

Article

BookGPT: A General Framework for Book Recommendation Empowered by Large Language Model

Zhiyu Li ¹, Yanfang Chen ^{2,*}, Xuan Zhang ³ and Xun Liang ³¹ Institute for Advanced Algorithms Research, Shanghai 200232, China² Libraries, Renmin University of China, Beijing 100872, China³ School of Information, Renmin University of China, Beijing 100872, China

* Correspondence: cyf@ruc.edu.cn

Abstract: With the continuous development and change exhibited by large language model (LLM) technology, represented by generative pretrained transformers (GPTs), many classic scenarios in various fields have re-emerged with new opportunities. This paper takes ChatGPT as the modeling object, incorporates LLM technology into the typical book resource understanding and recommendation scenario for the first time, and puts it into practice. By building a ChatGPT-like book recommendation system (BookGPT) framework based on ChatGPT, this paper attempts to apply ChatGPT to recommendation modeling for three typical tasks: book rating recommendation, user rating recommendation, and the book summary recommendation; it also explores the feasibility of LLM technology in book recommendation scenarios. At the same time, based on different evaluation schemes for book recommendation tasks and the existing classic recommendation models, this paper discusses the advantages and disadvantages of the BookGPT in book recommendation scenarios and analyzes the opportunities and improvement directions for subsequent LLMs in these scenarios. The experimental research shows the following: (1) The BookGPT can achieve good recommendation results in existing classic book recommendation tasks. Especially in cases containing less information about the target object to be recommended, such as zero-shot or one-shot learning tasks, the performance of the BookGPT is close to or even better than that of the current classic book recommendation algorithms, and this method has great potential for improvement. (2) In text generation tasks such as book summary recommendation, the recommendation effect of the BookGPT model is better than that of the manual editing process of Douban Reading, and it can even perform personalized interpretable content recommendations based on readers' attribute and identity information, making it more persuasive than interpretable one-size-fits-all recommendation models. Finally, we have open-sourced the relevant datasets and experimental codes, hoping that the exploratory program proposed in this paper can inspire the development of more LLMs to expand their applications and theoretical research prospects in the field of book recommendation and general recommendation tasks.

Keywords: book recommendation; large language model; general recommendation

Citation: Li, Z.; Chen, Y.; Zhang, X.; Liang, X. BookGPT: A General Framework for Book Recommendation Empowered by Large Language Model. *Electronics* **2023**, *12*, 4654. <https://doi.org/10.3390/electronics12224654>

Academic Editors: Yong Zheng, Peng Liu and Lemei Zhang

Received: 9 September 2023

Revised: 19 October 2023

Accepted: 10 November 2023

Published: 15 November 2023



Copyright: © 2023 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Book understanding and personalized recommendation (BUPR) are vital applications in the realm of library and information science (LIS). Within the BUPR context, we typically address three subproblems, including (1) recommending suitable books based on users' interests and preferences, (2) predicting a new book's popularity to determine its procurement, and (3) providing interpretable recommendations to different users to enhance user adoption rates. Generally, we must model the interactions between readers and books, account for readers' basic attributes, and consider books' fundamental attributes, among other feature and attribute data. We also utilize various machine learning methods to train and optimize unique recommendation models for each subtask, thereby improving the final recommendation effectiveness. However, as the scenario grows more complex

and the volume of data increases, it becomes challenging to fulfill diverse application and recommendation needs. Is it feasible to establish a unified personalized recommendation framework capable of solving all fundamental problems in the BUPR context with merely a handful of task-relevant training examples?

In recent years, the field of natural language processing (NLP) has witnessed significant advancements, with substantial changes in both model parameter scale and training data richness. For instance, in early December 2022, OpenAI unveiled a GPT-3.5-based chatbot [1] named Chat Generative Pretrained Transformer (ChatGPT) [2]. This chatbot utilizes large-scale pre-trained language models and is finely tuned for efficient natural language comprehension and logical reasoning in multi-turn conversations. Specifically, it can execute a range of NLP tasks, such as aiding with code writing, summarizing documents, and continuing novel writing. Following its launch, this model has ignited considerable industry discussion.

Figure 1 illustrates the trend of the Baidu search index for ChatGPT from its initial release in early December 2022 until now, with key events annotated. The chart reveals that during the initial release phase (November 2022–February 2023), due to the model's performance and user interface being less than perfect, the overall popularity of ChatGPT remained relatively low. However, starting in February, with OpenAI's model iteration and the news of several significant events, such as ChatGPT passing the Google engineer interview [3] and the broad participation of multiple companies, ChatGPT began a flourishing journey of research and application. By the end of April, the overall daily search volume for ChatGPT on Baidu reached 90,000. Moreover, ChatGPT's technical foundations and its large language models (LLMs) have also garnered substantial attention from the academic community [1].

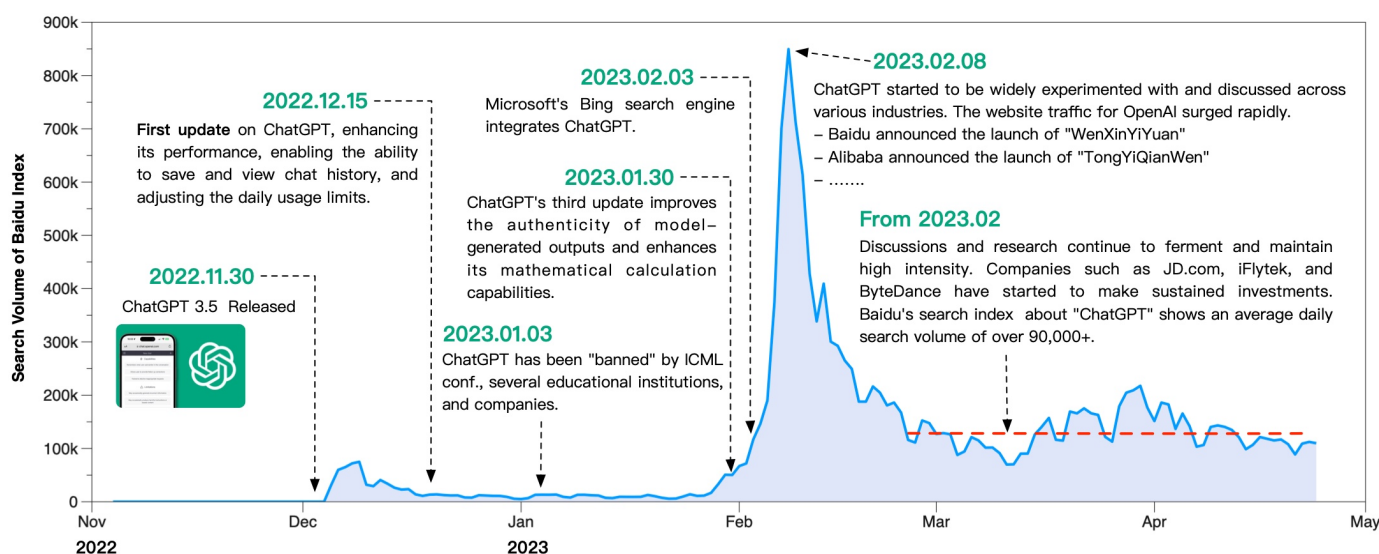


Figure 1. ChatGPT's search volume in the Baidu Index from November 2022 to April 2023.

In the LIS field, there has been a growing body of research focusing on the theoretical aspects related to models akin to ChatGPT. For instance, some studies have delved into the technical ethics and risks associated with ChatGPT-like models [4–6], and others have explored how to better harness ChatGPT in LIS [7,8]. However, these studies have primarily centered on the theoretical influences and application analyses of ChatGPT without conducting extensive practical experiments or testing within actual LIS scenarios. There is a need to investigate and research comprehensive experimental verification of whether ChatGPT models can be employed to construct recommendation frameworks to address recommendation issues in LIS and how these models perform in comparison to traditional recommendation models in these scenarios, along with their pros and cons. Hence, this paper zeroes in on the BUPR scenario to empirically examine the feasibility of

using ChatGPT models in book recommendation contexts and assesses the performance of ChatGPT models through experimental comparisons on BUPR tasks.

The main contributions to this article include the following: (1) This article applies LLMs such as ChatGPT to build a unified recommendation system framework for the classic BUPR task in the LIS field to explore the possibility of applying LLMs in LIS. (2) Based on the BUPR scenario, we discuss construction ideas and prompt engineering methods for three subtasks: the book rating task, the user book rating preference recommendation task, and the book summary recommendation task. Two different prompt modeling methods, zero-shot modeling and few-shot modeling, are verified in terms of their feasibility through empirical research. (3) Finally, we open-sourced the data and testing schemes involved in the experimental process to facilitate further research and discussion on the corresponding issues.

