

Evaluating ChatGPT as a Recommender System: A Rigorous Approach

Dario Di Palma^{1*}, Giovanni Maria Biancofiore^{1*},
Vito Walter Anelli^{1*}, Fedelucio Narducci¹, Tommaso Di Noia¹,
Eugenio Di Sciascio¹

¹Department of Electrical and Information Engineering Polytechnic
University of Bari, Via Edoardo Orabona 4, Bari, 70126, Italy.

*Corresponding author(s). E-mail(s): name.surname@poliba.it;

Abstract

In recent times, the popularity of large language models in the field of artificial intelligence has soared. These models demonstrate impressive abilities to comprehend and respond effectively to natural language requests, significantly contributing to various natural language-related tasks. The prompt-based learning approach, the emerging strategy of repurposing pretrained models without the need for additional training, has played a crucial role in enabling the application of general-purpose language models to specific tasks with minimal resources. By leveraging this approach, the full potential of large language models is unlocked, leading to improved sentence generation precision and generalization. Consequently, research communities are deeply exploring the capabilities of large language models across various dedicated tasks, culminating in the highly acclaimed ChatGPT.

Despite extensive research on large language models, their potential in the context of a recommendation scenario remains relatively underexplored. This study seeks to address this gap by investigating ChatGPT's capabilities as a zero-shot recommender system. Specifically, we aim to evaluate its ability to utilize user preferences to provide helpful recommendations, rerank existing recommendation lists, leverage information from similar users, and effectively handle cold start situations. To assess ChatGPT's performance, we conduct a comprehensive experimental evaluation using three datasets (MovieLens Small, Last.FM, and Facebook Book).

In this evaluation, we compare ChatGPT's performance against standard recommendation algorithms, representing the current state-of-the-art in the field. Furthermore, we compare ChatGPT's performance with other Large Language

Models, including GPT-3.5 and PaLM-2, for the recommendation task. To measure the effectiveness of the recommendations, we employ widely-used evaluation metrics, such as Mean Average Precision (MAP), Recall, Precision, F1, normalized Discounted Cumulative Gain (nDCG), Item Coverage, Expected Popularity Complement (EPC), Average Coverage of Long Tail (ACLT), Average Recommendation Popularity (ARP), and Popularity-based Ranking-based Equal Opportunity (PopREO).

Through a meticulous exploration of ChatGPT’s abilities in the scenario of recommender systems, our study seeks to enrich the growing body of research concerning the versatility and potential applications of large language models. The code used for conducting the experiments in this study is publicly accessible*.

Keywords: ChatGPT, Recommender Systems, Evaluation, Zero-Shot

1 Introduction

The increasing growth of social network applications and digital platforms highlights the vital role of information sharing and retrieval in our everyday lives. With human and corporate activities generating a vast amount of data, particularly in text format, the web has become a treasure trove of valuable information. Needs, opinions, and knowledge are quickly and efficiently expressed through natural language (NL) sentences. To effectively manage and automatically analyze and understand the abundance of this textual content, natural language processing (NLP) algorithms are essential [1]. In such a way, automatic systems can interface with users to understand their intents and offer personalized services. A prime example of the applications of NLP is seen in information filtering systems. These systems aim to alleviate the well-known information overload problem impacting users’ digital experience: sifting through a massive amount of data to find valuable information [2]. By implementing NLP algorithms, these systems can efficiently assist users in finding relevant information and items in the vast sea of data available.

Recently, researchers have found that incorporating interactive systems to engage users during their research leads to more accurate outcomes [3]. This explains the great diffusion of conversational agents (CAs) such as Amazon Alexa, Google Assistant, Microsoft Cortana, and Apple Siri [4]. Language models (LMs) play a crucial role with CAs and have garnered significant attention. An LM is a method for estimating the probability distribution of linguistic units such as words, sentences, and entire documents [5]. The natural progression of LMs led to the development of large language models (LLMs), which are pre-trained on vast unlabelled corpora. These LLMs exhibit remarkable adaptation capabilities in various downstream tasks [6].

Numerous CAs based on LLMs have emerged, aiming for better performance, wider scope, and reduced harm. Notable examples include BARD¹, Vicuna [7], and Alpaca [8], each implementing unique features like wider accessibility for Alpaca and

*<https://github.com/sisinflab/Recommender-ChatGPT>

¹<https://bard.google.com/>

cost-effectiveness for Vicuna. In this context, ChatGPT² gained particular attention. ChatGPT is a conversational agent fine-tuned from the GPT-3.5 pre-trained generative transformer [6] and continuously updated through reinforcement learning with human feedback [9]. With its robust architecture and vast knowledge from extensive training, ChatGPT delivers highly relevant and informative answers to users engaging with it, enriched with convincing explanations and supporting facts. As a result, researchers have taken a keen interest in exploring ChatGPT’s potential in various applications [10], tailoring it for diverse and specific tasks [11]. In particular, there has been a notable increase in interest from researchers in exploring its potential for the recommendation task [12, 13]. However, most of these studies have predominantly focused on ChatGPT’s fairness or have only built initial Recommender Systems (RSs) based on its LLM, rather than conducting a thorough and comprehensive evaluation of its effectiveness. Consequently, there is a need for further exploration and in-depth evaluations of ChatGPT’s capabilities as a recommender system to fully understand its potential and applicability in this specific scenario.

In this work, we promote a holistic analysis of ChatGPT’s ability to naturally act as an RS with a meticulous and reproducible investigation. We conduct a meticulous and replicable investigation by designing a well-structured experimental setup. Through this approach, we aim to identify and highlight the fundamental characteristics of an efficient RS. Since ChatGPT is a black-box model subject to continuous changes, this work reports empirical results on ChatGPT-3.5 evaluations on May 2023. In detail, we evaluate its ability to exploit user preferences in offering practical recommendations benefiting from similar users’ information.

Our main contributions are multifaceted. Firstly, we develop a rigorous pipeline of prompt-based experimental settings, ensuring a fair comparison and accurate positioning of ChatGPT alongside state-of-the-art RS baselines. Secondly, we unveil the natural capabilities of ChatGPT in recommending items to users based on their preferences, showcasing diverse and interesting behaviour across different domains, such as movies, music, and books.

In pursuit of these goals, we address the following research questions:

- **RQ1:** Is ChatGPT able to recommend items with a quality comparable to the state-of-the-art recommendation models? This question is further divided into sub-questions: **RQ1a:** How much is ChatGPT accurate compared with the state-of-the-art?, **RQ1b** How much is ChatGPT proposing diverse and novel recommendations compared with the state-of-the-art?, **RQ1c** How much is ChatGPT biased compared with the state-of-the-art?, **RQ1d** Which type of recommender system is ChatGPT more similar to?)
- **RQ2:** Is ChatGPT able to exploit user preferences to re-rank a recommendation list?
- **RQ3:** Does the substantial amount of knowledge utilized to train ChatGPT compensate for the absence of a complete user history in a cold-start scenario?

In our investigation, we deliberately avoid applying the Prompt Engineering approach. Instead, we design a unique prompt tailored to each experimental setup. This choice

²<https://openai.com/blog/chatgpt/>

allows us to isolate the effect of prompt engineering from recommendation performance, thereby establishing a lower bound for further investigations. We employ ChatGPT using a Zero-Shot approach. Zero-shot refers to the ability of an LLM to perform tasks different from the one learned during its training phase.

Our primary focus is to uncover the inherent capabilities of the standard ChatGPT³ as a recommender system. We recognize the potential impact of diverse prompts on the outcomes, and we acknowledge the necessity for a dedicated and targeted investigation, which exceeds the scope of our current study. Our primary objective is to comprehensively evaluate vanilla ChatGPT’s capabilities as a recommender system. Instead of presenting an optimized version with best performance, we focus on assessing its inherent abilities in this context.

We validate our study with an exhaustive analysis that engages three diverse datasets (i.e. MovieLens [15], Last.FM [16], and Facebook Book⁴) and a broad spectrum of baseline recommender algorithms and LLMs. To ensure reproducibility, we utilize the Elliot framework [17] to compute baseline performances, and we provide our code as an accessible GitHub repository⁵.

To the best of our knowledge, this study marks the first endeavour to analyze zero-shot ChatGPT in comparison to LLMs and specialized Recommender Systems without employing prompt engineering or in-context examples, making our investigation unique in its approach.

2 Related Work

Pretrained Foundation Models (PFMs) are powerful general models effectively studied and exploited in various fields such as natural language processing, computer vision, and graph learning [18]. PFMs use large amounts of data and can be fine-tuned in several downstream applications. ChatGPT is an excellent example of a PFM application and is fine-tuned from the generative pre-trained Transformer GPT-3.5 using Reinforcement Learning with Human Feedback approach [9, 19], which has become a promising way to align Large Language Models (LLMs) with a human’s intent [20].

Since re-training such models requires enormous computational resources, prompt learning [21] is a helpful technique for adapting pre-trained models without the cost of a fine-tuning procedure.

Accordingly, it emerged that the great potential of a pre-trained Transformer is to perform novel tasks for which it was not targeted during training [22] through prompts especially tailored. In detail, prompt learning relies on a suite of appropriate prompts, either hard text templates [6], or soft continuous embeddings [23], to reformulate the downstream tasks. Kojima et al. [24] give additional insights on the advantages that refined prompting approaches bring to the solution of such downstream tasks through LLMs. The recommendation is one of those downstream tasks, but the investigations

³With the term ‘standard (or vanilla) ChatGPT’, we refer to a model that operates without employing further artefacts that improve its performance. Such artefacts include prompt learning approaches, where prompts are a set of instructions that are learned to customise, enhance, or refine the capabilities of an LLM [14].

⁴<https://2015.eswc-conferences.org/program/semwebeval.html>

⁵<https://github.com/sisinflab/Recommender-ChatGPT>

are at a very early stage (this is the motivation because most cited works are in a pre-printing version).

Li et al. [22] propose two prompt learning approaches to exploit the rich knowledge contained in pre-trained language models (i.e., GPT-2) for recommendation explanation generation. Extensive experiments demonstrate the effectiveness of their approach in generating high-quality explanations as measured by text quality and explainability metrics.

GPT-2 is also leveraged in [25] for building a recommender system that uses prompts to reformulate the session-based recommendation task to a multi-token cloze task. The method is evaluated on a movie recommendation dataset in zero-shot and fine-tuned settings with limited training data. In the zero-shot setting, the Pretrained Language Model (PLM)-based method outperforms a random recommendation baseline, but under-performs traditional recommender systems such as GRU4Rec [26].

GPT-2 is also the basis of GPT4Rec [27], a flexible framework that generates hypothetical “search queries” given item titles in a user’s history and then retrieves items for recommendation by searching these queries. GPT4Rec combines GPT-2, to learn both item and user embeddings in the language space and the BM25 search engine to retrieve items for recommendation.

Similar to GPT-3, M6 [28] is an existing large-scale industrial pretrained language model on which M6-Rec [28] is based. The M6-Rec framework unifies various tasks in an industrial recommender system. It is able to perform retrieval, ranking, zero-shot recommendation, explanation generation, personalised content creation, and conversational recommendation by representing user behaviour data as plain texts and converting the tasks to either language understanding or generation. The authors verify the M6-Rec’s ability to perform zero-shot ranking on three datasets of different domains, and they demonstrate that it can match the performance of a traditional ID-based ranker trained on a million samples.

Another unified PLM-based framework which integrates the item recommendation into the generation process is RecInDial [29]. RecInDial finetunes the powerful PLMs like DialoGPT [30] together with a Relational Graph Convolutional Network (RGCN) to encode the node representation of an item-oriented knowledge graph. Besides, the authors design a vocabulary pointer mechanism to unify the response generation and item recommendation into the existing PLMs. Extensive experiments on the Conversational Recommender System (CRS) benchmark dataset REDIAL show that RecInDial significantly outperforms the state-of-the-art methods.

Wang and Lim [31] propose a strategy that incorporates a 3-step prompting that guides GPT-3 to carry subtasks that capture the user’s preferences, select representative previously watched movies, and recommend a ranked list of 10 movies. The proposed approach is evaluated on MovieLens 100K dataset, and it shows strong zero-shot performance, even outperforming some strong sequential recommendation models trained on the entire training dataset.

P5 (Pretrain, Personalized Prompt, and Predict Paradigm) [32], proposed by Geng et al., demonstrates that it is possible to learn multiple recommendation-related tasks by formulating these problems as prompt-based natural language tasks, where

user-item information and corresponding features are integrated with personalised prompt templates as model inputs. P5 is able to make predictions in a zero-shot or few-shot manner and largely reduces the necessity for extensive fine-tuning. In the same vein, Personalized prompt-based recommendation (PPR) [33] implements a prompt generator based on user profiles for cold-start recommendations. Evaluations on three large-scale datasets in few-shot and zero-shot scenarios confirm that PPR significantly improves both scenarios.

Gao et al. [34] are the first to use ChatGPT to improve recommendations given by a traditional recommender system. They exploit the conversations with ChatGPT to inject the users' preferences in order to refine recommendations generated by an existing recommender system. They do not investigate the main capabilities of vanilla ChatGPT and its behaviour in a recommendation scenario against sota RSs.

As demonstrated by the highly recent literature, there is great excitement in the research community to investigate the potential of PFM and LLM to support recommendation tasks. However, to the best of our knowledge, this is the first extensive investigation comparing zero-shot prompt-based recommendations, precisely the capabilities of vanilla ChatGPT, to other LLMs (e.g. PaLM2 [35, 36] or GPT-3.5 [6]) and traditional content-based/collaborative filtering approaches with the twofold aim of assessing the most similar paradigm and the effectiveness in cold-start settings without making use of prompt engineering.

3 ChatGPT as a Recommender

Chat Generative Pretrained Transformer (ChatGPT) is an advanced generative pre-trained Large Language Model (LLM) developed by OpenAI. It is fine-tuned from the generative pre-trained transformer GPT-3.5 [6]. The development process involves three crucial steps: supervised model training dialogue, enhanced training optimization, and model fine-tuning. Through these steps, ChatGPT emulates human language learning and knowledge acquisition while also aligning with user intent. The model aims to encompass both explicit intentions, such as following instructions, as well as implicit intentions, such as maintaining truthfulness and avoiding bias or any other harmful behavior [20].

