

A Survey on Large Language Models for Recommendation

Likang Wu^{1,2*}, Zhi Zheng^{1,2*}, Zhaopeng Qiu^{2*}, Hao Wang^{1†}, Hongchao Gu¹, Tingjia Shen¹, Chuan Qin², Chen Zhu², Hengshu Zhu^{2†}, Qi Liu¹, Hui Xiong^{3†}, Enhong Chen^{1†}

¹University of Science and Technology of China, ²Career Science Lab, BOSS Zhipin, ³Hong Kong University of Science and Technology (Guangzhou)

{wulk,zhengzhi97,hcgu,jts_stj}@mail.ustc.edu.cn,

{zhpengqiu,chuanqin0426,zc3930155,zhuhengshu}@gmail.com,

{wanghao3,qiliuql,cheneh}@ustc.edu.cn, xionghui@ust.hk

Abstract

Large Language Models (LLMs) have emerged as powerful tools in the field of Natural Language Processing (NLP) and have recently gained significant attention in the domain of Recommendation Systems (RS). These models, trained on massive amounts of data using self-supervised learning, have demonstrated remarkable success in learning universal representations and have the potential to enhance various aspects of recommendation systems by some effective transfer techniques such as fine-tuning and prompt tuning, and so on. The crucial aspect of harnessing the power of language models in enhancing recommendation quality is the utilization of their high-quality representations of textual features and their extensive coverage of external knowledge to establish correlations between items and users. To provide a comprehensive understanding of the existing LLM-based recommendation systems, this survey presents a taxonomy that categorizes these models into two major paradigms, respectively Discriminative LLM for Recommendation (DLLM4Rec) and Generative LLM for Recommendation (GLLM4Rec), with the latter being systematically sorted out for the first time. Furthermore, we systematically review and analyze existing LLM-based recommendation systems within each paradigm, providing insights into their methodologies, techniques, and performance. Additionally, we identify key challenges and several valuable findings to provide researchers and practitioners with inspiration.

1 Introduction

Recommendation systems play a critical role in assisting users in finding relevant and personalized items or content. With the emergence of Large Language Models (LLMs) in Natural Language Processing (NLP), there has been a growing interest in harnessing the power of these models to enhance recommendation systems.

*Equal Contribution.

†Corresponding Author.

The key advantage of incorporating LLMs into recommendation systems lies in their ability to extract high-quality representations of textual features and leverage the extensive external knowledge encoded within them [Liu *et al.*, 2023b]. And this survey views LLM as the Transformer-based model with a large number of parameters, trained on massive datasets using self/semi-supervised learning techniques, e.g., BERT, GPT series, PaLM series, etc[‡]. Different from traditional recommendation systems, the LLM-based models excel in capturing contextual information, comprehending user queries, item descriptions, and other textual data more effectively [Geng *et al.*, 2022]. By understanding the context, LLM-based RS can improve the accuracy and relevance of recommendations, leading to enhanced user satisfaction. Meanwhile, facing the common data sparsity issue of limited historical interactions [Da'u and Salim, 2020], LLMs also bring new possibilities to recommendation systems through zero/few-shot recommendation capabilities [Sileo *et al.*, 2022]. These models can generalize to unseen candidates due to the extensive pre-training with factual information, domain expertise, and common-sense reasoning, enabling them to provide reasonable recommendations even without prior exposure to specific items or users.

The aforementioned strategies are already well-applied in discriminative models. However, with the evolution of AI learning paradigms, generative language models have started to gain prominence [Zhao *et al.*, 2023]. A prime example of this is the emergence of ChatGPT and other comparable models, which have significantly disrupted human life and work patterns. Furthermore, the fusion of generative models with recommendation systems offers the potential for even more innovative and practical applications. For instance, the interpretability of recommendations can be improved, as LLM-based systems are able to provide explanations based on their language generation capabilities [Gao *et al.*, 2023], helping users understand the factors influencing the recommendations. Moreover, generative language models enable more personalized and context-aware recommendations, such as users' customizable prompts [Li *et al.*, 2023] in the chat-based recommendation system, enhancing user engagement and satisfaction with the diversity of results.

Motivated by the remarkable effectiveness of the afo-

LLM definition

why
LLMs
are
advantageous

better for
cold starts

LLMs can
explain
recommendations

context-based
recommendations
are possible

[‡]https://en.wikipedia.org/wiki/Large_language_model

mentioned paradigms in solving data sparsity and efficiency issues, the adaptation of language modeling paradigms for recommendation has emerged as a promising direction in both academia and industry, significantly advancing the state-of-the-art in the research of recommendation systems. So far, there are a few studies that review relevant papers in this domain [Zeng *et al.*, 2021; Liu *et al.*, 2023b]. Zeng *et al.* (2021) summarizes some research on the pre-training of recommendation models and discusses knowledge transfer methods between different domains. Liu *et al.* (2023b) proposes an orthogonal taxonomy to divide existing pre-trained language model-based recommendation systems w.r.t. their training strategies and objectives, analyzes and summarizes the connection between pre-trained language model-based training paradigms and different input data types. However, both of these surveys primarily focus on the transfer of training techniques and strategies in pretraining language models, rather than exploring the potential of language models and their capabilities, i.e., LLM-based way. Additionally, they lack a comprehensive overview of the recent advancements and systematic introductions of generative large language models in the recommendation field. To address this issue, we delve into LLM-based recommendation systems, categorizing them into **discriminative LLMs for recommendation** and **generative LLMs for recommendation**, and the focus of our review is on the latter. To the best of our knowledge, our survey is the first work that concludes an up-to-date and comprehensive review of generative large language models for recommendation systems. The main contributions of our survey are summarized as follows: **discriminative/generative**

- We present a systematic survey of the current state of LLM-based recommendation systems, focusing on **expanding the capacity of language models**. By analyzing the existing methods, we provide a systematic overview of related advancements and applications.
- To the best of our knowledge, our survey is the first comprehensive and up-to-date review specifically dedicated to generative large language models for recommendation systems.
- Our survey critically analyzes the **advantages, disadvantages, and limitations of existing methods**. We identify **key challenges** faced by LLM-based recommendation systems and propose valuable findings that can inspire further research in this potential field.

