

GenRec: Large Language Model for Generative Recommendation

Jianchao Ji
Rutgers University
New Brunswick, NJ, US
jianchao.ji@rutgers.edu

Zelong Li
Rutgers University
New Brunswick, NJ, US
zelong.li@rutgers.edu

Shuyuan Xu
Rutgers University
New Brunswick, NJ, US
shuyuan.xu@rutgers.edu

Wenyue Hua
Rutgers University
New Brunswick, NJ, US
wenyue.hua@rutgers.edu

Yingqiang Ge
Rutgers University
New Brunswick, NJ, US
yingqiang.ge@rutgers.edu

Juntao Tan
Rutgers University
New Brunswick, NJ, US
juntao.tan@rutgers.edu

Yongfeng Zhang
Rutgers University
New Brunswick, NJ, US
yongfeng.zhang@rutgers.edu

ABSTRACT

In recent years, large language models (LLM) have emerged as powerful tools for diverse natural language processing tasks. However, their potential for recommender systems under the generative recommendation paradigm remains relatively unexplored. This paper presents an innovative approach to recommendation systems using large language models (LLMs) based on text data. In this paper, we present a novel LLM for generative recommendation (GenRec) that utilized the expressive power of LLM to directly generate the target item to recommend, rather than calculating ranking score for each candidate item one by one as in traditional discriminative recommendation. GenRec uses LLM’s understanding ability to interpret context, learn user preferences, and generate relevant recommendation. Our proposed approach leverages the vast knowledge encoded in large language models to accomplish recommendation tasks. We first we formulate specialized prompts to enhance the ability of LLM to comprehend recommendation tasks. Subsequently, we use these prompts to fine-tune the LLaMA backbone LLM on a dataset of user-item interactions, represented by textual data, to capture user preferences and item characteristics. Our research underscores the potential of LLM-based generative recommendation in revolutionizing the domain of recommendation systems and offers a foundational framework for future explorations in this field. We conduct extensive experiments on benchmark datasets, and the experiments shows that our GenRec has significant better results on large dataset. Code and data are open-sourced at <https://github.com/rutgerswiselab/GenRec>.

KEYWORDS

Large Language Model; Recommender Systems; Natural Language Processing; Generative Recommendation

1 INTRODUCTION

Large Language Models (LLMs) have made a particularly significant milestone in this technological evolution. These LLMs, designed to understand and generate human-like text, have revolutionized numerous applications, from search engines to chatbots, and have facilitated more natural and intuitive interactions between humans and machines. This paper seeks to explore a relatively new and promising application of these models in the recommendation systems.

Recommendation systems have become an integral part of our digital experience. They are the unseen force guiding us through

vast amounts of data, suggesting relevant products on e-commerce websites, recommending movies on streaming platforms, and even proposing what news to read or videos to watch. The primary aim of these systems is to predict the individual user preferences and enhance user experience and engagement.

Traditionally, recommendation systems have been built around methods such as collaborative filtering [5, 6, 14], content-based filtering [16, 18], and hybrid approaches [1, 11]. Collaborative filtering leverages user-item interactions, making suggestions based on patterns found in the behavior of similar users or items. On the other hand, content-based filtering uses item features to recommend similar items to those a user has previously interacted with. Hybrid methods attempt to combine the strengths of these two approaches to overcome their respective limitations.

Despite the progress made with these traditional techniques, there still have some significant challenges. For instance, collaborative filtering struggles with the cold start problem, where it fails to provide accurate recommendations for new users or items due to lack of historical interaction data. Both content-based filtering hard to handle the issue of data sparsity, given that most users interact with only a small fraction of the total items available. Additionally, because of the computational complexity of processing large interaction matrices, these models often struggle to scale effectively with the growth of users and items.

The integration of text-based LLMs into recommendation systems presents an exciting opportunity to address these challenges [3]. These models can learn and understand complex patterns in human language, which allows for a more nuanced interpretation of user preferences and a more sophisticated generation of recommendations. However, a significant number of the prevailing recommendation models are trained using user and item indexes. This approach leads to the lack of text-based information in the dataset, including details like item titles and category information.

In this paper, we propose a novel large language model for generative recommendation (GenRec). One of the primary benefits of the GenRec model is that it capitalizes on the rich, descriptive information inherently contained within the item names, which often contain features that can be semantically analyzed, enabling a better understanding of the item’s potential relevance to the user. This could potentially provide more accurate and personalized recommendations, thereby enhancing the overall user experience.

