

# Large Language Models for Data Annotation: A Survey

Anonymous ACL submission

## Abstract

Data annotation generally refers to the labeling or generating of raw data with relevant information, which could be used for improving the efficacy of machine learning models. The process, however, is labor-intensive and costly. The emergence of advanced Large Language Models (LLMs), exemplified by GPT-4, presents an unprecedented opportunity to automate the complicated process of data annotation. While existing surveys have extensively covered LLM architecture, training, and general applications, we uniquely focus on their specific utility for data annotation. This survey contributes to three core aspects: LLM-Based Annotation Generation, LLM-Generated Annotations Assessment, and LLM-Generated Annotations Utilization. Furthermore, this survey includes an in-depth taxonomy of data types that LLMs can annotate, a comprehensive review of learning strategies for models utilizing LLM-generated annotations, and a detailed discussion of the primary challenges and limitations associated with using LLMs for data annotation. Serving as a key guide, this survey aims to assist researchers and practitioners in exploring the potential of the latest LLMs for data annotation, thereby fostering future advancements in this critical field.

## 1 Introduction

In the complex realm of machine learning and natural language processing (NLP), data annotation stands out as a critical yet challenging task, extending beyond simple label attachment to encompass a diverse array of fundamental or auxiliary information. This detailed process typically involves ❶ categorizing raw data with class or task labels for basic classification, ❷ adding intermediate labels for contextual depth (Yu et al., 2022), ❸ assigning confidence scores to assess annotation reliability (Lin et al., 2022), ❹ applying alignment or preference labels to tailor outputs to specific criteria or user needs, ❺ annotating entity relationships

to understand how entities within a dataset interact with each other (Wadhwa et al., 2023), ❻ marking semantic roles to define the underlying roles that entities play in a sentence (Larionov et al., 2019), or ❼ tagging temporal sequences to capture the order of events or actions (Yu et al., 2023).

Despite its wide applications, data annotation poses significant challenges for current machine learning models due to the complexity, subjectivity, and diversity of data. This process requires domain expertise and is resource-intensive, particularly when manually labeling large datasets. Advanced LLMs such as GPT-4 (OpenAI, 2023), Gemini (Team et al., 2023), and LLaMA-2 (Touvron et al., 2023b) offer a promising opportunity to revolutionize data annotation. LLMs serve as more than just tools but play a crucial role in improving the effectiveness and precision of data annotation. Their ability to automate annotation tasks (A, 2022), ensure consistency across large volumes of data (Hou et al., 2023), and adapt through fine-tuning or prompting for specific domains (Song et al., 2023), significantly mitigates the challenges encountered with traditional annotation methods, setting a new standard for what is achievable in the realm of NLP. This survey delves into the nuances of using LLMs for data annotation, exploring methodologies, utilizing strategies, and associated challenges in this transformative approach. Through this exploration, we aim to shed light on the motivations behind embracing LLMs as catalysts for redefining the landscape of data annotation in machine learning and NLP. We explore the utilization of LLMs for data annotation in this survey, making four main contributions:

- **LLM-Based Annotation Generation:** We dive into the process of generating annotations for various data types, including instruction & response, rationale, pairwise feedback, textual feedback, and other domain-specific data. Additionally, we discuss the criteria (e.g., diversity and quality) in

the annotation process.

- **Assessing LLM-Generated Annotations:** We explore various methods for assessing the quality of annotations and strategies for selecting high-quality annotations from numerous options.
- **LLM-Generated Annotations Utilization:** We investigate the methodologies at different stages, including supervised fine-tuning, alignment tuning, and inference time, to train machine learning models based on LLM-generated annotations.
- **Social Impact and Future Work:** We discuss issues ranging from ethical dilemmas, such as bias and implications, to technical limitations, including hallucination and efficiency in LLM-generated annotations.

Focusing on this underrepresented aspect of LLM application, the survey aims to serve as a valuable guide for academics and practitioners who intend to deploy LLMs for annotation purposes. Note that in this survey, we primarily focus on pure language models and do not extensively cover recently emerging multimodal LLMs, such as LLaVA (Liu et al., 2023b). Figure 1 illustrates the general structure of this survey. Additionally, a list of potential tools for utilizing LLMs for annotation is included in Appendix A, along with explanatory examples.

**Differences from Other LLM-related Surveys.** While existing surveys in the NLP domain extensively cover architectural nuances (Zhao et al., 2023a), training methodologies (Liu et al., 2023d), and evaluation protocols (Chang et al., 2023) associated with LLMs, their main focus lies on the capabilities of models for specific end tasks such as machine translation (Min et al., 2021), alignment (Wang et al., 2023g), code generation (Zan et al., 2023), and medical analysis (Thirunavukarasu et al., 2023). In contrast, this survey distinguishes itself by focusing primarily on the application of these potent next-generation LLMs to the intricate realm of data annotation, a domain that is crucial yet underexplored.

## 2 Preliminaries

In this section, we delve into our approach to the annotation process. We introduce two core models: an annotator model, denoted as  $\mathcal{A}$ , which maps input data to annotations, and a task learner, represented as  $\mathcal{L}$ , that utilizes or learns from these annotated data to accomplish specific tasks. Our primary focus is on utilizing advanced LLMs like GPT-4 (OpenAI, 2023) and LLaMA (Touvron et al., 2023a) as annotators ( $\mathcal{A}$ ), while the task learner ( $\mathcal{L}$ )

can be another large model (Chiang et al., 2023a) or a less complex one such as BERT (Devlin et al., 2018), which utilizes these annotated data to perform designated tasks. LLM-generated annotations encompass categorical labels and enhance raw data points with a comprehensive array of auxiliary signals. These annotations, including confidence scores, contextual details, and other metadata, extend beyond traditional categorical labels.

## 3 LLM-Based Annotation Generation

The emergence of LLMs has sparked significant interest in their capacity for high-quality, context-sensitive data annotation. This section discusses various kinds of annotations produced via LLMs.

### 3.1 Instruction & Response

Instruction and response are the two fundamental components that constitute a dataset for LLM fine-tuning and in-context learning (ICL). Previous NLP datasets (Li et al., 2017; Wang et al., 2018; Ouyang et al., 2022) mainly rely on human annotators to construct. Recently, with the advent of LLMs, automatic and generative methods (Meng et al., 2022; Ye et al., 2022a,b; Wang et al., 2024c) have gained more focus in data annotation.

**Instruction Diversity.** The diversity of instruction has been proven crucial for LLM learning (Li et al., 2023e; Song et al., 2024b,a). Recent studies have explored various methods to diversify and augment instructions in the original datasets. For example, Yoo et al. (2021) enhance data diversity by mixing two different samples to create a new one. Wang et al. (2022b) use a few manually-written seed instructions and iteratively augment them with a generate-then-filter pipeline. Additionally, Meng et al. (2023); Wang et al. (2023f) train an instruction generation model in the original dataset to augment the diversity of instruction. Gupta et al. (2023) employ a multi-step prompting method to first generate task descriptions, which are then used as instance seeds to guide LLMs in instruction generation. To obtain informative and diverse examples, Wang et al. (2023c) propose an explain-then-generate pipeline with LLMs for iterative data synthesis. Besides, Li et al. (2023a) paraphrase the given sample multiple times to help LLMs understand them from different perspectives. Köksal et al. suggest a clustering-based data selection method to ensure diversity in the initial seed data for augmentation. Recently, Yu et al. (2024) introduce AttrPrompt as an effective way to balance

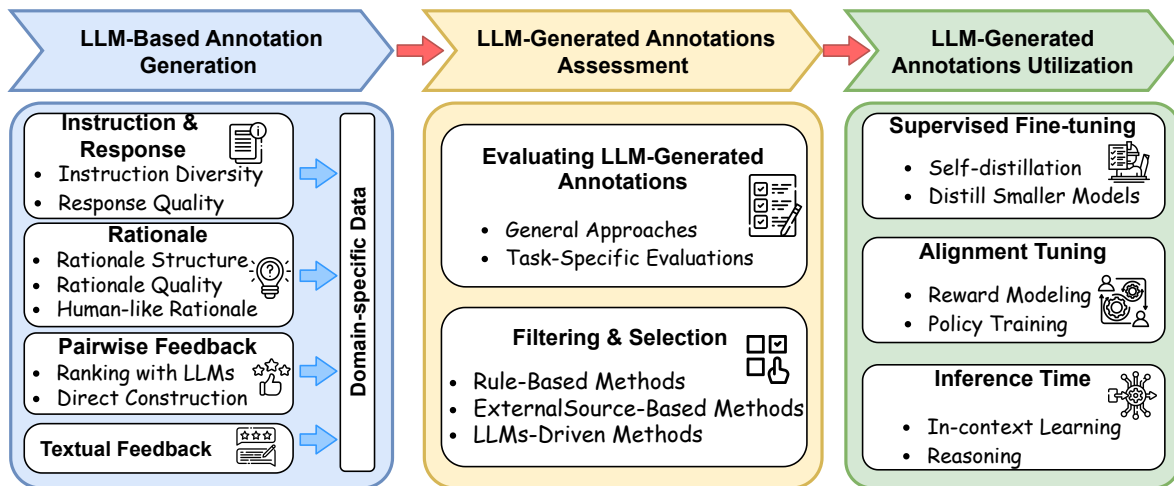


Figure 1: The proposed taxonomy of existing research on LLM for data annotation.

185 diversity and cost in LLM-based data annotation.

186 **Response Quality.** High-quality responses are essential for effective fine-tuning and ICL (Luo et al., 2024). To improve the quality of the generated response, Zhang and Yang (2023a) frame the response generation as reading comprehension tasks and create detailed prompts for LLMs. Huang et al. (2023) adopt self-consistency (Wang et al., 2022b) in response generation, selecting from the candidate response with the highest confidence score. Furthermore, Yang et al. (2024b) propose self-distill and augment the instruction tuning dataset by rewriting the original responses. Pang et al. (2024b) conduct social simulations to ensure high-quality, human-valued responses from LLMs. Moreover, Liu et al. (2024) introduce a multi-step prompting including question analysis, answer guidance and safe answer production in their response generation pipeline. Guo et al. (2024a) enhance the LLMs outputs' quality by implementing retrieval-augmented ICL and providing LLMs with relevant documents. To ensure LLMs provide responses aligned with human values, Sun et al. (2024) and Wang et al. (2024a) conduct principle-driven prompting, guiding LLMs with well-crafted and detailed principles.

### 210 3.2 Rationale

211 The rationale reflects the detailed thought process and reasoning pathway an individual follows when solving a given question, being considered valuable auxiliary information for the final answer prediction. In early studies (Ling et al., 2017; Cobbe et al., 2021; Wei et al., 2022), the rationale in each dataset was annotated by human experts, significantly limiting its availability and scalability. Kojima et al. (2022) initially confirm the efficacy of

220 the chain-of-thought (CoT) approach in LLMs and 221 boosting LLMs' reasoning through the integration 222 of self-generated rationales.

223 **Rationale Structure.** Following Kojima et al. (2022), there is a notable interest in abstracting the reasoning process of LLMs into diverse structures and format, including trees (Hao et al., 2023; Yao et al., 2024), graphs (Besta et al., 2024; Yao et al., 2023), tables (Wang et al., 2024d), programs (Chen et al., 2023e), recursion (Qi et al., 2023), and concepts (Tan et al., 2023).

231 **Rationale Quality.** To produce high-quality and 232 fine-grained rationale, diverse methodologies have 233 been employed. Wang et al. (2022a) prompt frozen 234 LLMs to produce choice-specific rationales to elu- 235 cidate each choice in a sample. Wang et al. (2023b) 236 employ contrastive decoding to foster more plau- 237 sible rationales, taking into account gold-standard 238 answers. Liu et al. (2023a) curate meticulously 239 designed prompts to derive high-quality rationales 240 from GPT-4 and construct a logical CoT instruc- 241 tion tuning dataset. For attaining fine-grained ra- 242 tionales, Shridhar et al. (2023) introduce Socratic 243 CoT by decomposing the original question into a 244 series of subquestion-solution pairs and generat- 245 ing CoT for them separately. Additionally, Kang 246 et al. (2024) propose a neural reranker to acquire 247 supplementary relevant documents for rationale 248 generation in knowledge-intensive reasoning tasks.

249 **Human-like Rationale.** Another intriguing avenue 250 in synthesized rationale delves into making the rea- 251 soning process more human-like. Many studies em- 252 ulate human diverse thinking in problem-solving, 253 sampling multiple reasoning pathways for a given 254 question (Gao et al., 2021; Wang et al., 2022b; 255 Chen et al., 2023f; Liu et al., 2023c). Subsequent

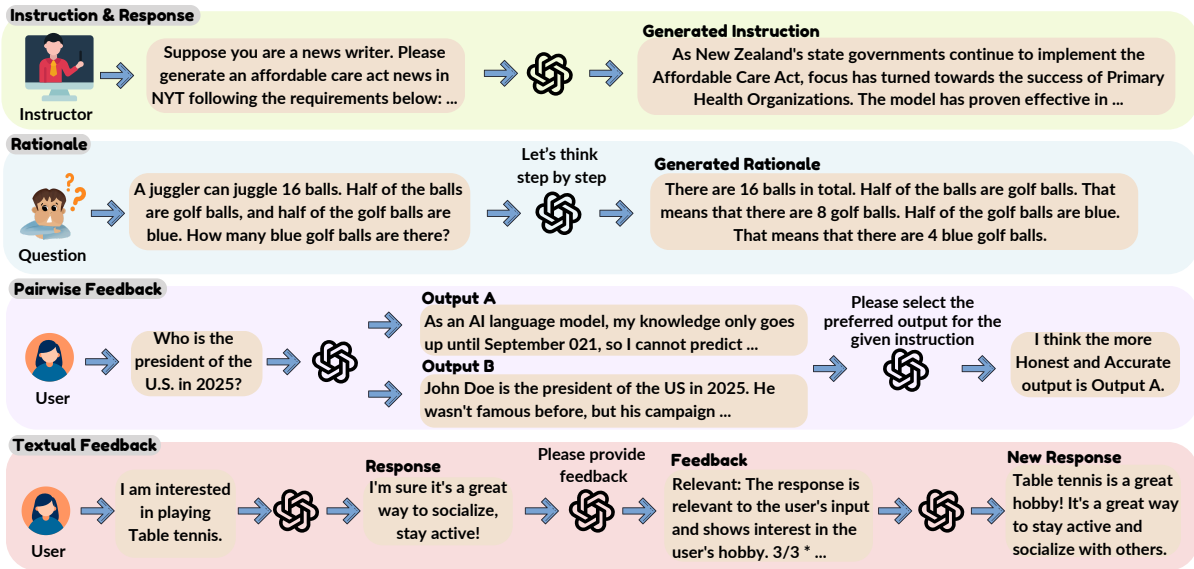


Figure 2: The examples for LLM-based annotation generation.

studies (Tong et al., 2023; Balepur et al., 2023; Ma and Du, 2023) explore the elimination reasoning in LLMs, checking each reasoning pathway reversely and removing the incorrect candidates. Moreover, various works (Yin et al., 2023; Liang et al., 2023; Xu et al., 2023d; Liu et al., 2023e) explore the peer collaboration and debate among individual LLMs to capture human-like discussions as rationales.

### 3.3 Pairwise Feedback

While high-quality human feedback is proven to be effective in aligning LLMs’ values and preferences with us humans, recent advancements aim to automate this pairwise feedback mechanism.

**Ranking with LLMs.** One technique is to sample multiple responses and have the LLM rank these candidates based on various criteria (Bai et al., 2022; Lee et al., 2023b; Yuan et al., 2024). Sun et al. (2023b) sample two responses from the initial policy model and use the model to select the preferred response based on a human-written principle (Sun et al., 2024). Zhang et al. (2024a) propose a self-evaluation mechanism, generating questions for each response and measuring factuality by the LLM’s confidence in the answers. To improve synthetic data quality, Pace et al. (2024) combine the Best-of-N and Worst-of-N sampling strategies and introduce the West-of-N approach. They constructed data pairs by identifying the best- and worst-scored responses according to a pre-trained preference model. In robotics, Zeng et al. (2024) iteratively update the reward function with the self-ranked responses from LLMs, enhancing learning efficiency without human supervision.

**Direct Construction.** Another effort towards

automatic pairwise feedback generation involves directly generating responses of various qualities (Feng et al., 2024; Lee et al., 2024a). To accomplish this, they typically have to make various assumptions when determining the factors influencing response quality. For example, Kim et al. (2023b) assume larger LLM with more shots will give better responses and produce synthetic pairs based on this. Tong et al. (2024b) follow the rule of thumb that the supervised fine-tuning model will perform better than its unfinetuned base model. Adhere to this criterion, they start with a few seed data, iteratively training the model and synthesizing comparison data pairs. Yang et al. (2023c) create quality differences by prompting LLMs to either follow or violate given principles. To measure the response quality more subjectively, Xu et al. (2023c) introduce multiple LLMs and utilize benchmark scores to define superiority.

### 3.4 Textual Feedback

Textual feedback (Pan et al., 2024) generated by LLMs typically highlights the shortcomings of the current output or suggests specific improvements, thus offering rich and valuable information for polishing or evaluating the generated response. Many existing works tailor appropriate prompts and instruct LLMs to generate such informative feedback in various tasks, including question answering (Madaan et al., 2024; Shinn et al., 2024), machine translation (Chen et al., 2023c; Raunak et al., 2023) and hallucination detection (Yang et al., 2023d; Manakul et al., 2023). Some investigations have explored leveraging debate and peer review as feedback to enhance LLMs’ reasoning (Du et al.,

2023a; Xu et al., 2023d; Cohen et al., 2023; Fu et al., 2023) and evaluation (Li et al., 2023d; Chu et al., 2024b; Ning et al., 2024) capabilities. Additionally, efforts have been made to analyze reasons for undesired or incorrect responses produced by LLMs, thus facilitating reflection and learning from their previous mistakes (Wang and Li, 2023; An et al., 2023; Chen et al., 2023a; Tong et al., 2024a).

