby converting semantic understanding into personalized recommendation decisions.

Optimization and Adaptation Mechanisms Module: This module aims to precisely tailor general-purpose large language models to specific recommendation scenarios and objectives. It primarily integrates three strategies: fine-tuning, prompt learning, and instruction tuning. Through parameter updates, prompt engineering, or instruction-based supervised training, it guides the model to capture domain-specific patterns and user preferences, continuously enhancing recommendation accuracy, diversity, and adaptability to new tasks and cold-start scenarios.

Output and Interaction Module: This module transforms the model's internal inference results into final user-perceivable outputs, typically presenting recommended items and their explanatory rationale in natural language. Simultaneously, it supports multi-turn dialogue interactions, dynamically adjusting subsequent recommendation directions based on real-time user feedback. This forms a closed-loop interaction cycle of "recommendation-feedback-optimization," thereby enhancing the system's interpretability, interactivity, and user experience.

## 3. Advanced Application of Large Language Models in Recommendation Systems

The following section delves into the advanced application of large language models in recommendation systems, with a particular

focus on their diverse adaptation strategies, which primarily include fine-tuning, prompt learning, and instruction tuning.

### 3.1 Fine-Tuning-Based Recommendation Methods

Fine-tuning pre-trained language models is widely adopted in natural language processing, extending to recommendation systems and beyond. The core concept involves training a model initially trained on large-scale text data to adapt its parameters for specific tasks or domains. For recommendation tasks, models can be fine-tuned using recommendation datasets—such as user interaction histories and item descriptions—to better capture user-item relationships, enabling more personalized and precise recommendations.

For sequential recommendation tasks, research has designed a deep bidirectional self-attention model based on the BERT architecture. This model mitigates information leakage by introducing a cloze task during BERT fine-tuning, enabling context-conditional prediction of randomly masked items in user behavior sequences. This enhances sequence modeling and recommendation effectiveness. In e-commerce scenarios, another approach proposes an end-to-end multi-task learning framework. It enhances domain semantic understanding through domain-adaptive fine-tuning of BERT and synergistically optimizes multiple objectives via multi-task training, thereby improving recommendation performance. However, most such methods still rely on traditional Top-K recommendation strategies, limiting flexibility in generating results. The introduction of generative language models offers novel solutions to this challenge. For instance, one model based on the GPT-2 architecture employs singular value decomposition (SVD) to perform tokenization on item IDs, decomposing them into multiple sub-identifiers to address the challenge of an item ID space vastly exceeding the vocabulary size typical for language models. Concurrently, this model adopts a Next-K generative recommendation strategy, progressively generating interdependent recommendation lists during inference. This enables more flexible and coherent generation of complex recommendation results.

In summary, fine-tuning pre-trained language models and integrating them into recommendation systems not only leverages powerful external knowledge and personalized user preferences to enhance recommendation accuracy but also enables cold-start capabilities for new items with limited historical data. This enriches the functionality and efficacy of recommendation systems, delivering diverse and personalized recommendation experiences to users while opening greater prospects and opportunities for the recommendation field.

### 3.2 Prompt Learning-Based Recommendation Methods

Large language models have also demonstrated exceptional capabilities in zero-shot or few-shot recommendation tasks. Some researchers have attempted to introduce prompt learning templates and their model generalization capabilities (leveraging model prior knowledge to enhance recommendation performance on specific tasks with minimal or no additional data) into recommendation systems to optimize model recommendation outcomes. Prompt learning, an NLP strategy recognized as an efficient personalized recommendation approach, enables the incorporation of additional prompt information without significantly altering the structure or parameters of pre-trained language models. This transforms downstream tasks into text generation tasks, significantly optimizing model performance for specific objectives and ensuring outputs align closely with task goals.

