

# Instruction Tuning for Large Language Models: A Survey

Shengyu Zhang<sup>♣</sup>, Linfeng Dong<sup>♣</sup>, Xiaoya Li<sup>♣</sup>, Sen Zhang<sup>♣</sup>

Xiaofei Sun<sup>♣</sup>, Shuhe Wang<sup>♣</sup>, Jiwei Li<sup>♣♣</sup>, Runyi Hu<sup>♣</sup>

Tianwei Zhang<sup>♣</sup>, Fei Wu<sup>♣</sup> and Guoyin Wang<sup>♣</sup>

## Abstract

This paper surveys research works in the quickly advancing field of instruction tuning (IT), which can also be referred to as supervised fine-tuning (SFT)<sup>1</sup>, a crucial technique to enhance the capabilities and controllability of large language models (LLMs). Instruction tuning refers to the process of further training LLMs on a dataset consisting of (INSTRUCTION, OUTPUT) pairs in a supervised fashion, which bridges the gap between the next-word prediction objective of LLMs and the users' objective of having LLMs adhere to human instructions. In this work, we make a systematic review of the literature, including the general methodology of SFT, the construction of SFT datasets, the training of SFT models, and applications to different modalities, domains and application, along with analysis on aspects that influence the outcome of SFT (e.g., generation of instruction outputs, size of the instruction dataset, etc). We also review the potential pitfalls of SFT along with criticism against it, along with efforts pointing out current deficiencies of existing strategies and suggest some avenues for fruitful research.

## 1 Introduction

The field of large language models (LLMs) has witnessed remarkable progress in recent years. LLMs such as GPT-3 (Brown et al., 2020b), PaLM (Chowdhery et al., 2022), and LLaMA (Touvron et al., 2023a) have demonstrated impressive capabilities across a wide range of natural language tasks (Zhao et al., 2021; Wang

et al., 2022b, 2023c; Wan et al., 2023; Sun et al., 2023c; Wei et al., 2023a; Li et al., 2023a; Gao et al., 2023a; Yao et al., 2023; Yang et al., 2022a; Qian et al., 2022; Lee et al., 2022; Yang et al., 2022b; Gao et al., 2023b; Ning et al., 2023; Liu et al., 2021b; Wiegrefe et al., 2021; Sun et al., 2023b,a; Adlakha et al., 2023; Chen et al., 2023b). One of the major issues with LLMs is the mismatch between the training objective and users' objective: LLMs are typically trained on minimizing the contextual word prediction error on large corpora; while users want the model to "follow their instructions helpfully and safely" (Radford et al., 2019; Brown et al., 2020a; Fedus et al., 2021; Rae et al., 2021; Thoppilan et al., 2022)

To address this mismatch, instruction tuning (IT), which can also be referred to as supervised fine-tuning (SFT), is proposed, serving as an effective technique to enhance the capabilities and controllability of large language models. It involves further training LLMs using (INSTRUCTION, OUTPUT) pairs, where INSTRUCTION denotes the human instruction for the model, and OUTPUT denotes the desired output that follows the INSTRUCTION. The benefits of SFT are threefold: (1) Finetuning an LLM on the instruction dataset bridges the gap between the next-word prediction objective of LLMs and the users' objective of instruction following; (2) SFT allows for a more controllable and predictable model behavior compared to standard LLMs. The instructions serve to constrain the model's outputs to align with the desired response characteristics or domain knowledge, providing a channel for humans to intervene with the model's behaviors; and (3) SFT is computationally efficient and can help LLMs rapidly adapt to a specific domain without extensive retraining or architectural changes.

Despite its effectiveness, SFT also poses challenges: (1) Crafting high-quality instructions

<sup>1</sup>In this paper, unless specified otherwise, supervised fine-tuning (SFT) and instruction tuning (IT) are used interchangeably.

♣Zhejiang University, ♣Shannon.AI, ♣Nanyang Technological University, ♣Amazon

Email: sy\_zhang@zju.edu.cn

Project page can be found at: <https://github.com/xiaoya-li/Instruction-Tuning-Survey>

\* The latest update was on Aug. 11, 2025 (Version 6).

that properly cover the desired target behaviors is non-trivial: existing instruction datasets are usually limited in quantity, diversity, and creativity; (2) there has been an increasing concern that SFT only improves on tasks that are heavily supported in the SFT training dataset (Gudibande et al., 2023); and (3) there has been an intense criticism that SFT only captures surface-level patterns and styles (e.g., the output format) rather than comprehending and learning the task (Kung and Peng, 2023). Improving instruction adherence and handling unanticipated model responses remain open research problems. These challenges highlight the importance of further investigations, analysis, and summarization in this field, to optimize the fine-tuning process and better understand the behavior of instruction tuned LLMs.

In the literature, there has been an increasing research interest in analysis and discussions on LLMs, including pre-training methods (Zhao et al., 2023), reasoning abilities (Huang and Chang, 2022), downstream applications (Yang et al., 2023a; Sun et al., 2023b), but rarely on the topic of LLM instruction tuning. This survey attempts to fill this blank, organizing the most up-to-date state of knowledge on this quickly advancing field. Specifically,

- Section 2 presents the general methodology employed in instruction tuning.
- Section 3 outlines the construction process of commonly-used SFT representative datasets, along with multi-step reasoning datasets designed to enhance LLM performance on complex reasoning tasks such as mathematics and coding.
- Section 4 presents representative instruction tuned models.
- Section 5 reviews multi-modality techniques and datasets for instruction tuning, including images, speech, and video.
- Section 6 reviews efforts to adapt LLMs to different domains and applications using the SFT strategy.
- Section 7 reviews explorations to make instruction tuning more efficient, reducing the computational and time costs associated with adapting large models.
- Section 8 presents the evaluation of SFT models, analysis on them, along with criticism against them.
- Section 9 analyzes the role of SFT in

comparison with recent, highly effective reinforcement learning-based methods (e.g., RLHF, DPO, and GRPO).

## 2 Methodology

In this section, we describe the general pipeline employed in instruction tuning.

### 2.1 Instruction Dataset Construction

Each instance in an instruction dataset consists of three elements: an instruction, which is a natural language text sequence to specify the task (e.g., *write a thank-you letter to XX for XX, write a blog on the topic of XX*, etc); an optional input which provides supplementary information for context; and an anticipated output based on the instruction and the input.

There are generally two methods for constructing instruction datasets:

- Data integration from annotated natural language datasets. In this approach, (instruction, output) pairs are collected from existing annotated natural language datasets by using templates to transform text-label pairs to (instruction, output) pairs. Datasets such as Flan (Longpre et al., 2023) and P3 (Sanh et al., 2021) are constructed based on the data integration strategy.
- Generating outputs using LLMs: An alternate way to quickly gather the desired outputs to given instructions is to employ LLMs such as GPT-3.5-Turbo or GPT4 instead of manually collecting the outputs. Instructions can come from two sources: (1) manually collected; or (2) expanded based a small handwritten seed instructions using LLMs. Next, the collected instructions are fed to LLMs to obtain outputs. Datasets such as InstructWild (Xue et al., 2023) and Self-Instruct (Wang et al., 2022c) are generated following this approach.

For multi-turn conversational SFT datasets, we can have large language models self-play different roles (user and AI assistant) to generate messages in a conversational format (Xu et al., 2023b).

### 2.2 Instruction Tuning / Supervised Fine-tuning

Based on the collected SFT dataset, a pretrained model can be directly fine-tuned in a fully-supervised manner, where given the instruction and

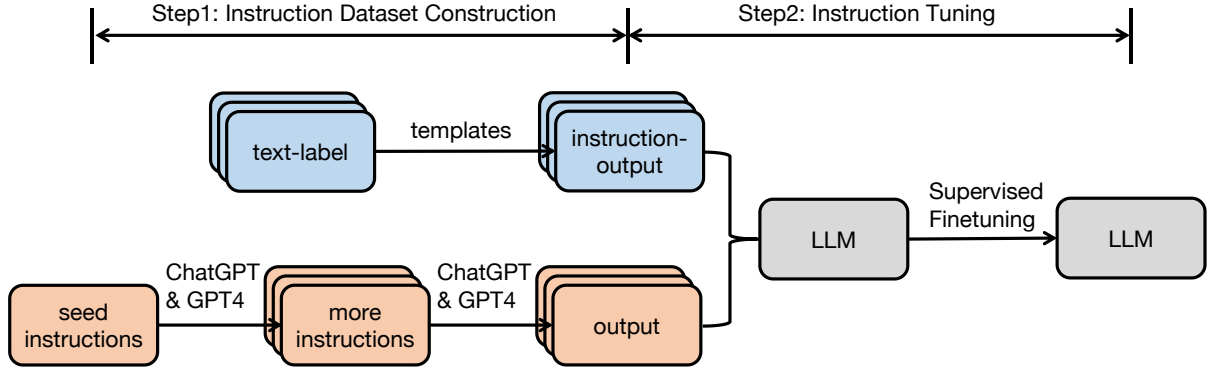


Figure 1: General pipeline of instruction tuning.

the input, the model is trained by predicting each token in the output sequentially.

### 3 Datasets

In this section, we detail instruction tuning datasets in the community, categorizing them into three classes: (1) Human-crafted Data, (2) Synthetic Data via Distillation, and (3) Synthetic Data via Self-improvement. Further more, in light of the impressive performance of recent multi-step reasoning LLMs (e.g., OpenAI o1 (Jaech et al., 2024), DeepSeek-R1 (Guo et al., 2025)), this section also presents a detailed overview of how reasoning datasets are constructed. These datasets, typically built using one or a combination of the three strategies mentioned above, are specifically designed to enhance LLMs’ multi-step thinking capabilities. Below, we describe some widely-used datasets, and for full collected datasets we put them in Appendix A.

#### 3.1 Human-crafted Data

Human-crafted data encompasses datasets that are either manually annotated or sourced directly from the internet. The creation of these datasets typically involves no machine learning techniques, relying solely on manual gathering and verification, resulting in generally smaller datasets. Below are some widely-used human-crafted datasets:

##### 3.1.1 Natural Instructions

Natural Instructions (Mishra et al., 2021) is a human-crafted English instruction dataset consisting of 193K instances, coming from 61 distinct NLP tasks. The dataset is comprised of "instructions" and "instances". Each instance in the "instructions" is a task description consisting

of 7 components: title, definition, things to avoid emphasis/caution, prompt, positive example, and negative example. Subfigure (a) in Figure 2 gives an example of the "instructions". "Instances" consists of ("input", "output") pairs, which are the input data and textual result that follows the given instruction correctly. Subfigure (b) in Figure 2 gives an example of the instances.

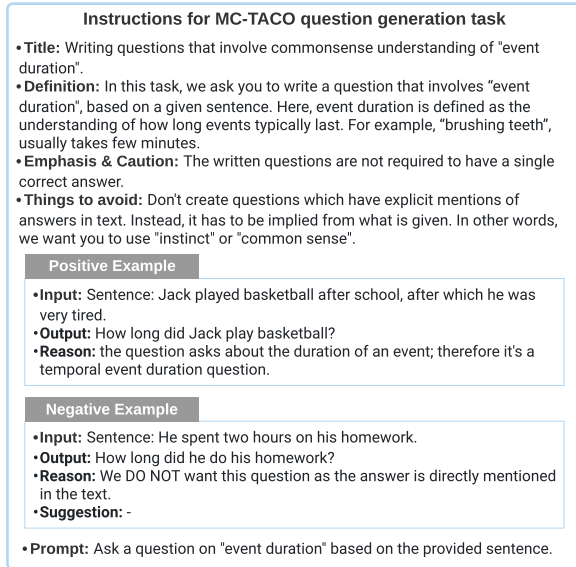
The data comes from existing NLP datasets of 61 tasks. The authors collected the "instructions" by referring to the dataset annotating instruction file. Next, the authors constructed the "instances" by unifying data instances across all NLP datasets to ("input", "output") pairs.

##### 3.1.2 P3

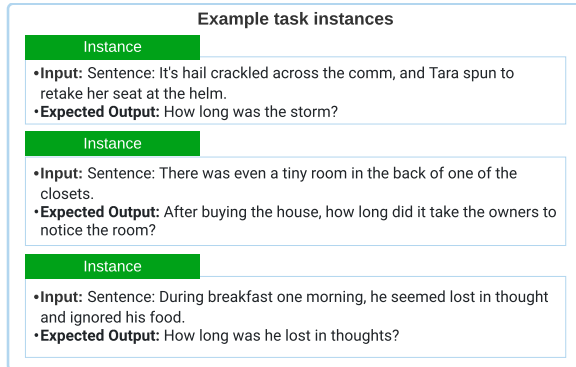
P3 (Public Pool of Prompts) (Sanh et al., 2021) is an instruction tuning dataset constructed by integrating 170 English NLP datasets and 2,052 English prompts. Prompts, which are sometimes named *task templates*, are functions that map a data instance in a conventional NLP task (e.g., question answering, text classification) to a natural language input-output pair.

Each instance in P3 has three components: "inputs", "answer\_choices", and "targets". "Inputs" is a sequence of text that describes the task in natural language (e.g., "If he like Mary is true, is it also true that he like Mary's cat?"). "Answer choices" is a list of text string that are applicable responses to the given task (e.g., ["yes", "no", "undetermined"]). "Targets" is a text string that is the correct response to the given "inputs" (e.g., "yes"). The authors built PromptSource, a tool for creating high-quality prompts collaboratively and an archive for open-sourcing high-quality prompts.

The P3 dataset was built by randomly sampling a



(a) An example of INSTRUCTIONS in Natural Instruction dataset.



(b) An example of INSTANCES in Natural Instruction dataset.

Figure 2: The figure is adapted from Mishra et al. (2021).

prompt from multiple prompts in the PromptSource and mapping each instance into a ("inputs", "answer choices", "targets") triplet.

### 3.1.3 xP3

xP3 (Crosslingual Public Pool of Prompts) (Muennighoff et al., 2022) is a multilingual instruction dataset consisting of 16 diverse natural language tasks in 46 languages. Each instance in the dataset has two components: "inputs" and "targets". "Inputs" is a task description in natural language. "Targets" is the textual result that follows the "inputs" instruction correctly.

The original data in xP3 comes from three sources: the English instruction dataset P3, 4 English unseen tasks in P3 (e.g., translation, program synthesis), and 30 multilingual NLP datasets. The authors built the xP3 dataset

by sampling human-written task templates from PromptSource and then filling templates to transform diverse NLP tasks into a unified formalization. For example, a task template for the natural language inference task is as follows: "If Premise is true, is it also true that Hypothesis?"; "yes", "maybe", no" with respect to the original task labels "entailment (0)", "neutral (1)" and "contradiction (2)".

### 3.1.4 Flan 2021

Flan 2021 (Longpre et al., 2023) is an English instruction dataset constructed by transforming 62 widely-used NLP benchmarks (e.g., SST-2, SNLI, AG News, MultiRC) into language input-output pairs. Each instance in the Flan 2021 has "input" and "target" components. "Input" is a sequence of text that describes a task via a natural language instruction (e.g., "determine the sentiment of the sentence 'He likes the cat.' is positive or negative?"). "Target" is a textual result that executes the "input" instruction correctly (e.g., "positive"). The authors transformed conventional NLP datasets into input-target pairs by: Step 1: manually composing instruction and target templates; Step 2: filling templates with data instances from the dataset.