### 3.5 Other Domain-specific Data

Distilling multi-round conversations from LLMs presents a highly cost-effective approach for constructing high-quality dialogue datasets (Kim et al., 2023a; Xu et al., 2023b; Chen et al., 2023b; Li et al., 2024d) or enhancing existing ones (Zheng et al., 2023; Chen et al., 2022; Zhou et al., 2022a). In graph and tabular data, several studies prompt LLMs to contextualize these structural data (Xi-ang et al., 2022; Kim et al., 2023a; Li et al., 2024b; Ronzano and Nanavati, 2024) or distill structural insights from raw text (Bi et al., 2024; Li et al., 2024c; Ding et al., 2024; Xiong et al., 2024; Tuozzo, 2022). Moreover, LLMs have also been widely adopted in the research of robotics and agents, serving as proficient data annotators to generate plans (Huang et al., 2022; Brohan et al., 2023; Rana et al., 2023; Singh et al., 2023; Lin et al., 2023a), simulation tasks (Wang et al., 2023a; Ha et al., 2023) and supervised signal (Kwon et al., 2022; Du et al., 2023b). Besides, LLMs are acting as efficient data annotators in various artificial intelligence domains, including multi-modal (Li et al., 2023f; Yin et al., 2024; Chen et al., 2024a), recommendation system (Acharya et al., 2023; Shen et al., 2024; Wei et al., 2024; Zhang et al., 2024b), information extraction (Josifoski et al., 2023; Jeronimo et al., 2023; Li et al., 2024a; Ma et al., 2024; Bonn et al., 2024) and etc (Chu et al., 2024a; Bhattacharjee et al., 2024; Martorana et al., 2024).

## 4 LLM-Generated Annotations Assessment

Effective evaluation of annotations generated by LLMs is crucial to fully harness their potential. This section focuses on two main aspects:

### 4.1 Evaluating LLM-Generated Annotations

This subsection explores various methods for assessing annotation quality, ranging from human-led to automated approaches.

**General Approaches:** Research has investigated diverse methods for evaluating LLM annotations.

The “Turking Test” by Efrat and Levy (2020), evaluates LLMs’ adherence to data annotation guidelines, with human annotators comparing LLM outputs against benchmarks like SNLI (Bowman et al., 2015), SQuAD (Rajpurkar et al., 2016), and NewsQA (Trischler et al., 2016). Similarly, Honovich et al. (2022) manually examined the originality, accuracy, and variety of datasets created by LLMs, focusing on their response to instructions. Additionally, studies such as by Alizadeh et al. (2023) measure the performance of open-source LLMs against human-annotated labels in tasks like relevance and topic detection.

**Task-Specific Evaluations:** Methodologies vary by application. For instance, in knowledge graph enhancement, token ranking metrics assess LLM contributions in fact completion. Additionally, evaluations of counterfactual generation often utilize diversity metrics like Self-BLEU (Chen et al., 2023g), while code generation relies on metrics such as Pass@k (Nijkamp et al., 2022). In scenarios requiring extensive datasets, the quality of LLM-generated annotations is compared to gold standard labels within a small, labeled subset (Zhao et al., 2021; Agrawal et al., 2022; He et al., 2023).

### 4.2 Filtering & Selection

Selecting high-quality annotations from numerous options is crucial. In this section, we categorize the filtering and selection methods for LLM-generated data into three types: rule-based filtering, external source utilization, and LLMs-driven selection.

**Rule-Based Methods.** Rule-based methods follow various heuristic assumptions concerning sample length (Li et al., 2023f; Kim et al., 2023a), keyword occurrence (Kim et al., 2023b; Zheng et al., 2023) and specific patterns (Zhang and Yang, 2023a; Guo et al., 2024a; Ding et al., 2024) to filter low-quality or undesired synthetic data points. Zheng et al. (2023); Kim et al. (2023a) establish thresholds for the number of rounds in generated conversations to guarantee each synthetic dialogue is informative enough. Ho et al. (2023); Kang et al. (2024) employ ground truth parsing to filter out incorrect CoT rationales within each candidate reasoning sample. To encourage diversity among the generated data points, Wang et al. (2022b); Lee et al. (2023a); Ding et al. (2024) utilize semantic similarity metrics to identify and remove redundant samples.

**External-Source-Based Methods.** There are also many works that depend on the external source’s feedback to clean and refine synthetic

424 datasets (Kim et al., 2023a). With a pre-trained  
425 reward model, Gulcehre et al. (2023); Dong et al.  
426 (2023) augment the original dataset only with sam-  
427 ples that obtain high reward values. When dis-  
428 tillling smaller models, Lin et al. (2023b); Wang  
429 et al. (2024c) meticulously select appropriate data  
430 through the feedback from the student models.  
431 Other approaches (Chen et al., 2023g; Zheng et al.,  
432 2023) utilize pre-trained classification models to  
433 discern between target and unwanted data points.  
434 **LLMs-Driven Methods.** The versatility of LLMs  
435 has invoked interest in leveraging LLMs them-  
436 selves to do data selection. Some approaches use  
437 signals or features produced by LLMs, such as  
438 perplexity score (Wang et al., 2023f), confidence  
439 levels (Wang et al., 2022b; Huang et al., 2023),  
440 and logits (Pace et al., 2024), as criteria for con-  
441 structing data selectors. Others directly prompt  
442 the LLMs for this task. For instance, Lu et al.  
443 (2023) query the target LLM to assess the quality  
444 of generated samples. Kim et al. (2023a) leverage  
445 ChatGPT to determine if the social commonsense  
446 knowledge is appropriately conveyed in the syn-  
447 thetic dialogues. Additionally, there are also works  
448 that adopt the LLMs to rank multiple candidate an-  
449 notations and utilize the top ones in the subsequent  
450 stages (Jeronymo et al., 2023; Li et al., 2024c). In  
451 pairwise feedback synthesis, Tong et al. (2024b)  
452 task the base LLM with judging whether one re-  
453 sponse genuinely surpasses another.

## 454 5 LLM-Generated Annotations 455 Utilization

456 LLM-generated annotations provide a valuable re-  
457 source of labeled data for NLP models in different  
458 stages. Hereby we explore the methods for utiliz-  
459 ing and learning with LLM-Generated Annotations.

### 460 5.1 Supervised Fine-tuning 461

462 Supervised fine-tuning can effectively enhance  
463 models’ specific capabilities or knowledge. In this  
464 section, we discuss the utilization of generated an-  
465 notation for supervised fine-tuning.

466 **Self-distillation.** Huang et al. (2023) first propose  
467 the concept of self-improve that utilizes LLMs as  
468 both data annotators and learnable models and it-  
469 eratively fine-tune LLMs in their self-annotated  
470 data. Wang et al. (2023e) also tune a GPT3 in  
471 the instruction tuning dataset to improve its zero-  
472 shot generalization capability. To foster LLMs’  
473 evolution, Lu et al. (2023) iteratively fine-tune the  
474 LLMs in self-refined synthetic responses. To miti-

475 gate the distribution gap between task datasets and  
476 the LLMs, Yang et al. (2024b) use self-distillation  
477 which guides fine-tuning with a distilled dataset  
478 generated by the model itself. Both Chen et al.  
479 (2024b) and Cheng et al. (2024) introduce a self-  
480 play mechanism, where the LLM refines its capa-  
481 bility by playing against instances of itself.

482 **Distill Smaller Models.** For efficiency issues,  
483 many studies aim to use the data generated by a  
484 large and powerful LLM to train a flexible and  
485 affordable smaller model. For a better instruction-  
486 following ability, many medium and small-sized  
487 LLMs are trained on the synthetic dataset pro-  
488 duced by larger LLMs (Taori et al., 2023; Chiang  
489 et al., 2023b; Xu et al., 2023a). In classification  
490 tasks, Meng et al. (2022, 2023); Wang et al. (2023d)  
491 augment the original datasets and train smaller bidi-  
492 rectional attention models on them. To foster mod-  
493 els’ reasoning ability, many studies tune smaller  
494 models with synthetic rationales collected from  
495 LLMs (Wang et al., 2022a; Shridhar et al., 2023;  
496 Liu et al., 2023a; Kang et al., 2024). Other task-  
497 specific capabilities distillation from LLMs include  
498 dialogue generation (Xu et al., 2023b), informa-  
499 tion extraction (Josifoski et al., 2023; Jeronymo  
500 et al., 2023) and code generation (Chaudhary,  
501 2023; Roziere et al., 2023). Moreover, LLMs  
502 have been proven to follow a scaling law in terms  
503 of their knowledge capacity. Therefore, there is  
504 also a growing interest in distilling vertical and  
505 domain-specific knowledge from LLMs, including  
506 medicine (Zhang et al., 2023; Xiong et al., 2023),  
507 finance (Zhang and Yang, 2023b) and science (Luo  
508 et al., 2023; Zhao et al., 2024), to smaller models.

### 509 5.2 Alignment Tuning

510 Alignment tuning methods, like RLHF (Ouyang  
511 et al., 2022), aim to align the output of LLMs with  
512 human intentions, ensuring they are helpful, ethical,  
513 and reliable. Synthetic data produced by LLMs are  
514 widely adopted in these alignment approaches for  
515 reward modeling and policy training.

516 **Reward Modeling.** LLMs-generated annotations  
517 can be used to train or refine the reward model  
518 for better alignment. Xu et al. (2023c) propose  
519 a data curriculum method that leverages the pair-  
520 wise feedback from LLMs to calculate the sample  
521 difficulty level and smooth LLMs’ learning from  
522 simple ones to hard ones. Kim et al. (2023b) de-  
523 sign reward model guided self-play to iteratively  
524 improve the reward model with synthesized data  
525 generated by the policy model. Pace et al. (2024)

propose to maximize the probability of correctly labeling a pair of on-policy responses to a given query according to the base preference model. In robotics, Zeng et al. (2024) learns a reward function from scratch using the LLMs’ feedback. With synthetic data pair, Sun et al. (2023b) train an instructable reward model to generate reward scores based on arbitrary human-defined principles.

**Policy Training.** While many direct alignment methods (Rafailov et al., 2024; Zhao et al., 2023b) have emerged recently, some works directly explore the use of annotated feedback for policy training. One common strategy is to directly apply DPO with the synthetic pairwise feedback produced by LLMs (Yuan et al., 2024; Zhang et al., 2024a; Lee et al., 2024b; Tong et al., 2024b; Lee et al., 2024a; Guo et al., 2024b). Besides, Gulcehre et al. (2023); Dong et al. (2023) leverage a pre-trained reward model to filter low-quality synthetic data and iteratively tune LLMs with growing datasets. Wang et al. (2024a) propose a bootstrapping self-alignment method to repeatedly utilize the synthetic data. Liu et al. (2024) introduce the Mixture of insightful Experts (MoTE) architecture, which applies the mixture of experts to enhance each component of the synthetic response, markedly increasing alignment efficiency. With the reasoning pairwise feedback generated by LLM itself, Pang et al. (2024a) use a modified DPO loss with an additional negative log-likelihood term to tune the LLM.

### 5.3 Inference

**In-context Learning.** In-context Learning (ICL) consists of three components: a task description (or prompt), several in-context samples (or demonstration), and the test case that needs to be inferred. Current studies have applied the annotations and data generated by LLMs in all these components for refining or augmenting. Zhou et al. (2022b) first showed that with a well-designed pipeline, LLMs can be human-level prompt engineers to generate accurate task descriptions. Following them, Yang et al. (2023b); Li et al. conduct augmentation and expansion to the original task prompt, making it more detailed for LLMs to follow. Demonstration augmentation (Kim et al., 2022; Li et al., 2023c; Chen et al., 2023d; He et al., 2024) is another useful skill to enrich and diversify the provided demonstrations, especially when the labeled data is limited. For the test sample, one augmentation method is to leverage LLMs to rephrase it once (Deng et al., 2023) or multiple times (Li et al., 2023a; Yang

et al., 2024a). Other works study how to polish the original test sample (Xi et al., 2023) or decompose it into several sub-questions (Wang et al., 2024b).

**Reasoning.** Reasoning plays a crucial role in enhancing the quality and accuracy of the content generated by LLMs. One efficient manner to boost LLMs’ reasoning with self-generated annotation is to provide the generated rationale directly before outputting the final answer/ response (Kojima et al., 2022). To improve LLMs’ performance with multiple reasoning pathways, majority voting (Wang et al., 2022b; Chen et al., 2023f) and elimination (Tong et al., 2023; Balepur et al., 2023; Ma and Du, 2023) are adopted to decide the final answer among several possible candidates. Post-hoc editing and refining (Madaan et al., 2024; Tong et al., 2024a) is another well-studied direction to utilize textual feedback and analysis for improving LLMs’ reasoning capabilities. Additionally, utilization of LLMs-generated annotations sometimes requires additional domain tools. For example, Chen et al. (2023e) use a program interpreter in program-of-thought (PoT) to execute the generated program and convert it to a specific answer. Besta et al. (2024) design a prompter to Build a prompt to be sent to the LLM and a parser to extract information from LLM thought. In tree-of-thought (ToT), Hao et al. (2023); Yao et al. (2024) build an additional state evaluator by designing specific prompts and repurposing the base LLM.

## 6 Societal Impact and Future Work

In this section, we outline LLM annotation challenges, including societal implications, technical concerns, and bias propagation.

### 6.1 Ethics Consideration

One critical concern of LLM-generated annotations is the ethics consideration, especially in high-stakes decision-making tasks like finance (Yang et al., 2023a), jurisprudence (Cui et al., 2023), and healthcare (Eloundou et al., 2023). Despite the efficiency of LLM annotation, the lack of human insight may lead to biased and unfair results (Wu et al., 2023; Abid et al., 2021; Cheng et al., 2021; Li et al., 2023g). Moreover, LLMs make human annotator roles redundant, potentially increasing social disparities (Dillion et al., 2023). Future studies should harmonize technological advancements with societal consequences, including considering social implications, ensuring ethical use, promoting fairness, and maintaining transparency.

627 **6.2 Challenges and Future Work**

628 **Model Collapse.** Model collapse refers to the gradual performance decrease of an LLM trained on the outputs of other LLMs (Sun et al., 2023a; Gunasekar et al., 2023; Hsieh et al., 2023; Honovich et al., 2022; Chiang et al., 2023a; Geng et al., 2023). It is unavoidable since LLM-generated data is occupying the information ecosystem. The imitation model often replicates stylistic elements without achieving the factual precision of superior models (Gudibande et al., 2023; Shumailov et al., 2023). This divergence is caused by *statistical approximation error* from limited sample sizes and *functional approximation error* from constrained model capacity. Both errors tend to amplify through successive training cycles (Alemohammad et al., 2023).

643 **Potential Solution.** It is important to ensure that the training data is diverse and high-quality, with a significant proportion of human-generated content. Gerstgrasser et al. (2024) avoid model collapse by accumulating real and machine-generated data. This method maintains data diversity, preventing performance degradation across different LLMs.

650 **Hallucinations.** Hallucinations in LLMs significantly undermine the integrity and reliability of their generated annotations (Alkaissi and McFarlane, 2023; Azamfirei et al., 2023; Chaudhary et al., 2024). Hulledinated outputs detached from factual information can cause the proliferation of misinformation (Jiang et al., 2024; Chen and Shu, 2023). Addressing hallucinations requires refining the training process and implementing validation mechanisms for annotations through automated and manual verification (Liao and Vaughan, 2023; Pan et al., 2023; Bian et al., 2023). Moreover, the inherent opacity of LLMs complicates efforts to investigate the causes of hallucinations.

664 **Potential Solution.** Yang et al. (2023d) addresses hallucinations in LLMs with the Reverse Validation method, detecting hallucinations at the passage level by constructing a query from the response and checking for a match within the LLM’s internal knowledge. Bertaglia et al. (2023) uses Chain-of-Thought (CoT) prompting and explanation generation, where CoT prompting produces explanations for predictions, ensuring logical and verifiable outputs. Li et al. (2023b) proposes the CoAnnotating framework, which uses uncertainty-guided work allocation between humans and LLMs, applying self-evaluation and entropy metrics to assess reliability and distribute tasks effectively.

**Efficiency of LLMs.** Efficiency in LLMs is crucial due to their growing size and complexity, which demand substantial computational resources (Wong et al., 2024). Efficient models reduce inference latency, vital for real-time applications, lower energy consumption for sustainable AI practices, and cut operational costs in cloud environments, making AI more cost-effective for researchers. Efficiency techniques for LLMs, such as pruning, compression, and distillation, are critical for deploying these models in resource-constrained environments.

**Potential Solution.** Pruning is an efficient technique to reduce the number of parameters in an LLM. For example, Ma et al. (2023) selectively removes redundant neurons based on gradient information while preserving most of the LLM’s capability. Mixture of Experts (MoE) is another promising technique that leverages a set of expert sub-models, where only a subset of these experts is activated for any given input (Artetxe et al., 2021). Researchers also adopt LLM Quantization to reduce the precision of the numbers used to represent a model’s parameters (Xiao et al., 2023). Instead of using 32-bit floating-point numbers, a quantized model might use 16-bit floats, 8-bit integers, or even lower precision. These techniques can be combined with each other to achieve further efficiencies.