We present experimental results to demonstrate the efficacy of our proposed method and compare its performance with other LLM recommendation models. The overarching aim of this paper is not

only to present our findings but also to inspire further research in this area. By highlighting the potential of LLMs in enhancing generative recommender systems, we hope to encourage a more widespread adoption of these models and stimulate further innovations in this field.

The key contributions of this paper can be summarized as follows:

- We highlight the promising paradigm of generative recommendation, which directly generates the target item to recommend, rather than traditional discriminative recommendation, which has to calculate a ranking score for each candidate item one by one and then sorts them for deciding which to recommend.
- We introduce a novel approach, GenRec, to enhance the generative recommendation performance by incorporating the textual information into the model.
- We also illustrate the efficacy of GenRec on practical recommendation tasks, underscoring its prospective abilities for a wider scope of applications.

In the following parts of the paper, we will discuss the related work in Section 2, introduce the proposed model in Section 3, analyze the experimental results in Section 4, and provide the conclusions as well as future work in Section 5.

2 RELATED WORK

2.1 Collaborative Filtering and Content-Based Recommendation Systems

Collaborative Filtering (CF) models are based on the concept of user-item interactions. Traditional CF models, such as the matrix factorization model [8], focus on latent factor modeling of user-item interaction matrices. More recent advancements, like NeuMF [5], have combined the merits of matrix factorization and neural networks to better capture complex user-item relationships.

On the other hand, Content-Based Recommendation systems rely on the features of items to make recommendations. Early works involved simple keyword matching [2] or cosine similarity based on TF-IDF vectors [13]. More advanced methods have started to exploit deep learning techniques, like CNN [10] and RNN [15], for extracting high-level features from item content.

2.2 Large Language Models for Recommendation

The use of large language models for recommendation systems has gained significant attention recently. These models exhibit great potential in the understanding and modeling of user-item interactions, exploiting rich semantics and long-range dependencies present in user activity data.

The pioneering work of P5 [3] illustrated the feasibility of formulating recommendation as a natural language task. P5 [3] fines the widely-used open-source T5 model [12] to create a unified system capable of handling various tasks. These tasks include not only recommendation ranking and retrieval but also complex functions like summary explanation. This innovative approach highlighted the versatility of large language models in handling multi-task learning in the recommendation context. However, the potential

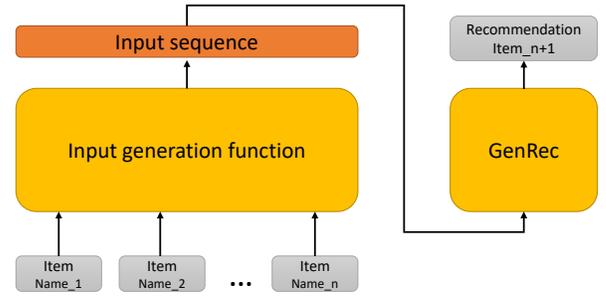


Figure 1: An illustration of GenRec. Our model will generate a input sequence based on the interaction history. Then the model will predict the next item the user may interact with.

of large language models to understand and generate text-based recommendations has not been fully explored.

In this paper, we propose a novel approach to text-based generative recommendation, leveraging the latest advances in large language models. We aim to address some of the limitations of previous works and push the boundaries of what is possible in the realm of recommendation systems.

3 METHOD

The architecture of the proposed framework is illustrated in Figure 1. Given a user’s item interaction sequence, the large language model for generative recommendation (GenRec) will format the item names with a prompt. This reformatted sequence is subsequently employed to fine-tune a Large Language Model (LLM). The adjusted LLM can then predict subsequent items the user is likely to interact with. In our paper, we select the LLaMA [17] language model as the backbone. However, our framework retains flexibility, allowing for seamless integration with any other LLM, thus broadening its potential usability and adaptability.

3.1 Sequence Generation

The initial component of GenRec is a generative function, tasked with producing various sequences that encapsulate user interests. To enhance the model’s comprehension of the recommendation task, we have devised multiple prompts that facilitate sequence generation. Take Figure 2 as an example, we use the user’s movie watching history as the training data and use this information to format the training sequence. The sequence consist of three part, instruction input and output. The instruction element outlines the specific task of movie recommendation, for which we have created several directives to enhance the LLM’s comprehension of the ongoing recommendation task. The input represents the history of the user’s interactions, excluding the most recent instance. And the output is the latest interaction in this record. The primary task for the LLM here is to predict this final interaction accurately.