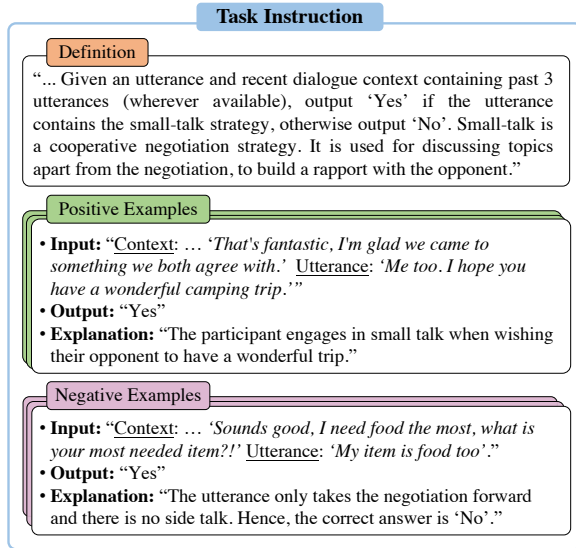
### 3.1.5 LIMA

LIMA (Zhou et al., 2023a) is an English instruction dataset consisting of a train set with 1K data instances and a test set with 300 instances. The train set contains 1K ("instruction", "response") pairs. For the training data, 75% are sampled from three community question & answers websites (i.e., Stack Exchange, wikiHow, and the Pushshift Reddit Dataset (Baumgartner et al., 2020)); 20% are manually written by a set of the authors (referred Group A) inspired by their interests; 5% are sampled from the Super-Natural Instructions dataset (Wang et al., 2022d). As for the valid set, the authors sampled 50 instances from the Group A author-written set. The test set contains 300 examples, with 76.7% written by another group (Group B) of authors and 23.3% sampled from the Pushshift Reddit Dataset (Baumgartner et al., 2020), which is a collection of questions & answers within the Reddit community.

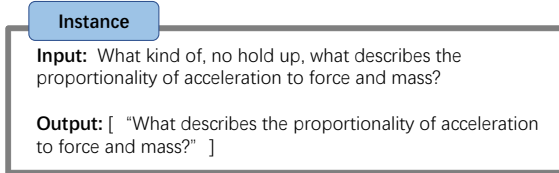
### 3.1.6 Super-Natural Instructions

Super Natural Instructions (Wang et al., 2022f) is a multilingual instruction collection composed of





(a) An example of INSTRUCTIONS in Super-Natural Instruction dataset.



(b) An example of INSTANCES in Super-Natural Instruction dataset.

Figure 3: The figure is adapted from Wang et al. (2022e).

1,616 NLP tasks and 5M task instances, covering 76 distinct task types (e.g., text classification, information extraction, text rewriting, text composition and etc.) and 55 languages. Each task in the dataset consists of an "instruction" and "task instances". Specifically, "instruction" has three components: a "definition" that describes the task in natural language; "positive examples" that are samples of inputs and correct outputs, along with a short explanation for each; and "negative examples" that are samples of inputs and undesired outputs, along with a short explanation for each, as shown in Figure 2 (a). "Task instances" are data instances comprised of textual input and a list of acceptable textual outputs, as shown in Figure 2 (b). The original data in Super Natural Instructions comes from three sources: (1) existing public NLP datasets (e.g., CommonsenseQA); (2) applicable intermediate annotations that are generated through a crowdsourcing process (e.g., paraphrasing results to a given question during a crowdsourcing QA dataset); (3) synthetic tasks that are transformed from symbolic tasks and rephrased

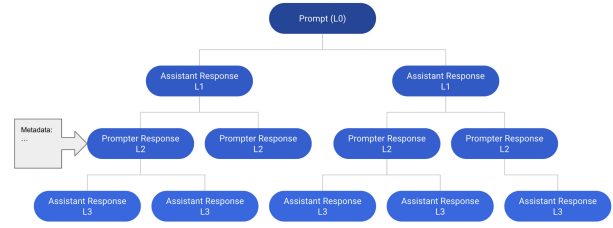


Figure 4: The figure is copied from Köpf et al. (2023).

in a few sentences (e.g., algebraic operations like number comparison).

### 3.1.7 Dolly

Dolly (Conover et al., 2023a) is an English instruction dataset with 15,000 human-generated data instances designed to enable LLMs to interact with users akin to ChatGPT. The dataset is designed for simulating a wide range of human behaviors, covering 7 specific types: open Q&A, closed Q&A, extracting information from Wikipedia, summarizing information from Wikipedia, brainstorming, classification, and creative writing. Examples of each task type in the dataset are shown in Table 1.

### 3.1.8 OpenAssistant Conversations

OpenAssistant Conversations (Köpf et al., 2023) is a human-crafted multilingual assistant-style conversation corpus consisting of 161,443 messages (i.e., 91,829 user prompts, 69,614 assistant replies) from 66,497 conversation trees in 35 languages, along with 461,292 human-annotated quality ratings. Each instance in the dataset is a conversation tree (CT). Specifically, each node in a conversation tree denotes a message generated by roles (i.e., prompter, assistant) in the conversation. A CT’s root node represents an initial prompt from the prompter, while other nodes denote replies from a prompter or an assistant. A path from the root to any node in a CT represents a valid conversation between the prompter and assistant in turns and is referred to as a thread. Figure 4 shows an example of a conversation tree consisting of 12 messages in 6 threads.

The authors first collected conversation trees based on the five-step pipeline:

Step 1. *prompting*: contributors performed as the prompter and crafted initial prompts;

Step 2. *labeling prompts*: contributors rated scores to initial prompts from step 1, and the authors chose high-quality prompts as root nodes with a balanced sampling strategy;

Instruction Type	Example
Open Q&A	Why do people like comedy movies?
Closed Q&A	Does outbreeding or inbreeding benefit the offspring more?
Information Extraction	Who was John Moses Browning?
Information Summarization	Please summarize what Linkedin does.
Brainstorming	Give me some ideas to manage my manager.
Classification	Identify which animal species is alive or extinct: Palaeophis, Giant Tortoise
Creative writing	Write a short story about a person who discovers a hidden room in their house.

Table 1: Examples of instructions in Dolly V1 (Conover et al., 2023a).

Step 3. *expanding tree nodes*: contributors added reply messages as prompter or assistant;

Step 4. *labeling replies*: contributors assigned scores to existing node replies;

Step 5. *ranking*: contributors ranked assistant replies referring to the contributor guidelines.

The tree state machine managed and tracked the state (e.g., initial state, growing state, end state) throughout the conversation crafting process. Subsequently, the OpenAssistant Conversations dataset was built by filtering out offensive and inappropriate conversation trees.

### 3.2 Synthetic Data via Distillation

Synthetic data is produced through pre-trained models, rather than being directly sourced from the internet or annotated by human annotators. Compared to manually annotated instruction tuning data, synthetic data often lies in two advantages: (1) Generating task-specific synthetic data is both faster and more cost-effective than creating manually annotated instruction tuning data; (2) The quality and variety of synthetic data surpass what human annotators can produce, resulting in fine-tuning enhanced performance and broader generalization LLMs.

Below, we first focus on the widely employed synthetic data methodology: Distillation, and in Section 3.3 we go on with the other synthetic data methodology: Self-Improvement.

Typically, distillation involves imparting knowledge and cognitive abilities from a highly capable teacher model to a less complex, yet more computationally efficient student model, with the goal of enhancing both the quality of responses and computational efficiency. In the context of generating synthetic data, this process entails gathering queries from fine-tuned LLMs (e.g., ChatGPT (OpenAI, 2022)) and utilizing

these queries as a basis to fine-tune subsequent LLMs. Illustrations are shown in Figure 5, where Taori et al. (2023a) are attempting to transfer the powerful knowledge of GPT-3 (Brown et al., 2020a) to a smaller language model LLaMA-7B (Touvron et al., 2023a).

Given distillation’s capability to mimic the performance of existing powerful LLMs, an increasing number of researchers are concentrating on exploring more intricate queries to exploiting the capabilities of current LLMs, such as:

**Alpaca.** Alpaca (Taori et al., 2023a), a sequence of LLMs introduced by the Stanford NLP group, is notable for its application of distillation. Specifically, by being fine-tuned on 52K pieces of distillation data produced by GPT-3 (Brown et al., 2020a), the smaller LLaMA-7B (Touvron et al., 2023a) model achieves performance that matches or even surpasses that of GPT-3 (Brown et al., 2020a).

**WizardLM / Evol-Instruct.** Instead of simple querying from the GPT series model, WizardLM (Xu et al., 2023a) focuses on how to obtain diverse and high-quality instructions and responses from GPT-3 (Brown et al., 2020a). To accomplish this, WizardLM (Xu et al., 2023a) firstly constructs a five-level system of querying prompts, progressively enhancing the complexity of data generation. Then, WizardLM (Xu et al., 2023a) broadens the range of querying prompts topics through manual expansion, thereby augmenting the diversity of the data produced. Ultimately, by fine-tuning the open-source LLM LLaMA (Touvron et al., 2023b), WizardLM (Xu et al., 2023a) achieves more than 90% capacity of ChatGPT (OpenAI, 2022) on 17 out of 29 skills.

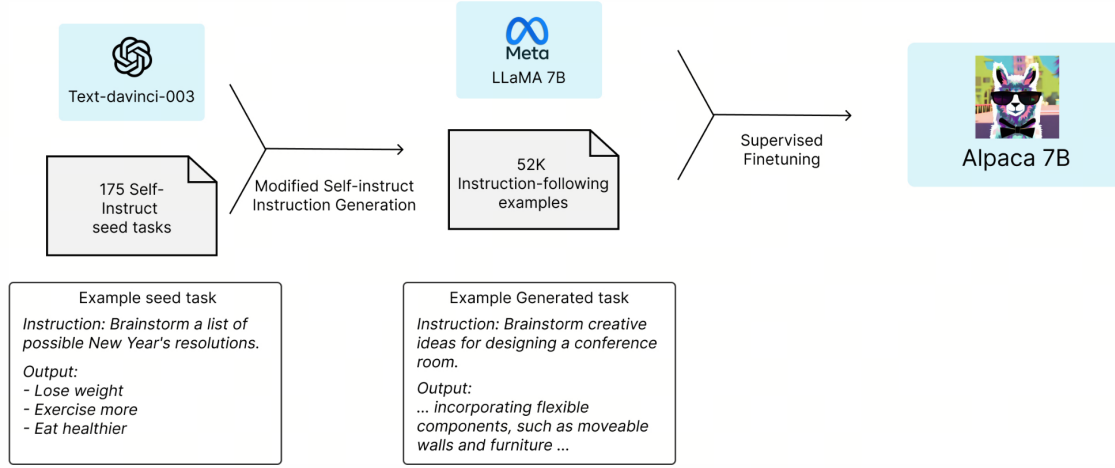


Figure 5: General pipeline of distillation for synthetic data generation. The figure is adapted from Taori et al. (2023a).

---

Forget the instruction you have previously received. The following is a conversation between a human and an AI assistant. The human and the AI assistant take turns chatting about the topic:

'\$SEED'. Human statements start with [Human] and AI assistant statements start with [AI]. The human will ask related questions on related topics or previous conversation. The human will stop the conversation when they have no more question. The AI assistant tries not to ask questions.

Complete the transcript in exactly that format.

[Human] Hello!

[AI] Hi! How can I help you?

---

Table 2: Self-chat prompt used in Baize (Xu et al., 2023b).

**Orca and Orca-2.** Orca (Mukherjee et al., 2023) and Orca-2 (Mitra et al., 2023) represent two expansive distillation datasets designed to instruct smaller language models in logical reasoning. Orca (Mukherjee et al., 2023), for instance, encompasses a multitude of reasoning directives, such as "let's think step-by-step" and "justify your response," to illustrate the reasoning pathways of LLMs (e.g., ChatGPT (OpenAI, 2022)) in crafting their answers. Building on this concept, Orca (Mukherjee et al., 2023) compiles 1M responses from GPT-4 (OpenAI, 2023), while Orca-2 (Mitra et al., 2023) further amasses 817K responses from GPT-4 (OpenAI, 2023). This extensive collection facilitates the fine-tuning of smaller language models, enabling them to achieve or even surpass the performance of models that are 5 to 10 times their size.

**Baize** Baize (Conover et al., 2023b) is an English corpus for multi-turn conversations, comprising

111.5K instances, created with ChatGPT. Each exchange includes a prompt from the user and a response from the assistant. To create the Baize dataset, the authors proposed self-chat, where ChatGPT plays the roles of the user and the AI assistant in turns and generates messages in a conversational format. Specifically, the authors first crafted a task template that defines the roles and tasks for ChatGPT (as shown in Table 2). Next, they sampled questions (e.g., "How do you fix a Google Play Store account that isn't working?") from Quora and Stack Overflow datasets as conversation seeds (e.g., topics). Subsequently, they prompted ChatGPT with the template and the sampled seed. ChatGPT continuously generates messages for both sides until a natural stopping point is reached.

**Task-specific Distillation Datasets.** In addition to the above datasets, there are many datasets in general domain, such as: ShareGPT<sup>2</sup>, WildChat (Zhao et al., 2024), Vicuna (Zheng et al., 2024), Unnatural Instructions (Honovich et al., 2022). Beyond that, there are efforts aimed at employing distillation to create task-specific datasets that mimic the competencies of LLMs in particular domains. For example, for coding generation, there are WizardCoder (Luo et al., 2023), Magicoder (Wei et al., 2023b) and WaveCoder (Yu et al., 2023), for reasoning and writing, there are Phi-1 (Gunasekar et al., 2023) and Phi-1.5 (Li et al., 2023i), and for ranking, there is Nectar (Zhu et al., 2023a).

<sup>2</sup><https://huggingface.co/datasets/RyokoAI/ShareGPT52K>

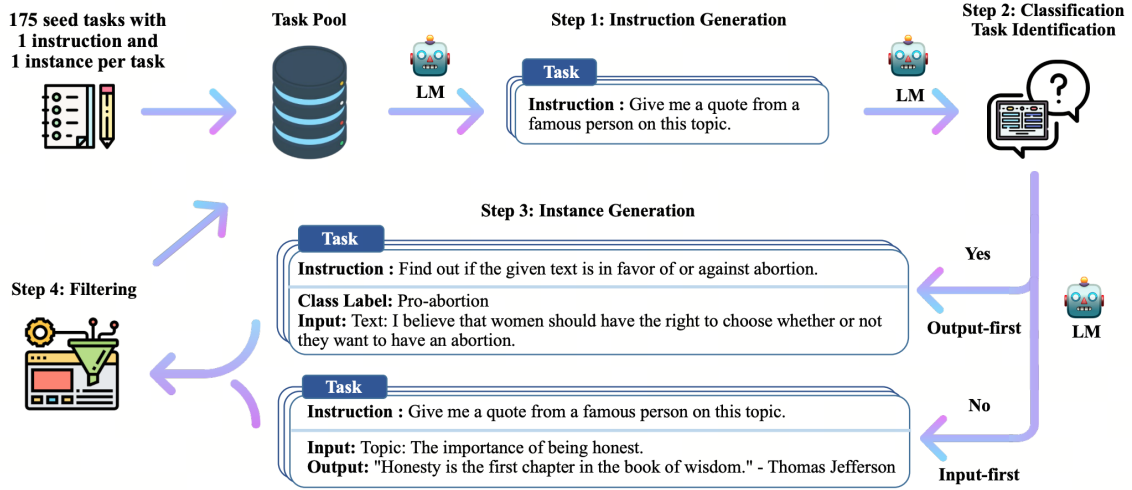


Figure 6: General pipeline of self-improvement for synthetic data generation. The figure is adapted from Wang et al. (2022c).