With its language understanding and text-generation capabilities, ChatGPT is trained on massive data to possess a broad knowledge. It can engage in continuous multi-round conversations based on contexts and has a particular writing ability to support various tasks such as art creation, technical transfer, office learning, and logical reasoning [37].

As an autoregressive generative language model, ChatGPT processes input sentences (questions/requests) along with the preceding exchange of messages, treated as a sequence of words. It computes the response by generating a series of terms, identifying the most probable ones to form a coherent reply. The output is constructed word by word, with each word selection depending on the context and the preceding words, ensuring a coherent and contextually relevant answer.

Formally, given the sequence of terms as input $\mathbf{x} = [x_1, x_2, \dots, x_n]$, the output $\mathbf{y} = [y_1, y_2, \dots, y_m]$ results from:

$$p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^m p(y_i|\mathbf{x}, y_{k<i})$$

where $k = 0$ identifies the prediction case of the first term of the answer y_1 , which takes support only from the input words \mathbf{x} [38]. Hence, the vast knowledge learned by ChatGPT during its training dictates the probability computed at inference time, making it highly efficient in solving the intended task.

Furthermore, ChatGPT is explicitly designed to excel in user-oriented tasks, leading us to presume that its vanilla Large Language Model (LLM) is naturally more inclined to address user-related assignments, such as item recommendation, compared to other state-of-the-art LLMs. However, it is crucial to recognize that the primary goal of ChatGPT is not solely to discern individual user preferences and provide tailored item recommendations. Instead, its broader objective lies in understanding and generating human-like text across a wide range of tasks and contexts. While ChatGPT may demonstrate proficiency in user-related tasks due to its user-centric design, it is not exclusively optimized for such specific functions. Its capabilities extend beyond user recommendations and encompass various language tasks, making it a versatile language model suitable for diverse applications.

Therefore, we adopt the prompt learning paradigm, which allows us not to specialise its vanilla model behaviour (i.e., fine-tuning) but yet to target the scenario we want to explore. Specifically, the prompt learning paradigm, also known as prompting, is a method of conditioning LLMs through well-designed sentences to reach goals different from the ones perceived during training [14].

Given a pre-trained and fixed LLM \mathcal{L} , a test set of input X and output Y to assess a specific task, the prompting approach transforms each sample of X through a template \mathcal{T} , which will guide \mathcal{L} in generating the predictions \hat{Y} close to Y . The template \mathcal{T} assumes the following form:

$$\mathcal{T} = \textit{prefix}[X]\textit{suffix}[\hat{Y}]$$

where \hat{Y} is a slot later filled by the \mathcal{L} predictions, while *prefix* and *suffix* identify some text specifically designed to guide \mathcal{L} in solving the task of interest. This paradigm constitutes the foundations of Prompt Engineering, which targets the prompts' search problem in optimising the downstream tasks where a given LLM is applied.

To this purpose, White et al. [39] defined a comprehensive catalogue of prompt patterns to enhance prompt engineering for ChatGPT, and we find the “*person pattern*” particularly fits our purposes (cf. Section 4). Nevertheless, modelling a prompt to grant ChatGPT to reach high performances in the recommendation task is out of the scope of this study.

In our experimental setup, we adopt a single, straightforward prompt to evaluate the inherent capabilities of ChatGPT in performing recommendations⁶. By doing so, we avoid any biased influence that may arise from complex prompt designs. This configuration highlights the potential for future applications of ChatGPT as a Recommender System (RS) in subsequent analyses.

Furthermore, according to Adomavicius and Tuzhilin [40], the recommendation problem can be defined as the task of maximising a utility function. In this context, we extend their definition by considering the utility function u serving as a measure of how useful an item s is to a recommendation context c , represented as $u : C \times S \rightarrow R$, where R is a totally ordered set. Here, the recommendation context c is intended as an information collection composed of the user profile, i.e. the list of her preferred items S_{user} , and those recommended by the system at recommendation time S_{rec} , s.t. $C = \{S_{\text{user}}, S_{\text{rec}}\}$. The primary objective is to select an item $s' \in S$ for each recommendation context $c \in C$ that maximises their utility. More formally:

$$\forall c \in C, s'_c = \arg \max_{s \in S} u(c, s).$$

The item thus recommended will increase the set of suggested items S_{rec} , hence expanding the recommendation context c , to calculate the subsequent suggestions.

In contrast, a language model estimates the probability distribution of a sentence, which is represented as a sequence of words, symbols, tokens, or token sequences. This estimation involves modeling the probability of the next word given the preceding words [41].

Formally, considering a text sequence of length T with tokens x_1, x_2, \dots, x_T forming a sentence. The language model aims to compute the probability $P(x_1, x_2, \dots, x_T)$ of observing all these tokens in the given order, namely estimating the entire sequence's joint probability. Applying the chain rule, this probability is formally defined as:

$$P(x_1, x_2, \dots, x_T) = \prod_{t=1}^T P(x_t | x_1, \dots, x_{t-1})$$

Specifically, an ideal language model would possess the ability to independently generate natural text by sequentially selecting individual tokens [42]:

$$x_t \sim P(x_t | x_{t-1}, \dots, x_1)$$

The intriguing parallel between modelling a Recommender System (RS) and a generative language system becomes evident. Both systems learn to predict the most probable next item (in RS) or word (in language models) based on the preceding context of items (RS) or words (language models) they encounter. This intriguing similarity, coupled with ChatGPT's training phase dedicated to comprehending and addressing human intents and requests, has been a driving force behind our decision to investigate ChatGPT's potential as a Recommender System.

⁶As previously stated in the introduction, creating an optimized prompt to maximize ChatGPT's performance in the recommendation task falls outside the scope of this study.

To reach this goal, here we outline a prompt examples we submit to ChatGPT:

1. *prefix* = “Given a user, as a recommender system, provide recommendations”,
X = “The user {*user_id*} likes the following items: {*item_list*”,
suffix = “Give me back 50 recommendations”.
2. *prefix* = “Given a user, as a recommender system, provide recommendations”,
X = “The user {*user_id*} likes the following items: {*item_list*”,
suffix = “Re-rank me the following list: {*list_to_be_re-ranked*”.

It is important to mention that the items in our experiment can encompass movies, books, or songs, depending on the specific dataset being used. However, we deliberately avoid explicitly indicating the recommendation domain in the prompt. The *item_list* is individually computed for each user present in the dataset. Please refer to the forthcoming section for more comprehensive insights and detailed information.

4 Experimental Setting

The primary objective of this study is to assess the core capabilities of the vanilla ChatGPT model as a recommender system. We design four experimental setups to achieve this goal, defining crucial elements to enable ChatGPT to play as a recommender system. We can effectively address the research questions guiding this analysis by comparing the outcomes with state-of-the-art baselines across different application domains.

In each scenario, we utilize the “*person pattern*” prompt as introduced by White et al. [39] to guide ChatGPT in performing recommendations. This pattern is designed to provide the Large Language Model (LLM) with a “*persona*” to aid in selecting the types of output to generate and the specific details to focus on. Referring to the previous prompt formalism, an example of a persona prompt might be *prefix* = “Act like a *P*”, *suffix* = “give me the results the *P* will return”, and *X* will be filled with the data on which the LLM should operate, simulating the persona *P*.

Practically, one complete example of a prompt employed in our experiments is: “Given a user, as a recommender system, provide recommendations. The user 1 likes the following movies: Mad Max: Fury Road (2015) 5/5, Whiplash (2014) 4/5, etc. Give me back 50 recommendations.”.

This prompt formulation aligns with the central objective of this study, which is to explore the capabilities of vanilla ChatGPT in emulating recommender systems. Our focus is not to identify the optimal ChatGPT configuration for the recommendation task. Instead, we craft a prompt that includes only essential information required to guide ChatGPT in tackling this task. As a result, the outcome of this experiment defines a lower bound, laying the foundation for future research to concentrate on the potential development of recommender systems based on ChatGPT.

Prompts enriched with more data, such as information about the type of recommender system to emulate, domain-specific details, etc. might yield significantly diverse results and performances. These variations necessitate dedicated investigations to effectively conduct prompt engineering, but this requires exploring prompts in the

latent space of the Large Language Model (LLM) and having direct access to the ChatGPT model is essential. However, currently, such access is limited to API calls, which poses a constraint in the prompt engineering process.

We have led the following experiments through the OpenAI ChatGPT3.5-turbo API⁷, setting the temperature parameter to zero to grant our work’s reproducibility (i.e. ChatGPT generates the same answers). Also, we adapt the persona pattern prompt to fit the tokens limit imposed by the exploited API (4096 tokens per message exchanged, including the ones that compose the ChatGPT response).

In ChatGPT, the training set was used only for generating the prompt for a given user, while the test set was to compute the metrics. Accordingly, we did not transfer to ChatGPT any other information, and we supposed that it acquired enough knowledge to generate recommendations during its initial training. An analogous approach was employed with the other LLMs under consideration. In contrast, the recommender baselines were trained using a conventional methodology, leveraging training and test sets through the recommendation framework Elliot [17].

In order to establish a correspondence between the recommendations produced by ChatGPT and the items in the test set, we employed the Damerau-Levenshtein distance algorithm (from the `difflib` Python library) on the titles. Remarkably, despite slight discrepancies between some generated titles and their corresponding item names in the dataset, all items were accurately identified, preventing any significant instances of hallucination. In other words, within the confined output of 50 elements, ChatGPT did not generate any hallucinated content, i.e., items that appear plausible but are incorrect.

We conduct our experiments through four configurations. In the first configuration, we quantitatively analyze the quality of recommendations provided by ChatGPT in an unrestricted scenario. Specifically, for each user on the selected datasets, we start a dialogue session with ChatGPT providing the set of items she liked in the past and requesting a list of the top 50 items it would recommend, sorted by relevance. In this setting, we refrain from influencing the output in any way through the prompt, allowing ChatGPT to autonomously comprehend the domain and offer its recommendations accordingly.

To enable ChatGPT to compute recommendations, we retrieve the names of each item, such as movie titles for MovieLens, authors’ names for Last.FM, and book titles for Facebook Book. This allows ChatGPT to access its knowledge for generating recommendations. As for the baselines, Collaborative Filtering (CF) RSs have information on users’ and items’ IDs, while Content-Based Filtering (CBF) RSs learn about the item content through their genres.

By comparing the metrics, the first configuration helps us understand ChatGPT’s ability to recommend items and determine which type of data it leverages, whether it’s collaborative, popularity-based, or content-based. Additionally, this analysis allows us to ascertain that ChatGPT does not generate hallucinated content while shedding light on which type of RS it bears greater similarity to.

The second and third configurations of the experiments focus on assessing ChatGPT’s proficiency in re-ranking a list of items using the user profile. In the second

⁷<https://platform.openai.com/docs/guides/chat>

configuration, the list of items is pre-determined and comprises the most popular items in the dataset. On the other hand, in the third configuration, the list of items is generated based on the preferences of each user’s nearest neighbours. The objective of these two configurations is to determine whether ChatGPT can proficiently utilize the user preferences to accurately re-rank the list of recommendations.

The fourth and final configuration of the experiments is designed to specifically evaluate ChatGPT’s performance with users in a cold start scenario. For this experiment, we identify users with the smallest number of interactions from each dataset and utilize them to assess ChatGPT’s performance compared with the baseline.

The derived considerations are presented and discussed in Section 5.

4.1 Dataset

To provide answers to our research questions, we used three state-of-the-art datasets, each one belonging to a different domain (Music, Books, and Movies). Their statistics are reported in Table 1. Hereby is a brief description of each dataset:

- The *MovieLens* dataset is widely exploited in the RS community [15]. Different versions are available online⁸, but the one used in our study was collected from the MovieLens website and contains ratings for movies on a 1-5 scale.
- The *Facebook Books* dataset has been released for the Linked Open Data challenge co-located with *ESWC 2015*⁹, and it refers to the book domain. Only implicit feedback is available here, but for each item, we exploit the items-feature mappings available at this link¹⁰ to retrieve the data about the book titles, genres, and authors.
- The *Last.FM* dataset corresponds to user-artist plays on Last.fm online music system released during *HETRec2011*¹¹ *Workshop* [16]. It contains social networking, tagging, artists, and music-listening information from a set of 2,000 users. As for Facebook Books, we fetch titles and genre information through the aforementioned mapping.

It is essential to note that due to the token limit of the ChatGPT API (i.e., 4,096 tokens), certain preprocessing steps were required for the MovieLens data. The objective was to handle users with interaction histories exceeding this limit effectively. Consequently, users with interaction histories surpassing a specific threshold (i.e., 230 interactions) were excluded from the dataset, resulting in the filtered MovieLens dataset, denoted as †. In the second and third experimental configurations, this threshold was further reduced (i.e., 200 interactions), leading to another MovieLens dataset, indicated as ‡. This adjustment was crucial to have a prompt for the entire list of items for re-ranking purposes.

⁸<https://grouplens.org/datasets/movielens/>

⁹<https://2015.eswc-conferences.org/program/semwebeval.html>

¹⁰<https://github.com/sisinflab/LinkedDatasets/>

¹¹<https://grouplens.org/datasets/hetrec-2011/>

Table 1: A comparative analysis of dataset characteristics before and after pre-processing, comprising interactions, user and item counts, dataset sparsity, and quantity of available content.