2 Modeling Paradigms and Taxonomy

inputs
outputs

The basic framework of all large language models is composed of several **transformer blocks**, e.g., GPT, PaLM, LLaMA, etc. The input of this architecture is generally composed of **token embeddings** or **position embeddings** and so on, while the **expected output embedding or tokens** can be obtained at the output module. Here, both the input and output data types are textual sequences. As shown in (1)-(3) in Figure 1, for the adaption of language models in recommendations, i.e., the modeling paradigm, existing work can be roughly divided into the following three categories:

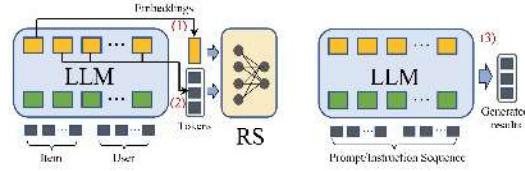


Figure 1: Three modeling paradigms of the research for large language models on recommendation systems.

- (1) **LLM Embeddings + RS**. This modeling paradigm views the language model as a **feature extractor**, which feeds the **features of items and users** into LLMs and outputs corresponding **embeddings**. A traditional RS model can utilize **knowledge-aware embeddings** for various recommendation tasks.
- (2) **LLM Tokens + RS**. Similar to the former method, this approach generates **tokens** based on the **inputted items** and **users' features**. The generated tokens capture **potential preferences through semantic mining**, which can be integrated into the decision-making process of a recommendation system.
- (3) **LLM as RS**. Different from (1) and (2), this paradigm aims to directly transfer pre-trained LLM into a powerful recommendation system. The **input sequence** usually consists of the **profile description**, **behavior prompt**, and **task instruction**. The output sequence is expected to offer **context and task instruction** a reasonable recommendation result.

In practical applications, the choice of language model significantly influences the design of modeling paradigms in recommendation systems. As shown in Figure 2, in this paper, we categorize existing works into two main categories, respectively discriminative LLMs and generative LLMs for recommendation. The taxonomy of LLMs for recommendation can be further subdivided based on the training manner, and the distinction among different manners is illustrated in Figure 3. Generally, discriminative language models are well-suited for embedding within the paradigm (1), while the response generation capability of generative language models further supports paradigms (2) or (3).

3 Discriminative LLMs for Recommendation

Indeed, so-called **discriminative language** models in the recommendation area mainly refer to those models of **BERT series** [Devlin *et al.*, 2019]. Due to the expertise of discriminative language models in natural language understanding tasks, they are often considered as **embedding backbones for downstream tasks**. This holds true for recommendation systems as well. Most existing works align the representations of pre-trained models like BERT with the **domain-specific data** through **fine-tuning**. Additionally, some research explores training strategies like **prompt tuning**. The representative approaches and common-used datasets are listed in Table 1 and Table 2.

returning
embeddings

semantic
mining

context and
task
instruction

discriminative:
embeddings
generative:
semantic mining/
task instruction

discriminative
BERT series
embedding backbones

specific data →
fine tuning
prompt + tuning

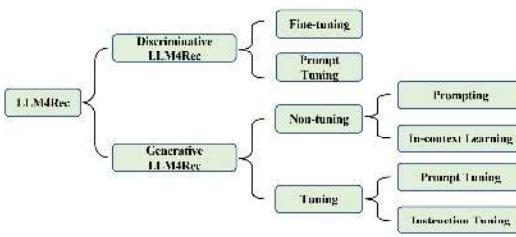


Figure 2: A taxonomy of the research for large language models on recommendation systems.

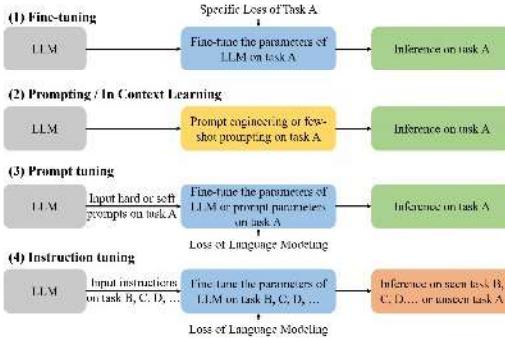


Figure 3: Detailed explanation of five different training (domain adaption) manners for LLM-based recommendations.

3.1 Fine-tuning

Fine-tuning pre-trained language models is a universal technique that has gained significant attention in various natural language processing (NLP) tasks, including recommendation systems. The idea behind fine-tuning is to take a **language model**, which has already learned rich linguistic representations from large-scale text data, and adapt it to a specific task or domain by further training it on task-specific data.

*fine tuning with specific data
parameters are updated*

The process of fine-tuning involves **initializing the pre-trained language model with its learned parameters and then training it on a recommendation-specific dataset**. This dataset typically includes **user-item interactions, textual descriptions of items, user profiles**, and other relevant contextual information. During fine-tuning, the model's parameters are **updated based on the task-specific data**, allowing it to adapt and specialize for recommendation tasks. The learning objectives can be different in the pre-training and fine-tuning stages.

Since the fine-tuning strategy is flexible, most bert-enhanced recommendation methods can be summarized into this track. For the basic representation task, Qiu *et al.* (2021) proposed a novel pre-training and fine-tuning-based approach U-BERT to learn users' representation, which leveraged content-rich domains to complement those users' feature with insufficient behavior data. A review co-matching layer is designed to capture implicit semantic interactions between the reviews of users and items. Similarly, in UserBERT [Wu *et al.*, 2021b], two self-supervision tasks are incorporated for

user model pre-training on unlabeled behavior data to empower user modeling. This model utilizes medium-hard contrastive learning, masked behavior prediction, and behavior sequence matching to train accurate user representation via captured inherent user interests and relatedness.