7 **Conclusion**

The exploration of LLMs for data annotation has revealed an exciting frontier in NLP, presenting novel solutions to longstanding challenges like data scarcity, and enhancing annotation quality and process efficiency. This survey meticulously reviews methodologies, applications, and hurdles associated with LLM employment, including detailed taxonomy from annotation generation to utilization. It evaluates the effects of LLM-generated annotations on training machine learning models while addressing both technical and ethical concerns like bias and societal ramifications. Highlighting our novel taxonomy of LLM methodologies, strategies for utilizing LLM-generated annotations, and a critical discussion on the challenges, this work aims to steer future progress in this crucial area. Additionally, we introduce a comprehensive categorization of techniques and compile extensive benchmark datasets to support ongoing research endeavors, concluding with an examination of persistent challenges and open questions, paving the way for future investigative pursuits in the domain.



## 728 Limitations

729 **Sampling Bias and Hallucination.** LLMs can display sampling bias, leading to incorrect or “hallucinated” data, impacting the reliability and quality of annotations for discriminative tasks.

730 **Social Bias and Ethical Dilemmas.** The inherent biases in training data can be perpetuated and amplified by LLMs, leading to ethical concerns and the propagation of social biases through annotated data. This is particularly problematic in tasks requiring fairness and impartiality.

731 **Dependence on High-Quality Data.** LLMs’ usefulness in generating annotations depends on large, high-quality datasets. But curating these datasets is labor-intensive, posing a scalability challenge for LLM-based annotation efforts.

732 **Complexity in Tuning and Prompt Engineering.** Successfully leveraging LLMs for data annotation requires sophisticated prompt engineering and fine-tuning techniques. This can pose a barrier to entry for practitioners and researchers without extensive expertise in NLP and machine learning.

733 **Generalization and Overfitting** While LLMs can be powerful tools for annotation, there’s a risk of overfitting to the training data, limiting their ability to generalize to unseen data or different contexts. This is a critical limitation for discriminative tasks where the goal is to develop models that perform well across diverse datasets and domains.

734 **Computational and Resource Requirements.** The training and deployment of state-of-the-art LLMs for data annotation require substantial computational resources, which may not be accessible to all researchers and organizations, thereby limiting widespread adoption.

## 763 References

764 Sujan Reddy A. 2022. [Automating human evaluation of dialogue systems](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*, pages 229–234, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

771 Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 298–306.

775 Bernardo Aceituno and Antoni Rosinol. 2022. [Stack ai: The middle-layer of ai](#).

Arkadeep Acharya, Brijraj Singh, and Naoyuki Onoe. 2023. Llm based generation of item-description for recommendation system. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 1204–1207.

Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. Large language models are zero-shot clinical information extractors. *arXiv preprint arXiv:2205.12689*.

Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Reza Babaei, Daniel LeJeune, Ali Siahkoochi, and Richard Baraniuk. 2023. [Self-consuming generative models go mad](#). *ArXiv*, abs/2307.01850.

Meysam Alizadeh, Maël Kubli, Zeynab Samei, Shirin Dehghani, Juan Diego Bermeo, Maria Kobobeynikova, and Fabrizio Gilardi. 2023. Open-source large language models outperform crowd workers and approach chatgpt in text-annotation tasks. *arXiv preprint arXiv:2307.02179*.

Hussam Alkaissi and Samy I McFarlane. 2023. Artificial hallucinations in chatgpt: implications in scientific writing. *Cureus*, 15(2).

Walid Amamou. 2021. [Ubiai: Text annotation tool](#).

Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, Jian-Guang Lou, and Weizhu Chen. 2023. Learning from mistakes makes llm better reasoner. *arXiv preprint arXiv:2310.20689*.

Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor Mihaylov, Myle Ott, Sam Shleifer, Xi Victoria Lin, Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru, et al. 2021. Efficient large scale language modeling with mixtures of experts. *arXiv preprint arXiv:2112.10684*.

Razvan Azamfirei, Sapna R Kudchadkar, and James Fackler. 2023. Large language models and the perils of their hallucinations. *Critical Care*, 27(1):1–2.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Nishant Balepur, Shramay Palta, and Rachel Rudinger. 2023. It’s not easy being wrong: Evaluating process of elimination reasoning in large language models. *arXiv preprint arXiv:2311.07532*.

Thales Bertaglia, Stefan Huber, Catalina Goanta, Gerassimos Spanakis, and Adriana Iamnitci. 2023. Closing the loop: Testing chatgpt to generate model explanations to improve human labelling of sponsored content on social media. In *World Conference on Explainable Artificial Intelligence*, pages 198–213. Springer.

831	Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 17682–17690.	Wenyong Huang, Zhenguo Li, et al. 2023a. Gaining wisdom from setbacks: Aligning large language models via mistake analysis. <i>arXiv preprint arXiv:2310.10477</i> .	885
832			886
833			887
834			888
835			
836		Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Seokhwan Kim, Andy Rosenbaum, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2023b. Places: Prompting language models for social conversation synthesis. In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 844–868.	889
837			890
838	Amrita Bhattacharjee, Raha Moraffah, Joshua Garland, and Huan Liu. 2024. Zero-shot llm-guided counterfactual generation for text. <i>arXiv preprint arXiv:2405.04793</i> .		891
839			892
840			893
841			894
842	Zhen Bi, Jing Chen, Yinuo Jiang, Feiyu Xiong, Wei Guo, Huajun Chen, and Ningyu Zhang. 2024. Codekgc: Code language model for generative knowledge graph construction. <i>ACM Transactions on Asian and Low-Resource Language Information Processing</i> , 23(3):1–16.	Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Andy Rosenbaum, Seokhwan Kim, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2022. Weakly supervised data augmentation through prompting for dialogue understanding. In <i>NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research</i> .	895
843			896
844			897
845			898
846			899
847			900
848	Ning Bian, Peilin Liu, Xianpei Han, Hongyu Lin, Yaojie Lu, Ben He, and Le Sun. 2023. A drop of ink may make a million think: The spread of false information in large language models. <i>arXiv preprint arXiv:2305.04812</i> .	Pinzhen Chen, Zhicheng Guo, Barry Haddow, and Kenneth Heafield. 2023c. Iterative translation refinement with large language models. <i>arXiv preprint arXiv:2306.03856</i> .	901
849			902
850			903
851			904
852			
853	Julia Bonn, Harish Tayyar Madabushi, Jena D Hwang, and Claire Bonial. 2024. Adjudicating llms as propbank annotators. <i>LREC-COLING 2024</i> , page 112.	Wei-Lin Chen, Cheng-Kuang Wu, Yun-Nung Chen, and Hsin-Hsi Chen. 2023d. Self-icl: Zero-shot in-context learning with self-generated demonstrations. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15651–15662.	905
854			906
855			907
856			908
857	Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. <i>arXiv preprint arXiv:1508.05326</i> .	Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2023e. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. <i>Transactions on Machine Learning Research</i> .	909
858			910
859			911
860	Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2023. Do as i can, not as i say: Grounding language in robotic affordances. In <i>Conference on robot learning</i> , pages 287–318. PMLR.	Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash, Charles Sutton, Xuezhi Wang, and Denny Zhou. 2023f. Universal self-consistency for large language model generation. <i>arXiv preprint arXiv:2311.17311</i> .	912
861			913
862			914
863			915
864			
865			916
866	Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. A survey on evaluation of large language models.	Yunkai Chen, Qimeng Wang, Shiwei Wu, Yan Gao, Tong Xu, and Yao Hu. 2024a. Tomgpt: Reliable text-only training approach for cost-effective multi-modal large language model. <i>ACM Transactions on Knowledge Discovery from Data</i> .	917
867			918
868			919
869			920
870			
871			921
872	Manav Chaudhary, Harshit Gupta, and Vasudeva Varma. 2024. Brainstorm@ ired at smm4h 2024: Leveraging translation and topical embeddings for annotation detection in tweets. <i>arXiv preprint arXiv:2405.11192</i> .	Zeming Chen, Qiyue Gao, Antoine Bosselut, Ashish Sabharwal, and Kyle Richardson. 2023g. Disco: Distilling counterfactuals with large language models. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5514–5528.	922
873			923
874			924
875			925
876	Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation. <i>Code alpaca: An instruction-following llama model for code generation</i> .	Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024b. Self-play fine-tuning converts weak language models to strong language models. <i>arXiv preprint arXiv:2401.01335</i> .	926
877			927
878			928
879			929
880	Canyu Chen and Kai Shu. 2023. Can llm-generated misinformation be detected? <i>arXiv preprint arXiv:2309.13788</i> .	Lu Cheng, Kush R Varshney, and Huan Liu. 2021. Socially responsible ai algorithms: Issues, purposes, and challenges. <i>Journal of Artificial Intelligence Research</i> , 71:1137–1181.	930
881			931
882			932
883	Kai Chen, Chunwei Wang, Kuo Yang, Jianhua Han, Lanqing Hong, Fei Mi, Hang Xu, Zhengying Liu,		933
884			934

940	Pengyu Cheng, Tianhao Hu, Han Xu, Zhisong Zhang,	for generative foundation model alignment. <i>arXiv</i>	994
941	Yong Dai, Lei Han, and Nan Du. 2024. Self-playing	<i>preprint arXiv:2304.06767.</i>	995
942	adversarial language game enhances llm reasoning.		
943	<i>arXiv preprint arXiv:2404.10642.</i>		
944	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,	Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-	996
945	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan	baum, and Igor Mordatch. 2023a. Improving fac-	997
946	Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion	tuality and reasoning in language models through	998
947	Stoica, and Eric P. Xing. 2023a. Vicuna: An open-	multiagent debate. <i>arXiv preprint arXiv:2305.14325.</i>	999
948	source chatbot impressing GPT-4 with 90%* chatgpt		
949	quality.		
950	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,	Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Co-	1000
951	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan	las, Trevor Darrell, Pieter Abbeel, Abhishek Gupta,	1001
952	Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al.	and Jacob Andreas. 2023b. Guiding pretraining in	1002
953	2023b. Vicuna: An open-source chatbot impressing	reinforcement learning with large language models.	1003
954	gpt-4 with 90%* chatgpt quality. <i>See https://vicuna.</i>	In <i>International Conference on Machine Learning</i> ,	1004
955	<i>lmsys.org (accessed 14 April 2023)</i> , 2(3):6.	pages 8657–8677. PMLR.	1005
956	Zhixuan Chu, Yan Wang, Longfei Li, Zhibo Wang,	Avia Efrat and Omer Levy. 2020. The turking test: Can	1006
957	Zhan Qin, and Kui Ren. 2024a. A causal explainable	language models understand instructions? <i>arXiv</i>	1007
958	guardrails for large language models. <i>arXiv preprint</i>	<i>preprint arXiv:2010.11982.</i>	1008
959	<i>arXiv:2405.04160.</i>		
960	Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and	Tyna Eloundou, Sam Manning, Pamela Mishkin, and	1009
961	Yiqun Liu. 2024b. Pre: A peer review based	Daniel Rock. 2023. Gpts are gpts: An early look at	1010
962	large language model evaluator. <i>arXiv preprint</i>	the labor market impact potential of large language	1011
963	<i>arXiv:2401.15641.</i>	models. <i>arXiv preprint arXiv:2303.10130.</i>	1012
964	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	Yunlong Feng, Yang Xu, Libo Qin, Yasheng Wang, and	1013
965	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	Wanxiang Che. 2024. Improving language model rea-	1014
966	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	soning with self-motivated learning. <i>arXiv preprint</i>	1015
967	Nakano, et al. 2021. Training verifiers to solve math	<i>arXiv:2404.07017.</i>	1016
968	word problems. <i>arXiv preprint arXiv:2110.14168.</i>		
969	Roi Cohen, May Hamri, Mor Geva, and Amir Globerson.	Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata.	1017
970	2023. Lm vs lm: Detecting factual errors via	2023. Improving language model negotiation with	1018
971	cross examination. <i>arXiv preprint arXiv:2305.13281.</i>	self-play and in-context learning from ai feedback.	1019
972	Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and	<i>arXiv preprint arXiv:2305.10142.</i>	1020
973	Li Yuan. 2023. Chatlaw: Open-source legal large	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021.	1021
974	language model with integrated external knowledge	Making pre-trained language models better few-shot	1022
975	bases. <i>arXiv preprint arXiv:2306.16092.</i>	learners. In <i>Proceedings of the 59th Annual Meet-</i>	1023
976	Yihe Deng, Weitong Zhang, Zixiang Chen, and Quan-	<i>ing of the Association for Computational Linguistics</i>	1024
977	quan Gu. 2023. Rephrase and respond: Let large	<i>and the 11th International Joint Conference on Natu-</i>	1025
978	language models ask better questions for themselves.	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	1026
979	<i>arXiv preprint arXiv:2311.04205.</i>	pages 3816–3830.	1027
980	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wal-	1028
981	Kristina Toutanova. 2018. Bert: Pre-training of deep	lace, Pieter Abbeel, Sergey Levine, and Dawn Song.	1029
982	bidirectional transformers for language understand-	2023. Koala: A dialogue model for academic re-	1030
983	ing. <i>arXiv preprint arXiv:1810.04805.</i>	search. <i>BAIR Blog.</i>	1031
984	Danica Dillion, Niket Tandon, Yuling Gu, and Kurt	Matthias Gerstgrasser, Rylan Schaeffer, Apratim Dey,	1032
985	Gray. 2023. Can ai language models replace human	Rafael Rafailov, Henry Sleight, John Hughes,	1033
986	participants? <i>Trends in Cognitive Sciences.</i>	Tomasz Korbak, Rajashree Agrawal, Dhruv Pai, An-	1034
987	Linyi Ding, Sizhe Zhou, Jinfeng Xiao, and Ji-	drey Gromov, et al. 2024. Is model collapse in-	1035
988	awei Han. 2024. Automated construction of	evitable? breaking the curse of recursion by ac-	1036
989	theme-specific knowledge graphs. <i>arXiv preprint</i>	cumulating real and synthetic data. <i>arXiv preprint</i>	1037
990	<i>arXiv:2404.19146.</i>	<i>arXiv:2404.01413.</i>	1038
991	Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan,	Arnav Gudibande, Eric Wallace, Charles Burton Snell,	1039
992	Shizhe Diao, Jipeng Zhang, Kashun Shum, and	Xinyang Geng, Hao Liu, P. Abbeel, Sergey Levine,	1040
993	Tong Zhang. 2023. Raft: Reward ranked finetuning	and Dawn Song. 2023. The false promise of imitating	1041
		proprietary llms. <i>ArXiv</i> , abs/2305.15717.	1042
		Caglar Gulcehre, Tom Le Paine, Srivatsan Srimi-	1043
		vasan, Ksenia Konyushkova, Lotte Weerts, Abhishek	1044
		Sharma, Aditya Siddhant, Alex Ahern, Miaosen	1045
		Wang, Chenjie Gu, et al. 2023. Reinforced self-	1046
		training (rest) for language modeling. <i>arXiv preprint</i>	1047
		<i>arXiv:2308.08998.</i>	1048

1049	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio Cesar Teodoro Mendes, Allison Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero C. Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, S. Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuan-Fang Li. 2023. Textbooks are all you need. <i>ArXiv</i> , abs/2306.11644.	Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. Unnatural instructions: Tuning language models with (almost) no human labor. <i>arXiv preprint arXiv:2212.09689</i> .	1102
1050			1103
1051			1104
1052			1105
1053		Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. 2023. Large language models are zero-shot rankers for recommender systems. <i>arXiv preprint arXiv:2305.08845</i> .	1106
1054			1107
1055			1108
1056			1109
1057	Hongyi Guo, Yuanshun Yao, Wei Shen, Jiaheng Wei, Xiaoying Zhang, Zhaoran Wang, and Yang Liu. 2024a. Human-instruction-free llm self-alignment with limited samples. <i>arXiv preprint arXiv:2401.06785</i> .		1110
1058			1111
1059		Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. <i>arXiv preprint arXiv:2305.02301</i> .	1112
1060			1113
1061	Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. 2024b. Direct language model alignment from online ai feedback. <i>arXiv preprint arXiv:2402.04792</i> .		1114
1062			1115
1063			1116
1064		Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023. Large language models can self-improve. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1051–1068.	1117
1065			1118
1066	Himanshu Gupta, Kevin Scaria, Ujjwala Anantheshwaran, Shreyas Verma, Mihir Parmar, Saurabh Arjun Sawant, Swaroop Mishra, and Chitta Baral. 2023. Targen: Targeted data generation with large language models. <i>arXiv preprint arXiv:2310.17876</i> .		1119
1067			1120
1068			1121
1069		Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In <i>International Conference on Machine Learning</i> , pages 9118–9147. PMLR.	1122
1070			1123
1071			1124
1072	Huy Ha, Pete Florence, and Shuran Song. 2023. Scaling up and distilling down: Language-guided robot skill acquisition. In <i>Conference on Robot Learning</i> , pages 3766–3777. PMLR.		1125
1073			1126
1074		Vitor Jeronimo, Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and Rodrigo Nogueira. 2023. Inpars-v2: Large language models as efficient dataset generators for information retrieval. <i>arXiv preprint arXiv:2301.01820</i> .	1127
1075			1128
1076	Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 8154–8173.		1129
1077			1130
1078			1131
1079		Bohan Jiang, Zhen Tan, Ayushi Nirmal, and Huan Liu. 2024. Disinformation detection: An evolving challenge in the age of llms. In <i>Proceedings of the 2024 SIAM International Conference on Data Mining (SDM)</i> , pages 427–435. SIAM.	1132
1080			1133
1081			1134
1082	Chase Harrison. 2022. <a href="#">Langchain</a> .		1135
1083			1136
1084	Wei He, Shichun Liu, Jun Zhao, Yiwen Ding, Yi Lu, Zhiheng Xi, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. Self-demos: Eliciting out-of-demonstration generalizability in large language models. <i>arXiv preprint arXiv:2404.00884</i> .	Martin Josifoski, Marija Sakota, Maxime Peyrard, and Robert West. 2023. Exploiting asymmetry for synthetic training data generation: Synthie and the case of information extraction. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1555–1574.	1137
1085			1138
1086			1139
1087			1140
1088	Xingwei He, Zhenghao Lin, Yeyun Gong, Alex Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming Yiu, Nan Duan, Weizhu Chen, et al. 2023. Annollm: Making large language models to be better crowdsourced annotators. <i>arXiv preprint arXiv:2303.16854</i> .		1141
1089			1142
1090		Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang. 2024. Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks. <i>Advances in Neural Information Processing Systems</i> , 36.	1143
1091			1144
1092			1145
1093	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023. Large language models are reasoning teachers. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14852–14882.		1146
1094			1147
1095		Hyuhng Joon Kim, Hyunsoo Cho, Junyeob Kim, Taeuk Kim, Kang Min Yoo, and Sang-goo Lee. 2022. Self-generated in-context learning: Leveraging autoregressive language models as a demonstration generator. <i>arXiv preprint arXiv:2206.08082</i> .	1148
1096			1149
1097			1150
1098	Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.		1151
1099			1152
1100		Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, et al. 2023a. Soda: Million-scale dialogue distillation with social commonsense contextualization. In <i>Proceedings of the</i>	1153
1101			1154
			1155
			1156
			1157