Refer to Figure 2 for an illustration. This figure represents how we utilize a user’s history of watched movies as interaction data. Given the prompt, "Based on the movie viewing habits, what is the most likely movie they will select to watch next?" and the provided input, we then allow GenRec to forecast the subsequent output.

Methods	MovieLens 25M				Amazon Toys			
	HR@5	NDCG@5	HR@10	NDCG@10	HR@5	NDCG@5	HR@10	NDCG@10
P5	0.0688	0.0464	0.1040	0.0577	0.0239	0.0145	0.0411	0.0201
GenRec	0.1034	0.0716	0.1311	0.0837	0.0190	0.0136	0.0251	0.0157

Table 1: Experimental results on Normalize Discounted Cumulative Gain (NDCG@k) and Hit Ratio (HR@k). Bold numbers represent best performance.

Interaction history: Pinocchio (1940), Legends of the Fall (1994), Once Were Warriors (1994), In the Name of the Father (1993), Shadowlands (1993), Heavenly Creatures (1994), Quiz Show (1994), In the Line of Fire (1993)

Recommendation Prompt Example:

Instruction: Given the movie viewing habits, what is the most probable movie they will choose to watch next?

input: Pinocchio (1940), Legends of the Fall (1994), Once Were Warriors (1994), In the Name of the Father (1993), Shadowlands (1993), Heavenly Creatures (1994), Quiz Show (1994)

output: In the Line of Fire (1993)

Figure 2: GenRec on recommendation task. Based on the interactive history, GenRec can convert them to a training sequence which consists of instruction, input and output.

3.2 Training Strategy

In this paper, we use the LLaMA model as the backbone for the training of GenRec. The LLaMA model is pre-trained on an expansive language corpus, offering a valuable resource for our intended purpose of efficiently capturing both user interests and item content information. However, it’s important to note that the memory requirements for GPU to fine-tune LLaMA, even the 7-billion parameter version, are pretty substantial.

To circumvent this challenge and conserve GPU memory, we adopt the LLaMA-LoRA architecture for fine-tuning and inference tasks within the scope of this study. By this measure, we have achieved a significant reduction in the GPU memory requirements. With this optimized approach, we can fine-tune the LLaMA-LoRA model on a single GPU with a memory capacity of 24GB.

However, in an effort to decrease the overall training time, we have employed a data parallel technique and leveraged multiple GPUs in the experiments. Further details about our experiments, including the implementation and results, will be shared in the following sections of this paper.

4 EXPERIMENTS

4.1 Dataset

We conduct extensive experiments on two real-world datasets from Amazon [9] and MovieLens [4], respectively, to evaluate the performance of our proposed GenRec approach on recommendation tasks. The Amazon datasets, which record user purchase histories across a diverse range of products, were sourced from the Amazon.com platform. MovieLens datasets comprise a large number of movie ratings and associated metadata, contributed by users of the

MovieLens website over various periods. The descriptive statistics of these datasets are depicted in Table 1 (see reference 2). For each user interaction sequence, the most recent item is used as the test data, the second-most recent is used as validation data, and the remaining is used for training.

Dataset	MovieLens 25M	Amazon Toys
#Users	162,541	19,412
#Items	62,423	11,924
#Interaction	25,000,095	2,252,771

Table 2: Basic statistics of the recommendation datasets.

4.2 Evaluation Metrics

In this paper, we evaluate the performance of the model using two widely used metrics : Hit Ratio (HR) and Normalized Discounted Cumulative Gain (NDCG). The HR metric indicates the percentage of items recommended by the model that match those in the ground truth data. On the other hand, NDCG is employed to assess the efficacy of the recommendations when they are ranked, factoring in the relevance of the suggested items. These metrics have found wide acceptance in the evaluation of recommendation systems due to their robustness and comprehensiveness.

4.3 Implement Details

The GenRec model was pretrained for 5 epochs using the AdamW optimization [7] on four NVIDIA RTX A5000 GPUs with a batch size of 128. The peak learning rate was set to 3×10^{-4} and the maximum input length was set to 256 tokens. A warm-up strategy was employed during training, where the learning rate was gradually increased over the first 1000 steps.

4.4 Baseline Methods

P5 [3]: The Pre-train, Personalized Prompt, and Predict Paradigm (P5) incorporates an array of templates for input and target sequences throughout the training process. This unique approach proficiently dissolves the boundaries between different tasks, promoting a more fluid and integrated training procedure. It has showcased noteworthy performance in the domain of sequential recommendation tasks, underlining its effectiveness and applicability.