The pre-trained BERT achieved outstanding breakthroughs in the ranking task as well. BECR [Yang *et al.*, 2022] proposed a lightweight composite re-ranking scheme that combined deep contextual token interactions and traditional lexical term-matching features at the same time. With a novel composite token encoding, BECR effectively approximates the query representations using pre-computable token embeddings based on uni-grams and skip-n-grams, allowing for a reasonable tradeoff between ad-hoc ranking relevance and efficiency. Besides, Wu *et al.* (2022) proposed an end-to-end multi-task learning framework for product ranking with fine-tuned domain-specific BERT to address the issue of vocabulary mismatch between queries and products. The authors utilized the mixture-of-experts layer and probability transfer between tasks to harness the abundant engagement data.

There are also many related studies in other specific tasks or scenarios, e.g., group recommendation [Zhang *et al.*, 2022], search/matching [Yao *et al.*, 2022], CTR prediction [Muhammed *et al.*, 2021]. Especially, the 'pre-train, fine-tuning' mechanism played an important role in several sequential or session-based recommendation systems, such as BERT4Rec [Sun *et al.*, 2019], RESETBERT4Rec [Zhao, 2022]. However, the above models only leveraged the advantages of the training strategy rather than expanding the large language model into the recommendation field, so it was not the focus of our discussion. The sequence representation learning model UniSRec [Liu *et al.*, 2022] developed a BERT-fine-tuned framework, which associated description text of items to learn transferable representations across different recommendation scenarios. For the content-based recommendation, especially news recommendation, NRMS [Wu *et al.*, 2021a], Tiny-NewsRec [Yu *et al.*, 2022], PREC [Liu *et al.*, 2022], exploited large language models to empower news recommendation via handling known domain shift problems or reducing transfer cost.

In summary, the integration of BERT fine-tuning into recommendation systems uses the powerful external knowledge and personalized user preference, which primarily aims to promote recommendation accuracy and simultaneously obtains a little cold-start handling capability for new items with limited historical data.

3.2 Prompt Tuning

Instead of adapting LLMs to different downstream recommendation tasks by designing specific objective functions, **prompt tuning** [Lester *et al.*, 2021] tries to align the tuning object of recommendation with pre-trained loss through **hard/soft prompts and label word verbalizer**. For example, Penha and Hauff (2020) leveraged BERT's Masked Language Modeling (MLM) head to uncover its understanding of item genres using cloze-style prompts. They further utilized BERT's Next Sentence Prediction (NSP) head and similarity (SIM) of representations to compare relevant and non-relevant search and recommendation query-document inputs.

BECR, ranking of results

*UniSRec
Tiny-NewsRec*

prompt tuning

*prompt
ensembling*

The experiment told that BERT, without any fine-tuning, can prioritize relevant items in the ranking process. Yang *et al.* (2021) developed a conversational recommendation system with prompts, where a BERT-based item encoder directly mapped the metadata of each item to an embedding. Recently, Prompt4NR [Zhang and Wang, 2023] pioneered the application of the prompt learning paradigm for news recommendation. This framework redefined the objective of predicting user clicks on candidate news as a cloze-style mask-prediction task. The experiments found that the performance of recommendation systems is noticeably enhanced through the utilization of multi-prompt ensembling, surpassing the results achieved with a single prompt on discrete and continuous templates. This highlights the effectiveness of prompt ensembling in combining multiple prompts to make more informed decisions.

*generative
recommendation
→ NLP*

4 Generative LLMs for Recommendation

Compared to discriminative models, generative models have better natural language generation capabilities. Therefore, unlike most discriminative model-based approaches that align the representation learned by LLMs to the recommendation domain, most generative model-based work translates recommendation tasks as natural language tasks, and then applies techniques such as in-context learning, prompt tuning, and instruction tuning to adapt LLMs to directly generate the recommendation results. Moreover, with the impressive capabilities demonstrated by ChatGPT, this type of work has received increased attention recently.

As shown in Figure 2, according to whether tuning parameters, these generative LLM-based approaches can be further subdivided into two paradigms: *non-tuning paradigm* and *tuning paradigm*. The following two sub-sections will address their details, respectively. The representative approaches and common-used datasets are also listed in Table 1 and Table 2.

4.1 Non-tuning Paradigm

*LLMs: success in
virgin tasks
+ same as a
recommendation
+ in-context
learning*

The LLMs have shown strong zero/few-shot abilities in many unseen tasks [Brown *et al.*, 2020; Ouyang *et al.*, 2022]. Hence, some recent works assume LLMs already have the recommendation abilities, and attempt to trigger these abilities by introducing specific prompts. They employ the recent practice of *instruction* and *In-Context Learning* [Brown *et al.*, 2020] to adopt the LLMs to recommendation tasks without tuning model parameters. According to whether the prompt includes the demonstration examples, the studies in this paradigm mainly belong to the following two categories: *prompting* and *in-context learning*.

*better
instructions*

Prompting
This category of works aims to design more suitable instructions and prompts to help LLMs better understand and solve the recommendation tasks. Liu *et al.* (2023a) systematically evaluated the performance of ChatGPT on five common recommendation tasks, i.e., *rating prediction*, *Sequential Recommendation*, *direct recommendation*, *explanation generation*, and *review summarization*. They proposed a general

recommendation prompt construction framework, which consists of: (1) *task description*, adapting recommendation tasks to natural language processing tasks; (2) *behavior injection*, incorporating user-item interaction to aid LLMs in capturing user preferences and needs; (3) *format indicator*, constraining the output format and making the recommendation results more comprehensible and assessable. Similarly, Dai *et al.* (2023) conducted an empirical analysis of ChatGPT's recommendation abilities on three common information retrieval tasks, including point-wise, pair-wise, and list-wise ranking. They proposed different prompts for different kinds of tasks and introduced the *role instructions* (such as *You are a news recommendation system now*) at the beginning of the prompts to enhance the domain adaption ability of ChatGPT.