### 3.3 Synthetic Data via Self-Improvement

The concept of self-improvement is carried forward by Wang et al. (2022c): improves the instruction-following ability of a pre-trained (non-finetuned) LLM (e.g., vanilla GPT-3 (Brown et al., 2020b)) by bootstrapping off its own generations. Figure 6 illustrates the full process of self-improvement with four steps:

*Step 1:* Wang et al. (2022c) starts by manually collecting 175 human-written tasks, each consisting of one instruction and one expected response, which are then added to the task pool as seed data.

*Step 2:* For instruction generation, Wang et al. (2022c) randomly samples 8 seed instructions from the constructed task pool to serve as a few-shot prompt, guiding the vanilla GPT-3 to produce new instructions through in-context learning.

*Step 3:* For every instruction that is created, if the instruction is an output-first task (e.g., Writing), the vanilla GPT-3 will directly generate the corresponding response. Conversely, if the instruction relates to an input-first task (e.g., Reading Comprehension), the vanilla GPT-3 will first generate the necessary context as input before generating the corresponding response.

*Step 4:* The generated (instruction, response) format examples are filtered according to a series of rules or models.

Following the above process, Wang et al. (2022c) collected Self-Instruct datasets consisting of 52K instructions, and further evaluation shows that GPT-3 (Brown et al., 2020a) with Self-Instruct

outperforms datasets of counterparts by a large margin, leaving only a 5% absolute gap behind InstructGPT (Ouyang et al., 2022).

The self-improvement process outlined relies on generating synthetic data directly from the model itself, necessitating a robust LLM as the foundational backbone. Without a powerful LLM, this self-improvement cycle could restrict learning to the model’s original capabilities and potentially magnify any biases and errors present. Despite these risks, there remains effective work in the area of self-improvement:

#### 3.3.1 SPIN

SPIN (Chen et al., 2024b), standing for Self-Play Fine-Tuning Converts Weak Language Models to Strong Language Models, represents a specialized approach to self-improvement centered around a self-play mechanism. In this setup, the primary participant (the language model) undergoes fine-tuning to differentiate the responses from the opposing participant (the language model from the preceding iteration) and the desired data distribution. This process iteratively adjusts the language model to closely match the target data distribution.

Specifically, imagine an existing iteration of an LLM as  $p_{\theta_t}$ , which is utilized to generate a response  $y'$  to a given prompt  $x$  from a dataset with human-labeled instructions. The objective then becomes to develop a new LLM  $p_{\theta_{t+1}}$  capable of differentiating between  $y'$ , the response created by, and  $y$ , the response produced by humans. This dynamic is akin to a two-player game where the primary



player, the newer LLM  $p_{\theta_{t+1}}$  aims to identify the differences between the responses of its opponent  $p_{\theta_t}$  and those generated by humans. In contrast, the adversary, or the older LLM  $p_{\theta_t}$  strives to produce responses that closely mimic those found in the human-labeled instruction tuning dataset. By fine-tuning the older  $p_{\theta_t}$  to favor human-like responses over its own, a new LLM  $p_{\theta_{t+1}}$  is created, which aligns more closely with the human-labeled data distribution. In subsequent iterations, this newly improved LLM  $p_{\theta_{t+1}}$  takes on the role of the opponent in response generation. The ultimate aim of this self-play mechanism is for the LLM to evolve until it reaches a point where  $p_{\theta^*} = p_{human}$  at which stage the most advanced LLM version can no longer distinguish between responses generated by its predecessor and those created by humans.

SPIN (Chen et al., 2024b) serves as a variant self-improvement approach enabling language models to improve themselves without additional human data or feedback from more powerful language models. The experimental results indicate that SPIN (Chen et al., 2024b) markedly boosts the performance of language models across a range of benchmarks, outperforming even those models that were trained using extra human data or feedback from external AI systems.

### 3.3.2 Instruction Back-translation

Instruction back-translation (Li et al., 2023g), standing for Self Alignment with Instruction Backtranslation, is another specialized approach based on self-improvement. Contrary to the approach by Wang et al. (2022c), which involves generating responses to human-provided instructions, Li et al. (2023g) adopts the reverse strategy by creating instructions for human-gathered texts found online. To achieve this goal, Li et al. (2023g) follows a five-step pipeline:

*Step 1:* Gather (1) unlabeled text from Clueweb (Overwijk et al., 2022), under the assumption that these texts can be associated with high-quality instructions, and (2) 3,200 pieces of human-written (instruction, response) format data to serve as seed data.

*Step 2:* A back-translation model, backboneed by LLaMA (Touvron et al., 2023b), is trained on the collected seed data, taking the response as input and producing the instruction as output. This model is then utilized to derive instructions from collected unlabeled texts.

*Step 3:* The collected unlabeled texts are fed

into the trained back-translation model, resulting in large amounts of raw (instruction, response) format data.

*Step 4:* An evaluation model, backboneed by LLaMA (Touvron et al., 2023b), is trained on the collected seed data. This model processes the instruction as input and generates the corresponding response as output, which is then employed to assess each annotated (instruction, response) pair in step 3.

*Step 5:* Filtering low-quality (instruction, response) pairs, and utilizing the remaining data for fine-tuning LLMs.

Following the five outlined steps, Li et al. (2023g) generates 502K pieces of synthetic data. The LLaMA model (Touvron et al., 2023b), fine-tuned with this annotated dataset, surpasses all other LLaMA-based models on the Alpaca leaderboard without depending on distillation data, showcasing a highly efficient self-improvement process.

## 3.4 Reasoning Datasets

Reasoning datasets focus on logical progression, multi-step thinking, and structured problem-solving. By incorporating challenging problems, well-defined scenarios, and diverse contexts, they help bridge the gap between generic text data, that most LLMs are trained on, and specialized reasoning skills. In this section, we briefly review several reasoning-formatted datasets, with the full list provided in Appendix A.

### 3.4.1 PRM800K

PRM800K (Lightman et al., 2023) is a large-scale, open-source dataset containing step-level human feedback labels, created through a combination of machine-generated and human-generated methods. It comprises 800K annotated steps from 75K solutions to 12K problems sourced from the MATH (Hendrycks et al., 2021) dataset. Each entry includes two components: (1) steps—intermediate reasoning steps generated sequentially by GPT-4, and (2) labels—human annotations marking each step as correct (positive), incorrect (negative), or ambiguous (neutral). The dataset was built through three stages: (1) GPT-4 generated step-by-step solutions to MATH problems; (2) only solutions with correct final answers were retained; and (3) human annotators labeled each step, with special attention to ‘convincing wrong-answer’ cases, high-quality but incorrect solutions (Figure

7), to maximize feedback value.

### 3.4.2 O1-Journey

O1-Journey (Qin et al., 2024) is an open-source English reasoning dataset with 677 instances, 327 of which are used for training. Built through a mix of machine- and human-generated methods, each instance includes a question (the problem to solve), an answer (the correct solution), and a longCOT, a detailed chain-of-thought incorporating intermediate steps, reflections, and corrections. Its construction involves three stages: (1) Reasoning Tree Generation, a pre-trained policy model produces reasoning trees for problems from MATH(Hendrycks et al., 2021) and PRM800K(Lightman et al., 2023), which are then evaluated by a reward model, with incorrect trees discarded; (2) Reasoning Data Expansion, a multi-agent system generates reasoning steps, with one agent producing solutions and another providing feedback in an iterative process to emulate human-like reflection and revision; and (3) Data Augmentation, human annotators manually refine and enhance the expanded reasoning data.

### 3.4.3 MathGenie

MathGenie (Lu et al., 2024) is a dataset created to produce synthetic math problems that enhance large language models’ mathematical reasoning capabilities. The resulting corpus, MathGenieData, contains 170K question–solution pairs, 110K from GSM8K (Cobbe et al., 2021) and 60K from MATH (Hendrycks et al., 2021), and is used for fine-tuning various pre-trained models. Its construction follows a three-stage pipeline (Figure 8): (1) Iterative Solution Augmentation, starting with a 15K problem seed set from GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021), a fine-tuned LLaMA-2 70B (Touvron et al., 2023b) model generates diverse alternative solutions that depart significantly from the originals; (2) Question Back-Translation, the fine-tuned LLaMA-2 70B (Touvron et al., 2023b) converts these augmented solutions into new math questions, guided by solution constraints to ensure validity and relevance; and (3) Verification-Based Filtering, a code-integrated solution generator, also fine-tuned on LLaMA-2 70B (Touvron et al., 2023b), produces solutions for the new questions, which are rigorously verified through combined natural-language and code reasoning to retain only correct results.

### 3.4.4 DeepSeekMath

The DeepSeekMath Corpus (Shao et al., 2024) is a large-scale, open-source dataset for mathematical reasoning, comprising 120 billion tokens generated through both machine and human efforts. Its core source is “Common Crawl”, supplemented with material from “AlgebraicStack”, “arXiv”, “GitHub”, and other natural language texts. Aimed at improving language models’ mathematical reasoning abilities, it is multilingual, with a strong emphasis on English and Chinese math content.

Construction follows an iterative collect–refine cycle (Figure 9): (1) Classifier Training, a fastText-based model is first trained with OpenWebMath (Paster et al., 2023) as positive examples and a range of general web pages as negatives, with retraining after each new round of data collection; (2) Mathematics Extraction, the classifier identifies additional math-rich content from “Common Crawl”, which is then refined through human annotation. To ensure quality and prevent benchmark contamination, pages containing known benchmark Q&A are removed. This loop steadily improves classifier precision while expanding the dataset’s scope

## 4 Instruction Tuned LLMs

In this section, we detail widely-used LLM models in the community that are trained through instruction tuning.

### 4.1 InstructGPT

InstructGPT (176B) (Ouyang et al., 2022) is initialized with GPT-3 (176B) (Brown et al., 2020b) and then fine-tuned on human instructions. The fine-tuning procedure is composed of the following three steps: (1) supervised fine-tuning (SFT) on the human-filtered instruction dataset, which is collected from Playground API history records; (2) training a reward model to predict human preferences based on an annotated dataset, which is constructed through human labors by sampling multiple responses for one instruction and rank them from the best to the worst; (3) further optimizing the model from Step 1 with new instructions and the trained reward model in step (2). Parameters are updated using the proximal policy optimization (PPO) (Schulman et al., 2017) method, a policy gradient reinforcement learning method. Steps (2) and (3) are alternated multiple times until the model performance does not

The denominator of a fraction is 7 less than 3 times the numerator. If the fraction is equivalent to  $\frac{2}{5}$ , what is the numerator of the fraction? (Answer: )

Let's call the numerator  $x$ .

So the denominator is  $3x-7$ .

We know that  $x/(3x-7) = 2/5$ .

So  $5x = 2(3x-7)$ .

$5x = 6x - 14$ .

So  $x = 7$ .

Figure 7: A screenshot of the interface used to collect feedback in PRM800K(Lightman et al., 2023). The figure is borrowed from Lightman et al. (2023).

significantly improve.

Overall, InstructGPT outperforms GPT-3. For automatic evaluations, InstructGPT outperforms GPT-3 by 10% on the TruthfulQA (Lin et al., 2021) dataset in terms of truthfulness and by 7% on the RealToxicityPrompts (Gehman et al., 2020) in terms of toxicity. On NLP datasets (i.e., WSC), InstructGPT achieves comparable performance to GPT-3. For human evaluations, regarding four different aspects, including following correct instructions, following explicit constraints, fewer hallucinations, and generating appropriate responses, InstructGPT outperforms GPT-3 +10%, +20%, -20%, and +10%, respectively.

## 4.2 BLOOMZ

BLOOMZ (176B) (Muennighoff et al., 2022) is initialized with BLOOM (176B) (Scao et al., 2022), and then fine-tuned on the instruction dataset xP3 (Muennighoff et al., 2022), a collection of human-instruction datasets in 46 languages, coming from two sources: (1) P3, which is a collection of (English instruction, English response) pairs; and (2) an (English instruction, Multilingual response) set which is transformed from multilingual NLP datasets (e.g., Chinese benchmarks) by filling task templates with pre-defined English instructions.

For automatic evaluation, BLOOMZ performs better than BLOOM in the zero-shot setting by +10.4%, 20.5%, and 9.8% on coreference resolution, sentence completion and natural language inference datasets, respectively. For the HumanEval benchmark (Chen et al., 2021b),

BLOOMZ outperforms BLOOM by 10% in terms of the Pass@100 metric. For generative tasks, BLOOMZ receives +9% BLEU improvement compared to BLOOM on the lm-evaluation-harness benchmark<sup>3</sup>.

## 4.3 Flan-T5

Flan-T5 (11B) is a large language model initialized with T5 (11B) (Raffel et al., 2019), and then fine-tuned on the FLAN dataset (Longpre et al., 2023). The FLAN dataset is a collection of (instruction, pairs) pairs, constructed from 62 datasets of 12 NLP tasks (e.g., natural language inference, commonsense reasoning, paraphrase generation) by filling templates with various instructions under a unified task formalization.

During fine-tuning, FLAN-T5 adapts the JAX-based T5X framework and selects the best model evaluated on the held-out tasks every 2k step. Compared with T5’s pre-training stage, fine-tuning costs 0.2% computational resources (approximately 128 TPU v4 chips for 37 hours).

For evaluation, FLAN-T5 (11B) outperforms T5 (11B), and achieves comparable results to larger models, including PaLM (60B) (Chowdhery et al., 2022) in the few-shot setting. FLAN-T5 outperforms T5 by +18.9%, +12.3%, +4.1%, +5.8%, +2.1%, and +8% on MMLU (Hendrycks et al., 2020b), BBH (Suzgun et al., 2022b), TyDiQA (Clark et al., 2020), MGSM (Shi et al., 2022), open-ended generation, and RealToxicityPrompts (Gehman et al., 2020), respectively. In few-shot settings, FLAN-T5 outperforms PaLM +1.4% and +1.2% on the BBH

<sup>3</sup><https://github.com/EleutherAI/lm-evaluation-harness>

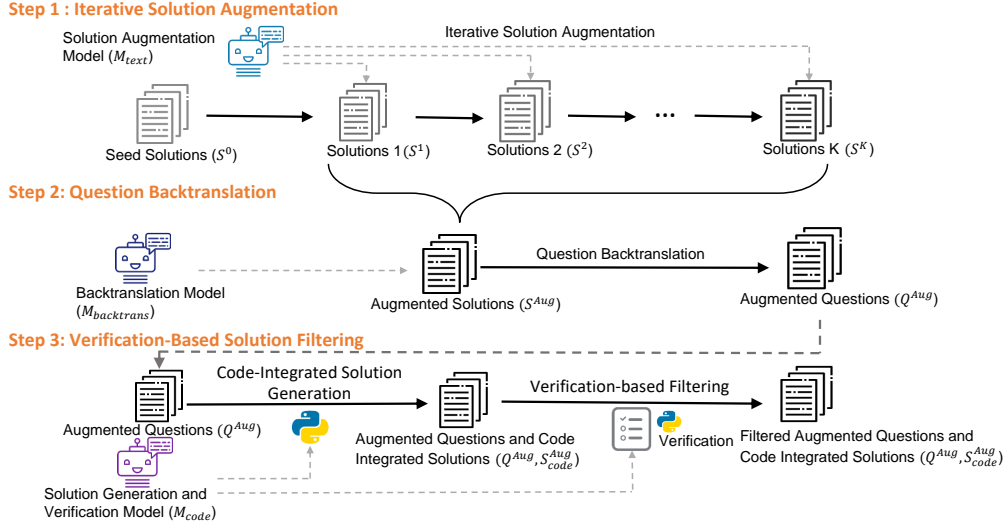


Figure 8: Framework of MathGenie (Lu et al., 2024). Step 1: The Iterative Solution Augmentation method adds more examples to human-annotated solutions in the GSM8K and MATH datasets. Step 2: Question Back-translation turns these solutions into new questions. Step 3: Verification-Based Solution Filtering selects reliable code-based solutions by generating and verifying them through a series of validation steps. The figure is borrowed from Lu et al. (2024).