Dataset	Interaction	Users	Items	Sparsity	Interaction	Users	Items	Sparsity	Content	
	<i>before pre-processing</i>				<i>after pre-processing</i>				<i>type</i>	<i>features</i>
MovieLens	100836	610	9724	98.30%	51576	603	1862	95.41%	genre	20
MovieLens †	100836	610	9724	98.30%	44309	603	1862	96.05%	genre	20
MovieLens ‡	100836	610	9724	98.30%	42456	603	1861	96.22%	genre	20
LAST.FM	60872	1883	5280	99.39%	38733	1797	1104	98.05%	genre	9748
FB Books	18978	1398	2,933	99.53%	12496	1398	1979	99.55%	genre, author	1970

4.2 Metrics

To ensure a comprehensive evaluation that sheds light on ChatGPT’s behaviour in a recommendation scenario, we carefully select the following metrics implemented in the Elliot framework:

- *Accuracy metrics:* We utilize various metrics, including Hit Ratio (HR), Mean Average Precision (MAP), Recall, Precision, and F1, to robustly quantify the relevance of the recommended items. Furthermore, we employ the **normalized Discounted Cumulative Gain (nDCG)** metric to evaluate the quality of the ranking in the recommendation list. Higher scores in these metrics indicate better recommendations, while lower scores suggest the opposite.
- *Coverage and novelty metrics:* These metrics provide valuable insights into the selection of items from the catalogue and the novelty of the recommendations offered to the users. Specifically, **ItemCoverage** quantifies the extent to which items are recommended, encompassing the coverage of the entire catalogue (i.e., the fraction of all available items that can potentially be recommended). On the other hand, the **Gini** score assesses the distribution of items, shedding light on the diversity of recommendations. Additionally, the **Expected Popularity Complement (EPC)** metric measures the expected number of relevant recommended items that were not previously seen by the user, reflecting the system’s ability to introduce novelty in the recommendations. For all the metrics higher values indicate superior performance.
- *Bias metrics:* Through the use of these metrics, our goal is to uncover the extent to which the recommendations generated by the system are influenced by biases. The **Average Coverage of Long Tail (ACLT)** metric enables us to evaluate the exposure that long-tail items receive in the entire recommendation process, providing insights into what fraction of the long-tail items the recommender has successfully covered [43]. Higher values better is the recommendation. On the other hand, the **Average Recommendation Popularity (ARP)** metric assesses the average rating popularity of the recommended items across testing users [44], helping us understand the distribution of popularity among the recommendations. Furthermore, the **Popularity-based Ranking-based Equal Opportunity (PopREO)** metric measures the bias that items in one or more groups may face, particularly concerning their lower recommendation probabilities when users express interest in these items [45]. Moreover, lower values of ARP and PopREO signify less biased

recommendations. By utilizing these metrics, we aim to shed light on the presence and impact of biases in the recommendation process.

- *Similarity index*: These metrics are also selected to assess the degree of similarity between the recommendation lists produced by ChatGPT and the baselines. Specifically, the **Jaccard** index evaluates the intersection of two given sets, offering insights into the common items recommended by both ChatGPT and the baselines. On the other hand, the **Kendall** index examines the same intersection but also takes into account the item positions, providing a measure of how closely the rankings of recommended items align between ChatGPT and the baselines. By leveraging these metrics, we aim to understand how ChatGPT’s recommendation behaviour compares to various types of Recommender Systems. Greater values of these indices indicate a higher degree of similarity in the recommendations.

4.3 Baselines

This section provides an outline of the foundational baselines that we deem essential for positioning ChatGPT within the current state-of-the-art of RSs. For this purpose, we select Collaborative Filtering (CF) and Content-Based Filtering (CBF) RSs as baselines to cover all the possible behaviours that ChatGPT may assume. To ensure the reproducibility of our experiments, we employ the following state-of-the-art models provided within the Elliot framework:

- **Random**, a non-personalized algorithm that randomly recommends items according to a uniform distribution.
- **MostPopular**, a non-personalized algorithm that recommends the same items’ list to all the users. This list encloses items’ from the most to least popular. Since it is known that it shows good performance because of statistical biases in the data [46], it is an important baseline to compare against [47].
- **ItemKNN** [48] is an item-based implementation of the K-nearest neighbours algorithm that finds the K-nearest item neighbours based on a similarity function like the Cosine Vector Similarity [49] and Pearson Correlation [50]. The items in the neighbourhood are then used to predict a score for each user-item pair.
- **UserKNN** [51], a user-based implementation of the K-nearest neighbours algorithm, similar to ItemKNN.
- **RP3beta** [52] a simple graph-based method conceptually similar to the ItemKNN.
- **EASER** [53] a linear model which works like a shallow autoencoder.
- **AttributeItemKNN** [54] a model that represents each item as a vector of weights computed through a TF-IDF model.
- **AttributeUserKNN** [54] a model that represents the users as a vector of weights instead of items.
- **VSM** (Vector Space Model fed with Knowledge Graph information) [55], a simple Vector Space Model in which the relevance is computed as the normalized TF-IDF value for each pair (item, feature). Thus, the recommendation list is computed based on the user-item similarities.

To give a comprehensive picture of ChatGPT’s position in the state of the art, we also compare its performances with the following LLMs baselines, which solve the recommendation task in a way identical to the one designed for ChatGPT:

- GPT-3.5¹² denoted as text-davinci-003, a pure large language model created by OpenAI, distinct from ChatGPT due to not be chat-oriented.
- PaLM-2¹³, specifically text-bison@001, is a state-of-the-art language model developed by Google.

Other models were considered during our experimentation but eventually discarded:

- LLaMA [56], a large language model developed by Meta. In our study, we aimed to evaluate Llama models of diverse scales, including those with 7B and 13B parameters. However, during the testing phase, we encountered challenges in generating an ordered list of 50 recommendations, resulting in incomparable outcomes, as the recommendation list was not provided. Subsequently, we attempted to explore models with parameter sizes of 33B and 65B. Regrettably, our specific hardware infrastructure posed substantial limitations, making it unsuitable for effectively running these extensive-scale models. This issue underscores the critical challenges associated with hardware resources when dealing with such large-scale language models.
- Alpaca [8], a fine-tuned LLaMA 7B model developed by Stanford University. However, its current input capacity of 512 tokens proved insufficient to feed the model with the user history.

5 Results

In this section, we conduct an in-depth analysis of the experimental evaluation results to assess ChatGPT’s performance as a recommendation system (RS). The primary objective of this analysis is to ascertain whether ChatGPT’s recommended item lists exhibit a stronger alignment with collaborative filtering, content-based recommender, or a hybrid approach.

We carried out a series of experiments on three popular recommendation datasets: MovieLens, LAST.FM, and Facebook Books, considering accuracy, coverage, novelty and bias metrics. A complete description of the prompt, datasets, metrics, and baselines is detailed in Section 4.

The Section is organized around the three research questions defined in the Introduction. Specifically, we address the RQ1 in subsections 5.1, 5.2, 5.3, and 5.4, RQ2 in subsections 5.5, and RQ3 in subsection 5.6. The results are shown in Tables 2 - 10 comparing ChatGPT against the baselines.

Table 2: Experiment 1: A Comparative analysis of ChatGPT-3.5 with various baselines on nDCG, HR, and MAP with cutoff at 10, 20, and 50. The results are ordered by nDCG@10 with the best results highlighted in bold and the second-best result underlined. (A- is for Attribute).

MovieLens									
Model	cutoff 10			cutoff 20			cutoff 50		
	\uparrow nDCG@10	\uparrow HR@10	\uparrow MAP@10	\uparrow nDCG@20	\uparrow HR@20	\uparrow MAP@20	\uparrow nDCG@50	\uparrow HR@50	\uparrow MAP@50
UserKNN	0.32358	0.81711	0.32792	0.31996	0.87919	0.27625	0.35223	0.93624	0.2077
ItemKNN	<u>0.31702</u>	0.81711	0.31869	0.31409	0.87584	<u>0.27227</u>	0.34784	0.94463	0.20649
RP ³ β	0.3151	0.82215	0.31814	<u>0.31754</u>	0.89094	0.27204	0.35393	0.94463	<u>0.20681</u>
A-UserKNN	0.2494	0.73826	0.25819	0.24657	0.80369	0.21759	0.26944	0.87919	0.16451
EASE ^R	0.22447	0.68121	0.23369	0.22404	0.81208	0.19535	0.25142	0.90436	0.1493
MostPop	0.17172	0.6057	0.17735	0.16748	0.70302	0.14994	0.19229	0.8557	0.11676
ChatGPT-3.5	0.16927	0.58221	0.181	0.15198	0.66275	0.14202	0.14081	0.70638	0.09216
GPT-3.5	0.15374	0.58389	0.161	0.14542	0.67114	0.12917	0.14533	0.74161	0.08869
PaLM 2	0.14032	0.49123	0.14666	0.13376	0.54035	0.1171	0.13655	0.55263	0.07874
A-ItemKNN	0.0438	0.29866	0.04426	0.04613	0.4094	0.04064	0.06021	0.61913	0.03609
VSM	0.0248	0.16107	0.02483	0.02431	0.23993	0.02222	0.03107	0.41443	0.01956
Random	0.0146	0.11745	0.01468	0.01464	0.1896	0.01361	0.02044	0.34732	0.0126

LAST.FM									
Model	cutoff 10			cutoff 20			cutoff 50		
	\uparrow nDCG@10	\uparrow HR@10	\uparrow MAP@10	\uparrow nDCG@20	\uparrow HR@20	\uparrow MAP@20	\uparrow nDCG@50	\uparrow HR@50	\uparrow MAP@50
RP ³ β	0.254841	0.800779	0.245995	0.30932	0.90317	0.19756	0.37565	0.97496	0.13787
ItemKNN	0.254563	0.799666	0.245618	0.29996	0.89538	0.19482	0.36331	0.95715	0.13375
EASE ^R	0.246097	0.788536	0.241016	0.29226	0.8798	0.19093	0.35423	0.96105	0.1312
UserKNN	0.243746	0.764051	0.237955	0.28913	0.86811	0.18851	0.35096	0.95492	0.12921
A-UserKNN	0.243389	0.782972	0.241725	0.29139	0.86032	0.19169	0.3533	0.936	0.13193
VSM	0.217008	0.75626	0.210701	0.25972	0.85476	0.16735	0.31553	0.94769	0.11525
A-ItemKNN	0.208757	0.739566	0.19778	0.25262	0.84363	0.1604	0.31212	0.93267	0.1133
PaLM-2	0.207328	0.614252	0.200345	0.25065	0.64429	0.15921	0.31479	0.65607	0.11273
ChatGPT-3.5	0.205844	0.721758	0.199038	0.23263	0.79855	0.15075	0.24938	0.8453	0.09485
GPT-3.5	0.198919	0.703395	0.195589	0.22161	0.77017	0.14669	0.24436	0.81413	0.09258
MostPop	0.093312	0.426266	0.090605	0.11652	0.53534	0.07658	0.15742	0.74513	0.05724
Random	0.006189	0.055648	0.005231	0.00942	0.09905	0.0053	0.01708	0.22092	0.00515

Facebook Books									
Model	cutoff 10			cutoff 20			cutoff 50		
	\uparrow nDCG@10	\uparrow HR@10	\uparrow MAP@10	\uparrow nDCG@20	\uparrow HR@20	\uparrow MAP@20	\uparrow nDCG@50	\uparrow HR@50	\uparrow MAP@50
ChatGPT-3.5	0.05742	0.17267	0.02404	0.0714	0.25129	0.02001	0.08536	<u>0.31521</u>	0.01441
A-ItemKNN	<u>0.05034</u>	0.14254	<u>0.02241</u>	<u>0.05944</u>	0.19838	<u>0.01774</u>	0.07104	0.28655	0.01208
VSM	0.04592	<u>0.15136</u>	0.01828	0.05805	<u>0.22483</u>	0.01654	0.07117	0.31668	<u>0.01243</u>
GPT-3.5	0.04214	0.13226	0.0163	0.0556	0.19251	0.01448	<u>0.08189</u>	0.21969	0.01219
A-UserKNN	0.04196	0.13446	0.01815	0.05092	0.18957	0.01487	0.0668	0.3086	0.01097
PaLM-2	0.03456	0.09384	0.01417	0.04173	0.10672	0.01141	0.05903	0.11684	0.00869
RP ³ β	0.03097	0.09552	0.01326	0.03953	0.14842	0.01088	0.0501	0.23512	0.00814
UserKNN	0.03002	0.10213	0.01221	0.04016	0.16605	0.01084	0.054	0.26231	0.00862
ItemKNN	0.02873	0.08229	0.01221	0.03691	0.13299	0.01	0.0489	0.22116	0.00747
EASE ^R	0.02107	0.07568	0.00844	0.03181	0.14769	0.00802	0.05309	0.30639	0.00759
MostPop	0.00914	0.03233	0.00402	0.01112	0.04702	0.00336	0.03761	0.26672	0.00405
Random	0.00142	0.00588	0.00053	0.00275	0.01616	0.00062	0.00585	0.04262	0.00083

5.1 RQ1a: How much is ChatGPT accurate compared with the state-of-the-art?

MovieLens. Tables 2, and 3 reports the accuracy of the competing models. We can observe that the UserKNN and ItemKNN models exhibit superior accuracy performance, which aligns with recent research on reproducibility [57, 58]. We expected this outcome since the MovieLens dataset does not hold much movie content information

¹²<https://platform.openai.com/docs/models/gpt-3-5>

¹³<https://ai.google/discover/palm2/>

except for their genre, leading to lower performance for the Content-based Filtering (CBF) models.

Although ChatGPT-3.5 outperforms CBF and random models, it falls short of MostPop, primarily due to the high Popularity Bias present in the dataset, and it does not surpass the collaborative filtering models. However, ChatGPT outperforms traditional CBF models such as AttributeItemKNN (A-ItemKNN) and VSM, highlighting that its performance are not at random and could be an indication of its extensive learned knowledge of content information during the training phase.

While ChatGPT does not occupy the highest position in the ranking, its accuracy is within the same order of magnitude as the collaborative models, except for the top-3 UserKNN, ItemKNN, and $RP^3\beta$. Increasing the cutoff value to 20 and 50 gives us further insights into the models' behaviour. The results remain relatively consistent with the cutoff of 10, indicating that the recommendation list size does not particularly impact the top-performing models (UserKNN, ItemKNN, and $RP^3\beta$) and ChatGPT.

Surprisingly, with a cutoff set to 50, GPT-3.5 slightly outperforms ChatGPT in nDCG, HR, and the other accuracy metrics. This indicates that GPT-3.5 makes better use of the increasing context it obtains for each item in generating the recommended list. This result could be attributed to GPT-3.5's main objective of generating sentences with strong semantic consistency. Conversely, ChatGPT might compromise a bit on semantic consistency to achieve better conversational results prioritising the effectiveness of user interaction. However, ChatGPT still outperforms the other LLMs in terms of MAP. To conclude, PaLM performs the worst among the LLM, but better than CBF approaches. This may depend on its training data that considers movies, but maybe it could not include vast information on popular ones.

LAST.FM. Concerning the music domain, the results of the top-performing models, namely UserKNN, ItemKNN, and $RP^3\beta$, are aligned with those of the MovieLens Small dataset.