1158		2023 Conference on Empirical Methods in Natural Language Processing, pages 12930–12949.	
1159			
1160	Sungdong Kim, Sanghwan Bae, Jamin Shin, Soyoung Kang, Donghyun Kwak, Kang Yoo, and Minjoon Seo. 2023b. Aligning large language models through synthetic feedback. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13677–13700.		
1161			
1162			
1163			
1164			
1165			
1166	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.		
1167			
1168			
1169			
1170			
1171	Abdullatif Köksal, Timo Schick, Anna Korhonen, and Hinrich Schuetze. Longform: Effective instruction tuning with reverse instructions. In <i>ICLR 2024 Workshop on Navigating and Addressing Data Problems for Foundation Models</i> .		
1172			
1173			
1174			
1175			
1176	Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. 2022. Reward design with language models. In <i>The Eleventh International Conference on Learning Representations</i> .		
1177			
1178			
1179			
1180	Daniil Larionov, Artem Shelmanov, Elena Chistova, and Ivan Smirnov. 2019. Semantic role labeling with pre-trained language models for known and unknown predicates. In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)</i> , pages 619–628.		
1181			
1182			
1183			
1184			
1185			
1186	Dong-Ho Lee, Jay Pujara, Mohit Sewak, Ryen White, and Sujay Jauhar. 2023a. Making large language models better data creators. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15349–15360.		
1187			
1188			
1189			
1190			
1191	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023b. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. <i>arXiv preprint arXiv:2309.00267</i> .		
1192			
1193			
1194			
1195			
1196	Kyungjae Lee, Dasol Hwang, Sunghyun Park, Youngsoo Jang, and Moontae Lee. 2024a. Reinforcement learning from reflective feedback (rlrf): Aligning and improving llms via fine-grained self-reflection. <i>arXiv preprint arXiv:2403.14238</i> .		
1197			
1198			
1199			
1200			
1201	Sangkyu Lee, Sungdong Kim, Ashkan Yousefpour, Minjoon Seo, Kang Min Yoo, and Youngjae Yu. 2024b. Aligning large language models by on-policy self-judgment. <i>arXiv preprint arXiv:2402.11253</i> .		
1202			
1203			
1204			
1205	Dawei Li, William Hogan, and Jingbo Shang. 2024a. Read: Improving relation extraction from an adversarial perspective. <i>arXiv preprint arXiv:2404.02931</i> .		
1206			
1207			
1208	Dawei Li, Yaxuan Li, Dheeraj Mekala, Shuyao Li, Xueqi Wang, William Hogan, Jingbo Shang, et al. 2023a. Dail: Data augmentation for in-context learning via self-paraphrase. <i>arXiv preprint arXiv:2311.03319</i> .		
1209			
1210			
1211			
1212			
	Dawei Li, Zhen Tan, Tianlong Chen, and Huan Liu. 2024b. Contextualization distillation from large language model for knowledge graph completion. <i>arXiv preprint arXiv:2402.01729</i> .		1213 1214 1215 1216
	Dawei Li, Shu Yang, Zhen Tan, Jae Young Baik, Sunkwon Yun, Joseph Lee, Aaron Chacko, Bojian Hou, Duy Duong-Tran, Ying Ding, et al. 2024c. Dalk: Dynamic co-augmentation of llms and kg to answer alzheimer’s disease questions with scientific literature. <i>arXiv preprint arXiv:2405.04819</i> .		1217 1218 1219 1220 1221 1222
	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2024d. Camel: Communicative agents for "mind" exploration of large language model society. <i>Advances in Neural Information Processing Systems</i> , 36.		1223 1224 1225 1226 1227
	Minzhi Li, Taiwei Shi, Caleb Ziems, Min-Yen Kan, Nancy F Chen, Zhengyuan Liu, and Diyi Yang. 2023b. Coannotating: Uncertainty-guided work allocation between human and large language models for data annotation. <i>arXiv preprint arXiv:2310.15638</i> .		1228 1229 1230 1231 1232
	Rui Li, Guoyin Wang, and Jiwei Li. 2023c. Are human-generated demonstrations necessary for in-context learning? <i>arXiv preprint arXiv:2309.14681</i> .		1233 1234 1235
	Ruosen Li, Teerth Patel, and Xinya Du. 2023d. Prd: Peer rank and discussion improve large language model based evaluations. <i>arXiv preprint arXiv:2307.02762</i> .		1236 1237 1238 1239
	Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason E Weston, and Mike Lewis. 2023e. Self-alignment with instruction back-translation. In <i>The Twelfth International Conference on Learning Representations</i> .		1240 1241 1242 1243 1244
	Yanda Li, Chi Zhang, Gang Yu, Zhibin Wang, Bin Fu, Guosheng Lin, Chunhua Shen, Ling Chen, and Yunchao Wei. 2023f. Stablelava: Enhanced visual instruction tuning with synthesized image-dialogue data. <i>arXiv preprint arXiv:2308.10253</i> .		1245 1246 1247 1248 1249
	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. Dailydialog: A manually labelled multi-turn dialogue dataset. <i>arXiv preprint arXiv:1710.03957</i> .		1250 1251 1252 1253
	Yichuan Li, Kaize Ding, Jianling Wang, and Kyumin Lee. Empowering large language models for textual data augmentation.		1254 1255 1256
	Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying Wang. 2023g. A survey on fairness in large language models. <i>arXiv preprint arXiv:2308.10149</i> .		1257 1258 1259
	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujia Yang, Zhaopeng Tu, and Shuming Shi. 2023. Encouraging divergent thinking in large language models through multi-agent debate. <i>arXiv preprint arXiv:2305.19118</i> .		1260 1261 1262 1263 1264
	Q Vera Liao and Jennifer Wortman Vaughan. 2023. Ai transparency in the age of llms: A human-centered research roadmap. <i>arXiv preprint arXiv:2306.01941</i> .		1265 1266 1267

1268	Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bohg. 2023a. Text2motion: From natural language instructions to feasible plans. <i>Autonomous Robots</i> , 47(8):1345–1365.	1323
1269		1324
1270		1325
1271		1326
1272	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. <i>arXiv preprint arXiv:2205.14334</i> .	1327
1273		1328
1274		
1275	Yen-Ting Lin, Alexandros Papangelis, Seokhwan Kim, Sungjin Lee, Devamanyu Hazarika, Mahdi Namazifar, Di Jin, Yang Liu, and Dilek Hakkani-Tur. 2023b. Selective in-context data augmentation for intent detection using pointwise v-information. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1463–1476.	1329
1276		1330
1277		1331
1278		1332
1279		1333
1280		
1281		
1282		
1283	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 158–167.	1334
1284		1335
1285		1336
1286		1337
1287		
1288		
1289	Hanmeng Liu, Zhiyang Teng, Leyang Cui, Chaoli Zhang, Qiji Zhou, and Yue Zhang. 2023a. Logicot: Logical chain-of-thought instruction tuning. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	1338
1290		1339
1291		1340
1292		1341
1293		1342
1294	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. <i>arXiv preprint arXiv:2304.08485</i> .	1343
1295		1344
1296		
1297	Tengxiao Liu, Qipeng Guo, Yuqing Yang, Xiangkun Hu, Yue Zhang, Xipeng Qiu, and Zheng Zhang. 2023c. Plan, verify and switch: Integrated reasoning with diverse x-of-thoughts. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2807–2822.	1345
1298		1346
1299		1347
1300		1348
1301		
1302		
1303	Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2023d. Trustworthy llms: a survey and guideline for evaluating large language models’ alignment. <i>arXiv preprint arXiv:2308.05374</i> .	1349
1304		1350
1305		1351
1306		1352
1307		1353
1308		1354
1309	Zhili Liu, Yunhao Gou, Kai Chen, Lanqing Hong, Jiahui Gao, Fei Mi, Yu Zhang, Zhenguo Li, Xin Jiang, Qun Liu, et al. 2024. Mixture of insightful experts (mote): The synergy of thought chains and expert mixtures in self-alignment. <i>arXiv preprint arXiv:2405.00557</i> .	1355
1310		1356
1311		1357
1312		1358
1313		
1314	Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. 2023e. Dynamic llm-agent network: An llm-agent collaboration framework with agent team optimization. <i>arXiv preprint arXiv:2310.02170</i> .	1359
1315		1360
1316		1361
1317		1362
1318	Jianqiao Lu, Wanjun Zhong, Wenyong Huang, Yufei Wang, Fei Mi, Baojun Wang, Weichao Wang, Lifeng Shang, and Qun Liu. 2023. Self: Language-driven self-evolution for large language model. <i>arXiv preprint arXiv:2310.00533</i> .	1363
1319		1364
1320		1365
1321		1366
1322		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378

1379	large pre-trained language models: A survey. <i>ACM Computing Surveys</i> .	model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36.	1433
1380			1434
1381	Ines Montani and Matthew Honnibal. 2018. <a href="#">Prodigy: A new annotation tool for radically efficient machine teaching</a> . <i>Artificial Intelligence</i> , to appear.	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. <i>arXiv preprint arXiv:1606.05250</i> .	1435
1382			1436
1383			1437
1384	Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. Codegen: An open large language model for code with multi-turn program synthesis. <i>arXiv preprint arXiv:2203.13474</i> .	Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf. 2023. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. In <i>7th Annual Conference on Robot Learning</i> .	1438
1385			1439
1386			1440
1387			1441
1388			1442
1389	Kun-Peng Ning, Shuo Yang, Yu-Yang Liu, Jia-Yu Yao, Zhen-Hui Liu, Yu Wang, Ming Pang, and Li Yuan. 2024. Peer-review-in-llms: Automatic evaluation method for llms in open-environment. <i>arXiv preprint arXiv:2402.01830</i> .	Vikas Raunak, Amr Sharaf, Yiren Wang, Hany Awadalla, and Arul Menezes. 2023. Leveraging gpt-4 for automatic translation post-editing. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12009–12024.	1443
1390			1444
1391			1445
1392			1446
1393			1447
1394	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	Francesco Ronzano and Jay Nanavati. 2024. Towards ontology-enhanced representation learning for large language models. <i>arXiv preprint arXiv:2405.20527</i> .	1448
1395	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> .	1449
1396			1450
1397			1451
1398			1452
1399			1453
1400			1454
1401	Alizée Pace, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. 2024. West-of-n: Synthetic preference generation for improved reward modeling. <i>arXiv preprint arXiv:2401.12086</i> .	Xiaoteng Shen, Rui Zhang, Xiaoyan Zhao, Jieming Zhu, and Xi Xiao. 2024. Pmg: Personalized multimodal generation with large language models. In <i>Proceedings of the ACM on Web Conference 2024</i> , pages 3833–3843.	1455
1402			1456
1403			1457
1404			1458
1405	Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2024. Automatically correcting large language models: Surveying the landscape of diverse automated correction strategies. <i>Transactions of the Association for Computational Linguistics</i> , 12:484–506.	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36.	1459
1406			1460
1407			1461
1408			1462
1409			1463
1410			1464
1411	Yikang Pan, Liangming Pan, Wenhui Chen, Preslav Nakov, Min-Yen Kan, and William Yang Wang. 2023. On the risk of misinformation pollution with large language models. <i>arXiv preprint arXiv:2305.13661</i> .	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2023. Distilling reasoning capabilities into smaller language models. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 7059–7073.	1465
1412			1466
1413			1467
1414			1468
1415	Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. 2024a. Iterative reasoning preference optimization. <i>arXiv preprint arXiv:2404.19733</i> .	Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget. <i>ArXiv</i> , abs/2305.17493.	1469
1416			1470
1417			1471
1418			1472
1419	Xianghe Pang, Shuo Tang, Rui Ye, Yuxin Xiong, Bolun Zhang, Yanfeng Wang, and Siheng Chen. 2024b. Self-alignment of large language models via monopolylogue-based social scene simulation. <i>arXiv preprint arXiv:2402.05699</i> .	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2023. Prog-prompt: Generating situated robot task plans using large language models. In <i>2023 IEEE International Conference on Robotics and Automation (ICRA)</i> , pages 11523–11530. IEEE.	1473
1420			1474
1421			1475
1422			1476
1423			1477
1424	Jingyuan Qi, Zhiyang Xu, Ying Shen, Minqian Liu, Di Jin, Qifan Wang, and Lifu Huang. 2023. The art of socratic questioning: Recursive thinking with large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4177–4199.	Feifan Song, Bowen Yu, Hao Lang, Haiyang Yu, Fei Huang, Houfeng Wang, and Yongbin Li. 2024a. Scaling data diversity for fine-tuning language models in human alignment. In <i>Proceedings of the 2024 Joint</i>	1478
1425			1479
1426			1480
1427			1481
1428			1482
1429			1483
1430	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language		1484
1431			1485
1432			1486

1487			1541
1488		<i>International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 14358–14369.	1542
1489			1543
1490	Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2023. Preference ranking optimization for human alignment. <i>arXiv preprint arXiv:2306.17492</i> .		1544
1491			1545
1492			
1493			
1494	Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2024b. Preference ranking optimization for human alignment. In <i>Proceedings of the AACL Conference on Artificial Intelligence</i> , volume 38, pages 18990–18998.		
1495			
1496			
1497			
1498			
1499	Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. 2023a. Is chatgpt good at search? investigating large language models as re-ranking agent. <i>arXiv preprint arXiv:2304.09542</i> .		
1500			
1501			
1502			
1503			
1504	Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong Zhou, Zhenfang Chen, David D. Cox, Yiming Yang, and Chuang Gan. 2023b. <a href="#">Salmon: Self-alignment with instructable reward models</a> .		
1505			
1506			
1507			
1508	Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2024. Principle-driven self-alignment of language models from scratch with minimal human supervision. <i>Advances in Neural Information Processing Systems</i> , 36.		
1509			
1510			
1511			
1512			
1513			
1514	Zhen Tan, Lu Cheng, Song Wang, Yuan Bo, Jundong Li, and Huan Liu. 2023. Interpreting pretrained language models via concept bottlenecks. <i>arXiv preprint arXiv:2311.05014</i> .		
1515			
1516			
1517			
1518	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.		
1519			
1520			
1521			
1522	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .		
1523			
1524			
1525			
1526			
1527			
1528	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. <i>Nature Medicine</i> , pages 1–11.		
1529			
1530			
1531			
1532	Yongqi Tong, Dawei Li, Sizhe Wang, Yujia Wang, Fei Teng, and Jingbo Shang. 2024a. Can llms learn from previous mistakes? investigating llms’ errors to boost for reasoning. <i>arXiv preprint arXiv:2403.20046</i> .		
1533			
1534			
1535			
1536	Yongqi Tong, Sizhe Wang, Dawei Li, Yifan Wang, Simeng Han, Zi Lin, Chengsong Huang, Jiaxin Huang, and Jingbo Shang. 2024b. Optimizing language model’s reasoning abilities with weak supervision. <i>arXiv preprint arXiv:2405.04086</i> .		
1537			
1538			
1539			
1540			
	Yongqi Tong, Yifan Wang, Dawei Li, Sizhe Wang, Zi Lin, Simeng Han, and Jingbo Shang. 2023. Eliminating reasoning via inferring with planning: A new framework to guide llms’ non-linear thinking. <i>arXiv preprint arXiv:2310.12342</i> .		1541
			1542
			1543
			1544
			1545
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		1546
			1547
			1548
			1549
			1550
			1551
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		1552
			1553
			1554
			1555
			1556
			1557
	Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman, and Kaheer Suleman. 2016. Newsqa: A machine comprehension dataset. <i>arXiv preprint arXiv:1611.09830</i> .		1558
			1559
			1560
			1561
	Gabriele Tuozzo. 2022. Moving from tabular knowledge graph quality assessment to rdf triples leveraging chatgpt.		1562
			1563
			1564
	Somin Wadhwa, Silvio Amir, and Byron C Wallace. 2023. Revisiting relation extraction in the era of large language models. <i>arXiv preprint arXiv:2305.05003</i> .		1565
			1566
			1567
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. <i>arXiv preprint arXiv:1804.07461</i> .		1568
			1569
			1570
			1571
			1572
	Danqing Wang and Lei Li. 2023. Learning from mistakes via cooperative study assistant for large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10667–10685.		1573
			1574
			1575
			1576
			1577
	Haoyu Wang, Guozheng Ma, Ziqiao Meng, Zeyu Qin, Li Shen, Zhong Zhang, Bingzhe Wu, Liu Liu, Yatao Bian, Tingyang Xu, et al. 2024a. Step-on-feet tuning: Scaling self-alignment of llms via bootstrapping. <i>arXiv preprint arXiv:2402.07610</i> .		1578
			1579
			1580
			1581
			1582
	Hongru Wang, Boyang Xue, Baohang Zhou, Tianhua Zhang, Cunxiang Wang, Guanhua Chen, Huimin Wang, and Kam-fai Wong. 2024b. Self-dc: When to retrieve and when to generate? self divide-and-conquer for compositional unknown questions. <i>arXiv preprint arXiv:2402.13514</i> .		1583
			1584
			1585
			1586
			1587
			1588
	Lirui Wang, Yiyang Ling, Zhecheng Yuan, Mohit Shridhar, Chen Bao, Yuzhe Qin, Bailin Wang, Huazhe Xu, and Xiaolong Wang. 2023a. Gensim: Generating robotic simulation tasks via large language models. In <i>The Twelfth International Conference on Learning Representations</i> .		1589
			1590
			1591
			1592
			1593
			1594