The outline of this article is as follows: We commence by discussing the related work involved in this study from two perspectives: “Large Language Models (LLMs)” and “book recommendations”. Next, we provide an overview of the BookGPT model, followed by comprehensive formal definitions of the three book recommendation subtasks. We then propose the construction methods of prompt engineering, the verification methods applied to the output results, and the methods for evaluating recommendation effectiveness. Subsequently, we conduct a detailed experimental evaluation of the ChatGPT-like book recommendation system (BookGPT), which includes an analysis of the dataset, the design of an evaluation scheme, and a discussion of the experimental results. Finally, we summarize the research content and focus of this article and suggest future research directions in the field of book recommendations based on LLMs.

2. Related Works

2.1. Large Language Models

LLMs typically refer to NLP models with parameter sizes exceeding billions. Recently, research on LLMs has become an important frontier in the field of NLP, from the widely used and researched statistical language models (SLMs) [9–11], to neural language models (NLMs) based on neural networks for NLP [12,13], to pretrained language models (PLMs) [14–16], and finally to LLMs [17–20]. With the iteration of these models, NLP technology has exhibited typical characteristics: the model parameter scale is becoming larger, the context awareness of these models is strengthening, the durations of multiturn conversations are lengthening, and multiple modalities of interaction are being utilized.

An LLM is a neural network architecture model based on the transformer mechanism [21] that extracts and expresses natural language features by introducing multihead attention and stacking multiple layers of deep neural networks. Among the various types of available LLMs [17–20], the main differences lie in the sizes of their training corpora, model parameter sizes, and scaling sizes. Based on a well-designed prompt engineering strategy [22], LLMs trained on large-scale corpora can usually produce good dialogue results. The natural language understanding and response abilities exhibited by current LLM models are generally believed to be emergent abilities [23,24] resulting from the tremendous growth in the number of model parameters and the size of the training corpus; i.e., when the parameter scale exceeds a certain level, the developed model exhibits new abilities that are radically different from those of its previous levels, such as in-context learning (ICL) [25] and chain of thought (CoT) [26].

Popular versions of LLM models currently include the GPT3/4 series of models (released by OpenAI) [17], the LLaMA model (released by Meta), and the GLM130B model (released by Tsinghua University). However, due to the strong commercial promotion and good product design provided by Microsoft and OpenAI, the ChatGPT application built on top of the GPT3.5/4 series of models is being increasingly adopted and used by researchers and enterprises in various real-world scenarios, such as intelligent customer service [27], interactive translation [28], and personal assistants [29]. Considering the cost of experimentation, this paper’s LLMs are based on OpenAI’s GPT3.5 and use gpt-3.5-

turbo-0301 as the kernel model, using application programming interfaces (APIs) provided by OpenAI.

2.2. Personalized Book Recommendations

Book understanding and recommendation is a fundamental application problem in the field of LIS. With the continuous increase in the number of book resources, both the number of book types and the number of interactions with readers are rapidly increasing. Therefore, how to select suitable books from a massive candidate set for recommendation is a fundamental problem. Generally, we can use personalized recommendation models to solve this typical information overload problem. In existing recommendation systems, the basic definition of a user's preference probability for an item i can be represented by the following function:

$$y_{i \rightarrow u} = f(h_i, h_u) \in [0, 1]$$

where h_i and h_u represent the learned item feature representation and user feature representation, respectively, and $f(\cdot)$ represents the scoring function that matches the user and item features, such as the cosine similarity function and multilayer perceptron module. Therefore, the existing research on personalized book recommendations can be partly summarized as follows.

Collaborative Filtering (CF)-Based Methods. CF is a classic and practical book recommendation algorithm based on the similarity between readers or books. The basic idea is that if user 1 likes book A and user 2 also likes books A and B, it can be assumed that there is a certain similarity between user 1 and user 2. Therefore, the other items liked by user 1 can be recommended to user 2, or vice versa. CF models can be divided into two types: user-based CF [30,31] and item-based CF models [32,33]. CF models usually require a large amount of user behavior data for training and prediction, so they are difficult to use in cases with sparse data or cold-start situations.

Deep Learning (DL)-Based Methods. In recent years, an increasing number of DL-based book recommendation algorithms have been proposed to better address the issues affecting CF models. DL models can learn more complex feature representations from user interaction data, thereby improving their recommendation accuracy. Specifically, DL-based book recommendation algorithms can be divided into two types: matrix factorization-based models [34,35] and sequence-based models [36].

Graph Neural Network (GNN)-Based Methods. A GNN-based recommendation system is a recently proposed algorithm that combines graph theory and DL [37]. A GNN-based book recommendation system represents readers, books, and their interaction information as a graph and then uses the GNN model to learn and aggregate the node and edge features in the graph, obtaining higher-order feature representations for different recommendation tasks [38,39].

3. Methods

3.1. Overview

This paper proposes a BookGPT. By combining the existing LLMs with typical tasks found in book recommendation scenarios, this framework constructs appropriate prompt strategies based on different task features and combines data validation, backtracking, and retrying methods to explore the possibility of using LLMs in book recommendation scenarios. As shown in Figure 2, the BookGPT framework is divided into four modules: (1) book recommendation, task definition, and data preparation, (2) prompt engineering, (3) GPT-based interaction and response parsing, and (4) task evaluation.

3.1.1. Book Recommendation Task Definition and Data Preparation

The BookGPT framework is primarily designed for three typical book recommendation applications: book rating recommendations, user rating recommendations, and book content summary recommendations. These three scenarios correspond to different aspects of applications: selecting high-quality new books based on rating potential, personalized

recommendations based on reader preferences, and explainability recommendations based on book summaries. In these scenarios, the feature system can be categorized into three types: basic user attribute features, book resource attribute features, and user-book interaction behavior features. Depending on the application scenario and data enrichment circumstances, it's typically possible to construct recommendation strategies for zero-shot and few-shot settings. This approach can enhance the satisfaction derived from the resulting recommendations. The specific formalizations and definitions of these tasks are detailed more thoroughly in the section titled "Book Recommendation Task Definition".