4.5 Performance Comparison

As we can see in the Table 1, P5 has better performance on Amazon Toys datasets, while our GenRec has significant better performance on movielens 25M datasets. The possible reasons behind this differential performance could be attributed to the distinct nature of the

datasets. The MovieLens 25M dataset, unlike Amazon Toys datasets, contains a richer amount of interaction information, which provides a more robust understanding of the user's preferences and behavior, thus likely leading to more accurate recommendations.

Our GenRec model, designed to effectively capture both user interests and item content information and produce more accurate and relevant recommendations. On the other hand, P5, while robust in handling sequential data, might not be as adept in leveraging this additional interaction information, resulting in relatively lower performance on the MovieLens 25M dataset.

5 CONCLUSION

In conclusion, our work on the text-based Large Language Model for Generative Recommendation (GenRec) has revealed a novel and promising approach in the field of recommendation systems. By focusing on the semantic richness of item names as input, GenRec promises more personalized and contextually relevant recommendations. Our practical demonstrations highlight GenRec's efficacy and point towards its adaptability across a diverse range of applications. Furthermore, the flexibility of the GenRec framework facilitates integration with any Large Language Model, hence widening its sphere of potential utility.

In terms of future work, there are several directions to explore. We intend to refine the generation of sequences by developing more sophisticated prompts, which could further enhance the model's understanding of recommendation tasks. Additionally, we plan to extend our research to incorporate more complex user interaction data, such as ratings or reviews, which could provide deeper insights into user behavior and preferences. A further direction would be to test GenRec's performance with different Large Language Models, investigating the possible benefits and trade-offs.

Our research with GenRec thus far has shown significant promise, and we look forward to continuing to develop and refine this approach. We believe that with further investigation, GenRec could revolutionize the way recommendation systems operate, ultimately leading to more personalized and satisfying user experiences.

REFERENCES

- [1] Justin Basilico and Thomas Hofmann. 2004. Unifying collaborative and content-based filtering. In *Proceedings of the twenty-first international conference on Machine learning*. 9.
- [2] Gaurav Bhalotia, Arvind Hulgeri, Charuta Nakhe, Soumen Chakrabarti, and Shashank Sudarshan. 2002. Keyword searching and browsing in databases using BANKS. In *Proceedings 18th international conference on data engineering*. IEEE, 431–440.
- [3] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as language processing (rlp): A unified pretrain, personalized prompt & predict paradigm (p5). In *Proceedings of the 16th ACM Conference on Recommender Systems*. 299–315.
- [4] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)* 5, 4 (2015), 1–19.
- [5] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.
- [6] Joseph A Konstan, Bradley N Miller, David Maltz, Jonathan L Herlocker, Lee R Gordon, and John Riedl. 1997. Grouplens: Applying collaborative filtering to usenet news. *Commun. ACM* 40, 3 (1997), 77–87.
- [7] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Bkg6RiCqY7>
- [8] Andriy Mnih and Russ R Salakhutdinov. 2007. Probabilistic matrix factorization. *Advances in neural information processing systems* 20 (2007).
- [9] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*. 188–197.
- [10] Keiron O'Shea and Ryan Nash. 2015. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458* (2015).
- [11] Michael J Pazzani. 1999. A framework for collaborative, content-based and demographic filtering. *Artificial intelligence review* 13 (1999), 393–408.
- [12] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* 21, 1 (2020), 5485–5551.
- [13] Juan Ramos et al. 2003. Using tf-idf to determine word relevance in document queries. In *Proceedings of the first instructional conference on machine learning*, Vol. 242. Citeseer, 29–48.
- [14] J Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen. 2007. Collaborative filtering recommender systems. *The adaptive web: methods and strategies of web personalization* (2007), 291–324.
- [15] Alex Sherstinsky. 2020. Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. *Physica D: Nonlinear Phenomena* 404 (2020), 132306.
- [16] Jieun Son and Seoung Bum Kim. 2017. Content-based filtering for recommendation systems using multiattribute networks. *Expert Systems with Applications* 89 (2017), 404–412.
- [17] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [18] Robin Van Meteren and Maarten Van Someren. 2000. Using content-based filtering for recommendation. In *Proceedings of the machine learning in the new information age: MLnet/ECML2000 workshop*, Vol. 30. Barcelona, 47–56.