1595	PeiFeng Wang, Aaron Chan, Filip Ilievski, Muhao Chen, and Xiang Ren. 2022a. Pinto: Faithful language reasoning using prompt-generated rationales. In <i>The Eleventh International Conference on Learning Representations</i> .	1650
1596		1651
1597		1652
1598		
1599		
1600	Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023b. Scott: Self-consistent chain-of-thought distillation. <i>arXiv preprint arXiv:2305.01879</i> .	1653
1601		1654
1602		1655
1603		1656
1604		1657
1605		1658
1606	Ruida Wang, Wangchunshu Zhou, and Mrinmaya Sachan. 2023c. Let's synthesize step by step: Iterative dataset synthesis with large language models by extrapolating errors from small models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11817–11831.	1659
1607		1660
1608		1661
1609		1662
1610		1663
1611	Song Wang, Zhen Tan, Ruocheng Guo, and Jundong Li. 2023d. Noise-robust fine-tuning of pretrained language models via external guidance. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12528–12540.	1664
1612		1665
1613		1666
1614		1667
1615		
1616	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022b. Self-consistency improves chain of thought reasoning in language models. In <i>The Eleventh International Conference on Learning Representations</i> .	1668
1617		1669
1618		1670
1619		1671
1620		
1621	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023e. Self-instruct: Aligning language models with self-generated instructions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13484–13508.	1672
1622		1673
1623		1674
1624		1675
1625		1676
1626		
1627		
1628	Yue Wang, Haoke Zhang, Juntao Li, Jinxiong Chang, Qishen Zhang, Zhongyi Liu, Guannan Zhang, and Min Zhang. 2023f. Sass: Self-alignment with semi-supervised instruction data generation.	1677
1629		1678
1630		1679
1631		1680
1632		1681
1633	Yufei Wang, Wanjuan Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023g. Aligning large language models with human: A survey. <i>arXiv preprint arXiv:2307.12966</i> .	1682
1634		1683
1635		
1636		
1637	Zifeng Wang, Chun-Liang Li, Vincent Perot, Long T Le, Jin Miao, Zizhao Zhang, Chen-Yu Lee, and Tomas Pfister. 2024c. Codeclm: Aligning language models with tailored synthetic data. <i>arXiv preprint arXiv:2404.05875</i> .	1684
1638		1685
1639		1686
1640		1687
1641		1688
1642		
1643	Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, et al. 2024d. Chain-of-table: Evolving tables in the reasoning chain for table understanding. <i>arXiv preprint arXiv:2401.04398</i> .	1689
1644		1690
1645		1691
1646		1692
1647		1693
1648		
1649	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	1694
		1695
		1696
		1697
		1698
		1699
		1700
		1701
		1702
	Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. 2024. Llmrec: Large language models with graph augmentation for recommendation. In <i>Proceedings of the 17th ACM International Conference on Web Search and Data Mining</i> , pages 806–815.	1703
		1704
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. <a href="#">Huggingface's transformers: State-of-the-art natural language processing</a> .	
	Siu Ming Wong, Ho Leung, and Ka Yan Wong. 2024. Efficiency in language understanding and generation: An evaluation of four open-source large language models.	
	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. <i>arXiv preprint arXiv:2303.17564</i> .	
	Zhiheng Xi, Senjie Jin, Yuhao Zhou, Rui Zheng, Songyang Gao, Jia Liu, Tao Gui, Qi Zhang, and Xuan-Jing Huang. 2023. Self-polish: Enhance reasoning in large language models via problem refinement. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11383–11406.	
	Jiannan Xiang, Zhengzhong Liu, Yucheng Zhou, Eric Xing, and Zhiting Hu. 2022. Asdot: Any-shot data-to-text generation with pretrained language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 1886–1899.	
	Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In <i>International Conference on Machine Learning</i> , pages 38087–38099. PMLR.	
	Honglin Xiong, Sheng Wang, Yitao Zhu, Zihao Zhao, Yuxiao Liu, Linlin Huang, Qian Wang, and Dinggang Shen. 2023. Doctorglm: Fine-tuning your chinese doctor is not a herculean task. <i>arXiv preprint arXiv:2304.01097</i> .	
	Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. 2024. Large language models can learn temporal reasoning. <i>arXiv preprint arXiv:2401.06853</i> .	
	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin	

1705	Jiang. 2023a. Wizardlm: Empowering large language models to follow complex instructions. <i>arXiv preprint arXiv:2304.12244</i> .	1757
1706		1758
1707		1759
1708	Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023b. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6268–6278.	1760
1709		1761
1710		1762
1711		1763
1712		1764
1713		1765
1714	Canwen Xu, Corby Rosset, Luciano Del Corro, Shweti Mahajan, Julian McAuley, Jennifer Neville, Ahmed Hassan Awadallah, and Nikhil Rao. 2023c. Contrastive post-training large language models on data curriculum. <i>arXiv preprint arXiv:2310.02263</i> .	1766
1715		1767
1716		1768
1717		1769
1718		1770
1719	Zhenran Xu, Senbao Shi, Baotian Hu, Jindi Yu, Dongfang Li, Min Zhang, and Yuxiang Wu. 2023d. Towards reasoning in large language models via multi-agent peer review collaboration. <i>arXiv preprint arXiv:2311.08152</i> .	1771
1720		1772
1721		1773
1722		1774
1723		1775
1724	Adam Yang, Chen Chen, and Konstantinos Pitas. 2024a. Just rephrase it! uncertainty estimation in closed-source language models via multiple rephrased queries. <i>arXiv preprint arXiv:2405.13907</i> .	1776
1725		1777
1726		1778
1727		1779
1728	Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023a. Fingpt: Open-source financial large language models. <i>arXiv preprint arXiv:2306.06031</i> .	1780
1729		1781
1730		1782
1731	Jinghan Yang, Shuming Ma, and Furu Wei. 2023b. Auto-icl: In-context learning without human supervision. <i>arXiv preprint arXiv:2311.09263</i> .	1783
1732		1784
1733		1785
1734	Kevin Yang, Dan Klein, Asli Celikyilmaz, Nanyun Peng, and Yuandong Tian. 2023c. Rlcd: Reinforcement learning from contrastive distillation for lm alignment. In <i>The Twelfth International Conference on Learning Representations</i> .	1786
1735		1787
1736		1788
1737		1789
1738		1790
1739	Shiping Yang, Renliang Sun, and Xiaojun Wan. 2023d. A new benchmark and reverse validation method for passage-level hallucination detection. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 3898–3908.	1791
1740		1792
1741		1793
1742		1794
1743		1795
1744	Zhaorui Yang, Qian Liu, Tianyu Pang, Han Wang, Haozhe Feng, Minfeng Zhu, and Wei Chen. 2024b. Self-distillation bridges distribution gap in language model fine-tuning. <i>arXiv preprint arXiv:2402.13669</i> .	1796
1745		1797
1746		1798
1747		1799
1748	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	1800
1749		1801
1750		1802
1751		1803
1752		1804
1753	Yao Yao, Zuchao Li, and Hai Zhao. 2023. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. <i>arXiv preprint arXiv:2305.16582</i> .	1805
1754		1806
1755		1807
1756		1808
	Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. 2022a. Zeroshot: Efficient zero-shot learning via dataset generation. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11653–11669.	1809
		1810
		1811
	Jiacheng Ye, Jiahui Gao, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2022b. Progen: Progressive zero-shot dataset generation via in-context feedback. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 3671–3683.	1812
		1813
	Zhangyue Yin, Qiushi Sun, Cheng Chang, Qipeng Guo, Junqi Dai, Xuan-Jing Huang, and Xipeng Qiu. 2023. Exchange-of-thought: Enhancing large language model capabilities through cross-model communication. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15135–15153.	
	Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Xiaoshui Huang, Zhiyong Wang, Lu Sheng, Lei Bai, et al. 2024. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo Lee, and Woomyoung Park. 2021. Gpt3mix: Leveraging large-scale language models for text augmentation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2225–2239.	
	Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2022. Generate rather than retrieve: Large language models are strong context generators. <i>arXiv preprint arXiv:2209.10063</i> .	
	Xinli Yu, Zheng Chen, Yuan Ling, Shujing Dong, Zongyi Liu, and Yanbin Lu. 2023. Temporal data meets llm—explainable financial time series forecasting. <i>arXiv preprint arXiv:2306.11025</i> .	
	Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander J Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. 2024. Large language model as attributed training data generator: A tale of diversity and bias. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. <i>arXiv preprint arXiv:2401.10020</i> .	
	Daoguang Zan, Bei Chen, Fengji Zhang, Dianjie Lu, Bingchao Wu, Bei Guan, Wang Yongji, and Jian-Guang Lou. 2023. Large language models meet nl2code: A survey. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7443–7464.	

1814	Yuwei Zeng, Yao Mu, and Lin Shao. 2024. Learning reward for robot skills using large language models via self-alignment. <i>arXiv preprint arXiv:2405.07162</i> .	1868
1815		1869
1816		1870
1817	Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Guiming Chen, Jianquan Li, Xiangbo Wu, Zhang Zhiyi, Qingying Xiao, et al. 2023. Huatuogpt, towards taming language model to be a doctor. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10859–10885.	1871
1818		1872
1819		1873
1820		1874
1821		
1822		
1823	Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024a. Self-alignment for factuality: Mitigating hallucinations in llms via self-evaluation. <i>arXiv preprint arXiv:2402.09267</i> .	
1824		
1825		
1826		
1827		
1828	Xiaoyu Zhang, Yishan Li, Jiayin Wang, Bowen Sun, Weizhi Ma, Peijie Sun, and Min Zhang. 2024b. Large language models as evaluators for recommendation explanations. <i>arXiv preprint arXiv:2406.03248</i> .	
1829		
1830		
1831		
1832	Xuanyu Zhang and Qing Yang. 2023a. Self-qa: Unsupervised knowledge guided language model alignment. <i>arXiv preprint arXiv:2305.11952</i> .	
1833		
1834		
1835	Xuanyu Zhang and Qing Yang. 2023b. Xuanyuan 2.0: A large chinese financial chat model with hundreds of billions parameters. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management</i> , pages 4435–4439.	
1836		
1837		
1838		
1839		
1840	Haiteng Zhao, Shengchao Liu, Ma Chang, Hannan Xu, Jie Fu, Zhihong Deng, Lingpeng Kong, and Qi Liu. 2024. Gimlet: A unified graph-text model for instruction-based molecule zero-shot learning. <i>Advances in Neural Information Processing Systems</i> , 36.	
1841		
1842		
1843		
1844		
1845		
1846	Mengjie Zhao, Fei Mi, Yasheng Wang, Minglei Li, Xin Jiang, Qun Liu, and Hinrich Schütze. 2021. Lm-turk: Few-shot learners as crowdsourcing workers in a language-model-as-a-service framework. <i>arXiv preprint arXiv:2112.07522</i> .	
1847		
1848		
1849		
1850		
1851	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023a. <a href="#">A survey of large language models</a> .	
1852		
1853		
1854		
1855		
1856		
1857		
1858	Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. 2023b. Slic-hf: Sequence likelihood calibration with human feedback. <i>arXiv preprint arXiv:2305.10425</i> .	
1859		
1860		
1861		
1862	Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng Zhang, and Minlie Huang. 2023. Augesc: Dialogue augmentation with large language models for emotional support conversation. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1552–1568.	
1863		
1864		
1865		
1866		
1867		
	Pei Zhou, Hyundong Cho, Pegah Jandaghi, Dong-Ho Lee, Bill Yuchen Lin, Jay Pujara, and Xiang Ren. 2022a. Reflect, not reflex: Inference-based common ground improves dialogue response quality. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10450–10468.	1875
		1876
		1877
		1878
		1879
	<b>A LLM-assisted Tools and Software for Annotation</b>	1880
		1881
	LLM-assisted annotation tools and software are invaluable resources designed specifically to facilitate the annotation process for various NLP tasks. One of their primary attributes is an intuitive and user-friendly interface, allowing engineers and even non-technical annotators to easily work with complex textual data. These tools are built to support numerous annotation types, from simple binary labels to more intricate hierarchical structures. The main goal of these tools is to simplify the labeling process, enhance the quality of the labels, and boost overall productivity in data annotation.	1882
		1883
		1884
		1885
		1886
		1887
		1888
		1889
		1890
		1891
		1892
		1893
	Below, we will present a selection of the libraries and tools that support Large Language Models for the annotation process:	1894
		1895
		1896
	• <b>LangChain:</b> LangChain (Harrison, 2022) is an open-source library <sup>1</sup> that offers an array of tools designed to facilitate the construction of LLM-related pipelines and workflows. This library specifically provides large language models with agents in order to interact effectively with their environment as well as various external data sources. Therefore, providing dynamic and contextually appropriate responses that go beyond a single LLM call.	1897
		1898
		1899
		1900
		1901
		1902
		1903
		1904
		1905
		1906
	In terms of the annotation process, their power mostly lies in the facilitation of annotation through the creation of a modularized structure called <i>chain</i> . In the chaining technique, a complex problem is broken down into smaller sub-tasks. The results obtained from one or more steps are then aggregated and utilized as input prompts for subsequent actions in the chain.	1907
		1908
		1909
		1910
		1911
		1912
		1913
		1914
		1915

<sup>1</sup>As of now, available only in JavaScript/TypeScript and Python languages.

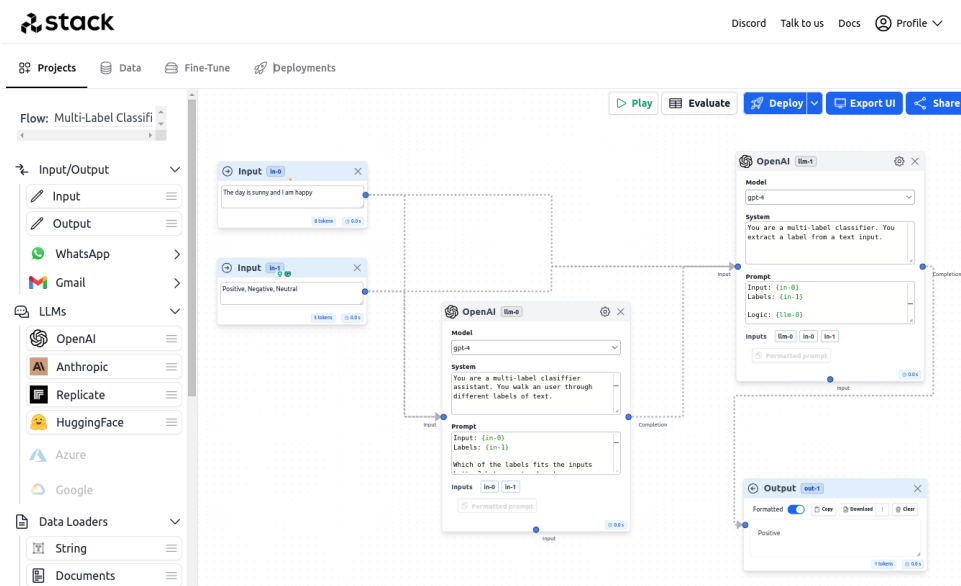


Figure 3: Stack AI dashboard. They provide a visual interface for users to design and track the AI workflow.

• **Stack AI:** Stack AI (Aceituno and Rosinol, 2022) is a paid service that offers an AI-powered data platform. It is designed explicitly for automating business processes allowing them to maximize efficiency. The essence of their platform lies in their ability to *visually* design, test, and deploy AI workflows through smooth integration of Large Language Models. Their user-friendly graphical interface (Figure 3) allows the users to create apps and workflows related to diverse tasks from content creation and data labeling to conversational AI apps and document processing. Moreover, Stack AI utilizes weakly supervised machine learning models to expedite the data preparation process.

• **UBIAI:** UBIAI (Amamou, 2021) is a paid annotation tool that offers multilingual cloud-based solutions and services in Natural Language Processing. The company aims to aid users in extracting valuable insights from unstructured documents. This tool not only provides a user interface that facilitates manual labeling but also offers several auto-labeling functionalities such as LLM-assisted zero- and few-shot labeling and model-assisted labeling. They also provide integration to various models on huggingface (Wolf et al., 2020) as well as an environment to fine-tune different models on the user’s labeled data.