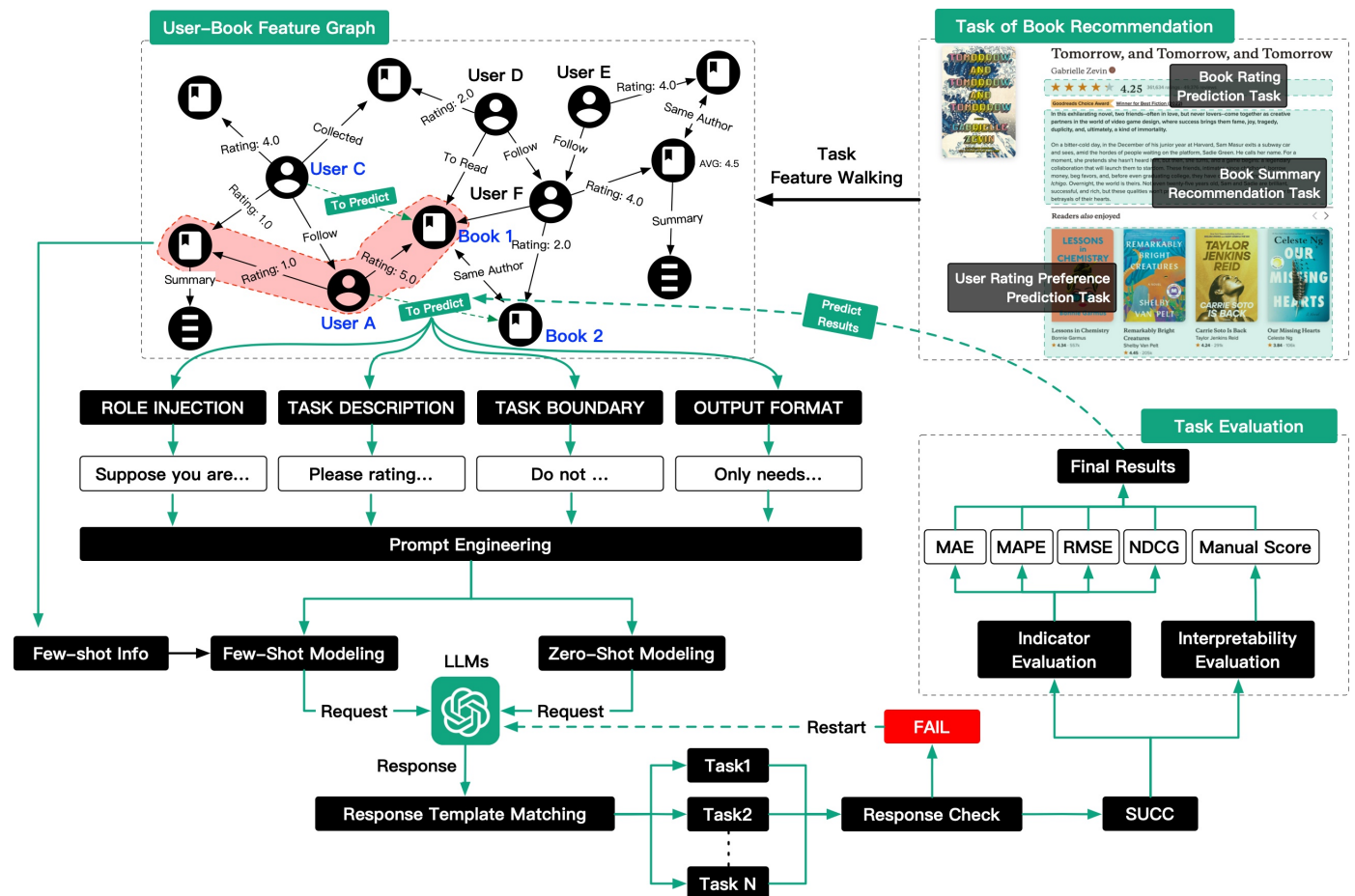


Figure 2. Framework of the BookGPT.

3.1.2. Prompt Engineering

Unlike traditional recommendation systems, the core recommendation module of the BookGPT is composed of an LLM; its recommendation process depends on the model's understanding and representation of natural language commands; and the output results are also highly flexible. Therefore, by designing appropriate prompt formats, the model's ability to understand tasks and the effectiveness of the final output results can be effectively enhanced [40]. Basic prompt engineering involves four core parts: an injected identity, a task description prompt, a task boundary prompt, and an output format prompt. At the same time, the prompt format of CoT [1] can be used to guide the model to solve complex tasks in a step-by-step manner and increase the modeling accuracy of the model. For a detailed analysis of the prompt engineering process of the BookGPT models, please refer to Section "Prompt Engineering for Book Recommendation".

3.1.3. GPT-Based Interaction and Response Parsing

The BookGPT model proposed in this article is tested and modeled using the ChatGPT (<https://platform.openai.com/docs/api-reference>, accessed on 15 May 2023) API provided

by OpenAI. The LLM used is the “gpt-3.5-turbo-0301” model, which is an optimized version based on GPT-3.5. It has faster response times and is approximately 10 times less expensive to call than the base model of GPT-3.5 with the same number of token requests. At the same time, due to the addition of many random factors in the response process of ChatGPT (to prevent the returned answers from being too convergent), it is necessary to perform a targeted formal verification for each returned result. If the obtained response does not satisfy the requirements of the established task, it is necessary to improve the prompt method or try to request it again. For a detailed description, please refer to Section “Output Verification and Task Restarting”.

3.1.4. Task Evaluation

As the BookGPT is the first framework to apply an LLM to the book recommendation system scenario, no evaluations of the application efficiency and feasibility of this scenario have been performed in previous studies. Therefore, we attempt to design two evaluation strategies, including a metric evaluation and an interpretability evaluation, to evaluate and discuss the performance achieved by the BookGPT on the three key tasks in the book recommendation scenario. Through detailed empirical research, we explore and analyze the advantages and problems of the BookGPT in the book recommendation scenario and conduct relevant discussions on subsequent research directions. The detailed analysis of this part can be found in Sections “Task Evaluation” and “Experiments”.

3.2. Book Recommendation Task Definition

3.2.1. Book Rating Recommendation

The book rating task is one of the fundamental tasks in book recommendation scenarios, especially those such as introducing new books and performing book evaluations. To effectively evaluate the ability of the recommendation system constructed by ChatGPT-like LLMs on this task, this paper verifies it through two modeling methods, zero-shot modeling and few-shot modeling, and measures the quality of the recommendation system by examining the difference between the rating results of the discriminant system and the actual rating results. The specific task definition is as follows.

Zero-Shot Modeling. Given a book name b_x and an author name without any background information, the system is required to output a rating result $R_{b_x} \in [0, 10]$ for the corresponding book, where a higher score indicates that the book is more recommended for reading.

Few-Shot Modeling. Compared to zero-shot testing, the essence of few-shot testing is to enhance the model’s understanding of the given task by providing partial sample information, with the hope of improving the final prediction performance. Therefore, in this paper, the few-shot testing case for book rating is defined as follows: given a list of books with their corresponding ratings as pairs for the same author, $P_{u_i} = (b_1, R_{b_1}), (b_2, R_{b_2}), \dots, (b_n, R_{b_n})$, a small portion of them (e.g., k) are selected as known input information, and the system is required to rate the remaining books from the same author. Finally, the system’s rating results are evaluated by the difference between the predicted ratings and the actual ratings of the remaining samples.

3.2.2. User Rating Preference Recommendation

The user rating recommendation task has a wider range of application scenarios than the book rating task, such as predicting the book preferences of users for e-commerce sales, predicting the interest levels of readers in book borrowing, predicting clicks, and predicting library browsing. This task typically uses historical interactions (clicks, browsing, borrowing, collecting, commenting, rating, etc.) between readers and books as feature data sources, combines them with basic user attributes and book attributes, and utilizes various machine learning models to build accurate recommendations. In this article, the specific task is defined as follows.

- **One-Shot User Preference Modeling.** Given a historical book behavior sample sequence (such as a rating sequence) for user u_i , $H_{u_i} = \{b_1, b_2, \dots, b_n\}$, the model is only provided with a single training sample as a hint or training set and is required to score the remaining samples in the behavior sequence. The final evaluation is based on the consistency between the model's scoring results and the original sample results.
- **N-Shot User Preference Modeling.** Given a historical book behavior sequence (such as a rating sequence) for user u_i , $H_{u_i} = \{b_1, b_2, \dots, b_n\}$, a certain proportion of the data is selected as the training set (or prompt set) from it. The model is required to score the remaining sequence based on the provided training set, and the final evaluation is the consistency between the model's scoring results and the original sample results.

3.2.3. Book Summary Recommendation

The task of book summary recommendation aims to automatically extract concise and accurate summary content from books, providing readers with a quick way to understand the main contents of the books. This task typically utilizes NLP techniques, including text summarization, text classification, information extraction, and other techniques, to achieve its goals. In practical applications, summaries can serve as important data sources for book recommendations, search result previewing, knowledge graph construction, and other areas. Therefore, in this section, we compare the summaries generated by the ChatGPT model with the standard summaries produced by humans and evaluate the effectiveness of the proposed recommendation system from the perspectives of interpretability and credibility. We answer two questions. (1) Can large-scale language models such as ChatGPT and WenXinYiYan (WenXin) [41] (released by Baidu) achieve better results than humans in book summary generation tasks? (2) Will ChatGPT and Wenxin exhibit different summarization abilities for different categories of literary genres, such as novels, essays, and poetry? The specific generation forms of the comparison task include the following.

- **Summaries Without Length Limitations.** This task generates summary recommendations based on specified author and book title information, with no limit imposed on the character length.
- **Summaries with Length Limitations.** For the reason that the summaries generated by LLMs usually contain more characters than those of humans, to further ensure the fairness of the comparison, the maximum number of characters that can be generated is further limited when generating an abstract with the model. Specifically, the model is required to generate summary recommendations based on the specified author and book title information, with a forced limit imposed on the maximum number of characters that can be generated; this limit is the same as the character count of the manual abstract provided for the same book.

3.3. Prompt Engineering for Book Recommendation

As far as we know, LLMs are typical generative language models, so the quality of the output contents of these models usually exhibits significant correlations with the input prompt contents. Therefore, in this section, we discuss how to design effective prompt content [22] for various types of book recommendation tasks to achieve improved recommendation efficiency. As shown in Figure 3, we provide prompt engineering examples for three typical book recommendation scenarios. Generally, prompt content typically includes four parts.

(1) **Role Injection Prompt.** This prompt is mainly used to indicate the role type represented by the LLM, guiding it to respond differently according to specific role types. As shown in Figure 4, in the book rating task, if no prompt is given for identity injection ("Assuming you are a professional book rating expert") and ChatGPT is instead directly asked to answer a task requirement such as "Please rate the book xxx", ChatGPT usually responds with a refusal to answer.