Regarding the performance of recommendation systems using content based on genres, notably, the available content for this particular dataset comprises a significantly larger magnitude with 9,748 genres, a substantial increase from the 20 genres found in the MovieLens dataset. Consequently, the Content-Based Filtering (CBF) models exhibit remarkable improvement compared to the MovieLens dataset, indicating the effectiveness of the approach. Furthermore, when comparing the performance of CBF models with ChatGPT, the former outperforms the latter, benefitting from its inherent recommendation-oriented nature.

Due to their performance similarity, defining distinct groups that separate model performance is intricate. However, if such groups exist, it is likely that $RP^3\beta$, ItemKNN, $EASE^R$, and UserKNN would be included in the same group, as they exhibit the highest comparable performance. Surprisingly, the same ranking holds for all cutoff values, except for $EASE^R$ which slightly outperforms ItemKNN at cutoff 20 for the HR.

Such a phenomenon may depend on the fewer items of the Last.FM dataset (1,104 artists compared to MovieLens' 1,862 movies), making the similarity computation

more challenging. ChatGPT and the other LLMs show relatively low performance, surpassing only MostPop and Random models.

The poor results on MostPop suggest that, in the Last.FM dataset, the popularity bias is not so evident, while ChatGPT’s underwhelming performance may have several underlying causes. The main cause might stem from the scarcity of music-related data or inadequate music review information in its training corpus. Nevertheless, the internal mechanisms of ChatGPT remain somewhat enigmatic, leaving room for speculation.

It is possible that the issue lies in the complexity of the recommendation task, and ChatGPT may not effectively utilise most of the content available to improve its performance. It is worth noting that all the LLMs demonstrate similar performance, which is still comparable to the state-of-the-art baselines. In particular, PaLM surpasses on average ChatGPT except for the HR metric, showing its ability to operate well in the music domain.

Facebook Books.

It is noteworthy that ChatGPT consistently outperforms other models for cutoff values of 10, 20, and in most cases, 50-item recommendations. ChatGPT not only outperforms but also significantly surpasses the other large language models by a considerable margin. This performance cannot be attributed to popularity, as MostPop ranks last as stated above. In fact, we observe that the MostPop model perform poorly, indicating a limited popularity bias in this dataset.

However, although also the other LLMs can take advantage of the large textual-content availability, the results demonstrate how they are less efficient than ChatGPT at exploiting it.

Furthermore, it is worth mentioning that incorporating content from curated sources such as Linked Data enhances the performance of content-based models significantly [55, 59, 60]. As a reflection, content-based models (A-ItemKNN, VSM, and A-UserKNN) outperform collaborative models ($RP^3\beta$, UserKNN, ItemKNN, and $EASE^R$), which could also be attributed to the limited and less significant collaborative information in the dataset. This result may be due to the datasets’ small number of users (1,398), which leads to a scarcity of identifiable patterns in the collaborative data. Some results support this intuition. Specifically, ItemKNN achieves an nDCG@10 of 0.317 and 0.254 in MovieLens and LAST.FM datasets, while it only reaches an nDCG@10 of 0.02873 in the Facebook Books dataset.

Conclusion. The answer to the research question “How much is ChatGPT accurate compared with the state-of-the-art?” is: ChatGPT solves the recommendation task remarkably accurately. Even in its vanilla version, thus without applying techniques to enhance its performance, ChatGPT reaches results comparable with the RSs state-of-the-art and distinguishes itself from other LLMs. In such a way, our empirical results open up new avenues of thinking and suggest that “user-oriented” LLMs like ChatGPT have the potential to revolutionise the recommendation scenario.

Table 3: Experiment 1: A Comparative analysis of ChatGPT-3.5 with various base-lines on Recall, Precision, and F1 with cutoff at 10, 20, and 50. The results are ordered by nDCG@10 with the best results highlighted in bold and the second-best result underlined. (A- is for Attribute).

MovieLens									
Model	cutoff 10			cutoff 20			cutoff 50		
	\uparrow Recall@10	\uparrow Precision@10	\uparrow F1@10	\uparrow Recall@20	\uparrow Precision@20	\uparrow F1@20	\uparrow Recall@50	\uparrow Precision@50	\uparrow F1@50
UserKNN	<u>0.18093</u>	0.26191	0.17867	0.26636	0.20201	0.18920	0.41324	0.13564	0.17283
ItemKNN	0.17844	<u>0.26174</u>	0.17692	0.26550	<u>0.20277</u>	<u>0.18993</u>	<u>0.41637</u>	0.13725	<u>0.17468</u>
RP ³ β	0.18096	<u>0.25872</u>	0.17877	0.27197	0.20294	0.19192	0.42799	<u>0.13648</u>	0.17540
A-UserKNN	0.12988	0.20302	0.13259	0.19473	0.16065	0.14451	0.31051	0.10909	0.13630
EASE [#]	0.11255	0.17836	0.11407	0.17673	0.14354	0.12841	0.29610	0.10191	0.12711
MostPop	0.08312	0.14094	0.08652	0.12777	0.11099	0.09609	0.23116	0.08272	0.10139
ChatGPT-3.5	0.08146	0.12970	0.08256	0.10772	0.08716	0.07839	0.13310	0.04084	0.05276
GPT-3.5	0.08248	0.11980	0.08147	0.11523	0.08565	0.08097	0.15627	0.04413	0.05875
PaLM-2	0.07665	0.10596	0.07522	0.10816	0.07509	0.07406	0.14961	0.03849	0.05348
A-ItemKNN	0.02268	0.03977	0.02398	0.03786	0.03557	0.02912	0.08390	0.03084	0.03705
VSM	0.01076	0.02198	0.01175	0.01702	0.01846	0.01417	0.03873	0.01708	0.01925
Random	0.00622	0.01393	0.00747	0.01150	0.01200	0.00960	0.02752	0.01208	0.01382

LAST.FM									
Model	cutoff 10			cutoff 20			cutoff 50		
	\uparrow Recall@10	\uparrow Precision@10	\uparrow F1@10	\uparrow Recall@20	\uparrow Precision@20	\uparrow F1@20	\uparrow Recall@50	\uparrow Precision@50	\uparrow F1@50
RP ³ β	0.24171	0.17969	0.20152	0.35299	0.13122	0.18814	0.52430	0.07817	0.13487
ItemKNN	<u>0.24045</u>	<u>0.17869</u>	<u>0.20034</u>	<u>0.33300</u>	<u>0.12379</u>	<u>0.17745</u>	<u>0.49627</u>	<u>0.07392</u>	<u>0.12755</u>
EASE [#]	0.22916	0.17368	0.19359	0.32317	0.12209	0.17446	0.48339	0.07279	0.12545
UserKNN	0.22619	0.17101	0.19069	0.31833	0.12006	0.17158	0.47879	0.07198	0.12406
A-UserKNN	0.22283	0.17173	0.19043	0.31880	0.12315	0.17525	0.47873	0.07314	0.12591
VSM	0.20207	0.15109	0.16911	0.28940	0.10785	0.15451	0.43500	0.06462	0.11150
A-ItemKNN	0.20141	0.14825	0.16679	0.29056	0.10735	0.15410	0.44441	0.06587	0.11369
PaLM-2	0.19223	0.14234	0.15998	0.28058	0.10418	0.14947	0.44421	0.06677	0.11512
ChatGPT-3.5	0.18529	0.13105	0.14978	0.23989	0.08561	0.12394	0.28202	0.04030	0.06985
GPT-3.5	0.17528	0.12794	0.14465	0.22189	0.08088	0.11657	0.27947	0.04078	0.07053
MostPop	0.08951	0.07401	0.07993	0.13528	0.05545	0.07782	0.24058	0.03855	0.06603
Random	0.00727	0.00556	0.00616	0.01401	0.00520	0.00744	0.03443	0.00502	0.00868

Facebook Books									
Model	cutoff 10			cutoff 20			cutoff 50		
	\uparrow Recall@10	\uparrow Precision@10	\uparrow F1@10	\uparrow Recall@20	\uparrow Precision@20	\uparrow F1@20	\uparrow Recall@50	\uparrow Precision@50	\uparrow F1@50
ChatGPT-3.5	0.08628	0.01837	0.02958	0.13031	0.01403	0.02497	<u>0.18504</u>	<u>0.00802</u>	<u>0.01527</u>
A-ItemKNN	0.07166	0.01594	0.02551	0.09993	0.01128	0.01999	0.14658	0.00660	0.01255
VSM	<u>0.07834</u>	<u>0.01712</u>	<u>0.02745</u>	<u>0.11510</u>	<u>0.01308</u>	<u>0.02318</u>	0.16771	0.00754	0.01434
GPT-3.5	0.06597	0.01440	0.02310	0.10967	0.01190	0.02116	0.22024	0.00891	0.01699
A-UserKNN	0.06264	0.01411	0.02254	0.09053	0.01040	0.01842	0.15381	0.00710	0.01349
PaLM-2	0.05172	0.01040	0.01680	0.07554	0.00750	0.01339	0.14775	0.00578	0.01103
RP ³ β	0.04695	0.00985	0.01593	0.07466	0.00790	0.01409	0.11620	0.00508	0.00968
UserKNN	0.04897	0.01051	0.01700	0.08110	0.00893	0.01589	0.13737	0.00607	0.01155
ItemKNN	0.04079	0.00860	0.01389	0.06718	0.00698	0.01247	0.11626	0.00489	0.00933
EASE [#]	0.03483	0.00794	0.01267	0.06890	0.00794	0.01405	0.15528	0.00710	0.01349
MostPop	0.01367	0.00331	0.00520	0.01947	0.00250	0.00437	0.12679	0.00605	0.01149
Random	0.00257	0.00059	0.00093	0.00664	0.00081	0.00142	0.01876	0.00087	0.00165

5.2 RQ1b: How much is ChatGPT proposing diverse and novel recommendations, compared with the state-of-the-art?

In this analysis, we aim to assess the extent of diversity and novelty in the recommendations generated by ChatGPT through the results collected in Table 4. This study, and consequently the Table, specifically examines i) the Item Coverage metric to quantify aggregate diversity by assessing the number of suggested items, ii) the Gini index to estimate the variation across the suggested lists, and iii) the EPC metric as an indicator of the novelty of the recommended items.

MovieLens. In this domain, the Random model obtains the highest value in terms of ItemCoverage, while the MostPop model performs poorly, aligning with their standard behaviours reported in the literature. Surprisingly, the EASE^R model exhibits a trend similar to MostPop, indicating an evident inefficiency when applied to datasets affected by popularity bias. The Attribute-based ItemKNN (A-ItemKNN) model emerges as the second-best performer, although its nDCG performance is underwhelming. Such an outcome derives from the model’s focus on recommending items that reflect the average content liked by the user rather than composing a highly relevant list (in this domain, popular items have a higher rank, while niche movies result in less pertinent recommendations).

In contrast, ItemKNN demonstrates an effective balance between leveraging diversity (ItemCoverage) and accuracy (nDCG), likely due to its reliance on collaborative information. However, ChatGPT’s performance remains mediocre despite its poor accuracy results, and this trend persists even when considering a cutoff of 20. Surprisingly, when the cutoff is raised to 50, all models, except ChatGPT, exhibit a significant coverage of items. This observation leads to an intriguing question: Could ChatGPT possess specialized knowledge about certain movies or reviews, making it less familiar with a broader range of movies beyond the 759 it is most familiar with? Notably, ChatGPT’s diversity is only half of what UserKNN achieves.

Furthermore, upon scrutinizing other large language models (LLMs) like PaLM-2 and GPT-3.5, it becomes apparent that they demonstrate comparable behaviour to ChatGPT. Consequently, the architectural characteristics of these models and their shared task of predicting the most probable token likely play a significant role in shaping this pattern. Such LLMs seem to prioritize a limited set of highly probable items, which could account for the similarity in their performance. This observation emphasizes the necessity of broadening our investigations to encompass a wider array of LLMs, aiming to understand their behaviour and potential limitations comprehensively.

LAST.FM. Similar outcomes are observed when examining the LAST.FM dataset. The Random model achieves the highest coverage, with AttributeItemKNN (A-ItemKNN) following closely behind. However, ItemKNN remains the best overall model. As for ChatGPT, it aligns with the top-performing models to some extent, but its coverage falls below expectations, reaching only 1,099 when evaluated with a cutoff of 50. Nevertheless, it appears that ChatGPT may have a higher upper limit for overall artists proposed compared to MovieLens, indicating a broader knowledge of the music domain and a less influential popularity bias.

Regarding the Gini index, a similar pattern is observed as in the previous dataset, although ChatGPT demonstrates relatively better performance. Additionally, ChatGPT exhibits remarkable efficiency in terms of EPC with a cutoff of 50, outperforming all other models except PaLM-2. This finding suggests that the one thousand items covered by ChatGPT may not necessarily align with the most popular ones, indicating its potential for promoting novelty.

It is worth noting that all LLM models demonstrate higher EPC values when evaluated with a cutoff of 50, emphasizing the significance of this particular cutoff in

Table 4: Experiment 1: A Diversity and Novelty Metrics comparative analysis of ChatGPT-3.5 with various baselines on ItemCoverage, Gini, and EPC with cutoff at 10, 20, and 50. The results are ordered by nDCG@10 with the best results highlighted in bold and the second-best result underlined. (A- is for Attribute).