In prompt-based recommendation approaches, prompts typically take forms such as task descriptions or example phrases. Integrating these into the input guides the model to generate expected outputs. For instance, one study proposed a real-time news recommendation framework that reframes the user click prediction task for candidate news items as a cloze-style masked word prediction problem. This approach designs a series of prompt templates—including discrete, continuous, and hybrid formats—to construct corresponding answer spaces, integrating predictions from multiple templates to enhance recommendation performance. Another study attempts to transform user profiles and historical interaction behaviors into prompt information, proposing a novel conversational recommendation paradigm based on large language models. This framework leverages conversational context to effectively learn user preferences and establish associations between users and products,

thereby enhancing the interactivity and interpretability of the recommendation process. Additionally, this approach supports cross-domain transfer of user preferences, making it applicable to multi-domain personalized recommendation scenarios. It also addresses the cold-start problem for new items through prompt information processing. Overall, prompt-learning-based large language model recommendation methods offer novel insights for recommendation system development while expanding the application scenarios of AI-generated content (AIGC) in the recommendation domain.

## 3.3 Instruction Tuning-Based Recommendation Methods

Instruction tuning is a flexible, versatile fine-tuning approach categorized as supervised fine-tuning. Its core advantage lies in its adaptability to task descriptions without reliance on specific labeled data, making it suitable for various natural language processing tasks under both supervised and unsupervised conditions. This approach enhances models' ability to capture and execute human intent, demonstrating exceptional performance in zero-shot scenarios. By integrating user task instructions with large language models, it further optimizes recommendation personalization and efficiency. When receiving natural language descriptions or task prompts from users, the model accurately parses user needs and adjusts recommendation content based on instructions to meet user expectations.

Research on instruction-tuned recommendation systems has yielded diverse innovative frameworks and methods. Representative work includes a universal paradigm unifying multiple recommendation tasks within a generative language architecture. This framework consolidates various recommendation tasks into a shared architecture based on conditional language generation, enhancing recommendation accuracy and diversity. Its core approach involves uniformly converting diverse data—such as user-item interactions, user descriptions, item attributes, and reviews—into natural language sequences for representation. This enables deeper modeling of semantic relationships between users and items, achieving more personalized recommendations. This universal recommendation framework acquires general-purpose language representations during pre-training, providing a shared foundational model for different recommendation tasks and facilitating knowledge transfer and reuse across tasks. Additionally, the framework supports zero-shot and few-shot learning capabilities. It dynamically generates personalized prompts based on user needs, enabling effective recommendations even when encountering new tasks or cold-start scenarios, thereby reducing reliance on fine-tuning.

On the other hand, some approaches focus on user-demand-driven instruction optimization strategies. By directly modifying open-source language models rather than invoking public interfaces, these methods enable flexible understanding and response to personalized instructions. This not only effectively mitigates data sparsity issues but also propels recommendation systems toward greater adaptability and interpretability. Collectively, these studies exemplify the paradigm shift in large language models within the recommendation domain—from foundational feature modeling toward generative, interactive intelligent services.

## 4. Analysis of the Characteristics of Llm-Based Recommendation Systems

### 4.1 Advantages of LLM-Based Recommendation Systems

1) Recommendation Accuracy

Traditional recommendation systems are directly constrained by data quality and feature engineering. Efficient feature engineering typically requires deep domain expertise, making data accuracy and completeness critical for generating suitable recommendations. In contrast, LLM-based recommendation systems leverage their profound understanding of natural language instructions to achieve superior accuracy, capturing user needs more precisely.

2) Handling Massive Data

Traditional recommendation systems rely heavily on robust engineering and resources, such as distributed computing and database optimization. Excessive data volumes degrade their performance and increase response times. In contrast, systems based on large language models excel at processing massive amounts of text and user-generated content. They can deeply analyze and understand user behavior, providing more precise suggestions from vast text repositories without constraints imposed by

data scale.

3) Cold Start Problem

Traditional recommendation systems often struggle with cold start problems due to their heavy reliance on existing data, making it difficult to provide effective recommendations for new users or new projects. Systems based on large language models mitigate this issue by parsing natural language user instructions or employing implicit feature learning. They can even deliver targeted suggestions in scenarios with limited sample data or insufficient historical information.