(2) **Task Description Prompt.** This prompt is used to provide information about the specific task that ChatGPT is being asked to perform. It typically includes details such as

the type of task, the goal of the task, and any relevant information about the input data or context. This prompt aids ChatGPT in comprehending its assignment, thus enabling it to generate more precise and relevant outputs. For instance, if the task was to recommend a book, the task description might encompass details about the user’s reading habits, the preferred genre of book, and any particular requirements or limitations that apply to the recommendation process.

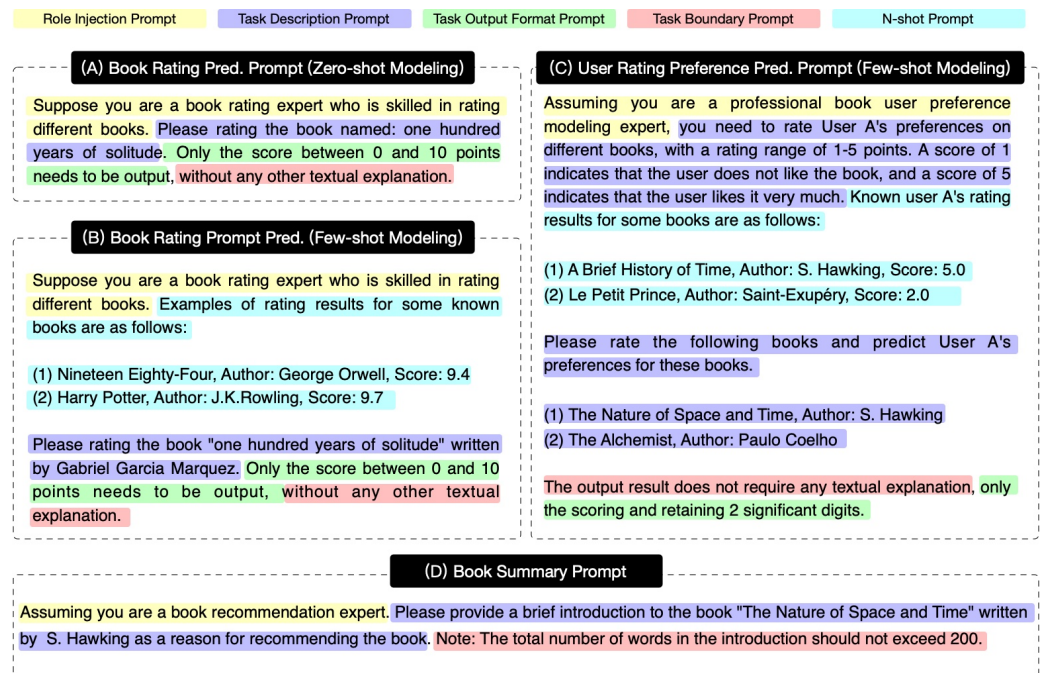


Figure 3. Prompt examples for the BookGPT.

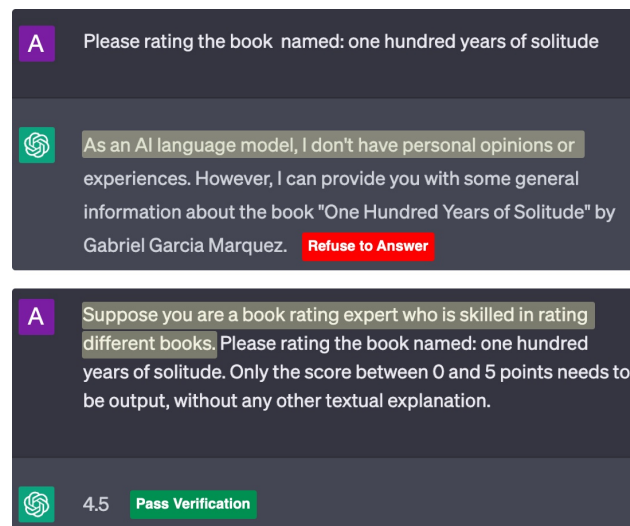


Figure 4. Example of role injection.

Indeed, due to existing constraints on model inference capabilities, task description prompts are often confined to a specific length. For instance, the current GPT-3.5 version imposes a limit of 4096 tokens for both model responses and prompt content to optimize inference outcomes (for more details, refer to <https://platform.openai.com/tokenizer>, accessed on 15 May 2023). Typically, content prompts derived from task examples (known as few-shot scenarios) act as supplementary training instances for the model. This additional information can improve the model’s comprehension and adaptation to the task at hand, thus enhancing prediction accuracy.

(3) Task Boundary Prompt. This part refers to the usage of prompts to set negative constraints on models akin to ChatGPT, instructing them on what to avoid in a given task. For instance, in a book rating task, if only identity injection and task description prompts are employed, the model might generate both a rating and an extensive explanatory text, which could complicate downstream applications. As such, it becomes necessary to explicitly define limitations and inform the model of the task boundaries; that is, a textual explanation is not required; only the rating result should be output. In such a scenario, the model will solely produce the corresponding rating score as needed.

(4) Task Output Format Prompt. After setting the role injection, task description, and boundary prompts in BookGPT, it's also necessary to instruct the model on the final output format. The primary benefits of this prompt include: (1) Enhancing the precision of the model's output. By defining the output format, we can ensure that the model's results meet the needs of downstream systems, circumventing mistakes and unnecessary extra processing steps resulting from incompatible data formats. (2) Increasing system maintainability. By explicitly stating the desired output format, we can avoid significant modifications and adjustments in downstream systems if the model's output format changes, promoting system maintainability and scalability. For instance, for a book rating prediction task, the output format should be restricted to a value with two decimal places. For a user preference estimation task, the output format should be confined to a Python list format.

3.4. Output Verification and Task Restarting

Through the design of a prompt engineering strategy, we can ensure that the output of the model satisfies the expected definition to some extent. However, since ChatGPT-type models are essentially natural language probability models, and because ChatGPT incorporates stochastic factors into its design to ensure the diversity of the generated results [1], it is possible that the model may produce different response results for the same input request. Therefore, in the end, we also need to recheck the legality of the generated content produced by ChatGPT-type models, i.e., perform a secondary verification of the critical output data format and the requirements.

In this module, we build independent validation functions for each type of book recommendation subtask. For example, for the book rating task, we need to check if the returned result is a numerical value. For the user rating recommendation task, we need to extract the rating values corresponding to the books in the returned result and check if the number of returned results matches the requested quantity. For the book summary task, we need to check if the length of the returned text satisfies the input requirements. If the result returned by ChatGPT does not meet the specified format requirements, we need to resend the request to the recommendation module until the maximum number of retries is reached (the maximum number of retries in this framework is set to 3).

3.5. Task Evaluation

To validate the performance of the BookGPT (including zero-shot and few-shot modeling), we evaluate the system from two aspects: a task metric evaluation and an interpretability evaluation.

3.5.1. Task Metric Evaluation

We evaluate the recommendation performance of the BookGPT on two subtasks: book rating recommendation and user rating preference recommendation. Specifically, we treat the book rating task as a regression model and evaluate the performance of different recommendation models in terms of the mean absolute error (MAE), mean absolute percentage error (MAPE), and root mean square error (RMSE) metrics. For the user rating preference recommendation task, which is treated as a sorted rating recommendation process, in addition to focusing on the MAE, MAPE, and RMSE, we further evaluate the performance of

different models in terms of the normalized discounted cumulative gain (NDCG) metric. The specific calculation methods for each metric are shown in Table 1.

Table 1. Task evaluation metrics.

Metrics	Equations	Concerns
MAE	$\frac{1}{n} \sum_{i=1}^n y_i - \tilde{y}_i $	absolute error
MAPE	$\frac{1}{n} \sum_{i=1}^n y_i - \tilde{y}_i / y_i $	percentage error
RMSE	$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2}$	divergence
NDCG@k	Details [42]	cumulative gain

3.5.2. Interpretability Evaluation

This part mainly evaluates the recommendation ability of the BookGPT in the book summary recommendation task. The goal of this task is to identify the key content of a book through its summary and arouse readers' interest in reading or purchasing it. Therefore, in this evaluation, we use manually generated summaries as the ground truth to evaluate the book summary generation performance of the BookGPT model when different books are tested. Specifically, since ChatGPT was trained on English corpora, we also introduce a large Chinese language model, Wenxin (released by Baidu Group), as a reference to test the performance of an LLM on Chinese summaries and compare it with ChatGPT. Finally, through the above two types of evaluations and empirical studies, we attempt to answer the following three questions.