MovieLens									
Model	cutoff 10			cutoff 20			cutoff 50		
	↑ItemCov@10	↑Gini@10	↑EPC@10	↑ItemCov@20	↑Gini@20	↑EPC@20	↑ItemCov@50	↑Gini@50	↑EPC@50
UserKNN	431	0.06005	<u>0.22631</u>	646	0.09096	<u>0.18938</u>	1055	0.15735	0.14056
ItemKNN	672	0.1025	0.23302	957	0.13176	0.19369	1427	0.19151	0.14349
RP ³ β	422	0.05512	0.22351	638	0.08305	0.18905	994	0.14218	<u>0.14078</u>
A-UserKNN	451	0.06181	0.1743	690	0.09503	0.14793	1166	0.16927	0.11061
EASE ^R	53	0.01054	0.14495	97	0.01876	0.12501	187	0.04176	0.09746
MostPop	40	0.0082	0.10873	64	0.01539	0.09289	124	0.03492	0.07523
ChatGPT-3.5	408	0.03708	0.12028	591	0.0624	0.09455	759	0.09774	0.08241
GPT-3.5	344	0.03725	0.10986	463	0.05127	0.08966	642	0.08087	0.0779
PaLM-2	277	0.02758	0.09883	385	0.03835	0.08343	472	0.0496	0.07695
A-ItemKNN	<u>1340</u>	<u>0.36863</u>	0.03798	<u>1654</u>	<u>0.42543</u>	0.03477	<u>1832</u>	<u>0.51702</u>	0.03048
VSM	385	0.05159	0.021	596	0.08068	0.01841	952	0.14054	0.01684
Random	1798	0.69364	0.01333	1859	0.78045	0.0118	1862	0.85914	0.01151

LAST.FM									
Model	cutoff 10			cutoff 20			cutoff 50		
	↑ItemCov@10	↑Gini@10	↑EPC@10	↑ItemCov@20	↑Gini@20	↑EPC@20	↑ItemCov@50	↑Gini@50	↑EPC@50
RP ³ β	761	0.11652	<u>0.18876</u>	1024	0.174	0.14953	1349	0.27484	0.1004
ItemKNN	1284	0.286471	0.192606	1438	0.34802	<u>0.14749</u>	<u>1504</u>	0.45097	0.09843
EASE ^R	567	0.079701	0.181923	824	0.12525	0.14097	1237	0.22145	0.09443
UserKNN	803	0.129114	0.182059	1096	0.18547	0.14057	1390	0.28436	0.09414
A-UserKNN	825	0.16048	0.184956	1134	0.23609	0.14442	1423	0.36406	0.09629
VSM	653	0.082577	0.159303	899	0.11602	0.12427	1210	0.17866	0.08338
A-ItemKNN	<u>1411</u>	<u>0.38553</u>	0.158457	<u>1487</u>	<u>0.45917</u>	0.12461	1507	<u>0.56704</u>	0.08471
PaLM-2	736	0.124381	0.161394	833	0.13252	0.1383	865	0.11606	0.12908
ChatGPT-3.5	833	0.151743	0.153509	1018	0.19144	0.11985	1099	0.22245	<u>0.10426</u>
GPT-3.5	654	0.106667	0.148682	771	0.12853	0.10962	895	0.1397	0.0854
MostPop	27	0.008103	0.066226	41	0.01568	0.05425	79	0.03702	0.04089
Random	1507	0.84102	0.004943	1507	0.88721	0.0048	1507	0.92653	0.00472

Facebook Books									
Model	cutoff 10			cutoff 20			cutoff 50		
	↑ItemCov@10	↑Gini@10	↑EPC@10	↑ItemCov@20	↑Gini@20	↑EPC@20	↑ItemCov@50	↑Gini@50	↑EPC@50
ChatGPT-3.5	689	0.0439	0.02102	876	0.05529	0.01707	1029	0.07221	0.01336
A-ItemKNN	1488	0.27354	<u>0.01917</u>	1732	0.33571	<u>0.01472</u>	1867	0.42203	0.00961
VSM	1382	0.24747	<u>0.01748</u>	1681	0.30591	0.0145	1856	0.39692	0.00968
GPT-3.5	505	0.02219	0.01548	631	0.02713	0.01329	725	0.02949	<u>0.0121</u>
A-UserKNN	809	0.07523	0.01586	1291	0.11023	0.01254	2010	0.19455	0.00904
PaLM-2	583	0.04371	0.01258	740	0.05053	0.01047	832	0.04926	0.01034
RP ³ β	1984	0.39997	0.01124	2150	0.46941	0.00933	<u>2222</u>	0.55403	0.00656
UserKNN	868	0.07124	0.01088	1377	0.10896	0.00957	<u>2067</u>	0.19848	0.00708
ItemKNN	<u>2139</u>	<u>0.52833</u>	0.01029	<u>2228</u>	<u>0.61147</u>	0.00854	2234	0.66477	0.00624
EASE ^R	252	0.01366	0.00757	478	0.02435	0.00757	1045	0.05554	0.00702
MostPop	17	0.0045	0.00311	28	0.00918	0.00257	61	0.0229	0.00506
Random	2230	0.77733	0.00054	2234	0.84123	0.0007	2234	0.89986	0.00081

encouraging novelty in the recommended items.

Facebook Books. In the domain of books, ChatGPT continues to exhibit outstanding performance, especially in terms of novelty, despite covering a relatively small number of books (1,029 with a cutoff of 50). Interestingly, other large language models (LLMs) also demonstrate a similar trend to ChatGPT, recommending novel and less popular items.

The exact reason behind ChatGPT’s exceptional performance in this context remains uncertain, giving rise to several plausible hypotheses. One possibility is that ChatGPT’s training data includes a substantial amount of book reviews, thereby

enhancing its capacity to provide diverse and innovative recommendations. Moreover, ChatGPT might effectively leverage collaborative information or thoroughly exploit the content of items to deliver well-targeted recommendations.

To validate these hypotheses and attain a deeper understanding of ChatGPT’s exceptional novelty in the Facebook Books dataset, subsequent experiments will thoroughly investigate and analyze ChatGPT’s performance.

Conclusion. In summary, the question: “How much is ChatGPT proposing diverse and novel recommendations, compared with the state-of-the-art?” has a context-dependent answer. Actually, the level of diversity and novelty exhibited by ChatGPT varies depending on the considered dataset and domain. In its vanilla version, ChatGPT tends to demonstrate lower diversity but higher novelty when dealing with Books, and it showcases good novelty in the realm of Music. These findings imply that ChatGPT and other LLMs have the potential to serve as effective RS with higher levels of novelty and satisfactory diversification.

5.3 RQ1c: How much is ChatGPT biased compared with the state-of-the-art?

Evaluating the presence of bias in ChatGPT and comparing its behaviour with other state-of-the-art models requires an analysis of several factors. Bias can manifest in various forms, including cultural, gender, racial, and ideological aspects, and it is essential to assess how each factor influences the recommendations. These biases can be ingrained in the training data, influenced by underlying algorithms, and reflected in the model’s ability to generate fair and unbiased recommendations.

Given the black-box nature of ChatGPT, which limits access to its internal workings, our evaluation focuses on examining the presence of biases in the recommendations it generates. This evaluation involves scrutinizing the recommendations for potential biases based on the identified factors. The results of this experiment are reported in Table 5.

MovieLens. As highlighted in the previous sections, the MovieLens dataset exhibits a significant popularity bias. To contextualize our discussion, we use MostPop as a reference point, given its ACLT value indicating the highest bias (0 value).

EASE^R demonstrates a biased trend similar to MostPop, which may be attributed to its simple yet effective architecture that proposes highly relevant recommendations. Although content-based models such as VSM and AttributeItemKNN exhibit different patterns with generally poor and unreliable performance in recommending relevant items. UserKNN reveals a high popularity bias trend, whereas ItemKNN is more biased on the long-tail items, further supporting previous findings on lower diversity but a higher novelty (refer to 5.2).

The evaluation of ChatGPT’s performance in the movies domain indicates that it outperforms most collaborative approaches (excluding ItemKNN) and other tested large language models (LLMs). However, when compared to content-based (CB) models, it does not achieve superior results. These observations point to the presence of

Table 5: Experiment 1: A Bias Metrics comparative analysis of ChatGPT-3.5 with various baselines on ACLT, ARP, and PopREO with cutoff at 10, 20, and 50. The results are ordered by nDCG@10 with the best results highlighted in bold and the second-best result underlined. (A- is for Attribute).

MovieLens									
Model	cutoff 10			cutoff 20			cutoff 50		
	↑ACLT@10	↓ARP@10	↓PopREO@10	↑ACLT@20	↓ARP@20	↓PopREO@20	↑ACLT@50	↓ARP@50	↓PopREO@50
UserKNN	0.05705	125.049	0.9958	0.22987	108.5662	0.98913	1.27685	86.85775	0.95677
ItemKNN	0.46477	103.1854	0.90628	1	95.56686	0.89775	3.25	82.11755	0.84785
RP ³ _β	0.21141	124.6554	0.96631	0.47819	109.3722	0.95455	1.54027	87.86866	0.91837
A-UserKNN	0.12248	125.7143	0.98918	0.36409	108.692	0.98295	1.73322	86.6751	0.97198
EASE ^R	0	175.4149	1	0	154.0077	1	0	121.6648	1
MostPop	0	182.5502	1	0	160.7216	1	0	127.0194	1
ChatGPT-3.5	0.41946	121.6523	0.92529	1.55201	99.77496	0.88917	3.98993	80.34624	0.84808
GPT-3.5	0.30705	120.7903	0.94715	0.84396	104.2706	0.92369	2.69966	85.00323	0.91203
PaLM-2	0.30877	104.7076	0.98709	0.86842	92.49537	0.96723	2.06842	81.30042	0.95882
A-ItemKNN	5.53523	<u>23.36779</u>	0.04259	10.98658	<u>23.55763</u>	<u>0.0469</u>	27.22987	<u>23.4457</u>	0.01073
VSM	3.66443	29.44077	0.31291	7.68456	27.65017	0.33438	20.75	26.11899	0.19908
Random	<u>5.44631</u>	22.52315	<u>0.11481</u>	<u>10.90268</u>	22.38414	0.03134	<u>27.0151</u>	22.25768	<u>0.11414</u>

LAST.FM									
Model	cutoff 10			cutoff 20			cutoff 50		
	↑ACLT@10	↓ARP@10	↓PopREO@10	↑ACLT@20	↓ARP@20	↓PopREO@20	↑ACLT@50	↓ARP@50	↓PopREO@50
RP ³ _β	1.160267	165.9932	0.608086	2.52977	137.5507	0.54217	7.24485	104.9925	0.45809
ItemKNN	2.169727	112.7532	0.407507	4.64162	95.52604	0.37312	13.09627	74.46029	0.29958
EASE ^R	0.254869	188.3455	0.846867	0.68447	155.7025	0.80695	3.18976	116.8123	0.68211
UserKNN	0.553144	159.6465	0.773285	1.40902	134.8585	0.71999	5.21536	105.6959	0.57606
A-UserKNN	1.074569	129.3922	0.752619	2.63717	108.5136	0.69619	8.97718	84.93656	0.53294
VSM	0.584864	177.9949	0.722211	1.40122	152.8326	0.65269	4.86311	120.894	0.55853
A-ItemKNN	<u>3.042849</u>	<u>87.86322</u>	<u>0.198491</u>	<u>6.55314</u>	<u>75.6611</u>	<u>0.16821</u>	<u>18.19811</u>	<u>61.60384</u>	<u>0.13711</u>
PaLM-2	0.869847	111.7505	0.586115	1.49293	104.1728	0.56617	1.86337	103.0726	0.56153
ChatGPT-3.5	1.047858	106.5101	0.516102	2.49193	93.3422	0.47508	4.86366	85.93759	0.44051
GPT-3.5	0.691152	114.9996	0.68311	2.15748	93.39119	0.61524	4.96828	78.49711	0.57589
MostPop	0	348.3308	1	0	293.6111	1	0	216.525	1
Random	5.716194	31.34936	0.184579	11.31942	31.6756	0.0374	28.16973	31.47403	0.02255

Facebook Books									
Model	cutoff 10			cutoff 20			cutoff 50		
	↑ACLT@10	↓ARP@10	↓PopREO@10	↑ACLT@20	↓ARP@20	↓PopREO@20	↑ACLT@50	↓ARP@50	↓PopREO@50
ChatGPT-3.5	1.73475	58.12016	0.43521	3.71932	48.14003	0.38361	8.38575	36.16197	0.34827
A-ItemKNN	5.84644	7.24181	0.07496	11.99265	7.27719	0.08573	30.80456	7.0872	0.08553
VSM	5.77663	7.36657	0.04653	11.63336	7.77355	0.09747	30.19251	7.53464	0.06996
GPT-3.5	1.05511	73.17358	0.65831	2.70977	61.16968	0.63209	5.26745	50.01659	0.55565
A-UserKNN	0.32697	68.63292	0.7953	0.99633	57.84971	0.79787	4.83615	41.7785	0.71704
PaLM-2	1.94664	55.75281	0.45066	2.96688	49.63106	0.49214	3.80405	46.27225	0.48654
RP ³ _β	4.43277	28.78332	0.41911	9.32917	20.36381	0.36144	24.99412	12.47866	0.3788
UserKNN	0.28582	76.73226	0.8727	0.90448	60.85636	0.8654	4.81999	42.12936	0.82569
ItemKNN	<u>5.96253</u>	25.08244	0.13088	11.53857	21.22469	0.23714	26.78913	16.63286	0.3259
EASE ^R	0.03233	120.4731	0.89825	0.09772	92.87744	0.94851	0.49302	62.04598	0.9463
MostPop	0	138.3632	1	0	108.5059	1	0	71.6974	1
Random	6.92799	5.78663	0.23049	13.79427	5.77384	0.54294	34.20353	5.76281	0.05755

bias in LLMs, which is further corroborated by the findings from the novelty measures analysis (refer to 5.2).

When considering cutoff 20 and cutoff 50, the models demonstrate similar behaviour, with Attribute-based ItemKNN emerging as the least biased model at cutoff 20. ChatGPT and other large language models (LLMs) demonstrate the news-worthy ability to avoid significant bias whenever they can recommend a longer list of items. Hence, they are able to recommend niche items to the users based on their content interests as soon as the most popular ones have already been proposed. This eventuality may imply that LLMs may combine “safe” recommendations based on the item’s popularity and “risky” ones.

LAST.FM. Our analysis in the music domain confirms a lower dataset bias compared to the MovieLens dataset. This conclusion is supported by the better performance of EASE^R, which is significantly different from the MostPop approach.

Moreover, ChatGPT demonstrates performance on par with the leading models in ACLT and ARP, outperforming all other models except ItemKNN in PopREO. This consistent trend persists when evaluating the same metrics at cutoff 20 and cutoff 50. Furthermore, as expected, the Random and AttributeItemKNN models produce high values –i.e. indicating a low bias– across all evaluations.

In addition, PaLM-2 and GPT-3.5 results follow ChatGPT, further supporting the assumption that an LLM fine-tuned on conversation reaches higher performance in user-related tasks.

Facebook Books. Similarly to the previous analysis, ChatGPT once again exhibits remarkable performance in the domain of books. Possible motivations for this impressive performance could be attributed to ChatGPT’s higher specialization in the book domain, but further exploration is needed to understand the underlying factors.

However, when comparing ChatGPT to other content-based models (A-ItemKNN, VSM, and A-UserKNN), it shows relatively worse results in terms of ACLT, indicating its higher susceptibility to popularity bias. This finding raises questions about how popularity bias affects ChatGPT’s book recommendations and warrants further investigation.