4) Recommendation Diversity

Traditional recommendation systems often exhibit limitations in diversity. Methods based on collaborative filtering can create "filter bubbles," leading to recommendations biased toward users' historical preferences. Content-based approaches primarily rely on item features, potentially resulting in monotonous recommendations. In contrast, methods leveraging large language models comprehensively analyze both user historical behavior and associated item information. Combined with natural language processing capabilities, they present users with a broader range of content, thereby enhancing the diversity of model recommendations.

5) User Behavior Analysis

Traditional recommendation systems typically rely on historical user behavior data—such as clicks, purchases, and ratings—to uncover user interests. In contrast, large language model-based systems leverage their exceptional contextual understanding to interpret natural language instructions beyond historical behavior, delivering deeper and richer recommendations.

## 4.2 Limitations of LLM-Based Recommendation Systems

1) Data Processing

Traditional recommendation systems primarily rely on structured data like user behavior and item attributes, which are relatively simple to process and manage. In contrast, systems based on large language models must handle not only structured data but also unstructured natural language text. This involves additional data processing steps such as text cleaning, tokenization, and vectorization, thereby imposing higher demands on data collection, processing, and storage.

2) Model Training and Deployment

Traditional recommendation systems feature streamlined training that can be completed on standard hardware. For instance, collaborative filtering's user-item interaction matrix computations can be deployed using conventional databases and servers. Systems based on large language models, however, have stringent training requirements, typically necessitating multiple GPUs or TPUs. To meet inference speed demands, cloud computing or specialized hardware may be required, resulting in higher costs and longer processing times.

3) Privacy and Security Concerns

Traditional recommendation systems typically handle smaller volumes of user data, often stored locally or anonymized, resulting in lower privacy and security risks. However, systems based on large language models, which require vast amounts of user data, may pose significant privacy and security risks.

## 5. Conclusion

This paper systematically investigates the application of large language models in recommendation systems. Research demonstrates that through strategies such as fine-tuning, prompt learning, and instruction tuning, large language models exhibit significant advantages in recommendation accuracy, diversity, and cold-start handling, while also achieving deep understanding of user intent. However, limitations remain in computational cost, data processing, and privacy security. Future research should focus on developing efficient lightweight technologies, constructing high-quality datasets, and enhancing model interpretability and security to advance the mature application of next-generation recommendation systems empowered by large language models.

## References

[1] Wu Z., Tang Y., Liu H. Survey of Personalized Learning Recommendation [J]. Journal of Frontiers of Computer Science and Technology, 2022, 16(1): 21-40.

[2] Brown T., Mann B., Ryder N., et al. Language Models are Few-shot Learners [J]. Advances in Neural Information Processing Systems, 2020, 159: 1877-1901.

[3] Stiennon N., Ouyang L., Wu J., et al. Learning to Summarize with Human Feedback [J]. Advances in Neural Information Processing Systems, 2020, 33:

3008-3021.

[4] Ding N., Qin Y., Yang G., et al. Parameter-efficient Fine-tuning of Large-scale Pre-trained Language Models [J]. Nature Machine Intelligence, 2023, 5(3): 220-235.

[5] Wu X., Magnani A., Chaidaroon S., et al. A Multi-Task Learning Framework for Product Ranking with BERT [C]. Proceedings of WWW2022, 2022: 493-501.

[6] Chen X., Zhang N., Xie X., et al. KnowPrompt: Knowledge-aware Prompt-tuning with Synergistic Optimization for Relation Extraction [C]. Proceedings of WWW2022, 2022:

2778-2788.

[7] Geng S., Liu S., Fu Z., et al. Recommendation as Language Processing (RLP): A Unified Pretrain, Personalized Prompt & Predict Paradigm (P5) [C]. Proceedings of the 16th ACM Conference on Recommender Systems (RecSys 2022), 2022: 299-315.

[8] Xiang W., Wang Z., Dai L., et al. ConnPrompt: Connective-cloze Prompt Learning for Implicit Discourse Relation Recognition [C]. Proceedings of the 29th International Conference on Computational Linguistics (COLING 2022), 2022: 902-911.