- Q1: What tasks in the book recommendation scenario are the BookGPT suitable for? How is its performance?
- Q2: Is there a significant difference between the final recommendation performances achieved with zero-shot modeling and few-shot modeling in the BookGPT?
- Q3: In the book recommendation scenario, what are the potential research directions concerning ChatGPT-like LLMs in the future? What problems can they solve?

4. Experiments

As shown in Table 2, the experimental datasets in this article include three types of data: book rating data, book summary data, and book-user interaction data.

Table 2. Datasets.

Datasets	Types	Amount	Additional Information
Douban Rating	scores	3228	title, author, comment
Douban Book Summary	summary text	50	title
Goodbook-10k	scores, interactions	10,000	label, metadata and tag

Douban Rating. This dataset was collected from the book rating channel of Douban (<https://book.douban.com/>, accessed on 10 May 2023), which includes four key fields: book title, author, rating, and number of comments. Due to the limited number of API requests for OpenAI, only books with more than 2000 ratings are selected as the test dataset, resulting in 3228 popular books for evaluation purposes. In this experiment, we evaluate the performance of the BookGPT under zero-shot learning and different levels of few-shot learning (1/3-shot learning and 2/3-shot learning), and the core observation indicators are the MAE, MAPE, and RMSE.

Douban Book Summary. To evaluate the model's ability to summarize and recommend book content, we selected 50 popular books from the Douban TOP250 book channel (<https://book.douban.com/top250>, accessed on 10 May 2023), including 4 categories of literary genres and their human-written summaries: 20 novels, 10 essays, 10 poems, and 10 dramas. The human-written summaries are provided by Douban's book editing experts and are used as the benchmark for the comparison. The summary results of the

BookGPT model and the Wenxin model in the “restricted size” and “free size” summary scenarios are used as comparison models.

Input Settings. For the BookGPT model, its input is a natural language description of the features, such as the book’s name, user preference score, the summary of the book (if available), etc. An example of the input for the BookGPT model can be observed in Figure 3. The input of the comparison model is usually the numerical or categorical features of the corresponding dataset.

Notably, to ensure the fairness of the evaluation results, we first randomly mix the human-written summaries and model summaries and present them in a random order to different annotators. Furthermore, 15 annotation participants are asked to rank the summaries, and each summary must be annotated for at least 3 min. If a summary is ranked higher, it is considered to have a stronger recommendation ability and more attractiveness. Finally, the results acquired from different annotators are processed to obtain the evaluation results for each model with respect to different literary genres of books and overall. The core observation indicators for this task are the summary evaluation score and average summary length. The calculation method for the final score of each model’s summary is as follows:

$$Score_{M_1} = \frac{1}{N} \sum_i^{rank} freq \times w_i$$

In the above formula, N is the total number of annotators in the test, and $freq$ is the number of times the corresponding model option appears in position i with a weight of w_i . For example, assuming that 15 people participate in annotation sorting, three options need to be sorted, with positions 1, 2, and 3 corresponding to scores of 3, 2, and 1, respectively. Furthermore, if model X’s sorted summary results for book Y are first place 10 times, second place 2 times, and third place 3 times, the comprehensive score of model X’s summary for book Y in this test is $(10 \times 3 + 2 \times 2 + 3 \times 1) / 15 = 2.47$ points.

GoodBook-10k. The GoodBook-10k dataset [43] was collected from the Goodreads (<https://goodreads.com/>, accessed on 10 May 2023) book review website, which is the largest online reading community in the world and is similar to Douban Reading in China. The GoodBook-10k dataset contains rating data for 10,000 popular books and 5.98 million users’ ratings, with fields including book ratings, user bookshelf labels, book metadata, and book tags. In this paper, we use it for the user rating task, which includes three forms at the prompt level: 1-shot, 10-shot, and 20-shot learning, where the model is provided with 1, 10, or 20 model rating records, respectively, and is required to predict the remaining records and provide user preference rankings. The benchmark models for this task include the BookGPT model proposed in this article, as well as four classic CF-based recommendation algorithm models for personalized recommendation scenarios: the matrix factorization model (FunkSVD) [44], the K-nearest neighbors (KNN; means) model [45], the SlopeOne [46] model, and the CoClustering [47] model. In terms of replicating the comparison model, we used the Surprise (<https://surprise.readthedocs.io/en/stable/>, accessed on 10 May 2023) package based on Python for the experiment. The evaluation metrics are the NDCG@5,10,15,20, MAE, MAPE, and RMSE.

5. Results

This section analyzes the performance achieved by the BookGPT model and the baseline models on different tasks, answering the questions raised above, namely, how does the BookGPT model perform on different tasks in the book recommendation scenario? Can few-shot learning improve the recommendation performance of the BookGPT model?

5.1. Book Rating Task

As shown in Table 3, the prediction results yielded by the BookGPT model in book rating tasks are analyzed under zero-shot modeling, 1/3-shot modeling, and 2/3-shot modeling. Overall, the BookGPT model exhibits a good regression prediction ability,

with MAPE values ranging from 8.8% to 5.4%, indicating that it performs well on book rating tasks. The overall absolute percentage error is within a small range of 10%, and even for 2/3-shot modeling, the MAPE can reach an estimated value of 5.4%.

Table 3. Results of the book rating task. The green numbers with ↓ represent the percentage decrease in the corresponding indicators compared to Zero-shot modeling.

BookGPT	MAE	MAPE	RMSE
Zero-shot modeling	0.682	0.088	0.886
1/3-shot modeling	0.441 (↓35.34%)	0.057 (↓35.23%)	0.558 (↓37.02%)
2/3-shot modeling	0.419 (↓38.56%)	0.054 (↓38.64%)	0.538 (↓39.28%)

Furthermore, based on the results of few-shot modeling with prompt enhancement, the accuracy of the estimated rating is significantly improved compared with that of zero-shot modeling. The model's MAE value decreases from 0.682 (zero-shot) to 0.441 (1/3-shot modeling) and 0.419 (2/3-shot modeling), representing 35.34% and 38.56% mean absolute error decreases compared with that of zero-shot modeling, respectively. In addition, the model's RMSE also decreases significantly, from 0.886 to 0.558 (1/3-shot modeling) and 0.538 (2/3-shot modeling), with a relative decrease in the optimal mean square error of 39.28% (2/3-shot modeling). This result indicates that few-shot modeling can significantly reduce the prediction error induced by zero-shot modeling in the book rating task by providing reference samples for the BookGPT model.

Finally, comparing the results of the BookGPT model with those of different levels of prompt-based few-shot learning, increasing the prompt size again (from 1/3 to 2/3 of the training set) leads to an additional improvement in the final performance. However, the improvement from zero-shot learning to 1/3-shot learning is larger than the improvement from 1/3- to 2/3-shot learning. This is because the 1/3 prompt size already provides a good information reference for the BookGPT model, and further increasing the prompt size brings limited information gain while increasing the model's reasoning overhead. Therefore, in practical applications, the appropriate prompt size can be selected based on the information gain inflection point through ablation experiments and in combination with the scenario's needs, thus effectively yielding improved few-shot learning performance.

5.2. User Rating Preference Recommendation Task

Figure 5 represents the NDCG evaluation results obtained for the user rating preference recommendation task.

First, overall, the MF (FunkSVD) model can achieve good results in different subtasks, especially in the single-sample scenario, where it performs best. The reason for this is that the recommendation strategy based on FunkSVD models the user-book rating interaction matrix through matrix factorization. The optimization goal of this modeling method is to make the residual between the user ratings and the rating product obtained by matrix multiplication as small as possible. Therefore, even with limited reference rating information provided by the user to be predicted (such as in 1-shot learning), FunkSVD can still achieve good results in terms of metrics such as the RMSE. However, as the number of effective prompt samples (features) for a single user increases, clustering models represented by KNN (means) begin to perform better. During the prediction process, KNN (means) relies on modeling the user's historical rating habits to generate the final estimation result, leading to a prediction accuracy increase with the increase in the number of prompt samples. It's noteworthy that, despite BookGPT not surpassing the comparative model in performance, our results illustrate its potential in recommendation scenarios as a language model. This is particularly evident when the number of available examples for training or prompts is limited. In such low-sample situations, BookGPT can leverage its strengths and achieve relatively competitive results. Furthermore, a comparison of Figure 5a–c, reveals that BookGPT's performance closely aligns with the comparison model when the number of prompt examples increases to 20. This observation suggests that the strategic

use of sample prompts can considerably enhance the effectiveness of models like BookGPT, thereby optimizing their recommendation performance.