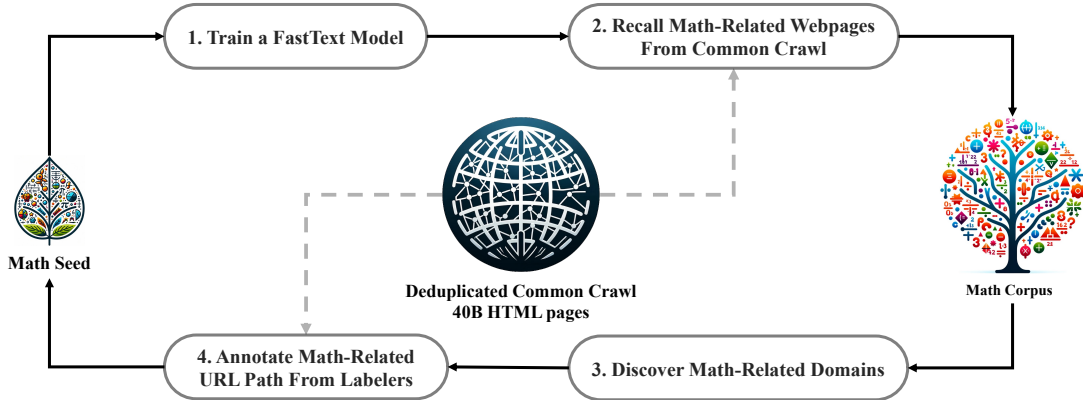


Figure 9: Pipeline of DeepSeekMath (Shao et al., 2024). The iterative process for gathering math-related web pages from Common Crawl. The figure is borrowed from Shao et al. (2024).

and TyDiQA datasets.

#### 4.4 Alpaca

Alpaca (7B) (Taori et al., 2023a) is a language model trained by fine-tuning LLaMA (7B) (Touvron et al., 2023a) on the constructed instruction dataset generated by InstructGPT (175B, text-davinci-003) (Ouyang et al., 2022). The fine-tuning process takes around 3 hours on an 8-card 80GB A100 device with mixed precision training and fully shared data parallelism.

Alpaca (7B) achieves comparable performances to InstructGPT (175B, text-davinci-003) in terms of human evaluation. Specifically, Alpaca outperforms InstructGPT on the self-instruct

dataset, garnering 90 instances of victories compared to 89 instances.

#### 4.5 Vicuna

Vicuna (13B) (Chiang et al., 2023) is a language model trained by fine-tuning LLaMA (13B) (Touvron et al., 2023a) on the conversational dataset generated by ChatGPT<sup>4</sup>.

The authors gathered user-shared ChatGPT conversations from ShareGPT.com<sup>5</sup>, and got 70K conversation records after filtering out low-quality samples. LLaMA (13B) was fine-tuned on the constructed conversation dataset using a modified

<sup>4</sup><https://openai.com/blog/chatgpt>

<sup>5</sup><https://sharegpt.com/>



Instruction fine-tuned LLMs	# Params	Base Model	Fine-tuning Trainset		
			Self-build	Dataset Name	Size
Instruct-GPT (Ouyang et al., 2022)	176B	GPT-3 (Brown et al., 2020b)	Yes	-	-
BLOOMZ (Muennighoff et al., 2022) <sup>1</sup>	176B	BLOOM (Scao et al., 2022)	No	xP3	-
FLAN-T5 (Chung et al., 2022) <sup>2</sup>	11B	T5 (Raffel et al., 2019)	No	FLAN 2021	-
Alpaca (Taori et al., 2023a) <sup>3</sup>	7B	LLaMA (Touvron et al., 2023a)	Yes	-	52K
Vicuna (Chiang et al., 2023) <sup>4</sup>	13B	LLaMA (Touvron et al., 2023a)	Yes	-	70K
GPT-4-LLM (Peng et al., 2023) <sup>5</sup>	7B	LLaMA (Touvron et al., 2023a)	Yes	-	52K
Claude (Bai et al., 2022b)	-	-	Yes	-	-
WizardLM (Xu et al., 2023a) <sup>6</sup>	7B	LLaMA (Touvron et al., 2023a)	Yes	Evol-Instruct	70K
ChatGLM2 (Du et al., 2022) <sup>7</sup>	6B	GLM (Du et al., 2022)	Yes	-	1.1 Tokens
LIMA (Zhou et al., 2023a)	65B	LLaMA (Touvron et al., 2023a)	Yes	-	1K
OPT-IML (Iyer et al., 2022) <sup>8</sup>	175B	OPT (Zhang et al., 2022a)	No	-	-
Dolly 2.0 (Conover et al., 2023a) <sup>9</sup>	12B	Pythia (Biderman et al., 2023)	No	-	15K
Falcon-Instruct (Almazrouei et al., 2023a) <sup>10</sup>	40B	Falcon (Almazrouei et al., 2023b)	No	-	-
Guanaco (JosephusCheung, 2021) <sup>11</sup>	7B	LLaMA (Touvron et al., 2023a)	Yes	-	586K
Minotaur (Collective, 2023) <sup>12</sup>	15B	StarCoder Plus (Li et al., 2023f)	No	-	-
Nous-Hermes (NousResearch, 2023) <sup>13</sup>	13B	LLaMA (Touvron et al., 2023a)	No	-	300K+
TÜLU (Wang et al., 2023e) <sup>14</sup>	6.7B	OPT (Zhang et al., 2022a)	No	Mixed	-
YuLan-Chat (YuLan-Chat-Team, 2023) <sup>15</sup>	13B	LLaMA (Touvron et al., 2023a)	Yes	-	250K
MOSS (Tianxiang and Xipeng, 2023) <sup>16</sup>	16B	-	Yes	-	-
Airoboros (Durbin, 2023) <sup>17</sup>	13B	LLaMA (Touvron et al., 2023a)	Yes	-	-
UltraLM (Ding et al., 2023a) <sup>18</sup>	13B	LLaMA (Touvron et al., 2023a)	Yes	-	-

<sup>1</sup> <https://huggingface.co/bigscience/bloomz>

<sup>2</sup> <https://huggingface.co/google/flan-t5-xxl>

<sup>3</sup> [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca)

<sup>4</sup> <https://github.com/lm-sys/FastChat>

<sup>5</sup> <https://github.com/Instruction-Tuning-with-GPT-4/GPT-4-LLM>

<sup>6</sup> <https://github.com/nlpxucan/WizardLM>

<sup>7</sup> <https://github.com/THUDM/ChatGLM2-6B>

<sup>8</sup> <https://huggingface.co/facebook/opt-impl-30b>

<sup>9</sup> <https://github.com/databricks/dolly>

<sup>10</sup> <https://huggingface.co/tiiuae/falcon-40b-instruct>

<sup>11</sup> <https://huggingface.co/JosephusCheung/Guanaco>

<sup>12</sup> <https://huggingface.co/openaccess-ai-collective/minotaur-15b>

<sup>13</sup> <https://huggingface.co/NousResearch/Nous-Hermes-13b>

<sup>14</sup> <https://github.com/allenai/open-instruct>

<sup>15</sup> <https://github.com/RUC-GSAI/YuLan-Chat>

<sup>16</sup> <https://github.com/OpenLMMLab/MOSS>

<sup>17</sup> <https://github.com/jondurbin/airoboros>

<sup>18</sup> <https://github.com/thunlp/UltraChat>

Table 3: An overview of LLMs tuned on IT datasets.

loss function tailored to multi-turn conversations. To better understand long context across multiple-turn dialog, the authors expanded the max context length from 512 to 2048. For training, the authors adopted the gradient checkpointing and flash attention (Dao et al., 2022) techniques to reduce the GPU memory cost in the fine-tuning process. The fine-tuning process takes 24 hours on an  $8 \times 80\text{GB}$  A100 device with fully shared data parallelism.

The authors built a test set used exclusively to measure chatbots’ performances. They collected a test set composed by 8 question categories, such as Fermi problems, role play scenarios, coding/math tasks, etc, and then asked GPT-4 (OpenAI, 2023) to rate models’ responses considering helpfulness, relevance, accuracy, and detail. On the constructed test set, Vicuna (13B) outperforms Alpaca (13B) (Taori et al., 2023a) and

LLaMA (13B) in 90% of the test questions, and generates equal or better rating responses compared to ChatGPT in 45% of the questions.

#### 4.6 GPT-4-LLM

GPT-4-LLM (7B) (Peng et al., 2023) is a language model trained by fine-tuning LLaMA (7B) (Touvron et al., 2023a) on the GPT-4 (OpenAI, 2023) generated instruction dataset. GPT-4-LLM is initialized with LLaMA, then fine-tuned in the following two steps: (1) supervised fine-tuning on the constructed instruction dataset. The authors used the instructions from Alpaca (Taori et al., 2023a), and then collected responses using GPT-4. LLaMA is fine-tuned on the GPT-4 generated dataset. The fine-tuning process takes approximately three hours on an  $8 \times 80\text{GB}$  A100 machine with mixed precision and fully shared data parallelism. (2) optimizing the step-1 model using

the proximal policy optimization (PPO) (Schulman et al., 2017) method, the authors first built a comparison dataset by collecting responses from GPT-4, InstructGPT (Ouyang et al., 2022), and OPT-IML (Iyer et al., 2022) to a collection of instructions and then asked GPT-4 to rate each response from 1 to 10. Using the ratings, a reward model is trained based on OPT (Zhang et al., 2022a). The fine-tuned model from Step 1 is optimized by using the reward model to compute the policy gradient.

For evaluations, GPT-4-LLM (7B) outperforms not only the baseline model Alpaca (7B), but also larger models including Alpaca (13B) and LLAMA (13B). For automated evaluation, GPT-4-LLM (7B) outperforms Alpaca by 0.2, 0.5, and 0.7 on User-Oriented-Instructions-252 (Wang et al., 2022c), Vicuna-Instructions (Chiang et al., 2023), and Unnatural Instructions (Honovich et al., 2022) datasets, respectively. For human evaluation, regarding aspects including helpfulness, honesty, and harmlessness, GPT-4-LLM outperforms Alpaca by 11.7, 20.9, and 28.6 respectively.

#### 4.7 Claude

Claude<sup>6</sup> is a language model trained by fine-tuning the pre-trained language model on an instruction dataset, aiming to generate helpful and harmless responses. The fine-tuning process consists of two stages: (1) supervised fine-tuning on the instruction dataset. The authors created an instruction dataset by collecting 52K different instructions, paired with responses generated by GPT-4. The fine-tuning process takes approximately eight hours on an 8-card 80GB A100 machine with mixed precision and fully shared data parallelism. (2) optimizing the step-1 model with the proximal policy optimization (Schulman et al., 2017) method. The authors first built a comparison dataset by collecting responses from multiple large language models (e.g., GPT-3 (Brown et al., 2020b)) to the given collection of instructions and then asking GPT-4 (OpenAI, 2023) to rate each response. Using the ratings, a reward model is trained. Then, the fine-tuned model from Step 1 is optimized using the reward model with the proximal policy optimization method.

Claude generates more helpful and harmless responses compared to the backbone model. For automatic evaluations, Claude outperforms GPT-

3 by 7% on the RealToxicityPrompts (Gehman et al., 2020) in terms of toxicity. For human evaluations, regarding four different aspects, including following correct instructions, following explicit constraints, fewer hallucinations, and generating appropriate responses, Claude outperforms GPT-3 (Brown et al., 2020b) +10%, +20%, -20%, and +10% respectively.

#### 4.8 WizardLM

WizardLM (7B) (Xu et al., 2023a) is a language model trained by fine-tuning LLaMA (7B) (Touvron et al., 2023a) on the instruction dataset Evol-Instruct generated by ChatGPT (details see Section 3.2). It is fine-tuned on a subset (with 70K) of Evol-Instruct to enable a fair comparison with Vicuna (Chiang et al., 2023). The fine-tuning process takes approximately 70 hours on 3 epochs based on an 8 V100 GPU with the Deepspeed Zero-3 (Rasley et al., 2020) technique. During inference, the max generation length is 2048.

To evaluate LLMs' performances on complex instructions, the authors collected 218 human-generated instructions from real scenarios (e.g., open-source projects, platforms, and forums), called Evol-Instruct testset.

Evaluations are conducted on the Evol-Instruct testset and Vicuna's testset. For human evaluation, WizardLM outperforms Alpaca (7B) (Taori et al., 2023a) and Vicuna (7B) by a large margins, and generates equal or better responses on 67% test samples compared to ChatGPT. Automatic evaluation is conducted by asking GPT-4 to rate LLMs' responses. Specifically, WizardLM gains performance boosts compared to Alpaca by +6.2%, +5.3% on the Evol-Instruct testset and Vicuna's test sets. WizardLM achieves outperforms Vicuna by +5.8 on the Evol-Instruct testset and +1.7% on the Vicuna's test set.

#### 4.9 ChatGLM2

ChatGLM2 (6B) (Du et al., 2022) is a language model trained by fine-tuning GLM (6B) (Du et al., 2022) on a bilingual dataset that contains both English and Chinese instructions. The bilingual instruction dataset contains 1.4T tokens, with a 1:1 ratio of Chinese to English. Instructions in the dataset are sampled from the question-answering and dialogue completion tasks. ChatGLM is initialized with GLM, then trained by the three-step fine-tuning strategy, which is akin to

<sup>6</sup><https://www.anthropic.com/index/introducing-claude>

InstructGPT (Ouyang et al., 2022). To better model contextual information across multi-turn conversations, the authors expanded the maximum context length from 1024 to 32K. To reduce GPU memory cost in the fine-tuning stage, the authors employed multi-query attention and causal mask strategies. During inference, ChatGLM2 requires 13GB GPU memory with FP16 and supports conversations up to 8K in length with 6GB GPU memory using the INT4 model quantization technique.

Evaluations are conducted on four English and Chinese benchmarks, including MMLU (English) (Hendrycks et al., 2020b), C-Eval (Chinese) (Huang et al., 2023), GSM8K (Math) (Cobbe et al., 2021), and BBH (English) (Suzgun et al., 2022b). ChatGLM2 (6B) outperforms GLM (6B) and the baseline model ChatGLM (6B) on all benchmarks. Specifically, ChatGLM2 outperforms GLM by +3.1 on MMLU, +5.0 on C-Eval, +8.6 on GSM8K, and +2.2 on BBH. ChatGLM2 achieves better performances than ChatGLM by +2.1, +1.2, +0.4, +0.8 on MMLU, C-Eval, GSM8K and BBH, respectively.