*recommendation
prompt
construction*

*role instruction
↓
enhance domain
adaption*

*bootstrapping/
sequential
prompting*

*use cases:
news
+, movies*

*LLMs can control
the recommendation
system*

*addresses cold-start
existing items vs.
new recommendations*

items by AIGC models.

In summary, these studies utilize natural language prompts to activate the zero-shot capability of LLM in recommendation tasks, providing a low-cost and practical solution.

In-context Learning

In-context learning is a technique used by GPT-3 and other LLMs to quickly adapt to new tasks and information. With a few demonstration input-label pairs, they can predict the label for an unseen input without additional parameter updates [Dai *et al.*, 2022]. Hence, some works attempt to add demonstration examples in the prompt to make LLMs better understand the recommendation tasks. For sequential recommendation, Hou *et al.* (2023) introduced demonstration examples by augmenting the input interaction sequence itself. In detail, they paired the prefix of the input interaction sequence and the corresponding successor as examples. Liu *et al.* (2023a) and Dai *et al.* (2023) designed the demonstration example templates for various recommendation tasks and the experimental results also showed the in-context learning method will improve the recommendation abilities of LLMs on most tasks.

However, in comparison to prompting, only a few studies have explored the use of In-context Learning of Language Models (LLMs) in recommendation tasks. Numerous open questions remain, including the selection of demonstration examples and the influence of the number of demonstration examples on recommendation performance.

4.2 Tuning Paradigm

As we mentioned above, LLMs have strong zero/few-shot abilities, and their recommendation performance can significantly surpass random guessing with appropriate prompt design. However, it is not surprising that recommendation systems constructed in this manner fail to surpass the performance of recommendation models trained specifically for a given task on specific data. Therefore, many researchers aim to enhance the recommendation ability of LLMs by further fine-tuning or prompt learning. In this paper, following [Wei *et al.*, 2022], we categorize the paradigm of the tuning methods into two different types, respectively prompt tuning and instruction tuning. Specifically, under the prompt tuning paradigm, the parameters of LLMs or soft prompts are tuned on a specific task, e.g., rating prediction, while the LLMs are fine-tuned for multiple tasks with different types of instructions under the instruction tuning paradigm. Therefore, the LLMs can get better zero-shot abilities by instruction tuning. However, we propose that currently there is no clear delineation or universally agreed definition regarding these two fine-tuning paradigms.

Prompt Tuning

In this paradigm, LLMs typically take the user/item information as input, and output the user preference (e.g., like or unlike) or ratings (e.g., 1-5) for the items. For example, Kang *et al.* (2023) proposed to format the user historical interactions as prompts, where each interaction is represented by information about the item, and formulated the rating prediction task as two different tasks, respectively multi-class classification and regression. Kang *et al.* (2023) further investigated

various LLMs in different sizes, ranging from 250M to 540B parameters and evaluate their performance in zero-shot, few-shot, and fine-tuning scenarios, and found that FLAN-T5-XXL (11B) model with fine-tuning can achieve the best result. Bao *et al.* (2023) proposed TALLRec which is trained by two tuning stages. Specifically, TALLRec is first fine-tuned based on the self-instruct data by Alpaca [Taori *et al.*, 2023]. Then, TALLRec is further fine-tuned by recommendation tuning, where the input is the historical sequence of users and the output is the "yes or no" feedback. Zhang *et al.* (2021) proposed to format the recommendation task as a next-token-prediction problem, and evaluated the proposed method on a movie recommendation dataset in zero-shot and fine-tuned settings. Zhang *et al.* (2021) also observed that language models show clear linguistic biases for recommendations, furthermore, fine-tuned LLMs can learn how to recommend but linguistic biases still exist for bottom predictions.

In addition to directly fine-tuning the LLMs, some studies also proposed to utilize prompt learning to achieve better performance. For example, Wang *et al.* (2022) designed a unified conversational recommendation system named UniCRS based on knowledge-enhanced prompt learning. In this paper, the authors proposed to freeze the parameters of LLMs, and trained the soft prompts for response generation and item recommendation by prompt learning. Li *et al.* (2023) proposed to provide user-understandable explanations based on the generation ability of LLMs. The authors tried both discrete prompt learning and continuous prompt learning, and further proposed two training strategies, respectively sequential tuning and recommendation as regularization.

Instruction Tuning

In this paradigm, LLMs are fine-tuned for multiple tasks with different types of instructions. In this way, LLMs can better align with human intent and achieve better zero-shot ability. For example, Geng *et al.* (2022) proposed to fine-tune a T5 model on five different types of instructions, respectively sequential recommendation, rating prediction, explanation generation, review summarization, and direct recommendation. After the multitask instruction tuning on recommendation datasets, the model can achieve the capability of zero-shot generalization to unseen personalized prompts and new items. Similarly, Cui *et al.* (2022) proposed to fine-tune an M6 model on three types of tasks, respectively scoring tasks, generation tasks and retrieval tasks. Zhang *et al.* (2023) first designed a general instruction format from three types of key aspects, respectively preference, intention and task form. Then, Zhang *et al.* (2023) manually designed 39 instruction templates and automatically generated a large amount of user-personalized instruction data for instruction tuning on a 3B FLAN-T5-XL model. The experiment results demonstrated that this approach can outperform several competitive baselines including GPT-3.5.

best model

linguistic biases
help and hurt

freeze parameters
+ train soft prompts

] multi-task
instruction
tuning

providing
an example
in the
prompt

number/
selection of
demonstration
examples

soft
prompts
+
instruction
tuning

user/item
information
↓
rating
prediction
↓
user
preference

Table 1: A list of some representative LLM-based recommendation methods.