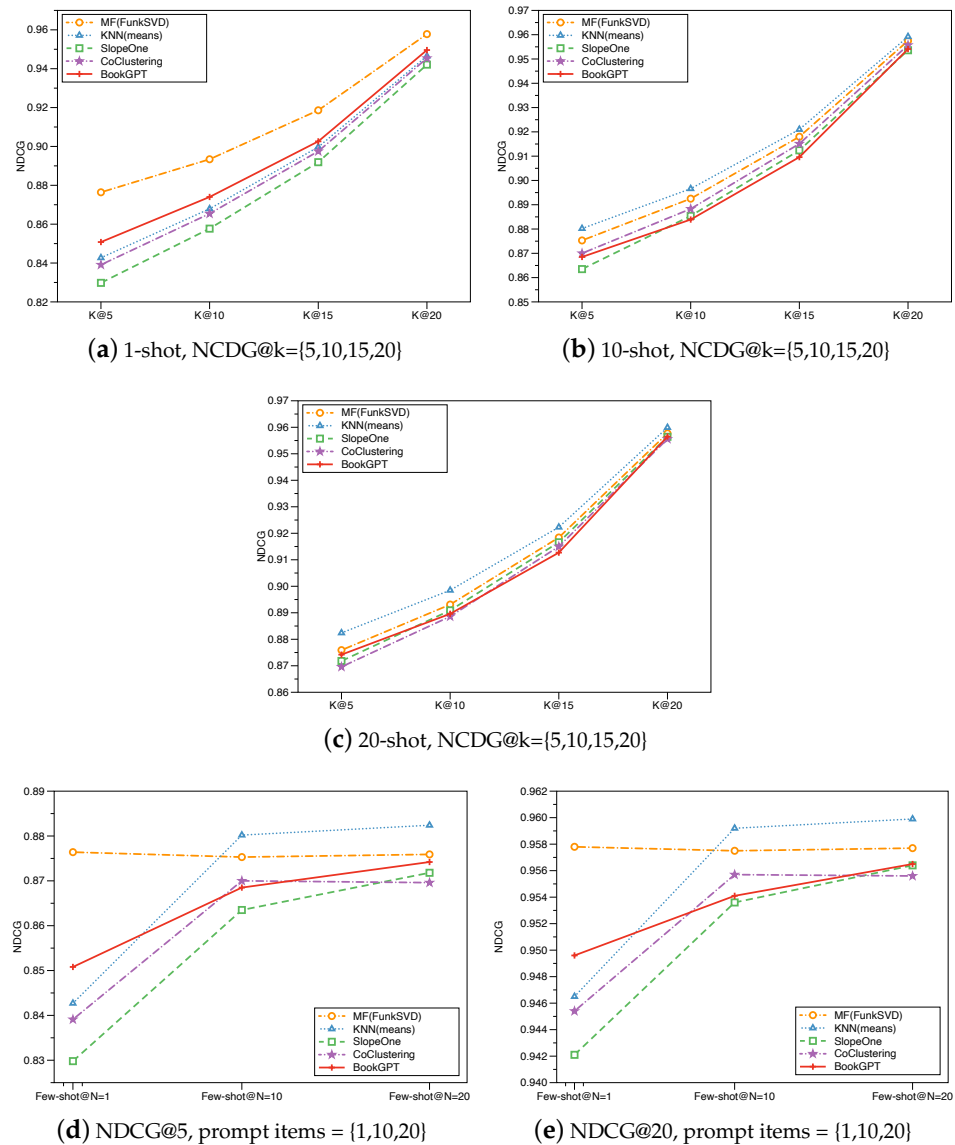


Figure 5. NDCG scores obtained in the user rating preference recommendation task. NDCG is a metric used to evaluate the performance of recommendation and information retrieval systems, considering both the relevance and ranking of recommended items. The value of NDCG ranges from 0 to 1, with 1 indicating optimal performance. Hence, a higher NDCG value signifies a more effective system.

Second, comparing the results based on different prompt sample quantities, the NDCG score of the FunkSVD model is not sensitive to the number of samples, and its performance remains consistent across the 1-shot, 10-shot, and 20-shot subtasks. However, other baseline models, such as KNN (means), SlopeOne, CoClustering, and the BookGPT model, exhibit significant changes in their NDCG scores as the number of prompt samples increases. This is because these models build a model of the user’s historical ratings during the result prediction process. For example, in the BookGPT model, during the rating process, the model uses the user’s historical book rating preferences as background knowledge to model the user’s interest preferences and makes comprehensive predictions for the newly predicted samples by referencing this knowledge. With an appropriate number of prompt

samples, the BookGPT model can typically learn real-time contextual knowledge and apply it to the prediction scenario. This capability is also known as in-context learning [1], which is one of the important foundational capabilities of LLMs. As shown in Table 4, by increasing the number of prompt samples, the BookGPT model achieves significant improvements in terms of its MAE, MAPE, and RMSE metrics in the 20-shot subtask, with error reductions of 21.67%, 16.53%, and 16.84%, respectively.

Table 4. The prediction errors induced in the user rating preference recommendation task.

Prompt Items	Models	Indicators		
		MAE	MAPE	RMSE
1-shot modeling	FunkSVD [44]	0.765	0.259	0.956
	KNN (means) [45]	0.865	0.283	1.153
	SlopeOne [46]	0.843	0.278	1.117
	CoClustering [47]	0.836	0.276	1.107
	BookGPT	1.075	0.297	1.342
10-shot modeling	FunkSVD [44]	0.733	0.249	0.917
	KNN (means) [45]	0.708	0.236	0.912
	SlopeOne [46]	0.741	0.242	0.965
	CoClustering [47]	0.737	0.243	0.961
	BookGPT	0.915	0.263	1.184
20-shot modeling	FunkSVD [44]	0.710	0.241	0.892
	KNN (means) [45]	0.676	0.226	0.879
	SlopeOne [46]	0.700	0.233	0.905
	CoClustering [47]	0.703	0.235	0.910
	BookGPT	0.842	0.248	1.116

Finally, when evaluating different types of metrics, it can be observed that in situations where the prompt samples are relatively few, such as in the 1-shot scenario, the BookGPT model performs well in terms of its ranking ability (NDCG), but it does not show any advantage in terms of error measurement (the MAE/MAPE/RMSE metrics) over the control models. Additionally, increasing the number of prompt samples for a single user from 10 shots to 20 shots does not result in a similar performance improvement in terms of the NDCG, including the BookGPT model. However, the error metric increases are still considerable. Therefore, if the application scenario emphasizes the absolute value preference of each user to be recommended, the effect can be improved by increasing the number of prompt samples for each user. If the focus is only on the relative ranking ability, performing modeling based on a small number of samples can satisfy the imposed requirements and further save inference resources.

5.3. Book Summary Recommendation Task

During the evaluation process of this task, to ensure the effectiveness of the results, all annotation processes are carried out anonymously with initial order randomization and cross-evaluation, requiring the average ranking results of each model to be obtained across multiple annotators. In addition, because LLMs such as ChatGPT and Wenxin are optimized based on human instructions, people tend to choose longer answers as high-quality answers during the optimization process, so if the summary size of the model is not controlled, the model tends to produce longer content recommendation results. To ensure fairness during the comparison with the human summaries of Douban, we add summary size restrictions to the prompt construction process, requiring that the summary recommendation results of the BookGPT and Wenxin be as close as possible to the number of words in human summaries to ensure the validity of the comparison. Table 5 shows the results of manual evaluations of the content summary recommendations provided by the ChatGPT-based BookGPT model and Wenxin for different book genres.

Table 5. The results obtained in the book summary recommendation task.

Size	Models	Genres									
		All		Novels		Poems		Essays		Dramas	
		Score	Length	Score	Length	Score	Length	Score	Length	Score	Length
Restricted size	Douban	2.05	300	2.07	323	2.09	360	2.14	255	1.88	237
	Wenxin [41]	2.35	251	2.15	278	2.49	291	2.34	204	2.65	206
	BookGPT	1.60	129	1.79	178	1.42	94	1.53	112	1.47	84
Free size	BookGPT	1.45	314	1.48	280	1.53	370	1.40	309	1.39	329
	Wenxin [41]	1.55	472	1.52	482	1.50	494	1.60	452	1.61	449

First, if we only compare the models based on the limited summary size, Wenxin achieves the best recommendation performance among all the compared models, both in the subgenre and overall tasks. Compared with the human-written summaries from Douban and the BookGPT model summaries, Wenxin achieves relative improvements of 14.97% and 47.25%, respectively.