This observation is further supported by ARP, where AttributeItemKNN and VSM exhibit significantly lower values. ChatGPT may still rely on collaborative data, although some content information is utilized in performing recommendations.

Moreover, when considering PopREO, AttributeItemKNN and VSM demonstrate a better tradeoff, while the other large language models (LLMs) exhibit behaviour similar to ChatGPT. There are marginal differences in ACLT and ARP, where PaLM outperforms ChatGPT and GPT-3.5. These differences may depend on PaLM’s lower specialization in the book domain or the lesser exploited collaborative data.

However, despite being influenced by popularity bias, ChatGPT performs comparably to other recommendation systems, and these trends remain consistent across various cutoff thresholds.

Conclusion. Finally, a noteworthy level of bias is evident when comparing ChatGPT to the state-of-the-art models. ChatGPT demonstrates varying degrees of popularity bias across different datasets and exhibits the behaviour of recommending popular items. Similar observations can be made for GPT-3.5 and PaLM-2. These findings highlight the need for intensive efforts to address and mitigate bias in the final recommendations generated by ChatGPT and other LLMs.

5.4 RQ1d: Which type of recommender system is ChatGPT more similar to?

The previous analysis positioned ChatGPT and other Large Language Models (LLMs) within the state-of-the-art. However, it remains unclear whether these models behave

Table 6: Experiment 1: A Comparative analysis of ChatGPT-3.5 with various baselines on Jaccard, and Kendall’s Tau with cutoff at 10, 20, and 50 on MovieLens recommendations. The results are ordered by Jaccard@10 with the best results highlighted in bold and the runner-up result underlined. (A- is for Attribute).

MovieLens						
Model	cutoff 10		cutoff 20		cutoff 50	
	↑Jaccard@10	↑Kendall@10	↑Jaccard@20	↑Kendall@20	↑Jaccard@50	↑Kendall@50
EASE ^R	0.16389	<u>0.03134</u>	0.17189	<u>0.04715</u>	0.16216	<u>0.05568</u>
MostPop	<u>0.14789</u>	0.02379	<u>0.16485</u>	0.07266	<u>0.15037</u>	0.06825
RP ³ β	0.12550	0.02255	0.11907	0.02506	0.10984	0.01834
A-UserKNN	0.11794	0.03581	0.10937	0.03928	0.10460	0.03892
UserKNN	0.11571	0.01052	0.10804	0.01596	0.09870	0.01997
ItemKNN	0.07449	0.00955	0.08655	0.01441	0.09594	0.01363
VSM	0.00604	-0.00212	0.00860	-0.00557	0.01735	-0.00589
A-ItemKNN	0.00515	0.00952	0.00898	0.02116	0.01843	0.00578
Random	0.00254	-0.03154	0.00419	-0.01742	0.00932	-0.01054

LAST.FM						
Model	cutoff 10		cutoff 20		cutoff 50	
	↑Jaccard@10	↑Kendall@10	↑Jaccard@20	↑Kendall@20	↑Jaccard@50	↑Kendall@50
RP ³ β	0.15970	<u>0.05588</u>	0.14719	<u>0.04396</u>	0.13815	0.02685
UserKNN	<u>0.15714</u>	0.04664	<u>0.13902</u>	0.03990	0.13058	0.02768
A-UserKNN	0.14621	0.04812	0.13314	0.03527	0.12743	0.02680
ItemKNN	0.13223	0.02428	0.12459	0.02256	0.11864	0.01783
EASE ^R	0.12914	0.06638	0.12848	0.05790	<u>0.13302</u>	<u>0.04070</u>
VSM	0.11455	0.03029	0.11043	0.03226	0.10994	0.02187
A-ItemKNN	0.10901	0.02546	0.10104	0.02454	0.09464	0.01773
MostPop	0.04193	0.02835	0.04571	0.02403	0.06711	0.04271
Random	0.00304	-0.00512	0.00518	-0.00171	0.01056	-0.00260

Facebook Books						
Model	cutoff 10		cutoff 20		cutoff 50	
	↑Jaccard@10	↑Kendall@10	↑Jaccard@20	↑Kendall@20	↑Jaccard@50	↑Kendall@50
MostPop	0.18824	<u>0.00730</u>	0.18235	-0.01066	0.15425	0.05266
EASE ^R	<u>0.17739</u>	-0.00829	<u>0.17619</u>	0.00409	<u>0.15064</u>	<u>0.05059</u>
A-UserKNN	0.10935	0.01376	0.11194	0.02083	0.10320	0.03981
UserKNN	0.09883	-0.00780	0.10382	<u>0.01166</u>	0.09736	0.04145
ItemKNN	0.03574	-0.00077	0.03735	0.00824	0.04029	-0.00221
RP ³ β	0.02770	-0.01885	0.02692	-0.00103	0.02692	0.01162
VSM	0.02365	-0.00362	0.02571	0.00314	0.03043	0.00760
A-ItemKNN	0.02324	-0.00783	0.02329	-0.00478	0.02694	0.00245
Random	0.00257	-0.00534	0.00442	0.00035	0.00919	-0.00364

as Content-based (CB) systems. In this section, we aim to address this uncertainty by analyzing the recommendations generated by ChatGPT and interpreting them in the context of state-of-the-art recommender models. The results of this evaluation are presented in Table 6, offering a comprehensive overview of the recommendation-set similarity (Jaccard index) and ranking similarity (Kendall’s Tau coefficient) between the recommendation lists of ChatGPT and those generated by other recommendation system (RS) algorithms.

MovieLens. The analysis of the Jaccard index uncovers intriguing findings concerning the similarity of ChatGPT to other recommender models on the MovieLens dataset. Notably, ChatGPT, EASE^R, and MostPop exhibit the highest similarity, indicating that ChatGPT tends to recommend popular items. However, going beyond these two models, ChatGPT’s behaviour aligns more closely with collaborative-based models like EASE^R, RP³ β , AttributeUserKNN (a hybrid model), and UserKNN.

The observed similarity in ChatGPT’s behaviour raises the question of whether it is primarily influenced by popularity bias or the collaborative information it gathers. It is worth noting that while ChatGPT suggests the same items for recommendation, their ranking differs. This disparity underscores ChatGPT’s unique position in the state-of-the-art, positioned between content-based (CB) and collaborative-based filtering (CBF) recommendation systems.

The analysis of the Kendall metric reveals interesting insights about ChatGPT’s behavior at different cutoffs. At cutoff 10, ChatGPT exhibits the highest similarity to AttributeUserKNN, indicating a notable alignment between the two models. However, this similarity decreases significantly at cutoffs 20 and 50. As a result, for the MovieLens dataset, ChatGPT primarily acts as a recommender system suggesting popular items while incorporating some content-based information, especially at cutoff 10. In contrast, at cutoffs 20 and 50, ChatGPT relies more on popularity and collaborative elements rather than content-based methods.

LAST.FM. The previous experiments conducted on the LAST.FM dataset revealed a consistent trend of low popularity bias for ChatGPT. This finding suggests that ChatGPT lacks specialized knowledge of the music domain, highlighting intriguing aspects of the LLM’s performance in the recommendation task.

A detailed analysis of both similarity metrics further confirms the validity of the earlier assumption. The results show that ChatGPT’s recommended lists have low similarity with the MostPop-generated lists but demonstrate high similarity with collaborative filtering models, i.e. RP³ β and UserKNN.

The observed low similarity with Content-based models provides additional evidence that ChatGPT has limited knowledge in the music domain. However, at cutoff 10, RP³ β , UserKNN, and AttributeUserKNN continue to be the most comparable models to ChatGPT. This finding implies that ChatGPT demonstrates greater similarity to hybrid and collaborative models rather than content-based ones, challenging the initial hypothesis that ChatGPT relies solely on content data.

This pattern remains consistent at cutoff 20 and 50, with the Kendall metric indicating that ChatGPT demonstrates purely collaborative behaviour at cutoff 50.

Facebook Books. As demonstrated in Section 5.3 (“How much is ChatGPT biased compared with the state-of-the-art?”), the analysis shows that ChatGPT is susceptible to popularity bias in the Facebook Books dataset. The Jaccard index analysis further validates this observation, indicating that ChatGPT exhibits the highest similarity to MostPop, followed by EASE^R(Jaccard @10/@20/@50 and Kendall’s Tau @50). This empirical evidence confirms our hypothesis that ChatGPT tends to prioritize collaborative-popularity data.

However, AttributeUserKNN stands out as the first non-collaborative model with similarities to ChatGPT, indicating a potential hybrid behaviour with occasional peaks of popularity influence. This pattern persists at cutoff 50, while at cutoff 20 (Kendall’s Tau), AttributeUserKNN emerges as the most similar model to ChatGPT.

This pattern remains consistent at cutoff 20 and 50, and with the Kendall metric.

Conclusion. These findings underscore ChatGPT’s tendency to align with hybrid and collaborative recommender models, showcasing its preference for a balanced approach rather than relying solely on content-based methods or popularity.

Specifically, in cases where content data is excessive or insufficient to make a recommendation, ChatGPT shows higher similarity with collaborative models, as demonstrated with Jaccard on MovieLens and LAST.FM.

This observation sheds new light on ChatGPT’s unique behaviour within the realm of recommendation systems, revealing results and features that have not been previously explored.

5.5 RQ2: Is ChatGPT able to exploit user preferences to re-rank a recommendation list?

In the previous experiments, ChatGPT has demonstrated the capability to recommend new items. Nonetheless, further investigation is required to assess its ability to understand user preferences and deliver personalised recommendations. To address this, we conducted two experiments:

- Experiment 1: Re-ranking a fixed list of the most popular items.
- Experiment 2: Dynamically generating personalised lists for each user starting from the preferences of their nearest neighbours.

The first task was to re-rank a list of items. Nevertheless, ChatGPT went beyond the assigned task. Indeed, ChatGPT replaced less relevant items from the fixed list by introducing new and different items. This behaviour was observed less frequently in the second experiment.

The accuracy, diversity, novelty, and bias evaluation results for both the re-ranking experiments are presented in Tables 7, 8, and 9.

Re-ranking a Most Popular list. To explore the re-ranking setting, we generated a fixed list of the fifty most popular items for each dataset. Subsequently, we have first evaluated the list in its original form across all users, and we have then incorporated it into ChatGPT with the user profile to generate and evaluate a personalised re-ranked list.

This investigation reveals that ChatGPT significantly improves overall recommendations through re-ranking, using the MostPop results as a baseline. Across all datasets and the majority of cutoff values, ChatGPT’s performance improvement is twofold: (i) a higher Gini index, suggesting its ability to diversify recommendations, while (ii) a higher EPC, indicating the successful introduction of novel items. However, it is

Table 7: Experiment 2 and 3 - A comparative analysis of ChatGPT-3.5 metrics after the re-rank on MostPop list, and personalized list with various baselines, with cutoff at 10, 20, and 50 on MovieLens. The best results are high-lighted in bold.

		MovieLens								
	model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACL \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
cutoff 10	MostPop	0.08055	0.03473	0.06812	0.00000	108.63205	1	26	0.00614	0.05819
	ChatGPT-3.5	0.14880	0.07446	0.12265	0.00336	153.13389	1	80	0.01202	0.09903
cutoff 20	MostPop	0.08946	0.07117	0.06980	0.00000	105.21535	1	47	0.01256	0.05713
	ChatGPT-3.5	0.13959	0.10780	0.09010	0.02181	133.63866	1	133	0.01880	0.08066
cutoff 50	MostPop	0.10983	0.15480	0.05054	0.00000	107.36097	1	50	0.02452	0.05832
	ChatGPT-3.5	0.13479	0.14450	0.04758	0.13423	108.35477	1	246	0.02758	0.07017

		LAST.FM								
	model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACL \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
cutoff 10	Nearest Neighbors	0.10568	0.06684	0.09983	1.51678	65.15302	0.72954	1146	0.28106	0.08347
	ChatGPT-3.5	0.20359	0.11696	0.16594	0.22483	119.27131	0.98089	533	0.06791	0.14229
cutoff 20	Nearest Neighbors	0.12622	0.13109	0.09740	2.96141	65.03809	0.73239	1346	0.30116	0.08229
	ChatGPT-3.5	0.20495	0.18005	0.12970	1.20973	92.73143	0.90861	1009	0.15534	0.12133
cutoff 50	Nearest Neighbors	0.19933	0.33021	0.09584	7.29027	64.85725	0.71754	1533	0.31651	0.08150
	ChatGPT-3.5	0.22823	0.28522	0.08399	5.58054	69.95444	0.73156	1501	0.29852	0.10086

Table 8: Experiment 2 and 3 - A comparative analysis of ChatGPT-3.5 metrics after the re-rank on MostPop list, and personalized list with various baselines, with cutoff at 10, 20, and 50 on LAST.FM. The best results are high-lighted in bold.

		LAST.FM								
	model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACL \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
cutoff 10	MostPop	0.07828	0.07915	0.06528	0.00000	325.16628	1.00000	24	0.00793	0.05614
	ChatGPT-3.5	0.13569	0.10910	0.08358	0.06845	220.38203	0.93419	338	0.02326	0.09810
cutoff 20	MostPop	0.09800	0.11889	0.04808	0.00000	272.35846	1.00000	37	0.01519	0.04574
	ChatGPT-3.5	0.15336	0.14449	0.05526	0.14190	205.41677	0.92894	443	0.02791	0.07379
cutoff 50	MostPop	0.13164	0.20401	0.03213	0.00000	211.94210	1.00000	50	0.03140	0.03893
	ChatGPT-3.5	0.16362	0.16974	0.02598	0.25598	203.68562	0.93633	564	0.03481	0.06444

		LAST.FM								
	model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACL \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
cutoff 10	Nearest Neighbors	0.07686	0.08915	0.06334	1.92453	86.58383	0.45409	1018	0.36418	0.05785
	ChatGPT-3.5	0.17413	0.16134	0.11820	1.03561	121.27462	0.62086	1151	0.24543	0.12858
cutoff 20	Nearest Neighbors	0.11450	0.16742	0.06024	3.90027	86.96334	0.35272	1220	0.39451	0.05605
	ChatGPT-3.5	0.20831	0.23071	0.08492	2.56149	106.10265	0.56087	1362	0.31457	0.10210
cutoff 50	Nearest Neighbors	0.19787	0.38450	0.05650	9.88949	86.56765	0.41512	1407	0.42510	0.05327
	ChatGPT-3.5	0.22778	0.28088	0.04067	4.98108	98.45203	0.50676	1457	0.37911	0.08836

important to note that this improvement varies for the Facebook Books dataset, with a cutoff of 50.