Adaption Way	Paper	Base Model	Recommendation Task	Modeling Paradigm	Source Code
Discriminative LLMs for Recommendation					
Fine-tuning	[Wu <i>et al.</i> , 2021a]	BERT/RoBERTa/UniLM	News Recommendation	LLM Embeddings + RS	https://shorturl.at/cimy7
	[Qiu <i>et al.</i> , 2021]	BERT	User Representation	LLM Embeddings + RS	N/A
	[Zhang <i>et al.</i> , 2022]	BERT	Group Recommendation	LLM as RS	N/A
	[Yao <i>et al.</i> , 2022]	BERT	Search/Matching	LLM Embeddings + RS	https://shorturl.at/suJ69
	[Muhammed <i>et al.</i> , 2021]	BERT	CTR Prediction	LLM Embeddings + RS	N/A
	[Xiao <i>et al.</i> , 2022]	BERT/RoBERTa	Conversational RS	LLM Embeddings + RS	https://shorturl.at/vSUZ8
Prompt Tuning	[Zhang and Wang, 2023]	BERT	Sequential Recommendation	LLM as RS	https://shorturl.at/ehOT0
	[Yang <i>et al.</i> , 2021]	DistilBERT/GPT-2	Conversational RS	LLM as RS	https://shorturl.at/gkuxz
	[Penha and Hauff, 2020]	BERT	Conversational RS	LLM as RS	https://shorturl.at/mqzEY
Generative LLMs for Recommendation					
Non-tuning	[Liu <i>et al.</i> , 2023c]	ChatGPT	News Recommendation	LLM Tokens - RS	https://shorturl.at/jkFST
	[Wang <i>et al.</i> , 2022]	DialoGPT/RoBERTa	Conversational RS	LLM Tokens - RS / LLM as RS	https://shorturl.at/sEU8
	[Silco <i>et al.</i> , 2022]	GPT-2	Sequential Recommendation	LLM as RS	https://shorturl.at/EJK29
	[Wang and Lim, 2023]	GPT-3.5	Sequential Recommendation	LLM Tokens - RS / LLM as RS	https://shorturl.at/qKU38
	[Gao <i>et al.</i> , 2023]	ChatGPT/GPT-3.5	Sequential Recommendation	LLM as RS	N/A
	[Wang <i>et al.</i> , 2023]	ChatGPT	Generative Recommendation	LLM as RS	https://shorturl.at/dBEP5
	[Hou <i>et al.</i> , 2023]	ChatGPT	Sequential Recommendation	LLM as RS	https://shorturl.at/KM056
	[Sun <i>et al.</i> , 2023]	ChatGPT/GPT-3.5	Passage Reranking	LLM as RS	https://shorturl.at/eATFY8
	[Li <i>et al.</i> , 2023a]	ChatGPT	Five Tasks	LLM as RS	N/A
	[Dai <i>et al.</i> , 2023]	ChatGPT/GPT-3.5	Sequential Recommendation	LLM as RS	https://shorturl.at/ijpT3
Tuning	[Zhang <i>et al.</i> , 2023]	FLAN T5	Three Tasks	LLM as RS	N/A
	[Kang <i>et al.</i> , 2023]	FLAN-T5/ChatGPT	Rating Prediction	LLM as RS	N/A
	[Bao <i>et al.</i> , 2023]	LLaMA-7B	Movie/Book RS	LLM as RS	https://shorturl.at/coET1
	[Li <i>et al.</i> , 2023]	GPT-2	Explainable RS	LLM as RS	https://shorturl.at/adT09
	[Geng <i>et al.</i> , 2022]	T5	Five Tasks	LLM as RS	https://shorturl.at/2CRY19
	[Cui <i>et al.</i> , 2022]	M6	Five Tasks	LLM as RS	N/A
				Zhang <i>et al.</i> , 2021	GPT-2
					Movie Recommendation
					LLM as RS
					N/A

els in recommendation systems, especially for generative language models. We have identified their potential to improve the performance of traditional recommendation models in specific tasks. However, it is necessary to note that the overall exploration in this field is still in the early stage. Researchers may find it challenging to determine the most worthwhile problems and pain points to investigate. To address this, we have summarized the common findings presented by numerous studies on large-scale model recommendations. These findings highlight certain technical challenges and present potential opportunities for further advancements in the field.

5.1 Model Bias

Position Bias. In the generative language modeling paradigm of recommendation systems, various information such as user behavior sequences and recommended candidates are input to the language model in the form of textual sequential descriptions, which can introduce some position biases inherent in the language model itself [Lu *et al.*, 2021]. For example, the order of candidates affects the ranking results of LLM-based recommendation models, i.e., LLM often prioritizes the items in the top order. And the model usually cannot capture the behavior order of the sequence well. Hou *et al.* (2022) used the random sampling-based bootstrapping to alleviate the position bias of candidates and emphasized the recently interacted items to enhance behavior order. However, these solutions are

not adaptive enough, and more robust learning strategies are needed in the future.

Popularity Bias. The ranking results of LLMs are influenced by the popularity levels of the candidates. Popular items, which are often extensively discussed and mentioned in the pre-training corpora of LLMs, tend to be ranked higher. Addressing this issue is challenging as it is closely tied to the composition of the pre-trained corpus.

Fairness Bias. Pre-trained language models have exhibited fairness issues related to sensitive attributes, which are influenced by the training data or the demographics of the individuals involved in certain task annotations [Ferrara, 2023]. These fairness concerns can result in models making recommendations that assume users belong to a specific group, potentially leading to controversial issues when deployed commercially. One example is the bias in recommendation results caused by gender or race. Addressing these fairness issues is crucial to ensure equitable and unbiased recommendations.

working against the recommendation of popular items

gender/race bias becomes an issue

using sufficient representation of data

sequential order in prompt has effects
random sampling bootstrapping

Table 2: A list of common datasets used in existing LLM-based recommendation methods.