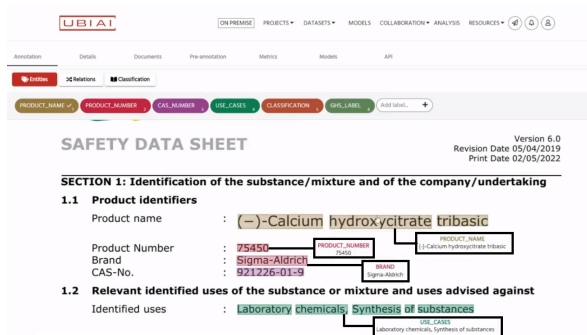


Figure 4: UBIAI annotation result on a pdf document. All the entities in the text of the document have been identified, annotated, and color-coded based on the type. This image has been borrowed from the videos provided in the UBIAI documentation (Amamou, 2021).

• **Prodigy:** Prodigy (Montani and Honnibal, 2018), designed by the creators of spaCy library (Honnibal and Montani, 2017), offers rule-based, statistical models, and LLM-assisted methods for annotation. This tool provides easy, flexible, and powerful annotation options such as named entity recognition, span categorization, and classification/labeling for different modalities including text, audio, and vision. Moreover, it can be easily integrated with large language models which are capable of zero- or few-shot learning, while also offering services and quantifiable methods for crafting prompts to address any noisy outcomes. This tool is not open-source.

1961 **B Acknowledgment of AI Assistance in**  
1962 **Writing and Revision**

1963 We utilized ChatGPT-4 for revising and enhancing  
1964 sections of this paper.

1965 **C Collections of Papers on LLM for Data**  
1966 **Annotation**

1967 This collection of tables provides a concise  
1968 overview of using Large Language Models (LLMs)  
1969 for data annotation, including state-of-the-art tech-  
1970 niques, methodologies, and practical applications.  
1971 Table 1 and Table 2 lists significant papers on LLM-  
1972 based data annotation, detailing their methods, core  
1973 technologies, publication venues, and links to re-  
1974 sources. Table 3 focuses on assessment and filter-  
1975 ing of LLM-generated annotations. Tables 4 ex-  
1976 plore strategies for learning with LLM-generated  
1977 annotations, covering supervised fine-tuning, align-  
1978 ment tuning and inference. Each table clearly out-  
1979 lines the data type, backbone, computational cost,  
1980 venues, and available resources, serving as a guide  
1981 to the latest in LLM-driven data annotation and  
1982 its implications for the future of automated data  
1983 processing and machine learning research.

Paper	Data Type	Backbone	Annotation Cost	Venue	Code/Data Link
<b>Instruction &amp; Response</b>					
GPT3Mix: Leveraging Large-scale Language Models for Text Augmentation <sup>[1]</sup>	Instruction	GPT-3	API Calling, 300 tokens per sample	EMNLP'21	<a href="#">Link</a>
SELF-INSTRUCT: Aligning Language Models with Self-Generated Instructions <sup>[2]</sup>	Instruction & Response	GPT-3	API Calling, \$600 for entire dataset	ACL'23	<a href="#">Link</a>
Tuning Language Models as Training Data Generators for Augmentation-Enhanced Few-Shot Learning <sup>[3]</sup>	Instruction	CTRL	Model Training, Nvidia A100 GPUs, 10 minutes per task	ICML'23	<a href="#">Link</a>
SASS: SELF-ALIGNMENT WITH SEMI-SUPERVISED INSTRUCTION DATA GENERATION <sup>[4]</sup>	Instruction	LLaMA	Model Training, Nvidia A100 GPUs	OpenReview'24	Not Available
DAIL: Data Augmentation for In-Context Learning via Self-Paraphrase <sup>[5]</sup>	Instruction	ChatGPT	API Calling	Arxiv'23	Not Available
LongForm: Effective Instruction Tuning with Reverse Instructions <sup>[6]</sup>	Instruction	GPT-3	PI Calling	ICLR'24	<a href="#">Link</a>
Large Language Model as Attributed Training Data Generator: A Tale of Diversity and Bias <sup>[7]</sup>	Instruction	ChatGPT	API Calling	NeurIPS'23	<a href="#">Link</a>
SELF-QA: Unsupervised Knowledge Guided Language Model Alignment <sup>[8]</sup>	Instruction & Response	BLOOM	Model Inference	Arxiv'23	Not Available
LARGE LANGUAGE MODELS CAN SELF-IMPROVE <sup>[9]</sup>	Response	PaLM-540B	Model Inference	EMNLP'23	Not Available
Self-Distillation Bridges Distribution Gap in Language Model Fine-Tuning <sup>[10]</sup>	Response	LLaMA-2	Model Inference	ACL'24	<a href="#">Link</a>
Mixture of insight/Tal Experts (MoTE): The Synergy of Thought Chains and Expert Mixtures in Self-Alignment <sup>[11]</sup>	Response	Alpaca	Model Inference	Arxiv'24	Not Available
Human-Instruction-Free LLM Self-Alignment with Limited Samples <sup>[12]</sup>	Instruction & Response	Multiple LLMs	Model Inference, single NVIDIA A100 80G GPU	Arxiv'24	Not Available
Principle-Driven Self-Alignment of Language Models from Scratch with Minimal Human Supervision <sup>[13]</sup>	Response	LLaMA	Model Inference	NeurIPS'23	<a href="#">Link</a>
Step-On-Feet Tuning: Scaling Self-Alignment of LLMs via Bootstrapping <sup>[14]</sup>	Response	LLaMA-2	Model Inference	Arxiv'24	Not Available
Assessing Empathy in Large Language Models with Real-World Physician-Patient Interactions <sup>[15]</sup>	Response	LLaMA	Model Inference	Arxiv'24	Not Available
<b>Rationale</b>					
Large Language Models are Zero-Shot Reasoners <sup>[16]</sup>	Rationale - CoT	Multiple LLMs	API Calling	NeurIPS'22	Not Available
Tree of Thoughts: Deliberate Problem Solving with Large Language Models <sup>[17]</sup>	Rationale - Tree	GPT-4	API Calling, \$0.74 per sample	NeurIPS'22	<a href="#">Link</a>
Reasoning with Language Model is Planning with World Model <sup>[18]</sup>	Rationale - Tree	LLaMA	Model Inference, 4x24 GB NVIDIA A5000 GPUs	EMNLP'23	<a href="#">Link</a>
Graph of Thoughts: Solving Elaborate Problems with Large Language Models <sup>[19]</sup>	Rationale - Graph	GPT-3.5	API Calling	AAAI'24	<a href="#">Link</a>
Beyond Chain-of-Thought, Effective Graph-of-Thought Reasoning in Language Models <sup>[20]</sup>	Rationale - Graph	GPT-3	API Calling	Arxiv'23	<a href="#">Link</a>
CHAIN-OF-TABLE: EVOLVING TABLES IN THE REASONING CHAIN FOR TABLE UNDERSTANDING <sup>[21]</sup>	Rationale - Table	Multiple LLMs	API Calling & Model Inference	ICLR'24	Not Available
Program of Thoughts Prompting: D disentangling Computation from Reasoning for Numerical Reasoning Tasks <sup>[22]</sup>	Rationale - Program	Multiple LLMs	API Calling & Model Inference	TMLR'23	Not Available
The Art of SOCRATIC QUESTIONING: Recursive Thinking with Large Language Models <sup>[23]</sup>	Rationale - Reversion	ChatGPT	API Calling, 9.22 calls per sample	EMNLP'23	<a href="#">Link</a>
Interpreting Pretrained Language Models via Concept Bottlenecks <sup>[24]</sup>	Rationale - Concept	ChatGPT	API Calling	PAKDD'24	<a href="#">Link</a>
PINTO: FAITHFUL LANGUAGE REASONING USING PROMPT-GENERATED RATIONALES <sup>[25]</sup>	Rationale - CoT	GPT-neox	Model Inference	ICLR'23	<a href="#">Link</a>
SCOTT: Self-Consistent Chain-of-Thought Distillation <sup>[26]</sup>	Rationale - CoT	GPT-neox	Model Inference	ACL'23	<a href="#">Link</a>
LogiCoT: Logical Chain-of-Thought Instruction Tuning <sup>[27]</sup>	Rationale - CoT	GPT-4	API Calling	EMNLP'23	Not Available
Distilling Reasoning Capabilities into Smaller Language Models <sup>[28]</sup>	Rationale - CoT	GPT-3	API Calling	ACL'23	Not Available
Knowledge-Augmented Reasoning Distillation for Small Language Models in Knowledge-Intensive Tasks <sup>[29]</sup>	Rationale - CoT	ChatGPT	API Calling	NeurIPS'23	<a href="#">Link</a>
Making Pre-trained Language Models Better Few-shot Learners <sup>[30]</sup>	Rationale - Diverse Thinking	GPT-3	API Calling	ACL'21	<a href="#">Link</a>
SELF-CONSISTENCY IMPROVES CHAIN OF THOUGHT REASONING IN LANGUAGE MODELS <sup>[31]</sup>	Rationale - Diverse Thinking	Multiple LLMs	API Calling & Model Inference	ICLR'23	Not Available
UNIVERSAL SELF-CONSISTENCY FOR LARGE LANGUAGE MODEL GENERATION <sup>[32]</sup>	Rationale - Diverse Thinking	Multiple LLMs	API Calling	Arxiv'23	Not Available
Plan, Verify and Switch: Integrated Reasoning with Diverse X-of-Thoughts <sup>[33]</sup>	Rationale - Diverse Thinking	ChatGPT	API Calling	EMNLP'23	<a href="#">Link</a>
Eliminating Reasoning via Inferring with Planning: A New Framework to Guide LLMs' Non-linear Thinking <sup>[34]</sup>	Rationale - Elimination	PaLM2	API Calling	Arxiv'23	Not Available
It's Not Easy Being Wrong: Large Language Models Struggle with Process of Elimination Reasoning <sup>[35]</sup>	Rationale - Elimination	Multiple LLMs	API Calling	ACL'24	<a href="#">Link</a>
POE: Process of Elimination for Multiple Choice Reasoning <sup>[36]</sup>	Rationale - Elimination	FLAN-T5	Model Inference	EMNLP'23	<a href="#">Link</a>
Exchange-of-Thought: Enhancing Large Language Model Capabilities through Cross-Model Communication <sup>[37]</sup>	Rationale - Collaboration	ChatGPT	API Calling	EMNLP'23	Not Available
Encouraging Divergent Thinking in Large Language Models through Multi-Agent Debate <sup>[38]</sup>	Rationale - Collaboration	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
Towards Reasoning in Large Language Models via Multi-Agent Peer Review Collaboration <sup>[39]</sup>	Rationale - Collaboration	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
DYNAMIC LLM-AGENT NETWORK: AN LLM-AGENT COLLABORATION FRAMEWORK WITH AGENT TEAM OPTIMIZATION <sup>[40]</sup>	Rationale - Collaboration	ChatGPT	API Calling, 16.5 calls per sample	Arxiv'23	<a href="#">Link</a>
<b>Pair-wise Feedback</b>					
Constitutional AI: Harmlessness from AI Feedback <sup>[41]</sup>	Pairwise Feedback	Multiple LLMs	Model Inference	Arxiv'22	<a href="#">Link</a>
RLAIF: Scaling Reinforcement Learning from Human Feedback with AI Feedback <sup>[42]</sup>	Pairwise Feedback	PaLM-2	Model Inference, \$0.67 per sample	Arxiv'23	Not Available
Self-Rewarding Language Models <sup>[43]</sup>	Pairwise Feedback	LLaMA-2	Model Inference	Arxiv'24	Not Available
SALMON: SELF-ALIGNMENT WITH INSTRUCTABLE REWARD MODELS <sup>[44]</sup>	Pairwise Feedback	LLaMA-2	Model Inference	ICLR'24	<a href="#">Link</a>
Self-Alignment for Factuality: Mitigating Hallucinations in LLMs via Self-Evaluation <sup>[45]</sup>	Pairwise Feedback	LLaMA	Model Inference	Arxiv'24	<a href="#">Link</a>
West-of-N: Synthetic Preference Generation for Improved Reward Modeling <sup>[46]</sup>	Pairwise Feedback	T5-XXL	Model Inference	Arxiv'24	Not Available
Learning Reward for Robot Skills Using Large Language Models via Self-Alignment <sup>[47]</sup>	Pairwise Feedback	ChatGPT	API Calling	ICML'24	<a href="#">Link</a>
Aligning Large Language Models through Synthetic Feedback <sup>[48]</sup>	Pairwise Feedback	LLaMA	Model Inference	EMNLP'23	<a href="#">Link</a>
Optimizing Language Model's Reasoning Abilities with Weak Supervision <sup>[49]</sup>	Pairwise Feedback	LLaMA	Model Inference	Arxiv'24	Not Available
RLCD: REINFORCEMENT LEARNING FROM CONTRASTIVE DISTILLATION FOR LM ALIGNMENT <sup>[50]</sup>	Pairwise Feedback	LLaMA	Model Inference	ICLR'24	<a href="#">Link</a>
Automatic Pair Construction for Contrastive Post-training <sup>[51]</sup>	Pairwise Feedback	LLaMA	Model Inference, 16 Nvidia V100 GPUs	NAACL'24	Not Available
Reinforcement Learning from Reflective Feedback (RLRF): Aligning and Improving LLMs via Fine-Grained Self-Reflection <sup>[52]</sup>	Pairwise Feedback	LLaMA-2	Model Inference, 16 Nvidia V100 GPUs	Arxiv'24	Not Available
Improving Language Model Reasoning with Self-motivated Learning <sup>[53]</sup>	Pairwise Feedback	LLaMA-2	Model Inference	LREC'24	Not Available

Note: <sup>[1]</sup>(Yoo et al., 2021); <sup>[2]</sup>(Wang et al., 2023e); <sup>[3]</sup>(Meng et al., 2023); <sup>[4]</sup>(Wang et al., 2023f); <sup>[5]</sup>(Li et al., 2023a); <sup>[6]</sup>(Köksal et al.); <sup>[7]</sup>(Yu et al., 2024); <sup>[8]</sup>(Zhang and Yang, 2023a); <sup>[9]</sup>(Huang et al., 2023); <sup>[10]</sup>(Yang et al., 2024b); <sup>[11]</sup>(Liu et al., 2024); <sup>[12]</sup>(Guo et al., 2024a); <sup>[13]</sup>(Sun et al., 2024); <sup>[14]</sup>(Wang et al., 2024a); <sup>[15]</sup>(Luo et al., 2024); <sup>[16]</sup>(Kojima et al., 2022); <sup>[17]</sup>(Yao et al., 2024); <sup>[18]</sup>(Hao et al., 2023); <sup>[19]</sup>(Besta et al., 2024); <sup>[20]</sup>(Yao et al., 2023); <sup>[21]</sup>(Wang et al., 2024d); <sup>[22]</sup>(Chen et al., 2023e); <sup>[23]</sup>(Qi et al., 2023); <sup>[24]</sup>(Tan et al., 2023); <sup>[25]</sup>(Wang et al., 2022a); <sup>[26]</sup>(Wang et al., 2023b); <sup>[27]</sup>(Liu et al., 2023a); <sup>[28]</sup>(Shridhar et al., 2023); <sup>[29]</sup>(Kang et al., 2024); <sup>[30]</sup>(Gao et al., 2021); <sup>[31]</sup>(Wang et al., 2022b); <sup>[32]</sup>(Chen et al., 2023f); <sup>[33]</sup>(Liu et al., 2023c); <sup>[34]</sup>(Tong et al., 2023); <sup>[35]</sup>(Balepur et al., 2023); <sup>[36]</sup>(Ma and Du, 2023); <sup>[37]</sup>(Yin et al., 2023); <sup>[38]</sup>(Liang et al., 2023); <sup>[39]</sup>(Xu et al., 2023d); <sup>[40]</sup>(Liu et al., 2023e); <sup>[41]</sup>(Bai et al., 2022); <sup>[42]</sup>(Lee et al., 2023b); <sup>[43]</sup>(Yuan et al., 2024); <sup>[44]</sup>(Sun et al., 2023b); <sup>[45]</sup>(Zhang et al., 2024a); <sup>[46]</sup>(Pace et al., 2024); <sup>[47]</sup>(Zeng et al., 2024); <sup>[48]</sup>(Kim et al., 2023b); <sup>[49]</sup>(Tong et al., 2024b); <sup>[50]</sup>(Yang et al., 2023c); <sup>[51]</sup>(Xu et al., 2023c); <sup>[52]</sup>(Lee et al., 2024a); <sup>[53]</sup>(Feng et al., 2024).

Table 1: A list of representative LLM-Based Annotation Generation (Instruction & Response, Rationale, Pairwise Feedback) papers with open-source code/data.