The enhanced ranked list suggests that ChatGPT effectively uses the provided user preferences to personalise and re-rank the recommendations. Additionally, it demonstrates its ability to introduce new items, as evidenced by a higher ItemCoverage metric value.

Re-rank on nearest neighbors' preferences. We conducted a second experiment to gain deeper insights into ChatGPT's ability to re-rank. Unlike the previous one, where a fixed list was used for all users, we generated personalized lists for each user and tasked ChatGPT with re-ranking them. We carefully curated these lists by selecting the five highest-rated items from each nearest neighbour, creating a tailored list of fifty items for each user. To ensure objectivity and minimize ranking bias based

Table 9: Experiment 2 and 3 - A comparative analysis of ChatGPT-3.5 metrics after the re-rank on MostPop list, and personalized list with various baselines, with cutoff at 10, 20, and 50 on Facebook Books. The best results are high-lighted in bold.

Facebook Books										
	model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACLT \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
cutoff 10	MostPop	0.00091	0.00257	0.00037	3.61866	74.86106	1.00000	17	0.00420	0.00018
	ChatGPT-3.5	0.01232	0.01954	0.00432	0.31845	90.87187	0.43095	331	0.01368	0.00452
cutoff 20	MostPop	0.00328	0.01016	0.00110	9.59589	45.57179	0.65500	27	0.00864	0.00071
	ChatGPT-3.5	0.01552	0.02965	0.00324	1.97024	70.71079	0.40133	430	0.01785	0.00374
cutoff 50	MostPop	0.01148	0.04449	0.00190	20.96988	33.07245	0.39975	49	0.02107	0.00150
	ChatGPT-3.5	0.01743	0.03734	0.00161	10.64360	54.75838	0.30657	531	0.02535	0.00291
	model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACLT \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
cutoff 10	Nearest Neighbors	0.01191	0.02132	0.00500	1.84203	36.50360	0.64570	1630	0.26445	0.00469
	ChatGPT-3.5	0.02993	0.04455	0.00975	0.50574	64.05953	0.72869	698	0.07020	0.01111
cutoff 20	Nearest Neighbors	0.01928	0.04483	0.00514	3.73475	36.74467	0.77165	1887	0.27770	0.00481
	ChatGPT-3.5	0.03390	0.05754	0.00621	1.69694	51.54417	0.73359	1237	0.14129	0.00811
cutoff 50	Nearest Neighbors	0.03017	0.08816	0.00397	7.43277	36.75897	0.77610	2033	0.28374	0.00475
	ChatGPT-3.5	0.03739	0.07111	0.00308	4.61472	42.91002	0.71285	1751	0.23747	0.00624

on neighbour similarity, we further introduced disorder by randomizing the item order in the list.

By using Nearest Neighbors and Cosine Similarity as a reference, the results obtained from ChatGPT show its ability to “understand” user preferences. Across all datasets and cutoff values, ChatGPT achieves higher values in nDCG and EPC, resulting in an improved ranking and an introduction of novel items in the recommended lists. Additionally, the model exhibits lower Item Coverage but higher Popularity Bias (PopREO), along with a decrease in the Gini value. These findings suggest that the re-ranked list generated by ChatGPT tends to shift towards the most popular items, limiting diversification. However, the higher performance in nDCG and the above observations reaffirm ChatGPT’s ability to offer personalized item suggestions to each user.

Conclusion. In summary, our exploration focuses on ChatGPT’s ability to utilize user profiles for re-ranking recommendations, leading to enhanced personalization. We investigate two scenarios: one involving a fixed list of the most popular items and the other based on preferences from nearest neighbours. The findings emphasize ChatGPT’s effectiveness in personalizing recommendations based on user profiles. Additionally, in scenarios with limited items (e.g., the first experimental scenario), ChatGPT showcases its reliance on its knowledge, offering the potential to address the Cold Start Problem.

5.6 RQ3: Does the substantial amount of knowledge utilized to train ChatGPT compensate for the absence of a complete user history in a cold-start scenario?

Firstly, we identify the cold start users by dividing the users in quartiles based on the number of their past interactions.

To investigate the performance of LLMs in cold-start scenarios, we employ a systematic two-step approach. Initially, we identify cold-start users by dividing them into quartiles based on their past interactions. Subsequently, we use the lower quartile as a

Table 10: Experiment 4: A Comparative analysis of ChatGPT-3.5 with various baselines and Datasets on Cold Start scenario with cutoff at 10. The results are ordered by nDCG with the best results high-lighted in bold and the runner-up result underlined. (A- is for Attribute).

MovieLens									
Model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACLT \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
PaLM-2	0.14032	0.07665	0.10596	0.30877	104.70760	0.98709	277	0.02758	0.09883
ChatGPT-3.5	<u>0.08719</u>	<u>0.09375</u>	<u>0.04211</u>	0.43421	129.05592	1.00000	211	0.02664	0.03797
GPT-3.5	0.08648	0.09868	<u>0.04474</u>	0.34211	132.23092	1.00000	211	0.03071	<u>0.03858</u>
RP $^3\beta$	0.03052	0.05000	<u>0.01974</u>	0.69737	67.09013	1.00000	518	0.15417	0.01539
ItemKNN	0.02710	0.03575	0.01645	1.51974	62.83618	1.00000	652	0.20589	0.01468
UserKNN	0.02329	0.03191	0.01447	0.99342	61.06711	0.61702	596	0.18760	0.01284
EASE ^R	0.02253	0.03388	0.01711	0.00000	94.42237	1.00000	142	0.03556	0.01287
MostPop	0.01803	0.02763	0.01184	0.00000	97.65724	1.00000	24	0.00763	0.00895
AttributeUserKNN	0.01376	0.01941	0.00921	1.24342	59.71184	1.00000	648	0.21535	0.00750
AttributeItemKNN	0.01201	0.01513	0.00592	<u>5.31579</u>	<u>22.73355</u>	<u>0.17391</u>	<u>957</u>	<u>0.37739</u>	0.00611
VSM	0.00568	0.00800	0.00395	4.69079	23.96316	0.00264	583	0.19462	0.00349
Random	0.00400	0.00647	0.00263	5.44737	21.70658	0.24752	1034	0.43162	0.00221
LAST.FM									
Model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACLT \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
ChatGPT-3.5	0.20311	0.21970	0.09859	1.92525	90.67174	0.31012	708	0.19659	0.11847
GPT-3.5	<u>0.17498</u>	0.17811	0.08101	1.33131	100.79920	0.42954	570	0.13548	<u>0.10251</u>
PaLM-2	0.17121	<u>0.19503</u>	<u>0.08948</u>	1.39056	93.32864	0.37234	588	0.13837	0.10235
UserKNN	0.03278	<u>0.04189</u>	<u>0.02000</u>	2.21616	75.06949	0.05075	1067	0.36214	0.02013
VSM	0.03030	0.04199	0.01960	1.90303	83.36566	0.08187	902	0.26529	0.01823
RP $^3\beta$	0.02892	0.04178	0.01980	2.68687	70.10788	0.11289	1090	0.38758	0.01764
ItemKNN	0.02769	0.03865	0.01798	4.14545	48.40222	<u>0.04236</u>	1319	0.56435	0.01689
AttributeItemKNN	0.02566	0.03475	0.01616	<u>4.93939</u>	<u>42.10768</u>	0.01351	<u>1373</u>	<u>0.60725</u>	0.01511
AttributeUserKNN	0.02169	0.03044	0.01495	4.21616	46.46586	0.26246	1121	0.40568	0.01344
EASE ^R	0.02139	0.02990	0.01414	0.56162	109.56162	0.53712	617	0.16076	0.01311
MostPop	0.01494	0.02175	0.01051	0.00000	146.27939	1.00000	21	0.00722	0.00889
Random	0.00392	0.00572	0.00242	5.59394	32.74949	0.12686	1467	0.69767	0.00220
Facebook Books									
Model	nDCG \uparrow	Recall \uparrow	Precision \uparrow	ACLT \uparrow	ARP \downarrow	PopREO \downarrow	ItemCoverage \uparrow	Gini \uparrow	EPC \uparrow
ChatGPT-3.5	0.04871	0.06809	0.01539	1.92245	55.01169	0.37060	635	0.04169	0.01858
PaLM-2	<u>0.03975</u>	0.05399	<u>0.01263</u>	2.38039	51.86428	0.34845	535	0.03969	<u>0.01551</u>
GPT-3.5	0.03689	<u>0.05440</u>	<u>0.01262</u>	1.30787	71.96433	0.62235	403	0.01773	0.01407
MostPop	0.01390	0.02276	0.00532	0.00000	40.99155	1.00000	18	0.00387	0.00530
EASE ^R	0.00918	0.01813	0.00417	0.15625	35.82002	0.87139	579	0.05465	0.00353
UserKNN	0.00755	0.01254	0.00278	1.29630	26.25428	1.00000	1363	0.19460	0.00284
AttributeUserKNN	0.00710	0.01427	0.00324	1.52083	25.70845	0.42765	1418	0.21258	0.00271
ItemKNN	0.00706	0.01177	0.00255	4.97454	13.45532	0.79446	<u>2017</u>	<u>0.47665</u>	0.00260
VSM	0.00447	0.00849	0.00185	6.21181	6.97500	<u>0.18142</u>	1340	0.28282	0.00168
RP $^3\beta$	0.00363	0.00752	0.00174	5.35995	5.71586	0.45979	1809	0.38695	0.00144
AttributeItemKNN	0.00362	0.00521	0.00127	6.29167	6.65231	0.05151	1400	0.26594	0.00147
Random	0.00094	0.00193	0.00046	6.81481	<u>5.95359</u>	0.41270	2194	0.61272	0.00037

filter on the test set, creating a subset of users representing the cold-start users. This enables us to evaluate all models and datasets under similar cold-start conditions, and the accuracy, diversity, novelty, and bias results are presented in Table 10.

In our evaluation, considering metrics such as nDCG, Recall, and Precision, LLMs consistently obtain the top three positions across all datasets. Notably, ChatGPT-3.5 excels in the music and books domains, while PaLM-2 demonstrates superior performance in the movie domain. However, these models do not exhibit notable results considering other metrics, with the exception of the novelty metric (EPC). Nevertheless, it is evident that ChatGPT (and other LLMs) achieve remarkable performance in generating recommendations even in cold-start situations, outperforming state-of-the-art models.

Conclusion. To summarize, ChatGPT and other Large Language Models can recommend items in cold-start scenarios. However, these models show limitations in Bias and Coverage metrics, while displaying strong performance in Novelty and Accuracy.

6 Conclusion

This study delves into the potential of ChatGPT as a recommender system, conducting an extensive experimental investigation on three datasets: MovieLens Small, Last.FM, and Facebook Book. We compare ChatGPT’s performance against various state-of-the-art recommender systems (RSs), such as UserKNN, ItemKNN, $RP^3\beta$, $EASE^R$, AttributeItemKNN, AttributeUserKNN, and VSM.

With the first experiment, we aim to understand how much ChatGPT is accurate in performing recommendations. Thus, we compared it with the state-of-the-art solutions, including other Large Language Models, and observed that it accurately solves the recommendation task. Even without any improvements through prompts engineering, ChatGPT produces results comparable to state-of-the-art Recommender Systems and stands out from other Large Language Models.

Then, our examination explored the diversity and novelty of ChatGPT in suggesting new items to users, and our findings have shown different conclusions depending on the domain. In its vanilla version, ChatGPT reached lower diversity but higher novelty when dealing with Books, while it showcased acceptable novelty in the Music domain.

Furthermore, our investigation aimed to reveal the potential bias inherent in ChatGPT’s recommendations. We proved that ChatGPT holds varying degrees of popularity bias across different datasets by recommending popular items. Also, we acknowledged an analogous biased behaviour in GPT-3.5 and PaLM-2.

Deepening the analysis of the recommendation abilities of ChatGPT, we further explored analogies with existing recommendation paradigms. By comparing the list of suggested items produced by ChatGPT with those from the baselines, we discovered a tendency to align with hybrid and collaborative recommenders. Thus, our results revealed that ChatGPT capabilities go beyond mere content information in selecting the items to suggest.

We were also interested in examining how ChatGPT leverages user profiles for re-ranking existing recommendations. We dissect two settings: the first with a fixed list of the most popular items and the second based on the user’s nearest neighbours. The findings demonstrated ChatGPT’s efficacy in further personalizing recommendations based on user profiles.

Finally, we investigate whether ChatGPT can address the Cold Start scenario. The experimental results proved that ChatGPT can still perform relevant suggestions, outperforming the other state-of-the-art recommenders.

In summary, our findings revealed that ChatGPT exhibited characteristics of a hybrid recommender system, leveraging collaborative and content-based information. Additionally, ChatGPT showcased domain-specific knowledge, particularly in books and movies domains. Moreover, it tended to recommend popular items and exhibited a remarkable proficiency in handling the cold start problem.

Despite its contributions, this study has certain limitations. Firstly, it does not delve into techniques for prompt-engineering, such as Chain-of-Thought or Tree-of-Thought, which could potentially enhance the quality of recommendations. Secondly, the rapid emergence of new LLMs has introduced additional baselines that were not considered in the current results. Lastly, the exploration of recommendations in scenarios with abundant user information was not feasible due to the context limits of the ChatGPT API.

Considering the promising outcomes of this study, our future research aims to delve deeper into ChatGPT’s performance in the recommendation task, with a focus on accurately designing a recommender framework that incorporates ChatGPT to enhance recommendation performance. To achieve this objective, we plan to explore potential improvements in the evaluation by investigating prompt engineering techniques and domain-specific fine-tuning methodologies.