Name	Scene	Tasks	Information	URL
Amazon Review	Commerce	Seq Rec/CF Rec	This is a large crawl of product reviews from Amazon. Ratings: 82.83 million, Users: 20.98 million, Items: 9.35 million, Timespan: May 1996 - July 2014	http://jmcauley.ucsd.edu/data/amazon/
Steam	Game	Seq Rec/CF Rec	Reviews represent a great opportunity to break down the satisfaction and dissatisfaction factors around games. Reviews: 7,793,069, Users: 2,567,538, Items: 15,474, Bundles: 615	https://eseweb.ucsd.edu/~jmcauley/datasets.html#steam
MovieLens	Movie	General	The dataset consists of 4 sub-datasets, which describe users' ratings to movies and free-text tagging activities from MovieLens, a movie recommendation service.	https://grouplens.org/datasets/movielens/
Yelp	Commerce	General	There are 6,990,280 reviews, 150,346 businesses, 200,100 pictures, 11 metropolitan areas, 908,915 tips by 1,987,897 users. Over 1.2 million business attributes like hours, parking, availability, etc.	https://www.yelp.com/dataset
Douban	Movie, Music, Book	Seq Rec/CF Rec	This dataset includes three domains, i.e., movie, music, and book, and different kinds of raw information, i.e., ratings, reviews, item details, user profiles, tags (labels), and date.	https://paperswithcode.com/dataset/douban
MIND	News	General	MIND contains about 160k English news articles and more than 15 million impression logs generated by 1 million users. Every news contains textual content including title, abstract, body, category, and entities.	https://microsoft.github.io/assets/doc/ACL2020_MIND.pdf
U-NEXT	Commerce	Conversation Rec	U-NEXT consists of 7,698 fine-grained annotated pre-sales dialogues, 333,879 user behaviors, and 332,148 product knowledge tuples,	https://github.com/LeeeoLiu/U-NEXT

tionally, it is critical to translate a user's heterogeneous behavior sequence (such as clicks, adding to cart, and purchases in the e-commerce domain) into natural language for preference modeling. ID-like features have been proven effective in traditional recommendation models, but incorporating them into prompts to improve personalized recommendation performance is also challenging.

Limited Context Length. The context length limitation of LLMs will constrain the length of users' behavioral sequences and the number of candidate items, resulting in sub-optimal performance [Zhang *et al.*, 2023]. Existing work has proposed some techniques to alleviate this problem, such as selecting representative items from user behavior sequence [Wang and Lim, 2023] and sliding window strategy for candidate list [Sun *et al.*, 2023].

5.3 Promising Ability

Zero/Few-shot Recommendation Ability. The experimental results on multiple domain datasets indicate that LLMs possess impressive zero/few-shot abilities in various recommendation tasks [Hou *et al.*, 2023; Liu *et al.*, 2023a; Dai *et al.*, 2023]. It is worth noting that few-shot learning, which is equivalent to in-context learning, does not change the parameters of LLMs. This suggests LLMs have the potential to mitigate the cold-start problem with limited data. However, there are still some open questions, such as the need for clearer guidance in selecting representative and effective demonstration examples for few-shot learning, as well as the need for experimental results across more domains to further support the conclusion regarding the zero/few-shot recommendation abilities.

Explainable Ability. Generative LLMs exhibit a remarkable ability for natural language generation. Thus, a natural thought is using LLMs to conduct explainable recommendation via text generation manner. Liu *et al.* (2023a) conduct a comparison experiment among ChatGPT and some baselines on explanation generation task. The results demonstrate that even without fine-tuning and under the in-context learning setting, ChatGPT still performs better than some supervised traditional methods. Moreover, according to human eval-

uation, ChatGPT's explanations are deemed even clearer and more reasonable than the ground truth. Encouraged by these exciting preliminary experimental results, the performance of fine-tuned LLMs in explainable recommendation is expected to be promising.

5.4 Evaluation Issues

Generation Controlling. As we mentioned before, many studies have employed large-scale models as recommendation systems by providing carefully designed instructions. For these LLMs, the output should strictly adhere to the given instruction format, such as providing binary responses (yes or no) or generating a ranked list. However, in practical applications, the output of LLMs may deviate from the desired output format. For instance, the model may produce responses in incorrect formats or even refuse to provide an answer [Dai *et al.*, 2023]. Therefore, addressing the challenge of ensuring better control over the output of LLMs is a pressing issue that needs to be resolved.

Evaluation Criteria. If the task performed by LLMs are standard recommendation tasks, such as rating prediction or item ranking, we can employ existing evaluation metrics for evaluation, e.g., NDCG, MSE, etc. However, LLMs also have strong generative capabilities, making them suitable for generative recommendation tasks [Wang *et al.*, 2023]. Following the generative recommendation paradigm, LLMs can generate items that have never appeared in the historical data and recommend them to users. In this scenario, evaluating the generative recommendation capability of LLMs remains an open question.

Datasets. Currently, most of the research in this area primarily tests the recommendation capability and zero/few-shot capability of LLMs using datasets like MovieLens, Amazon Books, and similar benchmarks. However, this may bring the following two potential issues. First, compared to real-world industrial datasets, these datasets are relatively small in scale and may not fully reflect the recommendation capability of LLMs. Second, the items in these datasets, such as movies and books, may have related information that appeared in the pre-training data of LLMs. This could introduce bias in eval-

output must be controlled more precisely

tougher to evaluate + evaluations

bias in recommendation
→ contains info about datasets

limit to number of candidate items

does not change the parameters
helps cold-start

explaining results

no
benchmark
for
evaluation

must
re-train for
specific
domains
large
computational
costs

uating the few-zero-shot learning capability of LLMs. Currently, we still lack a suitable benchmark for conducting a more comprehensive evaluation.