Paper	Data Type	Backbone	Annotation Cost	Venue	Code/Data Link
Textual Feedback					
SELF-REFINE: Iterative Refinement with Self-Feedback <sup>[1]</sup>	Textual Feedback	Multiple LLMs	API Calling	NeurIPS'23	Not Available
Reflexion: Language Agents with Verbal Reinforcement Learning <sup>[2]</sup>	Textual Feedback	GPT-3	API Calling	NeurIPS'23	<a href="#">Link</a>
Iterative Translation Refinement with Large Language Models <sup>[3]</sup>	Textual Feedback	GPT-3.5	API Calling	Arxiv'23	Not Available
Leveraging GPT-4 for Automatic Translation Post-Editing <sup>[4]</sup>	Textual Feedback	Multiple LLMs	API Calling	EMNLP'23	Not Available
A New Benchmark and Reverse Validation Method for Passage-level Hallucination Detection <sup>[5]</sup>	Textual Feedback	ChatGPT	API Calling	EMNLP'23	<a href="#">Link</a>
SELF-CHECKGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models <sup>[6]</sup>	Textual Feedback	Multiple LLMs	API Calling & Model Inference	EMNLP'23	<a href="#">Link</a>
Improving Factuality and Reasoning in Language Models through Multiagent Debate <sup>[7]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling		<a href="#">Link</a>
Towards Reasoning in Large Language Models via Multi-Agent Peer Review Collaboration <sup>[8]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling	Arxiv'23	<a href="#">Link</a>
LM vs LM: Detecting Factual Errors via Cross Examination <sup>[9]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling & Model Inference	EMNLP'23	Not Available
Improving Language Model Negotiation with Self-Play and In-Context Learning from AI Feedback <sup>[10]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling	Arxiv'23	<a href="#">Link</a>
PRD: Peer Rank and Discussion Improve Large Language Model based Evaluations <sup>[11]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling, \$0.14 per sample	Arxiv'23	<a href="#">Link</a>
PRE: A Peer Review Based Large Language Model Evaluator <sup>[12]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling	Arxiv'24	Not Available
PCO: Peer Review in LLMs based on the Consistency Optimization <sup>[13]</sup>	Textual Feedback - Peer Review	Multiple LLMs	API Calling & Model Inference	Arxiv'24	Not Available
Learning from Mistakes via Cooperative Study Assistant for Large Language Models <sup>[14]</sup>	Textual Feedback - Mistake	Multiple LLMs	Model Inference	EMNLP'23	<a href="#">Link</a>
Learning From Mistakes Makes LLM Better Reasoner <sup>[15]</sup>	Textual Feedback - Mistake	GPT-4	API Calling	Arxiv'23	<a href="#">Link</a>
GAINING WISDOM FROM SETBACKS: ALIGNING LARGE LANGUAGE MODELS VIA MISTAKE ANALYSIS <sup>[16]</sup>	Textual Feedback - Mistake	Multiple LLMs	API Calling & Modeling Inference	ICLR'24	Not Available
Can LLMs Learn from Previous Mistakes? Investigating LLMs' Errors to Boost for Reasoning <sup>[17]</sup>	Textual Feedback - Mistake	Multiple LLMs	API Calling & Modeling Inference	ACL'24	<a href="#">Link</a>
Other Domain-specific Data					
SODA: Million-scale Dialogue Distillation with Social Commonsense Contextualization <sup>[18]</sup>	Dialogue	GPT-3.5	API Calling, \$0.02 per dialogue	EMNLP'23	<a href="#">Link</a>
Baize: An Open-Source Chat Model with Parameter-Efficient Tuning on Self-Chat Data <sup>[19]</sup>	Dialogue	Alpaca	Model Inference	EMNLP'23	<a href="#">Link</a>
PLACES: Prompting Language Models for Social Conversation Synthesis <sup>[20]</sup>	Dialogue	Multiple LLMs	Model Inference	EACL'24	Not Available
CAMEL: Communicative Agents for "Mind" Exploration of Large Language Model Society <sup>[21]</sup>	Dialogue	ChatGPT	API Calling	NeurIPS'23	<a href="#">Link</a>
AUGESC: Dialogue Augmentation with Large Language Models for Emotional Support Conversation <sup>[22]</sup>	Dialogue	GPT-J	Model Inference	ACL'23	<a href="#">Link</a>
Weakly Supervised Data Augmentation Through Prompting for Dialogue Understanding <sup>[23]</sup>	Dialogue	GPT-J	Model Inference	NeurIPS'23	Not Available
Reflect, Not Reflex: Inference-Based Common Ground Improves Dialogue Response Quality <sup>[24]</sup>	Dialogue	GPT-3	API Calling	EMNLP'22	<a href="#">Link</a>
ASDOT: Any-Shot Data-to-Text Generation with Pretrained Language Models <sup>[25]</sup>	Context	GPT-3	API Calling, \$23 in total	EMNLP'22	<a href="#">Link</a>
Contextualization Distillation from Large Language Model for Knowledge Graph Completion <sup>[26]</sup>	Context	PaLM-2	API Calling	EACL'24	<a href="#">Link</a>
Towards Ontology-Enhanced Representation Learning for Large Language Models <sup>[27]</sup>	Context	ChatGPT	API Calling	Arxiv'24	<a href="#">Link</a>
DALK: Dynamic Co-Augmentation of LLMs and KG to answer Alzheimer's Disease Questions with Scientific Literature <sup>[28]</sup>	Graph	ChatGPT	API Calling	Arxiv'24	<a href="#">Link</a>
Automated Construction of Theme-specific Knowledge Graphs <sup>[29]</sup>	Graph	GPT-4	API Calling	Arxiv'24	Not Available
Large Language Models Can Learn Temporal Reasoning <sup>[30]</sup>	Graph	GPT-3.5	API Calling	ACL'24	<a href="#">Link</a>
Moving from Tabular Knowledge Graph Quality Assessment to RDF Triples Leveraging ChatGPT <sup>[31]</sup>	Graph	ChatGPT	API Calling	Arxiv'24	<a href="#">Link</a>
Language Models as Zero-Shot Planners: Extracting Actionable Knowledge for Embodied Agents <sup>[32]</sup>	Plan	GPT-3	API Calling	ICML'22	<a href="#">Link</a>
Do As I Can, Not As I Say: Grounding Language in Robotic Affordances <sup>[33]</sup>	Plan	Multiple LLMs	API Calling & Model Inference	CoRL'21	<a href="#">Link</a>
SayPlan: Grounding Large Language Models using 3D Scene Graphs for Scalable Robot Task Planning <sup>[34]</sup>	Plan	GPT-3.5	API Calling	CoRL'23	<a href="#">Link</a>
PROGPROMPT: Generating Situated Robot Task Plans using Large Language Models <sup>[35]</sup>	Plan	GPT-3	API Calling	ICRA'23	<a href="#">Link</a>
Text2Motion: From Natural Language Instructions to Feasible Plans <sup>[36]</sup>	Plan	GPT-3.5	API Calling	Autonomous Robots'23	<a href="#">Link</a>
GENSIM: GENERATING ROBOTIC SIMULATION TASKS VIA LARGE LANGUAGE MODELS <sup>[37]</sup>	Simulation Task	GPT-4	API Calling	ICML'24	<a href="#">Link</a>
Scaling Up and Distilling Down: Language-Guided Robot Skill Acquisition <sup>[38]</sup>	Simulation Task	Multiple LLMs	API Calling	CoRL'23	<a href="#">Link</a>
REWARD DESIGN WITH LANGUAGE MODELS <sup>[39]</sup>	Reward	GPT-3	API Calling	ICLR'23	<a href="#">Link</a>
Guiding Pretraining in Reinforcement Learning with Large Language Models <sup>[40]</sup>	Reward	GPT-3	API Calling, 0.02 second per call	ICML'23	Not Available
Enhanced Visual Instruction Tuning with Synthesized Image-Dialogue Data <sup>[41]</sup>	Visual Instruction Tuning Data	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
LAMM: Language-Assisted Multi-Modal Instruction-Tuning Dataset, Framework, and Benchmark <sup>[42]</sup>	Visual Instruction Tuning Data	GPT-4	API Calling	NeurIPS'23	<a href="#">Link</a>
TOMGPT: Reliable Text-Only Training Approach for Cost-Effective Multi-modal Large Language Model <sup>[43]</sup>	Context	ChatGPT	API Calling	TKDD'24	Not Available
LLM Based Generation of Item-Description for Recommendation System <sup>[44]</sup>	Item Description	Alpaca	Model Inference	RecSys'23	Not Available
PMG: Personalized Multimodal Generation with Large Language <sup>[45]</sup>	Context	Multiple LLMs	Model Inference	WWW'24	<a href="#">Link</a>
LLMRec: Large Language Models with Graph Augmentation for Recommendation <sup>[46]</sup>	Augmented Implicit Feedback	ChatGPT	API Calling, \$21.14	WSDM'24	<a href="#">Link</a>
Large Language Models as Evaluators for Recommendation Explanations <sup>[47]</sup>	Explanation	Multiple LLMs	API Calling & Model Inference, less than \$0.02 per sample	Arxiv'24	<a href="#">Link</a>
Exploiting Asymmetry for Synthetic Training Data Generation: SynthIE and the Case of Information Extraction <sup>[48]</sup>	IE Sample	GPT-3.5	API Calling, \$223.55 for entire dataset	EMNLP'23	<a href="#">Link</a>
InPars-v2: Large Language Models as Efficient Dataset Generators for Information Retrieval <sup>[49]</sup>	IE sample	GPT-J	Model Inference, 30 hours on an A100 GPU to generate 100k queries	Arxiv'23	<a href="#">Link</a>
READ: Improving Relation Extraction from an Adversarial Perspective <sup>[50]</sup>	IE Sample	ChatGPT	API Calling	NAACL'24	<a href="#">Link</a>
STAR: Boosting Low-Resource Information Extraction by Structure-to-Text Data Generation with Large Language Models <sup>[51]</sup>	IE Sample	Multiple LLMs	API Calling	AAAI'24	<a href="#">Link</a>
Adjudicating LLMs as PropBank Annotators <sup>[52]</sup>	IE Label	Multiple LLMs	API Calling	LREC'24	<a href="#">Link</a>
A Causal Explainable Guardrails for Large Language Models <sup>[53]</sup>	Representation	GPT-4	API Calling	Arxiv'24	Not Available
Zero-shot LLM-guided Counterfactual Generation for Text <sup>[54]</sup>	Context	Multiple LLMs	API Calling	Arxiv'24	Not Available
Text classification of column headers with a controlled vocabulary: leveraging LLMs for metadata enrichment <sup>[55]</sup>	Metadata	ChatGPT	API Calling	Arxiv'24	<a href="#">Link</a>

Note: <sup>[1]</sup>(Madaan et al., 2024); <sup>[2]</sup>(Shinn et al., 2024); <sup>[3]</sup>(Chen et al., 2023c); <sup>[4]</sup>(Raunak et al., 2023); <sup>[5]</sup>(Yang et al., 2023d); <sup>[6]</sup>(Manakul et al., 2023); <sup>[7]</sup>(Du et al., 2023a); <sup>[8]</sup>(Xu et al., 2023d); <sup>[9]</sup>(Cohen et al., 2023); <sup>[10]</sup>(Fu et al., 2023); <sup>[11]</sup>(Li et al., 2023d); <sup>[12]</sup>(Chu et al., 2024b); <sup>[13]</sup>(Ning et al., 2024); <sup>[14]</sup>(Wang and Li, 2023); <sup>[15]</sup>(An et al., 2023); <sup>[16]</sup>(Chen et al., 2023a); <sup>[17]</sup>(Tong et al., 2024a); <sup>[18]</sup>(Kim et al., 2023a); <sup>[19]</sup>(Xu et al., 2023b); <sup>[20]</sup>(Chen et al., 2023b); <sup>[21]</sup>(Li et al., 2024d); <sup>[22]</sup>(Zheng et al., 2023); <sup>[23]</sup>(Chen et al., 2022); <sup>[24]</sup>(Zhou et al., 2022a); <sup>[25]</sup>(Xiang et al., 2022); <sup>[26]</sup>(Li et al., 2024b); <sup>[27]</sup>(Ronzano and Nanavati, 2024); <sup>[28]</sup>(Li et al., 2024c); <sup>[29]</sup>(Ding et al., 2024); <sup>[30]</sup>(Xiong et al., 2024); <sup>[31]</sup>(Tuozzo, 2022); <sup>[32]</sup>(Huang et al., 2022); <sup>[33]</sup>(Brohan et al., 2023); <sup>[34]</sup>(Rana et al., 2023); <sup>[35]</sup>(Singh et al., 2023); <sup>[36]</sup>(Lin et al., 2023a); <sup>[37]</sup>(Wang et al., 2023a); <sup>[38]</sup>(Ha et al., 2023); <sup>[39]</sup>(Kwon et al., 2022); <sup>[40]</sup>(Du et al., 2023b); <sup>[41]</sup>(Li et al., 2023f); <sup>[42]</sup>(Yin et al., 2024); <sup>[43]</sup>(Chen et al., 2024a); <sup>[44]</sup>(Acharya et al., 2023); <sup>[45]</sup>(Shen et al., 2024); <sup>[46]</sup>(Wei et al., 2024); <sup>[47]</sup>(Zhang et al., 2024b); <sup>[48]</sup>(Josifoski et al., 2023); <sup>[49]</sup>(Jeronymo et al., 2023); <sup>[50]</sup>(Li et al., 2024a); <sup>[51]</sup>(Ma et al., 2024); <sup>[52]</sup>(Bonn et al., 2024); <sup>[53]</sup>(Chu et al., 2024a); <sup>[54]</sup>(Bhattacharjee et al., 2024); <sup>[55]</sup>(Martorana et al., 2024).

Table 2: A list of representative LLM-Based Annotation Generation (Textual Feedback, Other Domain-specific Data) papers with open-source code/data.

Paper	Data Type	Backbone	Annotation Cost	Venue	Code/Data Link
Filter & Selection					
Constitutional AI: Harmlessness from AI Feedback <sup>[1]</sup>	Pairwise Feedback	Multiple LLMs	Model Inference	Arxiv'22	<a href="#">Link</a>
SODA: Million-scale Dialogue Distillation with Social Commonsense Contextualization <sup>[2]</sup>	Dialogue	GPT-3.5	API Calling, \$0.02 per dialogue	EMNLP'23	<a href="#">Link</a>
Aligning Large Language Models through Synthetic Feedback <sup>[3]</sup>	Pairwise Feedback	LLaMA	Model Inference	EMNLP'23	<a href="#">Link</a>
AUGESC: Dialogue Augmentation with Large Language Models for Emotional Support Conversation <sup>[4]</sup>	Dialogue	GPT-J	Model Inference	ACL'23	<a href="#">Link</a>
SELF-QA: Unsupervised Knowledge Guided Language Model Alignment <sup>[5]</sup>	Instruction & Response	BLOOM	Model Inference	Arxiv'23	Not Available
Human-Instruction-Free LLM Self-Alignment with Limited Samples <sup>[6]</sup>	Instruction & Response	Multiple LLMs	Model Inference, single NVIDIA A100 80G GPU	Arxiv'24	Not Available
Automated Construction of Theme-specific Knowledge Graphs <sup>[7]</sup>	Graph	GPT-4	API Calling	Arxiv'24	Not Available
Large Language Models Are Reasoning Teachers <sup>[8]</sup>	CoT	GPT-3.5	API Calling	ACL'23	<a href="#">Link</a>
Knowledge-Augmented Reasoning Distillation for Small Language Models in Knowledge-Intensive Tasks <sup>[9]</sup>	Rationale - CoT	ChatGPT	API Calling	NeurIPS'23	<a href="#">Link</a>
SELF-CONSISTENCY IMPROVES CHAIN OF THOUGHT REASONING IN LANGUAGE MODELS <sup>[10]</sup>	Rationale - Diverse Thinking	Multiple LLMs	API Calling & Model Inference	ICLR'23	Not Available
Making Large Language Models Better Data Creators <sup>[11]</sup>	Instruction & Response	ChatGPT	API Calling	EMNLP'23	<a href="#">Link</a>
Automated Construction of Theme-specific Knowledge Graphs <sup>[12]</sup>	Graph	GPT-4	API Calling	Arxiv'24	Not Available
Reinforced Self-Training (ReST) for Language Modeling <sup>[13]</sup>	Response	Multiple LLMs	Model Inference	Arxiv'24	Not Available
RAFT: Reward rAnked FineTuning for Generative Foundation Model Alignment <sup>[14]</sup>	Response	LLaMA	Model Inference	TMLR	<a href="#">Link</a>
Selective In-Context Data Augmentation for Intent Detection using Pointwise V-Information <sup>[15]</sup>	Instruction	OPT	Model Inference	EACL'24	Not Available
CodeCLM: Aligning Language Models with Tailored Synthetic Data <sup>[16]</sup>	Instruction	LLaMA	Model Inference	NAACL'24	Not Available
DISCO: Distilling Counterfactuals with Large Language Models <sup>[17]</sup>	CoT	GPT-3	API Callin	ACL'23	<a href="#">Link</a>
LARGE LANGUAGE MODELS CAN SELF-IMPROVE <sup>[18]</sup>	Response	PaLM-540B	Model Inference	EMNLP'23	Not Available
West-of-N: Synthetic Preference Generation for Improved Reward Modeling <sup>[19]</sup>	Pairwise Feedback	TS-XXL	Model Inference	Arxiv'24	Not Available
SELF: SELF-EVOLUTION WITH LANGUAGE FEEDBACK <sup>[20]</sup>	Response	Multiple LLMs	Model Inference	Arxiv'23	Not Available
InPars-v2: Large Language Models as Efficient Dataset Generators for Information Retrieval <sup>[21]</sup>	IE sample	GPT-J	Model Inference, 30 hours on an A100 GPU to generate 100k queries	Arxiv'23	<a href="#">Link</a>
DALK: Dynamic Co-Augmentation of LLMs and KG to answer Alzheimer's Disease Questions with Scientific Literature <sup>[22]</sup>	Graph	ChatGPT	API Calling	Arxiv'24	<a href="#">Link</a>
Optimizing Language Model's Reasoning Abilities with Weak Supervision <sup>[23]</sup>	Pairwise Feedback	LLaMA	Model Inference	Arxiv'24	Not Available

Note: <sup>[1]</sup>(Bai et al., 2022); <sup>[2]</sup>(Kim et al., 2023a); <sup>[3]</sup>(Kim et al., 2023b); <sup>[4]</sup>(Zheng et al., 2023); <sup>[5]</sup>(Zhang and Yang, 2023a); <sup>[6]</sup>(Guo et al., 2024a); <sup>[7]</sup>(Ding et al., 2024); <sup>[8]</sup>(Ho et al., 2023); <sup>[9]</sup>(Kang et al., 2024); <sup>[10]</sup>(Wang et al., 2022b); <sup>[11]</sup>(Lee et al., 2023a); <sup>[12]</sup>(Ding et al., 2024); <sup>[13]</sup>(Gulcehre et al., 2023); <sup>[14]</sup>(Dong et al., 2023); <sup>[15]</sup>(Lin et al., 2023b); <sup>[16]</sup>(Wang et al., 2024c); <sup>[17]</sup>(Chen et al., 2023g); <sup>[18]</sup>(Huang et al., 2023); <sup>[19]</sup>(Pace et al., 2024); <sup>[20]</sup>(Lu et al., 2023); <sup>[21]</sup>(Jeronymo et al., 2023); <sup>[22]</sup>(Li et al., 2024c); <sup>[23]</sup>(Tong et al., 2024b).