References

- [1] Abdullah, M.H.A., Aziz, N., Abdulkadir, S.J., Alhussian, H.S.A., Talpur, N.: Systematic literature review of information extraction from textual data: Recent methods, applications, trends, and challenges. *IEEE Access* **11**, 10535–10562 (2023)
- [2] Burke, R.: *Recommender Systems: An Introduction*, by dietmar jannach, markus zanker, alexander felfernig, and gerhard friedrichcambridge university press, 2011, 336 pages. ISBN: 978-0-521-49336-9. *Int. J. Hum. Comput. Interact.* **28**(1), 72–73 (2012)
- [3] Biancofiore, G.M., Deldjoo, Y., Noia, T.D., Sciascio, E.D., Narducci, F.: Interactive question answering systems: Literature review. *CoRR* **abs/2209.01621** (2022)
- [4] Maarek, Y.: Alexa and her shopping journey. In: *CIKM*, p. 1. ACM, New York, NY, USA (2018)
- [5] Rosenfeld, R.: Two decades of statistical language modeling: where do we go from here? *Proc. IEEE* **88**(8), 1270–1278 (2000)
- [6] Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Nee-lakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners. In: *NeurIPS* (2020)
- [7] Zheng, L., Chiang, W., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E.P., Zhang, H., Gonzalez, J.E., Stoica, I.: Judging llm-as-a-judge with mt-bench and chatbot arena. *CoRR* **abs/2306.05685** (2023)

- [8] Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., Hashimoto, T.B.: Stanford Alpaca: An Instruction-following LLaMA model. GitHub (2023)
- [9] Christiano, P.F., Leike, J., Brown, T.B., Martic, M., Legg, S., Amodei, D.: Deep reinforcement learning from human preferences. In: NIPS, pp. 4299–4307 (2017)
- [10] Stokel-Walker, C., Van Noord, R.: What chatgpt and generative ai mean for science. *Nature* **614**(7947), 214–216 (2023)
- [11] Bang, Y., Cahyawijaya, S., Lee, N., Dai, W., Su, D., Wilie, B., Lovenia, H., Ji, Z., Yu, T., Chung, W., Do, Q.V., Xu, Y., Fung, P.: A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *CoRR* **abs/2302.04023** (2023)
- [12] Zhang, J., Bao, K., Zhang, Y., Wang, W., Feng, F., He, X.: Is chatgpt fair for recommendation? evaluating fairness in large language model recommendation. *CoRR* **abs/2305.07609** (2023)
- [13] Zhiyuli, A., Chen, Y., Zhang, X., Liang, X.: Bookgpt: A general framework for book recommendation empowered by large language model. *CoRR* **abs/2305.15673** (2023)
- [14] Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Comput. Surv.* **55**(9), 195–119535 (2023)
- [15] Harper, F.M., Konstan, J.A.: The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.* **5**(4), 19–11919 (2016)
- [16] Cantador, L., Brusilovsky, P., Kuflik, T.: Second workshop on information heterogeneity and fusion in recommender systems (hetrec2011). In: *RecSys*, pp. 387–388. ACM, New York, NY, USA (2011)
- [17] Anelli, V.W., Bellogín, A., Ferrara, A., Malitesta, D., Merra, F.A., Pomo, C., Donini, F.M., Noia, T.D.: Elliot: A comprehensive and rigorous framework for reproducible recommender systems evaluation. In: *SIGIR*, pp. 2405–2414. ACM, New York, NY, USA (2021)
- [18] Zhou, C., Li, Q., Li, C., Yu, J., Liu, Y., Wang, G., Zhang, K., Ji, C., Yan, Q., He, L., Peng, H., Li, J., Wu, J., Liu, Z., Xie, P., Xiong, C., Pei, J., Yu, P.S., Sun, L.: A comprehensive survey on pretrained foundation models: A history from BERT to chatgpt. *CoRR* **abs/2302.09419** (2023)
- [19] Stiennon, N., Ouyang, L., Wu, J., Ziegler, D.M., Lowe, R., Voss, C., Radford, A., Amodei, D., Christiano, P.F.: Learning to summarize with human feedback. In: *NeurIPS* (2020)

- [20] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P.F., Leike, J., Lowe, R.: Training language models to follow instructions with human feedback. In: NeurIPS (2022)
- [21] Han, X., Zhang, Z., Ding, N., Gu, Y., Liu, X., Huo, Y., Qiu, J., Yao, Y., Zhang, A., Zhang, L., Han, W., Huang, M., Jin, Q., Lan, Y., Liu, Y., Liu, Z., Lu, Z., Qiu, X., Song, R., Tang, J., Wen, J., Yuan, J., Zhao, W.X., Zhu, J.: Pre-trained models: Past, present and future. *AI Open* **2**, 225–250 (2021)
- [22] Li, L., Zhang, Y., Chen, L.: Personalized prompt learning for explainable recommendation. *CoRR* **abs/2202.07371** (2022)
- [23] Qin, G., Eisner, J.: Learning how to ask: Querying lms with mixtures of soft prompts. In: NAACL-HLT, pp. 5203–5212. Association for Computational Linguistics, Online (2021)
- [24] Kojima, T., Gu, S.S., Reid, M., Matsuo, Y., Iwasawa, Y.: Large language models are zero-shot reasoners. *CoRR* **abs/2205.11916** (2022)
- [25] Zhang, Y., DING, H., Shui, Z., Ma, Y., Zou, J., Deoras, A., Wang, H.: Language models as recommender systems: Evaluations and limitations. In: I (Still) Can’t Believe It’s Not Better! NeurIPS 2021 Workshop
- [26] Hidasi, B., Karatzoglou, A., Baltrunas, L., Tikk, D.: Session-based recommendations with recurrent neural networks. In: ICLR (Poster) (2016)
- [27] Li, J., Zhang, W., Wang, T., Xiong, G., Lu, A., Medioni, G.: Gpt4rec: A generative framework for personalized recommendation and user interests interpretation. *CoRR* **abs/2304.03879** (2023)
- [28] Lin, J., Men, R., Yang, A., Zhou, C., Zhang, Y., Wang, P., Zhou, J., Tang, J., Yang, H.: M6: multi-modality-to-multi-modality multitask mega-transformer for unified pretraining. In: KDD, pp. 3251–3261. ACM, New York, NY, USA (2021)
- [29] Wang, L., Hu, H., Sha, L., Xu, C., Jiang, D., Wong, K.: Recindial: A unified framework for conversational recommendation with pretrained language models. In: *ACL/IJCNLP* (1), pp. 489–500. Association for Computational Linguistics, Online (2022)
- [30] Zhang, Y., Sun, S., Galley, M., Chen, Y., Brockett, C., Gao, X., Gao, J., Liu, J., Dolan, B.: DIALOGPT : Large-scale generative pre-training for conversational response generation. In: *ACL (demo)*, pp. 270–278. Association for Computational Linguistics, Online (2020)
- [31] Wang, L., Lim, E.: Zero-shot next-item recommendation using large pretrained language models. *CoRR* **abs/2304.03153** (2023)

- [32] Geng, S., Liu, S., Fu, Z., Ge, Y., Zhang, Y.: Recommendation as language processing (RLP): A unified pretrain, personalized prompt & predict paradigm (P5). In: RecSys, pp. 299–315. ACM, New York, NY, USA (2022)
- [33] Wu, Y., Xie, R., Zhu, Y., Zhuang, F., Zhang, X., Lin, L., He, Q.: Personalized prompts for sequential recommendation. CoRR **abs/2205.09666** (2022)
- [34] Gao, Y., Sheng, T., Xiang, Y., Xiong, Y., Wang, H., Zhang, J.: Chat-rec: Towards interactive and explainable llms-augmented recommender system. CoRR **abs/2303.14524** (2023)
- [35] Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H.W., Sutton, C., Gehrmann, S., Schuh, P., Shi, K., Tsvyashchenko, S., Maynez, J., Rao, A., Barnes, P., Tay, Y., Shazeer, N., Prabhakaran, V., Reif, E., Du, N., Hutchinson, B., Pope, R., Bradbury, J., Austin, J., Isard, M., Gur-Ari, G., Yin, P., Duke, T., Levskaya, A., Ghemawat, S., Dev, S., Michalewski, H., Garcia, X., Misra, V., Robinson, K., Fedus, L., Zhou, D., Ippolito, D., Luan, D., Lim, H., Zoph, B., Spiridonov, A., Sepassi, R., Dohan, D., Agrawal, S., Omernick, M., Dai, A.M., Pillai, T.S., Pellat, M., Lewkowycz, A., Moreira, E., Child, R., Polozov, O., Lee, K., Zhou, Z., Wang, X., Saeta, B., Diaz, M., Firat, O., Catasta, M., Wei, J., Meier-Hellstern, K., Eck, D., Dean, J., Petrov, S., Fiedel, N.: Palm: Scaling language modeling with pathways. CoRR **abs/2204.02311** (2022)
- [36] Anil, R., Dai, A.M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., Chu, E., Clark, J.H., Shafey, L.E., Huang, Y., Meier-Hellstern, K., Mishra, G., Moreira, E., Omernick, M., Robinson, K., Ruder, S., Tay, Y., Xiao, K., Xu, Y., Zhang, Y., Ábrego, G.H., Ahn, J., Austin, J., Barham, P., Botha, J.A., Bradbury, J., Brahma, S., Brooks, K., Catasta, M., Cheng, Y., Cherry, C., Choquette-Choo, C.A., Chowdhery, A., Crepy, C., Dave, S., Dehghani, M., Dev, S., Devlin, J., Díaz, M., Du, N., Dyer, E., Feinberg, V., Feng, F., Fienber, V., Freitag, M., Garcia, X., Gehrmann, S., Gonzalez, L., al.: Palm 2 technical report. CoRR **abs/2305.10403** (2023)
- [37] Wang, F., Li, J., Qin, R., Zhu, J., Mo, H., Hu, B.: Chatgpt for computational social systems: From conversational applications to human-oriented operating systems. IEEE Trans. Comput. Soc. Syst. **10**(2), 414–425 (2023)
- [38] Liu, T., Jiang, Y.E., Monath, N., Cotterell, R., Sachan, M.: Autoregressive structured prediction with language models. In: EMNLP (Findings), pp. 993–1005. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (2022)
- [39] White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., Schmidt, D.C.: A prompt pattern catalog to enhance prompt engineering with chatgpt. CoRR **abs/2302.11382** (2023)

- [40] Adomavicius, G., Tuzhilin, A.: Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. Knowl. Data Eng.* **17**(6), 734–749 (2005)
- [41] Dauphin, Y.N., Fan, A., Auli, M., Grangier, D.: Language modeling with gated convolutional networks. In: *ICML. Proceedings of Machine Learning Research*, vol. 70, pp. 933–941. PMLR, ??? (2017)
- [42] Zhang, A., Lipton, Z.C., Li, M., Smola, A.J.: Dive into deep learning. *CoRR abs/2106.11342* (2021)
- [43] Abdollahpouri, H., Burke, R., Mobasher, B.: Managing popularity bias in recommender systems with personalized re-ranking. In: *FLAIRS*, pp. 413–418. AAAI Press, Palo Alto, California (2019)
- [44] Yin, H., Cui, B., Li, J., Yao, J., Chen, C.: Challenging the long tail recommendation. *Proc. VLDB Endow.* **5**(9), 896–907 (2012)
- [45] Zhu, Z., Wang, J., Caverlee, J.: Measuring and mitigating item under-recommendation bias in personalized ranking systems. In: *SIGIR*, pp. 449–458. ACM, New York, NY, USA (2020)
- [46] Bellogín, A., Castells, P., Cantador, I.: Statistical biases in information retrieval metrics for recommender systems. *Inf. Retr. J.* **20**(6), 606–634 (2017)
- [47] Cremonesi, P., Koren, Y., Turrin, R.: Performance of recommender algorithms on top-n recommendation tasks. In: *RecSys*, pp. 39–46. ACM, New York, NY, USA (2010)
- [48] Sarwar, B.M., Karypis, G., Konstan, J.A., Riedl, J.: Analysis of recommendation algorithms for e-commerce. In: *EC*, pp. 158–167. ACM, New York, NY, USA (2000)
- [49] Balabanovic, M., Shoham, Y.: Content-based, collaborative recommendation. *Commun. ACM* **40**(3), 66–72 (1997)
- [50] Herlocker, J.L., Konstan, J.A., Riedl, J.: An empirical analysis of design choices in neighborhood-based collaborative filtering algorithms. *Inf. Retr.* **5**(4), 287–310 (2002)
- [51] Breese, J.S., Heckerman, D., Kadie, C.M.: Empirical analysis of predictive algorithms for collaborative filtering. In: *UAI*, pp. 43–52. Morgan Kaufmann, ??? (1998)
- [52] Paudel, B., Christoffel, F., Newell, C., Bernstein, A.: Updatable, accurate, diverse, and scalable recommendations for interactive applications. *ACM Trans. Interact. Intell. Syst.* **7**(1), 1–1134 (2017)

- [53] Steck, H.: Embarrassingly shallow autoencoders for sparse data. In: WWW, pp. 3251–3257. ACM, New York, NY, USA (2019)
- [54] Gantner, Z., Rendle, S., Freudenthaler, C., Schmidt-Thieme, L.: Mymedialite: a free recommender system library. In: RecSys, pp. 305–308. ACM, New York, NY, USA (2011)
- [55] Noia, T.D., Mirizzi, R., Ostuni, V.C., Romito, D., Zanker, M.: Linked open data to support content-based recommender systems. In: I-SEMANTICS, pp. 1–8. ACM, New York, NY, USA (2012)
- [56] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G.: Llama: Open and efficient foundation language models. CoRR **abs/2302.13971** (2023)
- [57] Dacrema, M.F., Boglio, S., Cremonesi, P., Jannach, D.: A troubling analysis of reproducibility and progress in recommender systems research. ACM Trans. Inf. Syst. **39**(2), 20–12049 (2021)
- [58] Sun, Z., Yu, D., Fang, H., Yang, J., Qu, X., Zhang, J., Geng, C.: Are we evaluating rigorously? benchmarking recommendation for reproducible evaluation and fair comparison. In: RecSys, pp. 23–32. ACM, New York, NY, USA (2020)
- [59] Anelli, V.W., Noia, T.D., Sciascio, E.D., Ragone, A., Trotta, J.: Semantic interpretation of top-n recommendations. IEEE Trans. Knowl. Data Eng. **34**(5), 2416–2428 (2022)
- [60] Ferrara, A., Anelli, V.W., Mancino, A.C.M., Noia, T.D., Sciascio, E.D.: Kgflex: Efficient recommendation with sparse feature factorization and knowledge graphs. ACM Trans. Recomm. Syst. (2023)