In addition to the aforementioned prominent findings, there are also some limitations associated with the capabilities of large language models. For example, the challenge of knowledge forgetting may arise when training models for specific domain tasks or updating model knowledge [Jang *et al.*, 2022]. Another issue is the distinct performances caused by varying sizes of language model parameters, where using excessively large models would result in excessive computational costs for research and deployment in recommendation systems [Hou *et al.*, 2023]. These challenges also present valuable research opportunities in the field.

6 Conclusion

In this paper, we reviewed the research area of large language models (LLMs) for recommendation systems. We classified existing work into discriminative models and generative models, and then illustrated them in detail by the domain adaption manner. And in order to prevent conceptual confusion, we provided the definition and distinction of fine-tuning, prompting, prompt tuning, and instruction tuning in the LLM-based recommendation. To the best of our knowledge, our survey is the first systematic and up-to-date review specifically dedicated to generative LLMs for recommendation systems, which further summarized the common findings and challenges presented by numerous related studies. Therefore, this survey provided researchers with a valuable resource for gaining a comprehensive understanding of LLM recommendations and exploring potential research directions.

References

- [Bao *et al.*, 2023] Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. *CoRR*, abs/2305.00447, 2023.
- [Brown *et al.*, 2020] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- [Cui *et al.*, 2022] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingen Zhou, and Hongxia Yang. M6-rec: Generative pre-trained language models are open-ended recommender systems. *CoRR*, abs/2205.08084, 2022.
- [Dai *et al.*, 2022] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui, and Furu Wei. Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. *CoRR*, abs/2212.10559, 2022.
- [Dai *et al.*, 2023] Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. Uncovering chatgpt’s capabilities in recommender systems. *CoRR*, abs/2305.02182, 2023.
- [Da’u and Salim, 2020] Aminu Da’u and Naomie Salim. Recommendation system based on deep learning methods: a systematic review and new directions. *Artificial Intelligence Review*, 53(4):2709–2748, 2020.
- [Devlin *et al.*, 2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages 4171–4186. Association for Computational Linguistics, 2019.
- [Ferrara, 2023] Emilio Ferrara. Should chatgpt be biased? challenges and risks of bias in large language models. *arXiv preprint arXiv:2304.03738*, 2023.
- [Gao *et al.*, 2023] Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. Chat-rec: Towards interactive and explainable llms-augmented recommender system. *CoRR*, abs/2303.14524, 2023.
- [Geng *et al.*, 2022] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. Recommendation as language processing (RLP): A unified pretrain, personalized prompt & predict paradigm (P5). In *RecSys*, pages 299–315. ACM, 2022.
- [Hou *et al.*, 2022] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. Towards universal sequence representation learning for recommender systems. In *KDD*, pages 585–593. ACM, 2022.
- [Hou *et al.*, 2023] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian J. McAuley, and Wayne Xin Zhao. Large language models are zero-shot rankers for recommender systems. *CoRR*, abs/2305.08845, 2023.
- [Jang *et al.*, 2022] Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. Towards continual knowledge learning of language models. In *ICLR*, 2022.
- [Kang *et al.*, 2023] Wang-Cheng Kang, Jianmo Ni, Nikhil Mehta, Maheswaran Sathiamoorthy, Lichan Hong, Ed H. Chi, and Derek Zhiyuan Cheng. Do llms understand user preferences? evaluating llms on user rating prediction. *CoRR*, abs/2305.06474, 2023.
- [Lester *et al.*, 2021] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- [Li *et al.*, 2023] Lei Li, Yongfeng Zhang, and Li Chen. Personalized prompt learning for explainable recommendation. *ACM Transactions on Information Systems*, 41(4):1–26, 2023.
- [Liu *et al.*, 2022] Qijiong Liu, Jieming Zhu, Quanyu Dai, and Xiaoming Wu. Boosting deep CTR prediction with a plug-and-play pre-trainer for news recommendation. In *COLING*, pages 2823–2833, 2022.