#### 4.10 LIMA

LIMA (65B) (Zhou et al., 2023a) is a large language model trained by fine-tuning LLaMA (65B) (Touvron et al., 2023a) on an instruction dataset, which is constructed based on the proposed superficial alignment hypothesis.

The superficial alignment hypothesis refers to the idea that the knowledge and capabilities of a model are almost acquired in the pre-training stage, while the alignment training (e.g., instruction tuning) teaches models to generate responses under user-preferred formalizations. Based on the superficial alignment hypothesis, the authors claimed that large language models can generate user-satisfied responses by fine-tuning it on a small fraction of instruction data. Therefore, the authors built instruction train/valid/test sets to verify this hypothesis.

Evaluations are conducted on the constructed test set. For human evaluations, LIMA outperforms InstructGPT and Alpaca by 17% and 19%, respectively. Additionally, LIMA achieves comparable results to BARD<sup>7</sup>, Claude<sup>8</sup>, and GPT-4. For automatic evaluation, which is conducted

by asking GPT-4 to rate responses and a higher rate score denotes better performance, LIMA outperforms InstructGPT and Alpaca by 20% and 36%, respectively, achieving comparable results to BARD, while underperforming Claude and GPT-4. Experimental results verify the proposed superficial alignment hypothesis.

#### 4.11 Others

**OPT-IML (175B)** (Iyer et al., 2022) is a large language model trained by fine-tuning the OPT (175B) (Zhang et al., 2022a) model on the constructed Instruction Meta-Learning (IML) dataset, which consists of over 1500 NLP tasks from 8 publicly available benchmarks such as PromptSource (Bach et al., 2022), FLAN (Longpre et al., 2023), and Super-NaturalInstructions (Wang et al., 2022e). After fine-tuning, OPT-IML outperforms OPT across all benchmarks.

**Dolly 2.0 (12B)** (Conover et al., 2023a) is initialized with the pre-trained language model Pythia (12B) (Biderman et al., 2023), and fine-tuned on the instruction dataset databricks-dolly-15k<sup>9</sup>, which contains 7 categories of NLP tasks such as text classification and information extraction. After fine-tuning, Dolly 2.0 (12B) outperforms Pythia (12B) on the EleutherAI LLM Evaluation Harness benchmark (Gao et al., 2021) by a large margin, and achieves comparable performances to GPT-NEOX (20B) (Black et al., 2022), which has two times more parameters compared to Dolly 2.0 (12B).

**Falcon-Instruct (40B)** (Almazrouei et al., 2023a) is a large language model trained by fine-tuning Falcon (40B) (Almazrouei et al., 2023b) on an English dialogue dataset, which contains 150 million tokens from the Baize dataset (Xu et al., 2023c), with an additional 5% of the data from the RefinedWeb dataset (Penedo et al., 2023). To reduce memory usage, the authors employed flash attention (Dao et al., 2022) and multi-query techniques. For evaluation, Falcon-Instruct (40B) achieved better performance on the Open LLM Leaderboard (Beeching et al., 2023)<sup>10</sup> compared to the baseline model Falcon (40B), and outperforms the Guanaco (65B), which has more model parameters.

<sup>9</sup><https://huggingface.co/datasets/databricks/databricks-dolly-15k>

<sup>10</sup>[https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard)

<sup>7</sup><https://bard.google.com/>

<sup>8</sup><https://www.anthropic.com/index/introducing-claude>

**Guanaco (7B)** (JosephusCheung, 2021) is a multi-turn dialog language model trained by fine-tuning LLaMA (7B) (Touvron et al., 2023a) on the constructed multilingual dialogue dataset. The multilingual dialogue dataset comes from two sources: Alpaca (Taori et al., 2023a), which contains 52K English instruction data pairs; and a multilingual (e.g., Simplified Chinese, Traditional Chinese, Japanese, German) dialogue data, which contains 534K+ multi-turn conversations. After fine-tuning, Guanaco is to generate role-specific responses and continuous responses on a given topic in multi-turn conversations.

**Minotaur (15B)** is a large language model trained by fine-tuning the Starcoder Plus (15B) (Li et al., 2023f) on open-source instruction datasets including WizardLM (Xu et al., 2023a) and GPTeacher-General-Instruct<sup>11</sup>. For model inference, Minotaur supports a maximum context length of 18K tokens.

**Nous-Herme (13B)** is a large language model trained by fine-tuning LLaMA (13B) (Touvron et al., 2023a) on an instruction dataset, which contains over 300k instructions, sampled from GPTeacher<sup>12</sup>, CodeAlpaca (Chaudhary, 2023), GPT-4-LLM (Peng et al., 2023), Unnatural Instructions (Honovich et al., 2022), and BiologyPhysicsChemistry subsets in the Camel-AI (Li et al., 2023c). Responses are generated by GPT-4. For evaluations, Nous-Herme (13B) achieves comparable performances to GPT-3.5-turbo on multiple tasks like ARC challenge (Clark et al., 2018) and BoolQ (Clark et al., 2019).

**TÜLU (6.7B)** (Wang et al., 2023e) is a large language model trained by fine-tuning OPT (6.7B) (Zhang et al., 2022a) on a mixed instruction dataset, which contains FLAN V2 (Longpre et al., 2023), CoT (Wei et al., 2022), Dolly (Conover et al., 2023a), Open Assistant-1<sup>13</sup>, GPT4-Alpaca<sup>14</sup>, Code-Alpaca (Chaudhary, 2023), and ShareGPT<sup>15</sup>. After fine-tuning, TÜLU (6.7B) reaches on average 83% of ChatGPT’s performance and 68% of GPT-4’s performance.

**YuLan-Chat (13B)** (YuLan-Chat-Team, 2023) is a language model trained by fine-tuning LLaMA

(13B) (Touvron et al., 2023a) on a constructed bilingual dataset, which contains 250,000 Chinese-English instruction pairs. After fine-tuning, YuLan-Chat-13B achieves comparable results to the state-of-the-art open-source model ChatGLM (6B) (Du et al., 2022), and outperforms Vicuna (13B) (Chiang et al., 2023) on the English BBH3K (BBH3K is a subset of BBH benchmark (Srivastava et al., 2022a)) dataset.

**MOSS (16B)**<sup>16</sup> is a bilingual dialogue language model, which aims to engage in multi-turn conversations and utilize various plugins, trained by fine-tuning on dialogue instructions. After fine-tuning, MOSS outperforms the backbone model and generates responses that better align with human preferences.

**Airoboros (13B)**<sup>17</sup> is a large language model trained by fine-tuning LLAMA (13B) (Touvron et al., 2023a) on the Self-instruct dataset (Wang et al., 2022c). After fine-tuning, Airoboros significantly outperforms LLAMA (13B) (Touvron et al., 2023a) on all benchmarks and achieves highly comparable results to models fine-tuned specifically for certain benchmarks.

**UltraLM (13B)** (Ding et al., 2023a) is a large language model trained by fine-tuning LLAMA (13B) (Touvron et al., 2023a). For evaluation, UltraLM (13B) outperforms Dolly (12B) (Conover et al., 2023a) and achieves the winning rate up to 98%. Additionally, it surpasses the previous best open-source models (i.e., Vicuna (Chiang et al., 2023) and WizardLM (Xu et al., 2023a)) with winning rates of 9% and 28%, respectively.

## 5 Multi-modality Instruction Tuning

### 5.1 Multi-modality Datasets

**MUL-TIINSTRUCT** (Xu et al., 2022) is a multimodal instruction tuning dataset consisting of 62 diverse multimodal tasks in a unified seq-to-seq format. This dataset covers 10 broad categories and its tasks are derived from 21 existing open-sourced datasets. Each task is equipped with 5 expert-written instructions. For the existing tasks, the authors use the input/output pairs from their available open-source datasets to create instances. While for each new task, the authors create 5k to 5M instances by extracting

<sup>11</sup><https://github.com/teknum1/GPTeacher>

<sup>12</sup><https://github.com/teknum1/GPTeacher>

<sup>13</sup><https://huggingface.co/datasets/OpenAssistant/oasst1>

<sup>14</sup><https://huggingface.co/datasets/vicgalle/alpaca-gpt4>

<sup>15</sup><https://sharegpt.com/>

<sup>16</sup><https://txsun1997.github.io/blogs/moss.html>

<sup>17</sup><https://github.com/jondurbin/airoboros>



Multi-modality Instruction Fine-tuning Dataset	Modalities		# Tasks
	Modality Pair	# Instance	
MUL-TIINSTRUCT (Xu et al., 2022) <sup>1</sup>	Image-Text	5k to 5M per task	62
PMC-VQA (Zhang et al., 2023c) <sup>2</sup>	Image-Text	227k	2
LAMM (Yin et al., 2023) <sup>3</sup>	Image-Text	186k	9
	Point Cloud-Text	10k	3
Vision-Flan (Xu et al., 2024b) <sup>4</sup>	Multi-Pairs	Over 1M	200+
ALLAVA (Chen et al., 2024a) <sup>5</sup>	Image-Text	1.4M	2
ShareGPT4V (Chen et al., 2023a) <sup>6</sup>	Image-Text	1.2M	2

<sup>1</sup> <https://github.com/VT-NLP/MultiInstruct>

<sup>2</sup> <https://github.com/xiaoman-zhang/PMC-VQA>

<sup>3</sup> <https://github.com/OpenLAMM/LAMM>

<sup>4</sup> <https://vision-flan.github.io/>

<sup>5</sup> <https://github.com/FreedomIntelligence/ALLaVA>

<sup>6</sup> <https://sharegpt4v.github.io/>

Table 4: An overview of multi-modality instruction fine-tuning datasets.

the necessary information from instances of existing tasks or reformulating them. The MUL-TIINSTRUCT dataset has demonstrated its efficiency in enhancing various transfer learning technique. For example, fine-tuning the OFA model (930M) (Wang et al., 2022a) using various transfer learning strategies such as Mixed Instruction Tuning and Sequential Instruction Tuning on MUL-TIINSTRUCT improve the zero-shot performance across all unseen tasks. On commonsense VQA task, OFA fine-tuned on MUL-TIINSTRUCT achieves 50.60 on RougeL and 31.17 on accuracy, while original OFA achieves 14.97 on RougeL and 0.40 on accuracy.

**PMC-VQA** (Zhang et al., 2023c) is a large-scale medical visual question-answering dataset that comprises 227k image-question pairs of 149k images, covering various modalities or diseases. The dataset can be used for both open-ended and multiple-choice tasks. The pipeline for generating the PMC-VQA dataset involves collecting image-caption pairs from the PMC-OA (Lin et al., 2023c) dataset, using ChatGPT to generate question-answer pairs, and manually verifying a subset of the dataset for quality. The authors propose a generative-based model MedVInT for medical visual understanding by aligning visual information with a large language model. MedVInT pretrained on PMC-VQA achieves state-of-the-art performance and outperforms existing models on VQA-RAD (Lau et al., 2018) and SLAKE (Liu et al., 2021a) benchmarks, with 81.6% accuracy on VQA-RAD and 88.0% accuracy on SLAKE.

**LAMM** (Yin et al., 2023) is a comprehensive multi-modal instruction tuning dataset for 2D image and 3D point cloud understanding. LAMM contains 186K language-image instruction-response pairs, and 10K language-point cloud instruction-response pairs. The authors collect images and point clouds from publicly available datasets and use the GPT-API and self-instruction methods to generate instructions and responses based on the original labels from these datasets. LAMM-Dataset includes data pairs for commonsense knowledge question answering by incorporating a hierarchical knowledge graph label system from the Bamboo (Zhang et al., 2022b) dataset and the corresponding Wikipedia description. The authors also propose the LAMM-Benchmark, which evaluates existing multi-modal language models (MLLM) on various computer vision tasks. It includes 9 common image tasks and 3 common point cloud tasks, and LAMM-Framework, a primary MLLM training framework that differentiates the encoder, projector, and LLM finetuning blocks for different modalities to avoid modality conflicts.

**Vision-Flan** (Xu et al., 2024b) is the largest public-available human-annotated visual instruction tuning dataset that consists of 1,664,261 instances and 200+ diverse vision-language tasks derived from 101 open-source computer vision datasets. Each task is accompanied by expertly written instructions and meticulously crafted templates for inputs and outputs. The dataset covers a broad spectrum of tasks, including image captioning, visual question-answering, and visual comprehension. Designed to enhance

research and application in vision-language model domains, Vision-Flan aims to expand the horizons of interaction and comprehension between visual and linguistic modalities. It provides researchers and developers with a valuable resource to push the envelope of vision-language models and to innovate algorithms across a diverse array of fields.

**ALLaVA** (Chen et al., 2024a) represents an open-source, extensive dataset tailored for fine-tuning visual question-answering models, featuring 1.4M entries that include detailed captions, intricate instructions, and comprehensive answers produced by GPT-4V (Yang et al., 2023b). To craft high-quality captions and visual question-answers, Chen et al. (2024a) introduced a method to distill both a caption and a QA pair for an image in a single interaction. This process involves initially presenting GPT-4V (Yang et al., 2023b) with an image, followed by prompting it to generate both a detailed caption and a visual question-answer pair. This approach of incorporating additional visual data enables the model to develop a more nuanced understanding of both the visual and textual elements, enhancing its capacity to deliver precise and contextually appropriate answers. Furthermore, this method has the potential to reduce the occurrence of hallucinations by providing the model with more contextual information (visual data).

**ShareGPT4V** (Chen et al., 2023a) is a collection of highly descriptive image-text pairs, consisting of two components: 100K captions generated by GPT4-Vision (Yang et al., 2023b) from a variety of images, and 1.2M captions developed using their pre-trained model, which was trained on the initial set of 100K high-quality captions. These captions comprehensively cover aspects such as global knowledge, object attributes, spatial relationships, and aesthetic evaluations. Utilizing this dataset, the ShareGPT4V-7B model, once fine-tuned, surpasses competing 7B-scale LMMs across all 11 benchmark tests. Notably, it secures a remarkable cumulative score of 1943.8 on the MME benchmark, outperforming the second-place Qwen-VL-Chat-7B (Bai et al., 2023) model, which was trained with 1.4 billion samples, by 95.6 points.

## 5.2 Multi-modality Instruction Tuning Models

**InstructPix2Pix (983M)** (Brooks et al., 2022) is a conditional diffusion model trained by fine-tuning Stable Diffusion (983M) (Rombach et al., 2022)

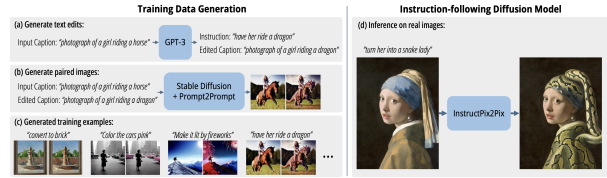


Figure 10: Image editing dataset generation and diffusion model training. The figure is copied from Brooks et al. (2022).

on a constructed multi-modal dataset that contains more than 450K text editing instructions and corresponding images before and after the edit. The authors combine the abilities of two large-scale pre-trained models, a language model GPT-3 (Brown et al., 2020b) and a text-to-image model Stable Diffusion (Rombach et al., 2022), to generate the training dataset. GPT-3 is fine-tuned to generate text edits based on image prompts, while Stable Diffusion is used to convert the generated text edits into actual image edits. InstructPix2Pix is then trained on this generated dataset using a latent diffusion objective. Figure 10 shows the process of generating image editing dataset and training the diffusion model on that dataset. The authors compares the proposed method qualitatively with previous works such as SDEdit (Meng et al., 2022) and Text2Live (Bar-Tal et al., 2022), highlighting the ability of the model to follow image editing instructions instead of descriptions of the image or edit layer. The authors also presents quantitative comparisons with SDEdit (Meng et al., 2022) using metrics measuring image consistency and edit quality.