Furthermore, although we limit and remind the models to pay attention to the number of characters during the prompt construction process, it is apparent that Wenxin and BookGPT have different understandings of the character limit requirement. During the actual testing process, we find that BookGPT is more conservative in terms of the character limit rule and usually strictly follows the limit requirement, while Wenxin tends to produce longer summaries. In terms of the average character lengths of the generated summaries, Wenxin exceeds BookGPT by 94.57%. Based on this result, we believe that one of the reasons for this is that Wenxin incorporates more Chinese language data into its training and fine-tuning processes and lacks intervention during rule-based prompt fine-tuning, making the model more inclined to produce longer results (with stronger expressions) for Chinese tasks. At the same time, this also indicates that Wenxin has a weaker sense of “rules”.

Afterwards, we remove the character limit during the prompt-building process, allowing the models to freely generate summary recommendations based on their own capabilities. As shown in Table 5, compared with the results obtained with limited character counts, the advantage of Wenxin over BookGPT is reduced in the free-scale evaluation, with a relative improvement of only 6.89% compared with BookGPT, as opposed to 47.25% under the character limit. This result also suggests that the advantage of Wenxin over BookGPT under the character limit may be due to the production of longer summaries. However, in terms of actual summary generation ability, Wenxin and BookGPT are relatively close.

From the performance results obtained for different genres, it can be seen that the improvement exhibited by Wenxin over the human summaries in the “poems” and “drama” genres is more significant than that in the “novels” and “essays” genres. Furthermore, if the summary sizes of the models are not limited, Wenxin and BookGPT perform similarly in the “novels” and “poems” genres, while Wenxin’s advantage is more obvious in the “essay” and “dramas” genres. However, regardless of the genre, the performance of Wenxin is better than that of the human summaries on Douban.

Overall, in the book content summary recommendation task, the BookGPT based on LLMs has certain advantages over the human-generated summaries on Douban and can provide relatively good improvements for different genres. However, we also discover some issues, such as “fantasizing” and “piecing together” in some summary content. For example, when BookGPT produces the summary of a book named “Demi-Gods and Semi-Devils”, it says, “The Legend of the Condor Heroes is one of Jin Yong’s representative works, telling the adventure story in the background of the prosperous Tang Dynasty and unfolding a multi-linear narrative centered on the protagonist Chen Jialuo with ups and downs”. For readers who are not familiar with this book, the result seems reasonable, but in reality, this summary not only describes the wrong dynasty of the story but also

uses an incorrect name for the protagonist. In contrast, Wenxin is correct in terms of all key information. Therefore, it can be seen that if one must achieve good fact-description results and accuracy in a specific scenario, it is usually necessary to further enhance the training process based on the corpus and prompt rules of that scenario. Otherwise, a model trained on a general language corpus may easily fail to ensure the factual correctness of the generated content.

6. Conclusions and Future Work

This paper introduces a book recommendation framework, BookGPT, leveraging Large Language Models (LLMs). This framework capitalizes on the comprehension and reasoning abilities of LLMs, applying them to the familiar context of book understanding and personalized recommendation within the Library and Information Science (LIS) field. We established a task definition, created a prompt engineering strategy, implemented an interactive querying method, and developed a result verification framework to explore the potential utility of LLMs in three typical subtasks of book recommendation: book rating, user-book preference recommendation, and book summary recommendation.

An extensive comparative analysis was conducted across a variety of prompt sample quantities, including zero-shot, one-shot, and few-shot modeling. Despite the current limitations of LLMs, our experimental results revealed that BookGPT showed promise in certain scenarios. Particularly, the model demonstrated impressive performance in the 1-shot scenario, indicating its substantial potential for handling tasks with limited sample availability. Additionally, as the number of prompt samples increases, the model's recommendation performance significantly improves.

While the overall performance of BookGPT in our experiments did not always outperform other models, it's crucial to note that the primary objective of this research was not merely to develop a high-performing model but to establish an adaptable framework for implementing LLMs in book recommendation systems. The limitations observed in BookGPT's performance can be viewed as opportunities for future enhancements rather than inherent flaws in the framework.

In the future, our objectives are to refine the BookGPT model within this framework, specifically focusing on areas identified as having potential for significant improvement. We aspire to delve more deeply into the untapped applications of LLMs within the LIS field, investigating innovative methods to augment efficiency. This exploration, in tandem with the insights obtained from this study, will pave the way for the advancement of more potent and versatile LLM-based recommendation systems. Our future work is aligned along three primary paths, each aiming to enhance a unique aspect of our current model: task-specific data fine-tuning, user feedback incorporation through multi-round conversation-based recommendations, and personalization of user information for explainable recommendations.

Optimization through task-specific data fine-tuning. The current BookGPT framework is built directly on pretrained LLMs such as ChatGPT and Wenxin, and its recommendation performance usually depends on the corpus during LLM training, with a focus on the model's generalization ability. It is not optimized for various proprietary scenarios in the LIS field. Therefore, an important research direction for the future is how to construct training data for fine-tuning specific domain scenarios to further leverage the knowledge and reasoning advantages of LLMs, improve the recommendation and prediction performance of the corresponding model, and even achieve better performance than that of the current state-of-the-art recommendation models in domain-specific scenarios.

Combining user feedback with multi-round conversation-based recommendations. In the current BookGPT recommendation paradigm, single-round offline recommendation is adopted, and no attention is paid to user feedback regarding the recommendation effect (such as clicking, pressing the "favourite" button, and borrowing behavior). Therefore, it is also worth exploring how to integrate different real-time user behaviors and interactions with the system into the recommendation paradigm and construct a multi-round conversation-based recommendation model. Through multi-round conversation-based

recommendation, not only can the contextual learning abilities of LLMs be maximized, but more training prompt language materials can be generated from the interactions, which can improve the model's training and fine-tuning results.

Incorporating personalized user information for explainable recommendations. For the reason that LLMs are developed and trained based on various natural language corpora; it is possible to incorporate more personalized user information into the recommendation results and express them in a more "natural" form rather than merely providing direct recommendations. For example, suppose that a student majoring in history wants to search for books on "recommendation algorithms". If the system can account for the reader's major background attributes when recommending books and explain why a certain book is recommended from a professional perspective, the common points or connections it may have with the reader's major attributes, or what issues need to be noted while reading, could this improve the user's acceptance rate? Therefore, this optimization strategy based on personalized and interpretable recommendations for users is also a very interesting research direction.

In summary, this paper aims to explore the possibility of applying LLMs in the LIS field through empirical research and evaluate their effectiveness in the typical book recommendation scenario. We hope that this study can inspire researchers to analyze more opportunities for applying LLMs to similar tasks and further improve their performance in existing scenarios.

Author Contributions: Methodology, Z.L.; Software, X.Z.; Writing—original draft, Z.L.; Writing—review & editing, Y.C. and X.L. All authors have read and agreed to the published version of the manuscript.

Funding: This research was supported by the Fundamental Research Funds for the Central Universities, and the Research Funds of Renmin University of China: 23XNQT24.

Data Availability Statement: This data can be found in <https://github.com/zhiyulee-RUC/bookgpt>.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Zhao, W.X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. A Survey of Large Language Models. *arXiv* **2023**, arXiv:2303.18223.
2. OpenAI. Introducing ChatGPT. Available online: <https://openai.com/blog/chatgpt> (accessed on 10 May 2023).
3. Dreibelbis, E. ChatGPT Passes Google Coding Interview for Level 3 Engineer with \$183K Salary. Available online: <http://985.so/mny2k> (accessed on 25 May 2023).
4. Lund, B.D.; Wang, T. Chatting about ChatGPT: How may AI and GPT impact academia and libraries? *Libr. Hi Tech News* **2023**, *40*, 26–29. [CrossRef]
5. Cox, C.; Tzoc, E. ChatGPT: Implications for academic libraries. *Coll. Res. Libr. News* **2023**, *84*, 99. [CrossRef]
6. Verma, M. Novel Study on AI-Based Chatbot (ChatGPT) Impacts on the Traditional Library Management. *Int. J. Trend Sci. Res. Dev.* **2023**, *7*, 961–964.
7. Panda, S.; Kaur, N. Exploring the viability of ChatGPT as an alternative to traditional chatbot systems in library and information centers. *Libr. Hi Tech News* **2023**, *40*, 22–25. [CrossRef]
8. Kirtania, D.K.; Patra, S.K. OpenAI ChatGPT Generated Content and Similarity Index: A study of selected terms from the Library & Information Science (LIS). *Qeios* **2023**. [CrossRef]
9. Jelinek, F. *Statistical Methods for Speech Recognition*; MIT Press: Cambridge, MA, USA, 1998.
10. Rosenfeld, R. Two decades of statistical language modeling: Where do we go from here? *Proc. IEEE* **2000**, *88*, 1270–1278. [CrossRef]
11. Liu, X.; Croft, W.B. Statistical Language Modeling. *Annu. Rev. Inf. Sci. Technol.* **2004**, *39*, 1. [CrossRef]
12. Bengio, Y.; Ducharme, R.; Vincent, P. A neural probabilistic language model. In Proceedings of the Advances in Neural Information Processing Systems, Denver, CO, USA, 20 June 2000; Volume 13.
13. Mikolov, T.; Karafiát, M.; Burget, L.; Cernocký, J.; Khudanpur, S. Recurrent neural network based language model. In Proceedings of the Interspeech, Makuhari, Japan, 26–30 September 2010; Volume 2, pp. 1045–1048.
14. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient estimation of word representations in vector space. *arXiv* **2013**, arXiv:1301.3781.
15. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv* **2019**, arXiv:1810.04805v2.