Table 3: A list of representative LLM-Generated Annotation Assessment papers with open-source code/data.



Paper	Data Type	Backbone	Annotation Cost	Venue	Code/Data Link
Supervised Fine-tuning					
LARGE LANGUAGE MODELS CAN SELF-IMPROVE <sup>[1]</sup>	Response	PaLM-540B	Model Inference	EMNLP'23	Not Available
SELF-INSTRUCT: Aligning Language Models with Self-Generated Instructions <sup>[2]</sup>	Instruction & Response	GPT-3	API Calling, \$600 for entire dataset	ACL'23	<a href="#">Link</a>
SELF: SELF-EVOLUTION WITH LANGUAGE FEEDBACK <sup>[3]</sup>	Response	Multiple LLMs	Model Inference	Arxiv'23	Not Available
Self-Distillation Bridges Gap in Language Model Fine-Tuning <sup>[4]</sup>	Response	LLaMA-2	Model Inference	ACL'24	<a href="#">Link</a>
Self-Play Fine-Tuning Converts Weak Language Models to Strong Language Models <sup>[5]</sup>	Response	zephyr	Model Inference	Arxiv'24	<a href="#">Link</a>
Self-playing Adversarial Language Game Enhances LLM Reasoning <sup>[6]</sup>	Response	Multiple LLMs	Model Inference	Arxiv'24	<a href="#">Link</a>
Stanford alpaca: An instruction-following llama model <sup>[7]</sup>	Response	GPT-3.5	API Calling	Arxiv'23	<a href="#">Link</a>
Vicuna: An open-source chatbot impressing gpt-4 with 90%+ chatgpt quality <sup>[8]</sup>	Response	GPT-4	API Calling	Arxiv'23	<a href="#">Link</a>
WizardLM: Empowering large language models to follow complex instructions <sup>[9]</sup>	Instruction	LLaMA	Model Inference	Arxiv'23	<a href="#">Link</a>
Generating training data with language models: Towards zero-shot language understanding <sup>[10]</sup>	Instruction	CTRL	Model Inference	NeurIPS	<a href="#">Link</a>
Tuning Language Models as Training Data Generators for Augmentation-Enhanced Few-Shot Learning <sup>[11]</sup>	Instruction	CTRL	Model Training	ICML'23	<a href="#">Link</a>
Noise-Robust Fine-Tuning of Pretrained Language Models via External Guidance <sup>[12]</sup>	Response	ChatGPT	API Calling	EMNLP'23	<a href="#">Link</a>
PINTO: FAITHFUL LANGUAGE REASONING USING PROMPT-GENERATED RATIONALES <sup>[13]</sup>	Rationale - CoT	GPT-neox	Model Inference	ICLR'23	<a href="#">Link</a>
Distilling Reasoning Capabilities into Smaller Language Models <sup>[14]</sup>	Rationale - CoT	GPT-3	API Calling	ACL'23	Not Available
LogCoT: Logical Chain-of-Thought Instruction Tuning <sup>[15]</sup>	Rationale - CoT	GPT-4	API Calling	EMNLP'23	Not Available
Knowledge-Augmented Reasoning Distillation for Small Language Models in Knowledge-Intensive Tasks <sup>[16]</sup>	Rationale - CoT	ChatGPT	API Calling	NeurIPS'23	<a href="#">Link</a>
Baize: An Open-Source Chat Model with Parameter-Efficient Tuning on Self-Chat Data <sup>[17]</sup>	Dialogue	Alpaca	Model Inference	EMNLP'23	<a href="#">Link</a>
Exploiting Asymmetry for Synthetic Training Data Generation: SynthE and the Case of Information Extraction <sup>[18]</sup>	IE Sample	GPT-3.5	API Calling, \$223.55 for entire dataset	EMNLP'23	<a href="#">Link</a>
InPars-v2: Large Language Models as Efficient Dataset Generators for Information Retrieval <sup>[19]</sup>	IE sample	GPT-J	Model Inference, 30 hours on an A100 GPU to generate 100k queries	Arxiv'23	<a href="#">Link</a>
Code alpaca: An instruction-following llama model for code generation <sup>[20]</sup>	Instruction & Response	Alpaca	Model Inference	Arxiv'23	<a href="#">Link</a>
Code llama: Open foundation models for code <sup>[21]</sup>	Instruction & Response	Multiple LLMs	Model Inference	Arxiv'23	<a href="#">Link</a>
HuatuGPT: Towards Taming Language Model to Be a Doctor <sup>[22]</sup>	Instruction & Response	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
Doctorglm: Fine-tuning your chinese doctor is not a herculean task <sup>[23]</sup>	Response	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
Xuanyuan 2.0: A large chinese financial chat model with hundreds of billions parameters <sup>[24]</sup>	Instruction & Response	BLOOM	Model Inference	CIKM'23	Not Available
Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct <sup>[25]</sup>	Pairwise Feedback	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
Gitmlt: A unified graph-text model for instruction-based molecule zero-shot learning <sup>[26]</sup>	Instruction	ChatGPT	API Calling	NeurIPS'23	<a href="#">Link</a>
Alignment Tuning					
Automatic Pair Construction for Contrastive Post-training <sup>[27]</sup>	Pairwise Feedback	LLaMA	Model Inference, 16 Nvidia V100 GPUs	NAACL'24	Not Available
Aligning Large Language Models through Synthetic Feedback <sup>[28]</sup>	Pairwise Feedback	LLaMA	Model Inference	EMNLP'23	<a href="#">Link</a>
West-of-N: Synthetic Preference Generation for Improved Reward Modeling <sup>[29]</sup>	Pairwise Feedback	TS-XXL	Model Inference	Arxiv'24	Not Available
Learning Reward for Robot Skills Using Large Language Models via Self-Alignment <sup>[30]</sup>	Pairwise Feedback	ChatGPT	API Calling	ICML'24	<a href="#">Link</a>
SALMON: SELF-ALIGNMENT WITH INSTRUCTABLE REWARD MODELS <sup>[31]</sup>	Pairwise Feedback	LLaMA-2	Model Inference	ICLR'24	<a href="#">Link</a>
Self-Rewarding Language Models <sup>[32]</sup>	Pairwise Feedback	LLaMA-2	Model Inference	Arxiv'24	Not Available
Self-Alignment for Factuality: Mitigating Hallucinations in LLMs via Self-Evaluation <sup>[33]</sup>	Pairwise Feedback	LLaMA	Model Inference	Arxiv'24	<a href="#">Link</a>
Aligning Large Language Models by On-Policy Self-Judgment <sup>[34]</sup>	Response	LLaMA-2	Model Inference	Arxiv'24	<a href="#">Link</a>
Optimizing Language Model's Reasoning Abilities with Weak Supervision <sup>[35]</sup>	Pairwise Feedback	LLaMA	Model Inference	Arxiv'24	Not Available
Reinforcement Learning from Reflective Feedback (RLRF): Aligning and Improving LLMs via Fine-Grained Self-Reflection <sup>[36]</sup>	Pairwise Feedback	LLaMA-2	Model Inference, 16 Nvidia V100 GPUs	Arxiv'24	Not Available
Direct language model alignment from online ai feedback <sup>[37]</sup>	Pairwise Feedback	PaLM-2	API Calling	Arxiv'24	Not Available
Reinforced Self-Training (ReST) for Language Modeling <sup>[38]</sup>	Response	Multiple LLMs	Model Inference	Arxiv'24	Not Available
RAFT: Reward-Anked FineTuning for Generative Foundation Model Alignment <sup>[39]</sup>	Response	LLaMA	Model Inference	TMLR	<a href="#">Link</a>
Step-On-Feet Tuning: Scaling Self-Alignment of LLMs via Bootstrapping <sup>[40]</sup>	Response	LLaMA-2	Model Inference	Arxiv'24	Not Available
Mixture of insightful Experts (MoTE): The Synergy of Thought Chains and Expert Mixtures in Self-Alignment <sup>[41]</sup>	Response	Alpaca	Model Inference	Arxiv'24	Not Available
Iterative reasoning preference optimization <sup>[42]</sup>	Pairwise Feedback	LLaMA-2	Model Inference	Arxiv'24	Not Available
Inference Time					
Large Language Models are Human-Level Prompt Engineers <sup>[43]</sup>	Instruction	GPT-3.5	API Calling	ICLR'23	<a href="#">Link</a>
Auto-ICL: In-Context Learning without Human Supervision <sup>[44]</sup>	Instruction	ChatGPT	API Calling	Arxiv'23	<a href="#">Link</a>
Empowering Large Language Models for Textual Data Augmentation <sup>[45]</sup>	Instruction	ChatGPT	API Calling	Arxiv'24	Not Available
Self-generated in-context learning: Leveraging auto-regressive language models as a demonstration generator <sup>[46]</sup>	Instruction	GPT-J	Model Inference	NAACL'22	<a href="#">Link</a>
Are Human-generated Demonstrations Necessary for In-context Learning? <sup>[47]</sup>	Instruction	Multiple LLMs	API Calling	Arxiv'23	<a href="#">Link</a>
Self-ICL: Zero-Shot In-Context Learning with Self-Generated Demonstrations <sup>[48]</sup>	Instruction	Multiple LLMs	API Calling	EMNLP'23	<a href="#">Link</a>
Self-Demos: Eliciting Out-of-Demonstration Generalizability in Large Language Models <sup>[49]</sup>	Instruction	ChatGPT	API Calling	NAACL'24	<a href="#">Link</a>
Rephrase and respond: Let large language models ask better questions for themselves <sup>[50]</sup>	Instruction	GPT-4	API Calling	Arxiv'23	<a href="#">Link</a>
DALL: Data Augmentation for In-Context Learning via Self-Paraphrase <sup>[51]</sup>	Instruction	ChatGPT	API Calling	Arxiv'23	Not Available
Just rephrase it! Uncertainty estimation in closed-source language models via multiple rephrased queries <sup>[52]</sup>	Instruction	Multiple LLMs	Model Inference	Arxiv'24	Not Available
Self-Polish: Enhance Reasoning in Large Language Models via Problem Refinement <sup>[53]</sup>	Instruction	GPT-3.5	API Calling	EMNLP'23	<a href="#">Link</a>
Self-DC: When to retrieve and when to generate? Self-Divide-and-Conquer for Compositional Unknown Questions <sup>[54]</sup>	Instruction	ChatGPT	API Calling	Arxiv'24	Not Available
Large Language Models are Zero-Shot Reasoners <sup>[55]</sup>	Rationale - CoT	Multiple LLMs	API Calling	NeurIPS'22	Not Available
SELF-CONSISTENCY IMPROVES CHAIN OF THOUGHT REASONING IN LANGUAGE MODELS <sup>[56]</sup>	Rationale - Diverse Thinking	Multiple LLMs	API Calling & Model Inference	ICLR'23	Not Available
UNIVERSAL SELF-CONSISTENCY FOR LARGE LANGUAGE MODEL GENERATION <sup>[57]</sup>	Rationale - Diverse Thinking	Multiple LLMs	API Calling	Arxiv'23	Not Available
Eliminating Reasoning via Inferring with Planning: A New Framework to Guide LLMs' Non-linear Thinking <sup>[58]</sup>	Rationale - Elimination	PaLM2	API Calling	Arxiv'23	Not Available
It's Not Easy Being Wrong: Large Language Models Struggle with Process of Elimination Reasoning <sup>[59]</sup>	Rationale - Elimination	Multiple LLMs	API Calling	ACL'24	<a href="#">Link</a>
POE: Process of Elimination for Multiple Choice Reasoning <sup>[60]</sup>	Rationale - Elimination	FLAN-T5	Model Inference	EMNLP'23	<a href="#">Link</a>
SELF-REFINE: Iterative Refinement with Self-Feedback <sup>[61]</sup>	Textual Feedback	Multiple LLMs	API Calling	NeurIPS'23	Not Available
Can LLMs Learn from Previous Mistakes? Investigating LLMs' Errors to Boost for Reasoning <sup>[62]</sup>	Textual Feedback - Mistake	Multiple LLMs	API Calling & Modeling Inference	ACL'24	<a href="#">Link</a>
Program of Thoughts Prompting: Disentangling Computation from Reasoning for Numerical Reasoning Tasks <sup>[63]</sup>	Rationale - Program	Multiple LLMs	API Calling & Model Inference	TMLR'23	Not Available
Graph of Thoughts: Solving Elaborate Problems with Large Language Models <sup>[64]</sup>	Rationale - Graph	GPT-3.5	API Calling	AAAI'24	<a href="#">Link</a>
Reasoning with Language Model is Planning with World Model <sup>[65]</sup>	Rationale - Tree	LLaMA	Model Inference, 4x24 GB NVIDIA A5000 GPUs	EMNLP'23	<a href="#">Link</a>

Note: <sup>[1]</sup>(Huang et al., 2023); <sup>[2]</sup>(Wang et al., 2023e); <sup>[3]</sup>(Lu et al., 2023); <sup>[4]</sup>(Yang et al., 2024b); <sup>[5]</sup>(Chen et al., 2024b); <sup>[6]</sup>(Cheng et al., 2024); <sup>[7]</sup>(Taori et al., 2023); <sup>[8]</sup>(Chiang et al., 2023a); <sup>[9]</sup>(Xu et al., 2023a); <sup>[10]</sup>(Meng et al., 2022); <sup>[11]</sup>(Meng et al., 2023); <sup>[12]</sup>(Wang et al., 2023d); <sup>[13]</sup>(Wang et al., 2022a); <sup>[14]</sup>(Shridhar et al., 2023); <sup>[15]</sup>(Liu et al., 2023a); <sup>[16]</sup>(Kang et al., 2024); <sup>[17]</sup>(Xu et al., 2023b); <sup>[18]</sup>(Josifoski et al., 2023); <sup>[19]</sup>(Jeronymo et al., 2023); <sup>[20]</sup>(Chaudhary, 2023); <sup>[21]</sup>(Roziere et al., 2023); <sup>[22]</sup>(Zhang et al., 2023); <sup>[23]</sup>(Xiong et al., 2023); <sup>[24]</sup>(Zhang and Yang, 2023b); <sup>[25]</sup>(Luo et al., 2023); <sup>[26]</sup>(Zhao et al., 2024); <sup>[27]</sup>(Xu et al., 2023c); <sup>[28]</sup>(Kim et al., 2023b); <sup>[29]</sup>(Pace et al., 2024); <sup>[30]</sup>(Zeng et al., 2024); <sup>[31]</sup>(Sun et al., 2023b); <sup>[32]</sup>(Yuan et al., 2024); <sup>[33]</sup>(Zhang et al., 2024a); <sup>[34]</sup>(Lee et al., 2024b); <sup>[35]</sup>(Tong et al., 2024b); <sup>[36]</sup>(Lee et al., 2024a); <sup>[37]</sup>(Guo et al., 2024b); <sup>[38]</sup>(Gulcehre et al., 2023); <sup>[39]</sup>(Dong et al., 2023); <sup>[40]</sup>(Wang et al., 2024a); <sup>[41]</sup>(Liu et al., 2024); <sup>[42]</sup>(Chen et al., 2023c); <sup>[43]</sup>(Zhou et al., 2022b); <sup>[44]</sup>(Yang et al., 2023b); <sup>[45]</sup>(Li et al.); <sup>[46]</sup>(Kim et al., 2022); <sup>[47]</sup>(Li et al., 2023c); <sup>[48]</sup>(Chen et al., 2023d); <sup>[49]</sup>(He et al., 2024); <sup>[50]</sup>(Deng et al., 2023); <sup>[51]</sup>(Li et al., 2023a); <sup>[52]</sup>(Yang et al., 2024a); <sup>[53]</sup>(Xi et al., 2023); <sup>[54]</sup>(Wang et al., 2024b); <sup>[55]</sup>(Kojima et al., 2022); <sup>[56]</sup>(Wang et al., 2022b); <sup>[57]</sup>(Chen et al., 2023f); <sup>[58]</sup>(Tong et al., 2023); <sup>[59]</sup>(Balepur et al., 2023); <sup>[60]</sup>(Ma and Du, 2023); <sup>[61]</sup>(Madaan et al., 2024); <sup>[62]</sup>(Tong et al., 2024a); <sup>[63]</sup>(Chen et al., 2023e); <sup>[64]</sup>(Besta et al., 2024); <sup>[65]</sup>(Hao et al., 2023).

Table 4: A list of representative LLM-Generated Annotation Utilization papers with open-source code/data.