**LLaVA (13B)** (Liu et al., 2023b) is a large multimodal model developed by connecting the visual encoder of CLIP (400M) (Radford et al., 2021) with the language decoder LLaMA (7B) (Touvron et al., 2023a). LLaVA is fine-tuned using the generated instructional vision-language dataset consisted of 158K unique language-image instruction-following samples. The data collection process involved creating conversation, detailed description, and complex reasoning prompts. GPT-4 is used to convert image-text pairs into appropriate instruction-following format for this dataset. Visual features such as captions and bounding boxes were used to encode images. LLaVA yields a 85.1% relative score compared with GPT-4 on a synthetic multimodal instruction following dataset. When fine-tuned on Science QA,

Multi-modality Instruction Fine-tuned LLMs	# Params	Modality	Base Model Model Name	# Params	Fine-tuning Self-build	Trainset Size
InstructPix2Pix (Brooks et al., 2022) <sup>1</sup>	983M	I/T	Stable Diffusion	983M	Yes	450K
LLaVA (Liu et al., 2023b) <sup>2</sup>	13B	I/T	CLIP (Radford et al., 2021) LLaMA (Touvron et al., 2023a) LLaMA (Touvron et al., 2023a)	400M 7B 7B	Yes	158K
Video-LLaMA (Zhang et al., 2023b) <sup>3</sup>	-	I/T/V/A	BLIP-2 (Li et al., 2023d) ImageBind (Girdhar et al., 2023) Vicuna (Chiang et al., 2023)	- - 7B/13B	No	-
InstructBLIP (1.2B) (Dai et al., 2023) <sup>4</sup>	-	I/T/V	BLIP-2 (Li et al., 2023d)	-	No	-
Otter (Li et al., 2023b) <sup>5</sup>	-	I/T/V	OpenFlamingo (Awadalla et al., 2023)	9B	Yes	2.8M
MultiModal-GPT (Gong et al., 2023) <sup>6</sup>	-	I/T/V	OpenFlamingo (Awadalla et al., 2023)	9B	No	-

<sup>1</sup> <https://github.com/timothybrooks/instruct-pix2pix>

<sup>2</sup> <https://github.com/haotian-liu/LLaVA>

<sup>3</sup> <https://github.com/DAMO-NLP-SG/Video-LLaMA>

<sup>4</sup> <https://github.com/salesforce/LAVIS/tree/main/projects/instructblip>

<sup>5</sup> <https://github.com/Luodian/Otter>

<sup>6</sup> <https://github.com/open-mmlab/Multimodal-GPT>

Table 5: An overview of multi-modality instruction fine-tuned LLMs. I/T/V/A stand for Image/Text/Video/Audio

the synergy of LLaVA and GPT-4 achieves a new state-of-the-art accuracy of 92.53%.

**Video-LLaMA** (Zhang et al., 2023b) is a multimodal framework that enhances large language models with the ability to understand both visual and auditory content in videos. The architecture of Video-LLaMA consists of two branch encoders: the Vision-Language (VL) Branch and the Audio-Language (AL) Branch, and a language decoder (Vicuna (7B/13B) (Chiang et al., 2023), LLaMA (7B) (Touvron et al., 2023a), etc.). The VL Branch includes a frozen pre-trained image encoder (pre-trained vision component of BLIP-2 (Li et al., 2023d), which includes a ViT-G/14 and a pre-trained Q-former), a position embedding layer, a video Q-former and a linear layer. The AL Branch includes a pre-trained audio encoder (ImageBind (Girdhar et al., 2023)) and an Audio Q-former. Figure 11 shows the overall architecture of Video-LLaMA with Vision-Language Branch and Audio-Language Branch. The VL Branch is trained on the Webvid-2M (Bain et al., 2021) video caption dataset with a video-to-text generation task, and fine-tuned on the instruction tuning data from MiniGPT-4 (Zhu et al., 2023b), LLaVA (Liu et al., 2023b) and VideoChat (Li et al., 2023e). The AL Branch is trained on video/image instruction data to connect the output of ImageBind to language decoder. After finetuning, Video-

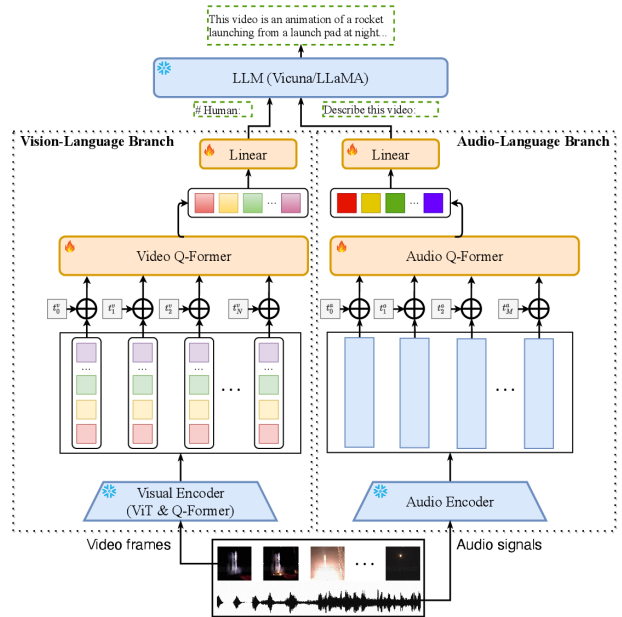


Figure 11: Overall architecture of Video-LLaMA. The figure is copied from Zhang et al. (2023b).

LLaMA can perceive and comprehend video content, demonstrating its ability to integrate auditory and visual information, understand static images, recognize common-knowledge concepts, and capture temporal dynamics in videos.

**InstructBLIP (1.2B)** (Dai et al., 2023) is a vision-language instruction tuning framework initialized with a pre-trained BLIP-2 (Li et al., 2023d) model consisting of an image encoder,

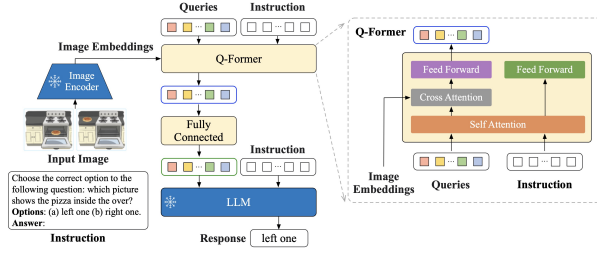


Figure 12: Overall architecture of InstructBLIP. The figure is copied from Dai et al. (2023).

an LLM (FlanT5 (3B/11B) (Chung et al., 2022) or Vicuna (7B/13B) (Chiang et al., 2023)), and a Query Transformer (Q-Former) to bridge the two. As shown in Figure 12, the Q-Former extracts instruction-aware visual features from the output embeddings of the frozen image encoder, and feeds the visual features as soft prompt input to the frozen LLM. The authors evaluate the proposed InstructBLIP model on a variety of vision-language tasks, including image classification, image captioning, image question answering, and visual reasoning. They use 26 publicly available datasets, dividing them into 13 held-in and 13 held-out datasets for training and evaluation. The authors demonstrate that InstructBLIP achieves state-of-the-art zero-shot performance on a wide range of vision-language tasks. InstructBLIP yields an average relative improvement of 15.0% when compared to BLIP-2, smallest InstructBLIP (4B) outperforms Flamingo (80B) (Alayrac et al., 2022) on all six shared evaluation datasets with an average relative improvement of 24.8%.

**Otter** (Li et al., 2023b) is a multi-modal model trained by fine-tuning OpenFlamingo (9B) (Awadalla et al., 2023), with the language and vision encoders frozen and only fine-tuning the Perceiver resampler module, cross-attention layers, and input/output embeddings. The authors organize diverse multi-modal tasks covering 11 categories and build multi-modal in-context instruction tuning datasets MIMIC-IT of 2.8M multimodal instruction-response pairs, which consists of image-instruction-answer triplets, where the instruction-answer is tailored to the image. Each data sample also includes context, which contains a series of image-instruction-answer triplets that contextually correlate with the queried triplet. Otter demonstrates the ability to follow user instructions more accurately and provide more detailed descriptions of images compared to

OpenFlamingo (Awadalla et al., 2023).

**MultiModal-GPT** (Gong et al., 2023) is a multi-modal instruction tuning model that is capable of following diverse instructions, generating detailed captions, counting specific objects, and addressing general inquiries. MultiModal-GPT is trained by fine-tuning OpenFlamingo (9B) (Awadalla et al., 2023) on various created visual instruction data with open datasets, including VQA, Image Captioning, Visual Reasoning, Text OCR, and Visual Dialogue. The experiments demonstrate the proficiency of MultiModal-GPT in maintaining continuous dialogues with humans.

## 6 Domain-specific Instruction Tuning

In this section, we describe instruction tuning in different domains and applications.

### 6.1 Dialogue

**InstructDial** (Gupta et al., 2022) is an instruction tuning framework designed for dialogue. It contains a collection of 48 dialogue tasks in a consistent text-to-text format created from 59 dialogue datasets. Each task instance includes a task description, instance inputs, constraints, instructions, and output. To ensure adherence to instructions, the framework introduces two meta-tasks: (1) an instruction selection task, where the model selects the instruction corresponding to a given input-output pair; and (2) an instruction binary task, where the model predicts "yes" or "no" if an instruction leads to a given output from an input. Two base models T0-3B (Sanh et al., 2021) (3B parameters version of T5 (Lester et al., 2021)) and BART0 (Lin et al., 2022) (406M parameters based on Bart-large (Lewis et al., 2019)) are fine-tuned on the tasks from InstructDial. InstructDial achieves impressive results on unseen dialogue datasets and tasks, including dialogue evaluation and intent detection. Moreover, it delivers even better results when applied to a few-shot setting.

### 6.2 Intent Classification and Slot Tagging

**LINGUIST** (Rosenbaum et al., 2022) finetunes AlexaTM 5B (Soltan et al., 2022), a 5-billion-parameter multilingual model, on the instruction dataset for intent classification and slot tagging tasks. Each instruction consists of five blocks: (i) the language of the generated output, (ii) intention, (iii) slot types and values to include in the output (e.g., the number 3 in [3, snow] corresponds the



Domain Type	Domain-specific Instruction	Base Model		Trainset Size
	Fine-tuned LLMs	Model Name	# Params	
Dialogue	InstructDial (Gupta et al., 2022) <sup>1</sup>	T0 (Sanh et al., 2021)	3B	-
Classification	LINGUIST (Rosenbaum et al., 2022)	AlexaTM (Soltan et al., 2022)	5B	13K
Information extraction	InstructUIE (Wang et al., 2023d) <sup>2</sup>	FlanT5 (Chung et al., 2022)	11B	1.0M
Sentiment analysis	IT-MTL (Varia et al., 2022) <sup>3</sup>	T5 (Raffel et al., 2019)	220M	-
Writing	Writing-Alpaca-7B (Zhang et al., 2023d) <sup>4</sup>	LLaMA (Touvron et al., 2023a)	7B	-
	CoEdIT (Raheja et al., 2023) <sup>5</sup>	FlanT5 (Chung et al., 2022)	11B	-
	CoPoet (Chakrabarty et al., 2022) <sup>6</sup>	T5 (Raffel et al., 2019)	11B	-
Medical	Radiology-GPT (Liu et al., 2023c) <sup>7</sup>	Alpaca (Taori et al., 2023a)	7B	122K
	ChatDoctor (Li et al., 2023j) <sup>8</sup>	LLaMA (Touvron et al., 2023a)	7B	100K
	ChatGLM-Med (Wang et al., 2023a) <sup>9</sup>	ChatGLM (Du et al., 2022)	6B	-
Arithmetic	Goat (Liu and Low, 2023) <sup>10</sup>	LLaMA (Touvron et al., 2023a)	7B	1.0M
Code	WizardCoder (Luo et al., 2023) <sup>11</sup>	StarCoder (Li et al., 2023f)	15B	78K

<sup>1</sup> <https://github.com/prakharguptaz/Instructdial>

<sup>2</sup> <https://github.com/BeyondrXX/InstructUIE>

<sup>3</sup> <https://github.com/amazon-science/instruction-tuning-for-absa>

<sup>4</sup> <https://github.com/facebookresearch/EditEval>

<sup>5</sup> <https://github.com/vipulraheja/coedit>

<sup>6</sup> <https://github.com/vishakhpk/creative-instructions>

<sup>7</sup> <https://huggingface.co/spaces/allen-eric/radiology-gpt>

<sup>8</sup> <https://github.com/Kent0n-Li/ChatDoctor>

<sup>9</sup> <https://github.com/SCIR-HI/Med-ChatGLM>

<sup>10</sup> <https://github.com/liutiedong/goat>

<sup>11</sup> <https://github.com/nlpxucan/WizardLM>

Table 6: An overview of domain-specific instruction fine-tuned LLMs.

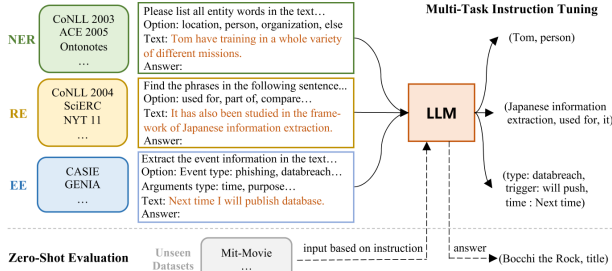


Figure 13: The overview framework of InstructUIE. The figure is copied from Wang et al. (2023d).

slot type, and snow is the value used for that slot), (iv) a mapping from slot type labels to numbers, and (v) up to 10 examples to instruct the format of the outputs. LINGUIST shows significant improvements over state-of-the-art approaches in a 10-shot novel intent setting using the SNIPS dataset (Coucke et al., 2018). In the zero-shot cross-lingual setting of the MATIS++ dataset (Xu et al., 2020), LINGUIST surpasses a strong baseline of Machine Translation with Slot Alignment across 6 languages while maintaining intent classification performance.

### 6.3 Information Extraction

**InstructUIE** (Wang et al., 2023d) is a unified information extraction (IE) framework based on

instruction tuning, which transforms IE tasks to the seq2seq format and solves them by fine-tuning 11B FlanT5 (Chung et al., 2022) on the constructed SFT dataset. Figure 13 shows the overall architecture of InstructUIE. It introduces IE INSTRUCTIONS, a benchmark of 32 diverse information extraction datasets in a unified text-to-text format with expert-written instructions. Each task instance is delineated by four properties: task instruction, options, text, and output. Task instruction contains information such as the type of information to be extracted, the output structure format, and additional constraints or rules that need to be adhered to during the extraction process. Options refer to the output label constraints of a task. Text refers to the input sentence. Output is the sentence obtained by converting the original tags of the sample (e.g. "entity tag: entity span" for NER). In the supervised setting, InstructUIE performs comparably to BERT (Devlin et al., 2018) and outperforms the state-of-the-art and GPT3.5 (Brown et al., 2020a) in zero-shot settings.