- [Liu *et al.*, 2023a] Junling Liu, Chao Liu, Renjie Lv, Kang Zhou, and Yan Zhang. Is chatgpt a good recommender? A preliminary study. *CoRR*, abs/2304.10149, 2023.
- [Liu *et al.*, 2023b] Peng Liu, Lemei Zhang, and Jon Atle Gulla. Pre-train, prompt and recommendation: A comprehensive survey of language modelling paradigm adaptations in recommender systems. *arXiv preprint arXiv:2302.03735*, 2023.
- [Liu *et al.*, 2023c] Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. A first look at llm-powered generative news recommendation. *CoRR*, abs/2305.06566, 2023.
- [Lu *et al.*, 2021] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*, 2021.
- [Muhammed *et al.*, 2021] Aashiq Muhammed, Iman Keivanloo, Sujan Perera, James Mracek, Yi Xu, Qingjun Cui, Santosh Rajagopalan, Belinda Zeng, and Trishul Chilimbi. Ctr-bert: Cost-effective knowledge distillation for billion-parameter teacher models. In *NeurIPS Efficient Natural Language and Speech Processing Workshop*, 2021.
- [Ouyang *et al.*, 2022] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Macdzie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022.
- [Penha and Hauff, 2020] Gustavo Penha and Claudia Hauff. What does BERT know about books, movies and music? probing BERT for conversational recommendation. In *RecSys*, pages 388–397. ACM, 2020.
- [Qiu *et al.*, 2021] Zhaopeng Qiu, Xian Wu, Jingyue Gao, and Wei Fan. U-BERT: pre-training user representations for improved recommendation. In *AAAI*, pages 4320–4327. AAAI Press, 2021.
- [Sileo *et al.*, 2022] Damien Sileo, Wout Vossen, and Robbe Raymaekers. Zero-shot recommendation as language modeling. In *ECAIR (2)*, volume 13186 of *Lecture Notes in Computer Science*, pages 223–230. Springer, 2022.
- [Sun *et al.*, 2019] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In *CIKM*, pages 1441–1450. ACM, 2019.
- [Sun *et al.*, 2023] Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. Is chatgpt good at search? investigating large language models as re-ranking agent. *CoRR*, abs/2304.09542, 2023.
- [Taori *et al.*, 2023] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. <https://github.com/tatsu-lab/stanford-alpaca>, 2023.
- [Wang and Lim, 2023] Lei Wang and Ee-Peng Lim. Zero-shot next-item recommendation using large pretrained language models. *CoRR*, abs/2304.03153, 2023.
- [Wang *et al.*, 2022] Xiaolei Wang, Kun Zhou, Ji-Rong Wen, and Wayne Xin Zhao. Towards unified conversational recommender systems via knowledge-enhanced prompt learning. In *KDD*, pages 1929–1937. ACM, 2022.
- [Wang *et al.*, 2023] Wenjie Wang, Xinyu Lin, Fuli Feng, Xiangnan He, and Tat-Seng Chua. Generative recommendation: Towards next-generation recommender paradigm. *CoRR*, abs/2304.03516, 2023.
- [Wei *et al.*, 2022] Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. In *ICLR*. OpenReview.net, 2022.
- [Wu *et al.*, 2021a] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. Empowering news recommendation with pre-trained language models. In *SIGIR*, pages 1652–1656. ACM, 2021.
- [Wu *et al.*, 2021b] Chuhan Wu, Fangzhao Wu, Yang Yu, Tao Qi, Yongfeng Huang, and Xing Xie. Userbert: Contrastive user model pre-training. *arXiv preprint arXiv:2109.01274*, 2021.
- [Wu *et al.*, 2022] Xuyang Wu, Alessandro Magnani, Suthee Chaidaroon, Ajit Puthenputhussery, Ciya Liao, and Yi Fang. A multi-task learning framework for product ranking with BERT. In *WWW*, pages 493–501. ACM, 2022.
- [Xiao *et al.*, 2022] Shitao Xiao, Zheng Liu, Yingxia Shao, Tao Di, Bhuvan Middha, Fangzhao Wu, and Xing Xie. Training large-scale news recommenders with pretrained language models in the loop. In *KDD*, pages 4215–4225. ACM, 2022.
- [Yang *et al.*, 2021] Bowen Yang, Cong Han, Yu Li, Lei Zuo, and Zhou Yu. Improving conversational recommendation systems’ quality with context-aware item meta information. *arXiv preprint arXiv:2112.08140*, 2021.
- [Yang *et al.*, 2022] Yingrui Yang, Yifan Qiao, Jinjin Shao, Xifeng Yan, and Tao Yang. Lightweight composite re-ranking for efficient keyword search with BERT. In *WSDM*, pages 1234–1244. ACM, 2022.
- [Yao *et al.*, 2022] Shaowei Yao, Jiwei Tan, Xi Chen, Juhao Zhang, Xiaoyi Zeng, and Keping Yang. Reprbert: Distilling BERT to an efficient representation-based relevance model for e-commerce. In *KDD*, pages 4363–4371. ACM, 2022.
- [Yu *et al.*, 2022] Yang Yu, Fangzhao Wu, Chuhan Wu, Jingwei Yi, and Qi Liu. Tiny-newsrec: Effective and efficient plm-based news recommendation. In *EMNLP*, pages 5478–5489. Association for Computational Linguistics, 2022.
- [Zeng *et al.*, 2021] Zheni Zeng, Chaojun Xiao, Yuan Yao, Ruobing Xie, Zhiyuan Liu, Fen Lin, Leyu Lin, and Maosong Sun. Knowledge transfer via pre-training for

- recommendation: A review and prospect. *Frontiers in big Data*, 4:602071, 2021.
- [Zhang and Wang, 2023] Zizhuo Zhang and Bang Wang. Prompt learning for news recommendation. *arXiv preprint arXiv:2304.05263*, 2023.
- [Zhang et al., 2021] Yuhui Zhang, Hao Ding, Zeren Shui, Yifei Ma, James Zou, Anoop Deoras, and Hao Wang. Language models as recommender systems: Evaluations and limitations. 2021.
- [Zhang et al., 2022] Song Zhang, Nan Zheng, and Danli Wang. GBERT: pre-training user representations for ephemeral group recommendation. In *CIKM*, pages 2631–2639. ACM, 2022.
- [Zhang et al., 2023] Junjie Zhang, Ruobing Xie, Yupeng Hou, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. Recommendation as instruction following: A large language model empowered recommendation approach. *CoRR*, abs/2305.07001, 2023.
- [Zhao et al., 2023] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [Zhao, 2022] Qihang Zhao. Resetbert4rec: A pre-training model integrating time and user historical behavior for sequential recommendation. In *SIGIR*, pages 1812–1816. ACM, 2022.

1. Why LMs are advantageous for recommendation systems
2. Discriminative LMs
 - a) Fine-tuning

a) fine-tuning

- i. WM embeddings+RS

ii. prompt-tuning

3. Generative WMs

a) Non-Tuning Paradigm

- i. WM Token+RS

- prompting
- in-context learning

b) Tuning Paradigm

- i. WM as RS

- prompt tuning
- instruction tuning

* findings (for generative WMs)

a) Model Biases

- i. Position Biases

- ii. Popularity Biases

- iii. Fairness Biases

b) recommendation prompt designing

- i. user-item representation

- ii. unified context length

c) Promoting fidelity

- i. zero/few shot recommendation ability

- ii. explainable ability

d) Evaluation issues

- i. generation controlling

- ii. evaluation criteria

- iii. datasets