16. Sarzynska-Wawer, J.; Wawer, A.; Pawlak, A.; Szymanowska, J.; Stefaniak, I.; Jarkiewicz, M.; Okruszek, L. Detecting formal thought disorder by deep contextualized word representations. *Psychiatry Res.* **2021**, *304*, 114135. [[CrossRef](#)] [[PubMed](#)]
17. Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. Language models are few-shot learners. *Adv. Neural Inf. Process. Syst.* **2020**, *33*, 1877–1901.
18. Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; Mishra, G.; Roberts, A.; Barham, P.; Chung, H.W.; Sutton, C.; Gehrmann, S.; et al. Palm: Scaling language modeling with pathways. *arXiv* **2022**, arXiv:2204.02311.
19. Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. Llama: Open and efficient foundation language models. *arXiv* **2023**, arXiv:2302.13971.
20. Zeng, A.; Liu, X.; Du, Z.; Wang, Z.; Lai, H.; Ding, M.; Yang, Z.; Xu, Y.; Zheng, W.; Xia, X.; et al. GLM-130B: An Open Bilingual Pre-trained Model. In Proceedings of the The Eleventh International Conference on Learning Representations (ICLR), Kigali, Rwanda, 1–5 May 2023.
21. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, Ł.; Polosukhin, I. Attention is all you need. In Proceedings of the Advances in Neural Information Processing Systems, Long Beach, CA, USA, 4–9 December 2017; Volume 30.
22. Saravia, E. Prompt Engineering Guide. 2022. Available online: <https://www.promptingguide.ai> (accessed on 15th May 2023).
23. Wei, J.; Tay, Y.; Bommasani, R.; Raffel, C.; Zoph, B.; Borgeaud, S.; Yogatama, D.; Bosma, M.; Zhou, D.; Metzler, D.; et al. Emergent abilities of large language models. *arXiv* **2022**, arXiv:2206.07682.
24. Jiang, H. A Latent Space Theory for Emergent Abilities in Large Language Models. *arXiv* **2023**, arXiv:2304.09960.
25. Lampinen, A.K.; Dasgupta, I.; Chan, S.C.; Matthewson, K.; Tessler, M.H.; Creswell, A.; McClelland, J.L.; Wang, J.X.; Hill, F. Can language models learn from explanations in context? *arXiv* **2022**, arXiv:2204.02329.
26. Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Chi, E.; Le, Q.; Zhou, D. Chain of thought prompting elicits reasoning in large language models. *Adv. Neural Inf. Process. Syst.* **2022**, *35*, 24824–24837.
27. George, A.S.; George, A.H. A review of ChatGPT AI’s impact on several business sectors. *Partners Univers. Int. Innov. J.* **2023**, *1*, 9–23.
28. Lu, Q.; Qiu, B.; Ding, L.; Xie, L.; Tao, D. Error Analysis Prompting Enables Human-Like Translation Evaluation in Large Language Models: A Case Study on ChatGPT. *arXiv* **2023**, arXiv:2303.13809.
29. Shafeeg, A.; Shazhaev, I.; Mihaylov, D.; Tularov, A.; Shazhaev, I. Voice Assistant Integrated with Chat GPT. *Indones. J. Comput. Sci.* **2023**, *12*. [[CrossRef](#)]
30. Tewari, A.S.; Priyanka, K. Book recommendation system based on collaborative filtering and association rule mining for college students. In Proceedings of the 2014 International Conference on Contemporary Computing and Informatics (IC3I), Mysore, India, 27–29 November 2014; pp. 135–138.
31. Bellogín, A.; Castells, P.; Cantador, I. Neighbor selection and weighting in user-based collaborative filtering: A performance prediction approach. *ACM Trans. Web (TWEB)* **2014**, *8*, 1–30. [[CrossRef](#)]
32. Linden, G.; Smith, B.; York, J. Amazon. com recommendations: Item-to-item collaborative filtering. *IEEE Internet Comput.* **2003**, *7*, 76–80. [[CrossRef](#)]
33. Rajpurkar, S.; Bhatt, D.; Malhotra, P.; Rajpurkar, M.; Bhatt, M. Book recommendation system. *Int. J. Innov. Res. Sci. Technol.* **2015**, *1*, 314–316.
34. Xu, C. A novel recommendation method based on social network using matrix factorization technique. *Inf. Process. Manag.* **2018**, *54*, 463–474. [[CrossRef](#)]
35. Dien, T.T.; Thanh-Hai, N.; Thai-Nghe, N. An approach for learning resource recommendation using deep matrix factorization. *J. Inf. Telecommun.* **2022**, *6*, 381–398. [[CrossRef](#)]
36. Lipton, Z.C.; Berkowitz, J.; Elkan, C. A critical review of recurrent neural networks for sequence learning. *arXiv* **2015**, arXiv:1506.00019.
37. Wu, S.; Sun, F.; Zhang, W.; Xie, X.; Cui, B. Graph neural networks in recommender systems: A survey. *ACM Comput. Surv.* **2022**, *55*, 1–37. [[CrossRef](#)]
38. Ma, C.; Ma, L.; Zhang, Y.; Sun, J.; Liu, X.; Coates, M. Memory augmented graph neural networks for sequential recommendation. In Proceedings of the AAAI Conference on Artificial Intelligence, New York, NY, USA, 7–12 February 2020; Volume 34, pp. 5045–5052.
39. Wang, X.; Huang, T.; Wang, D.; Yuan, Y.; Liu, Z.; He, X.; Chua, T.S. Learning intents behind interactions with knowledge graph for recommendation. In Proceedings of the Web Conference 2021, Ljubljana, Slovenia, 19–23 April 2021; pp. 878–887.
40. Liu, V.; Chilton, L.B. Design guidelines for prompt engineering text-to-image generative models. In Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, New Orleans, LA, USA, 30 April–5 May 2022; pp. 1–23.
41. Baidu. Wenxin Yiyao. Available online: <https://yiyao.baidu.com/welcome> (accessed on 15th May 2023).
42. Wang, Y.; Wang, L.; Li, Y.; He, D.; Liu, T.Y. A theoretical analysis of NDCG type ranking measures. In Proceedings of the Conference on Learning Theory, PMLR, Princeton, NJ, USA, 12–14 June 2013; pp. 25–54.
43. Zajac, Z. Goodbooks-10k: A new dataset for book recommendations. FastML. 2017. Available online: <http://fastml.com/goodbooks-10k> (accessed on 10th May 2023).
44. Mnih, A.; Salakhutdinov, R.R. Probabilistic matrix factorization. In Proceedings of the Advances in Neural Information Processing Systems, Vancouver, BC, Canada, 3–6 December 2007; Volume 20.

45. Koren, Y. Factor in the neighbors: Scalable and accurate collaborative filtering. *ACM Trans. Knowl. Discov. Data (TKDD)* **2010**, *4*, 1–24. [[CrossRef](#)]
46. Lemire, D.; Maclachlan, A. Slope one predictors for online rating-based collaborative filtering. In Proceedings of the 2005 SIAM International Conference on Data Mining, SIAM, Beach, CA, USA, 21–23 April 2005; pp. 471–475.
47. George, T.; Merugu, S. A scalable collaborative filtering framework based on co-clustering. In Proceedings of the Fifth IEEE International Conference on Data Mining (ICDM'05), Houston, TX, USA, 27–30 November 2005; p. 4.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.