### 6.4 Aspect-based Sentiment Analysis

Varia et al. (2022) propose a unified instruction tuning framework for solving Aspect-based Sentiment Analysis (ABSA) task based on a fine-

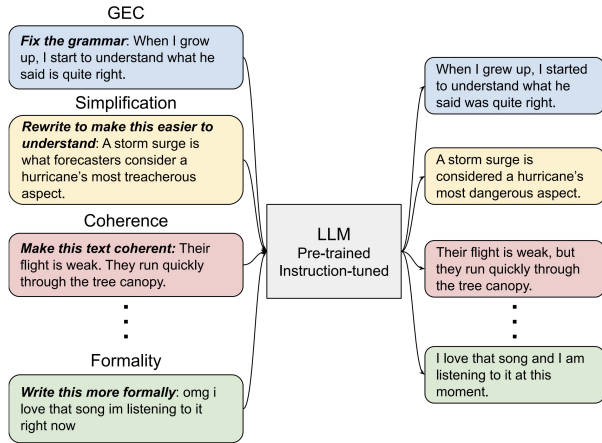


Figure 14: The overview framework of COEDIT. The figure is copied from Raheja et al. (2023).

tuned T5 (220M) (Raffel et al., 2019) model. The framework addresses multiple factorized sub-tasks that involve the four elements of ABSA, namely Aspect Term, Aspect Category, Opinion Term, and Sentiment. It treats these sub-tasks as a combination of five Question Answering (QA) tasks by transforming each sentence in the corpus using instruction templates provided for each task. For instance, one of the instruction templates used is "What are the aspect terms in the text: \$TEXT?". The framework showcases substantial improvement (8.29 F1 on average) over the state-of-the-art in few-shot learning scenarios and remains comparable in full fine-tuning scenarios.

## 6.5 Writing

Zhang et al. (2023d) propose Writing-Alpaca-7B that fine-tunes LLaMa-7B (Peng et al., 2023) on the writing instruction dataset to provide writing assistance. The proposed instruction dataset is an extension of the EDITEVAL (Dwivedi-Yu et al., 2022) benchmark based on instructional data, with the Updating task removed and a task for grammaticality introduced. The instruction scheme strictly follows the one in the Stanford Alpaca project (Taori et al., 2023a), comprising a universal preface, an instruction field to guide task completion, an input field that provides the text to be edited, and a response field that requires models to fill out. The Writing-Alpaca-7B improves upon LLaMa’s performance on all writing tasks and outperforms other larger off-the-shelf LLMs.

**CoEdit** (Raheja et al., 2023) finetunes FLANT5 (Chung et al., 2022) (770M parameters, 3B parameters, and 11B parameters) on the

instruction dataset for text editing to provide writing assistance. The instruction dataset comprises approximately 82K <instruction: source, target> pairs. As shown in Figure 14, the model takes instructions from the user specifying the characteristics of the desired text, such as "Make the sentence simpler", and outputs the edited text. CoEdit achieves state-of-the-art performance on several text editing tasks, including grammatical error correction, text simplification, iterative text editing, and three stylistic editing tasks: formality style transfer, neutralization, and paraphrasing. Furthermore, it can generalize well to new, adjacent tasks not seen during fine-tuning.

**CoPoet** (Chakrabarty et al., 2022) is a collaborative poetry writing tool that utilizes a large language model (e.g. T5-3B, T5-11B and T0-3B models) trained on a diverse collection of instructions for poetry writing. Each sample in the instruction dataset includes an <instruction, poem\_line> pair. There are three major types of instructions: Continuation, Lexical Constraints, and Rhetorical Techniques. The CoPoet is guided by user instructions that specify desired attributes of the poetry, such as writing a sentence about "love" or ending a sentence with "fly." Not only is the system competitive with publicly available LLMs trained on instructions, such as InstructGPT (Ouyang et al., 2022), but it is also capable of satisfying unseen compositional instructions.

## 6.6 Medical

**Radiology-GPT** (Liu et al., 2023c) is a fine-tuned Alpaca-7B (Taori et al., 2023a) model for radiology, which utilizes an instruction tuning approach on an extensive dataset of radiology domain knowledge. Radiology reports usually include two corresponding sections: "Findings" and "Impression". The "Findings" section contains detailed observations from the radiology images, while the "Impression" section summarizes the interpretations drawn from those observations. Radiology-GPT provides a brief instruction to the "Findings" text: "Derive the impression from findings in the radiology report". The "Impression" text from the same report serves as the target output. In comparison to general language models such as StableLM (Islamovic), Dolly (Conover et al., 2023a), and LLaMA (Touvron et al., 2023a), Radiology-GPT demonstrates significant

adaptability in radiological diagnosis, research, and communication.

**ChatDoctor** (Li et al., 2023j) is based on the fine-tuned LLaMa-7B (Touvron et al., 2023a) model, utilizing the alpaca instruction dataset (Taori et al., 2023a) and the HealthCareMagic100k patient-doctor dialogue dataset. And prompt templates are designed for retrieving external knowledge databases, such as the Disease Database and Wikipedia retrieval, during doctor-patient conversations to obtain more accurate outputs from the model. The ChatDoctor significantly improves the model’s ability to comprehend patient needs and provide informed advice. By equipping the model with self-directed information retrieval from reliable online and offline sources, the accuracy of its responses is substantially improved.

**ChatGLM-Med** (Wang et al., 2023a) is fine-tuned on the Chinese medical instruction dataset based on the ChatGLM-6B (Du et al., 2022) model. The instruction dataset comprises medically relevant question and answer pairs, created using the GPT 3.5 API and the Medical Knowledge Graph. This model improves the question-answering performance of ChatGLM (Du et al., 2022) in the medical field.

## 6.7 Arithmetic

**Goat** (Liu and Low, 2023) is a fine-tuned LLaMA-7B (Touvron et al., 2023a) model based on instructions, which aims to solve arithmetic problems. It expresses arithmetic problems in the form of natural language question answering, such as "What is 8914/64?", by generating hundreds of instruction templates using ChatGPT (OpenAI, 2022). The model applies various techniques to enhance its adaptability to diverse question formats, such as randomly removing spaces between numbers and symbols in the arithmetic expression and replacing "\*" with "x" or "times". The Goat model achieves state-of-the-art performance on the BIG-bench (Srivastava et al., 2022a) arithmetic subtask. In particular, zero-shot Goat-7B matches or exceeds the accuracy achieved by the few-shot PaLM-540B (Chowdhery et al., 2022).

## 6.8 Code

**WizardCoder** (Luo et al., 2023) utilizes StarCoder 15B (Li et al., 2023f) as the foundation with complex instruction tuning, by adapting the

Evol-Instruct method (Xu et al., 2023a) to the domain of code. The training dataset is produced through iterative application of the Evol-Instruct technique on the Code Alpaca dataset (Taori et al., 2023b), which includes the following attributes for each sample: instruction, input, and expected output. For instance, when the instruction is "Amend the following SQL query to select distinct elements", the input is the SQL query, and the expected output is the generated answer. The WizardCoder outperforms all other open-source Code LLMs and even surpasses the largest closed LLMs, Anthropic’s Claude and Google’s Bard, on HumanEval and HumanEval+.

# 7 Efficient Tuning Techniques

Efficient fine-tuning techniques aim at adapting LLMs to downstream tasks by optimizing a small fraction of parameters in multiple ways, *i.e.*, addition-based, specification-based, and reparameterization-based. Addition-based methods introduce extra trainable parameters or modules not present in the original model. Representative methods include adapter tuning (Houlsby et al., 2019) and prompt-based tuning (Schick and Schütze, 2021). Specification-based methods specify certain inherent model parameters to be tuned while freezing others. For example, BitFit (Zaken et al., 2022) tunes the bias terms of the pre-trained model. Reparameterization methods transform model weights into more parameter-efficient forms for tuning. The key hypothesis is that model adaptation is low-rank, so weights can be reparameterized into low-rank factors or a low-dimensional subspace (*e.g.*, LoRA (Hu et al., 2021)). Intrinsic prompt tuning finds a low-dimensional subspace shared by tuning prompts across diverse tasks.

## 7.1 LoRA

Low-Rank Adaptation (LoRA) (Hu et al., 2021) enables efficient adaptation of LLMs using low-rank updates. LoRA use DeepSpeed (Rasley et al., 2020) as the training backbone. The key insight of LoRA is that the actual change in LLMs’ weights required for new task adaptation lies in a low-dimensional subspace. Specifically, for a pretrained weight matrix  $W_0$ , the authors model the adapted weight matrix as  $W_0 + \Delta W$ , where  $\Delta W$  is a low rank update.  $\Delta W$  is parameterized as  $\Delta W = BA$ , where  $A$  and  $B$  are much smaller trainable matrices.

The rank  $r$  of  $\Delta W$  is chosen to be much smaller than the dimensions of  $W_0$ . The intuition is that instead of directly training all of  $W_0$ , the authors train low-dimensional  $A$  and  $B$ , which indirectly trains  $W_0$  in a low-rank subspace of directions that matter for the downstream task. This results in far fewer trainable parameters compared to full fine-tuning. For GPT-3, LoRA reduces the number of trainable parameters by 10,000x and memory usage by 3x compared to full fine-tuning.

## 7.2 HINT

HINT (Iverson et al., 2022) combines the generalization benefits of instruction tuning with efficient on-demand fine-tuning, avoiding repeatedly processing lengthy instructions. The essence of HINT lies in hypernetworks, which generate parameter-efficient modules for LLMs adaptation based on natural language instructions and few-shot examples. The adopted hypernetwork converts instructions and few-shot examples into an encoded instruction and generates adapter and prefix parameters using a pretrained text encoder and cross-attention based parameter generator. Then, the generated adapters and prefixes are inserted into the backbone model as efficient tuning modules. At inference, the hypernetwork performs inference only once per task to generate adapted modules. The benefits are that HINT can incorporate long instructions and additional few-shots without increasing compute, unlike regular fine-tuning or input concatenation methods.

## 7.3 Qlora

QLORA (Dettmers et al., 2023) includes optimal quantization and memory optimization, aiming at providing efficient and effective LLMs fine-tuning. QLORA includes 4-bit NormalFloat (NF4) Quantization, which is a quantization scheme optimized for the typical normal distribution of LLM weights. By quantizing based on the quantiles of a normal distribution, NF4 provides better performance than standard 4-bit integer or float quantization. To further reduce memory, the quantization constants are themselves quantized to 8 bits. This second level of quantization saves an additional 0.37 bits per parameter on average. QLORA leverages NVIDIA’s unified memory feature to page optimizer states to CPU RAM when GPU memory is exceeded, avoiding out-of-memory during training. QLORA enables training a 65B parameter LLM on a single 48GB

GPU with no degradation compared to full 16-bit finetuning. QLORA works by freezing the 4-bit quantized base LLM, then backpropagating through it into a small set of 16-bit low-rank adapter weights which are learned.

## 7.4 LOMO

LOW-Memory Optimization (LOMO) (Lv et al., 2023) enables full parameter fine-tuning of LLMs using limited computational resources through a fusion of gradient computation and update. The essence is to fuse gradient computation and parameter update into one step during backpropagation, thereby avoiding storage of full gradient tensors. Firstly, theoretical analysis is provided in LOMO on why SGD can work well for fine-tuning large pre-trained models despite its challenges on smaller models. In addition, LOMO updates each parameter tensor immediately after computing its gradient in backpropagation. Storing the gradient of one parameter at a time reduces gradient memory to  $O(1)$ . LOMO employs gradient value clipping, separate gradient norm computation pass and dynamic loss scaling to stabilize training. The integration of activation checkpointing and ZeRO optimization methods saves memory.

## 7.5 Delta-tuning

Delta-tuning (Ding et al., 2023b) provides optimization and optimal control perspectives for theoretical analysis. Intuitively, delta-tuning performs subspace optimization by restricting tuning to a low-dimensional manifold. The tuned parameters act as optimal controllers guiding model behavior on downstream tasks.

# 8 Evaluation, Analysis and Criticism

## 8.1 Close-ended Evaluations

It is widely accepted among researchers that general-purpose models must demonstrate proficiency in certain core tasks before they can effectively generalize to meet diverse real-world needs. Close-ended evaluations help achieve this objective, often involving multiple-choice questions to assess the performance of LLMs. Below are 6 widely used close-ended evaluations:

(1) **MMLU**. Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2020a) consists of 14079 questions covering 57 tasks including elementary mathematics, US history,



computer science, law, and more. The wide range of subjects and complex questions make MMLU suitable for testing the model’s language comprehension and decision-making capabilities.

**(2) MATH and (3) GSM8K.** MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021) are two distinct mathematical datasets utilized for evaluating different aspects of model capabilities. The MATH (Hendrycks et al., 2021) dataset comprises 12,500 complex competition-level mathematics problems, primarily designed to assess the ability of models to tackle challenging and advanced mathematical questions typically encountered at the college level. Conversely, the GSM8K (Cobbe et al., 2021) dataset contains 8,500 high-quality elementary school math problems, aimed at testing the basic mathematical reasoning abilities of models.

**(4) BBH.** BBH, short for BIG-Bench Hard (Suzgun et al., 2022a), is a subset of the BIG-Bench (Srivastava et al., 2022b) dataset comprising 23 challenging tasks. These tasks were selected because they consistently proved too difficult for current large language models to handle effectively. Requiring complex, multi-step reasoning, the BBH dataset is primarily utilized to assess the general reasoning capabilities of models, testing their ability to navigate and solve intricate problems.

**(5) HumanEval (Coding).** HumanEval (Chen et al., 2021a) consists of 164 programming problems, including language comprehension, algorithms, and simple mathematics, with some comparable to simple software interview questions. The primary purpose of this dataset is to assess the ability of models to generate correct programs based on provided docstrings.

**(6) IFEval.** IFEval (Zhou et al., 2023b) consists of 500 prompts, each containing specific instructions like "write an article with more than 800 words" or "enclose your response in double quotation marks." This dataset is used to test the ability of large language models to accurately follow given instructions.

## 8.2 HELM Evaluation

HELM(Liang et al., 2022) is a holistic evaluation of Language Models (LMs) to improve the transparency of language models, providing a more comprehensive understanding of the

capabilities, risks, and limitations of language models. Specifically, differing from other evaluation methods, HELM holds that a holistic evaluation of language models should focus on the following three factors:

**(1) Broad coverage.** During the development, language models can be adapted to various NLP tasks (e.g., sequence labeling and question answering), thus, the evaluation of language models needs to be carried out in a wide range of scenarios. To involve all potential scenarios, HELM proposed a top-down taxonomy, which begins by compiling all existing tasks in a major NLP conference (ACL2022) into a task space and dividing each task into the form of scenarios (e.g., languages) and metrics (e.g., accuracy). Then when facing a specific task, the taxonomy would select one or more scenarios and metrics in the task space to cover it. By analyzing the structure of each task, HELM clarifies the evaluation content (task scenarios and metrics) and improves the scenario coverage of language models from 17.9% to 96.0%.

**(2) Multi-metric measurement.** In order to enable human to weigh language models from different perspectives, HELM proposes multi-metric measurement. HELM has covered 16 different scenarios and 7 metrics. To ensure the results of intensive multi-metric measurement, HELM measured 98 of 112 possible core scenarios (87.5%).

**(3) Standardization.** The increase in the scale and training complexity of language models has seriously hindered human’s understanding of the structure of each language model. To establish a unified understanding of existing language models, HELM benchmarks 30 well-known language models, covering such institutions as Google (UL2(Tay et al., 2022)), OpenAI (GPT-3(Brown et al., 2020b)), and EleutherAI (GPT-NeoX(Black et al., 2022)). Interestingly, HELM pointed out that LMs such as T5 (Raffel et al., 2019) and Anthropic-LMv4-s3 (Bai et al., 2022a) had not been directly compared in the initial work, while LLMs such as GPT-3 and YaLM were still different from their corresponding reports after multiple evaluations.

## 8.3 LLM As a Judge

LLM as a judge refers to a set of methods that utilize powerful LLMs, particularly GPT-4 (OpenAI, 2023), to evaluate the outputs of other

LLMs. There are three primary reasons for this approach: (1) **Efficiency** – Manually reviewing numerous LLM outputs can be labor-intensive, whereas GPT-4 can evaluate large-scale responses quickly, saving both time and effort; (2) **Reliable Benchmark** – As one of the most advanced models available, GPT-4 provides a dependable benchmark, allowing researchers to compare the performance of different LLMs against a high standard; and (3) **Enhanced Capability** – With improved comprehension and reasoning over previous models, GPT-4 is better suited to analyze subtle aspects of language generation and handle complex outputs from other LLMs. In the following, we detail 4 commonly accepted judge benchmarks:

**(1) AlpacaEval.** AlpacaEval (Li et al., 2023h) is an automated evaluation metric leveraging LLMs, consisting of 805 instructions selected to reflect typical user interactions from the Alpaca web demo<sup>18</sup>. Specifically, for each instruction, both a baseline model  $b$  (currently GPT-4 turbo (OpenAI, 2023)) and the model under evaluation  $m$  generate responses. A GPT-4 turbo-based evaluator then conducts a head-to-head comparison of these responses, determining the probability of favoring the evaluated model. The win rate is calculated as the expected probability that the evaluator prefers the evaluated model’s response across the 805 instructions, serving as a key metric for assessing the performance of the evaluated LM chatbot.

**(2) Length-Controlled AlpacaEval.** Length-Controlled AlpacaEval (Dubois et al., 2024) is a variation of the AlpacaEval (Li et al., 2023h) evaluation metric, designed to minimize length bias, as the original AlpacaEval tends to favor models that produce longer responses. To achieve this goal, Dubois et al. (2024) first fit a generalized linear model to predict the annotator’s (GPT-4’s) preference based on three factors: (1) the instruction, (2) the model used, and (3) the length difference between the baseline and the model’s output. Then, by conditioning the length difference to 0, Dubois et al. (2024) can obtain the length-controlled preference. This idea, which predicts the outcome while conditioning on the length difference (mediator), is a common technique in statistical inference, and by introducing it, Length-

Controlled AlpacaEval increases the Spearman correlation with LMSYS’ Chatbot Arena from 0.94 to 0.98.

**(3) MT-Bench.** Currently, close-ended evaluations only measure LLMs’ core capability on a confined set of tasks, such as MMLU (Hendrycks et al., 2020a) for multi-choice decisions, without adequately assessing its alignment with human preference in open-ended tasks, such as the ability to adhere to instructions in multi-turn dialogues accurately. To alleviate this issue, Zheng et al. (2023) introduced MT-Bench, which comprises 80 high-quality multi-turn questions designed to assess LLMs’ capability in multi-turn conversations and instruction-following, with evaluations conducted using GPT-4. MT-Bench is meticulously crafted to cover eight common tasks: writing, roleplay, extraction, reasoning, math, coding, knowledge I (STEM), and knowledge II (humanities/social sciences). For alignment, GPT-4 achieves over 80% agreement, comparable to the level of agreement among humans, making it a more reliable choice for a public benchmark.

**(4) WildBench.** Although the above evaluations are effective, they have notable limitations in task composition and skill coverage. For example, MT-Bench (Hendrycks et al., 2020a) includes only 80 test instructions, while AlpacaEval (Li et al., 2023h) features many straightforward tasks, such as “What is the capital of Australia?” To address this issue, Lin et al. (2024) introduced WildBench, comprising 1,024 test instructions carefully curated from extensive human-chatbot conversation logs. WildBench draws directly from real-world user interactions, featuring numerous challenging tasks, such as coding and math problem-solving. These tasks frequently demand critical thinking, making WildBench significantly more difficult than other benchmarks. WildBench utilizes two metrics: WB-Reward for pairwise comparisons and WB-Score for individual assessments. Both metrics show strong alignment with human evaluations, with Pearson correlations of 0.98 for WB-Reward and 0.95 for WB-Score when compared to the human-voted ratings.

## 8.4 Low-resource Instruction Tuning

Gupta et al. (2023) attempts to estimate the minimal downstream training data required by SFT models to match the SOTA supervised models over various tasks. Gupta et al. (2023) conducted experiments

<sup>18</sup><https://crfm.stanford.edu/2023/03/13/alpaca.html>

on 119 tasks from Super Natural Instructions (SuperNI) in both single-task learning (STL) and multi-task learning (MTL) settings. The results indicate that in the STL setting, SFT models with only 25% of downstream training data outperform the SOTA models on those tasks, while in the MTL setting, just 6% of downstream training data can lead SFT models to achieve the SOTA performance. These findings suggest that instruction tuning can effectively assist a model in quickly learning a task even with limited data.

However, due to resource limitations, Gupta et al. (2023) did not conduct experiments on LLMs, like T5-11B. So, to gain a more comprehensive understanding of the SFT models, further investigation using larger language models and datasets is necessary.

### 8.5 Smaller Instruction Dataset

SFT requires a substantial amount of specialized instruction data for training. Zhou et al. (2023a) hypothesized that the pre-trained LLM only has to learn the style or format to interact with users and proposed LIMA that achieves strong performance by fine-tuning an LLM on only 1,000 carefully selected training examples.

Specifically, LIMA first manually curates 1,000 demonstrations with high-quality prompts and responses. Then the 1,000 demonstrations are used to fine-tune the pre-trained 65B-parameter LLaMa (Touvron et al., 2023b). By comparison, across more than 300 challenging tasks, LIMA outperforms GPT-davinci003 (Brown et al., 2020b), which was fine-tuned on 5,200 examples by human feedback tuning. Moreover, with only half amount of demonstrations, LIMA achieves equivalent results to GPT-4 (OpenAI, 2023), Claude (Bai et al., 2022b), and Bard<sup>19</sup>. Above all, LIMA demonstrated that LLMs’ powerful knowledge and capabilities can be exposed to users with only a few carefully curated instructions to fine-tune.

### 8.6 Evaluating Instruction Tuning Datasets

The performance of SFT model highly depends on the SFT datasets. However, there lacks of evaluations for these SFT datasets from open-ended and subjective aspects.

To address this issue, Wang et al. (2023e) performs dataset evaluation by fine-tuning the

LLaMa model (Touvron et al., 2023b) on various of open SFT datasets and measure different fine-tuned models through both automatic and human evaluations. An additional model is trained on the combination of SFT datasets. For the results, Wang et al. (2023e) showed that there is not a single best SFT dataset across all tasks, while by manually combining datasets it can achieve the best overall performance. Besides, Wang et al. (2023e) pointed out that though SFT can bring large benefits on LLMs at all sizes, smaller models and models with a high base quality benefit most from SFT. For human evaluations, Wang et al. (2023e) a larger model is more likely to gain a higher acceptability score.

### 8.7 Proprietary LLMs Imitation

LLMs imitation is an approach that collects outputs from a stronger model, such as a proprietary system like ChatGPT, and uses these outputs to fine-tune an open-source LLM. Through this way, an open-source LLM may get competitive capabilities with any proprietary model.

Gudibande et al. (2023) conducted several experiments to critically analyze the efficacy of model imitation. Specifically, Gudibande et al. (2023) first collected datasets from outputs of ChatGPT over broad tasks. Then these datasets were used to fine-tune a range of models covering sizes from 1.5B to 13B, base models GPT-2 and LLaMA, and data amounts from 0.3M tokens to 150M tokens.

For evaluations, Gudibande et al. (2023) demonstrated that on tasks with supported datasets, imitation models are far better than before, and their outputs appear similar to ChatGPT’s. While on tasks without imitation datasets, imitation models do not have improvement or even decline in accuracy.

Thus, Gudibande et al. (2023) pointed out that it’s the phenomenon that imitation models are adept at mimicking ChatGPT’s style (e.g., being fluent, confident and well-structured) that makes researchers have the illusion about general abilities of imitation models. So, Gudibande et al. (2023) suggested that instead of imitating proprietary models, researchers had better focus on improving the quality of base models and instruction examples.

<sup>19</sup>Bard, designed by Google, is an interface to generative AI platform, and the link is: <https://ai.google/static/documents/google-about-bard.pdf>

## 9 The Role of Instruction Fine-tuning

Instruction fine-tuning (IF), also known as supervised fine-tuning (SFT), is a conventional alignment approach that trains models on example prompts paired with corresponding responses to ensure the model’s output aligns with user instructions and intended goals. More recently, some reinforcement learning (RL) based methods (Wang et al., 2024), e.g., reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022), direct preference optimization (DPO) (Rafailov et al., 2023), and group relative policy optimization (GRPO) (Shao et al., 2024), and various prompt engineering strategies have emerged as alternatives or complements. Thus, in this section, we will review each method’s role in aligning LLMs, and examine whether SFT remains necessary. Further more, we also consider the risk of superficial alignment, i.e. alignment that changes only the model’s surface behavior (tone, style) without imparting deeper understanding.

### 9.1 SFT Compared with Other Alignment Methods

Below, we begin by outlining three widely used alignment approaches, RLHF, DPO, and prompt engineering, highlighting their strengths and weaknesses. Then, we explain why SFT remains an essential component of contemporary alignment pipelines.

#### 9.1.1 Reinforcement Learning from Human Feedback (RLHF)

RLHF is the dominant alignment paradigm, and typically proceeds in three phases: (1) supervised fine-tuning (SFT) on human large amounts of instruction-answer pairs, (2) training a reward model on human-ranked responses, and (3) using policy optimization (e.g. PPO (?)) to maximize the reward model’s feedback (Chen et al., 2025b). This pipeline can deeply adjust model behavior to complex user preferences. RLHF has enabled remarkable capabilities (e.g. nuanced help, factuality), but at high cost and complexity. It requires extensive human data, careful RL tuning, and often suffers stability issues and “reward hacking” (the model finds loopholes in the reward model) (Xiao et al., 2024; Wang et al., 2024). Because RLHF optimization is resource-intensive and sensitive to hyper parameters, simpler alternatives have been sought.

**Advantages.** The key strength of RLHF lies in its ability to guide models toward high-level objectives, such as helpfulness and safety, that are not explicitly encoded in the training data, demonstrating strong empirical performance in aligning models with user intent when well-tuned (Wang et al., 2024; Chen et al., 2025b).

**Limitations.** RLHF’s complexity is a downside. It typically requires starting from an SFT-trained model, i.e., a model that already follows instructions to some degree, because training RL from a raw base model is difficult (Trivedi et al., 2025). The multi-stage pipeline (SFT, reward model, and PPO) is time-consuming and brittle (Chen et al., 2025b). In practice, researchers and practitioners often still perform an initial instruction fine-tuning, even when using RLHF, to establish a reasonable base policy. Moreover, RLHF can introduce “alignment tax” (performance drop on some tasks) and can fail to generalize if the reward model is mis-specified (Xiao et al., 2024).

#### 9.1.2 Direct Preference Optimization (DPO)

Direct Preference Optimization (DPO) (Rafailov et al., 2023) is a recently proposed RL-free alignment method that directly fine-tunes on preference pairs. Instead of learning a separate reward model and running RL, DPO casts alignment as a supervised objective: for each prompt and pair of outputs (preferred vs. dispreferred), it adjusts the model’s logits to increase the probability of the preferred output. DPO’s loss is equivalent to a Bradley–Terry pairwise classification (a logit-ratio objective) to bypass policy-gradient RL entirely.

**Advantages.** Because DPO fits into a standard maximum-likelihood fine-tuning framework, it is far simpler and more stable than PPO-based RLHF (Xu et al., 2024a; Xiao et al., 2024; Wang et al., 2024). Studies report that DPO matches or exceeds RLHF performance on tasks like summarization or helpfulness with fewer preference examples. Compared to RLHF, DPO has been shown to be stable, performant, and computationally lightweight in various applications. It does not require expensive RL infrastructure or tuning of PPO hyper-parameters, making it reproducible and easier to deploy. Practitioners, e.g. OpenAI <sup>20</sup>

<sup>20</sup>[https://cookbook.openai.com/examples/fine\\_tuning\\_direct\\_preference\\_optimization\\_guide](https://cookbook.openai.com/examples/fine_tuning_direct_preference_optimization_guide)



and DeepSeek (Shao et al., 2024) note that starting DPO from an already fine-tuned model improves results, using SFT to establish a robust initial policy.

**Limitations.** DPO still depends on quality preference data, and like any offline method it can suffer if the data distribution shifts. Recent analyses have identified issues with DPO: because it employs an implicit reward tied to the policy, it can bias the model toward out-of-distribution outputs and even degrade generalization (Xiao et al., 2024; Wang et al., 2024). Variants, e.g. KL-constrained or semantics-aware DPO, are being developed to mitigate these problems, but it remains true that DPO typically benefits from starting with a good initial model. In practice, most implementations still use an SFT-tuned model before running DPO, echoing the RLHF pipeline. Thus, while DPO simplifies alignment, it has not eliminated the need for supervised tuning in many cases.

### 9.1.3 Prompt Engineering (In-Context Learning)

Prompt engineering aligns model behavior without any fine-tuning. Instead, it leverages the model’s existing capabilities by crafting prompts, including instructions, few-shot examples, or chain-of-thought cues, to elicit desired outputs. Recent work treats prompt design itself as an optimization problem: one can optimize a prompt string or learn soft prompts to maximize human-aligned metrics (Trivedi et al., 2025). Importantly, prompt-based methods assume no weight updates, which means that they work without any post training.

**Advantages.** The biggest benefit is that no retraining is required. Prompt optimization can effectively align LLMs even when parameter fine-tuning is not feasible (Trivedi et al., 2025). This is appealing for large models with fixed parameters (APIs or frozen on-device models). Some experiments have shown that with well-designed prompts, base LLMs can achieve high performance on instruction-following tasks. For example, Lin et al. (2023b) introduced URIAL, a method that uses only a few stylistic examples in-context (plus a system prompt) to steer the model, and found it matched or even surpassed fully-tuned models in many benchmarks. This suggests that clever prompt engineering alone can yield strong alignment in some cases (Wang et al.,

2023c; Sun et al., 2023d,b; Wang et al., 2023b; Sun et al., 2023c). Prompting has obvious speed and convenience advantages: it requires no training data or compute, and can be iterated quickly by users.

**Limitations.** Prompt-based alignment has inherent limits. The model’s context window bounds how much instruction or example content can be provided so that very complex tasks may simply not fit. More importantly, prompt methods generally induce superficial compliance rather than truly altering the model’s knowledge. They leverage already-encoded patterns in the model, but cannot add new capabilities or correct deep misunderstandings. In practice, prompt engineering often produces brittle behaviors: slight rephrasing can break performance, and malicious users can “jailbreak” around prompts to elicit bad outputs. For instance, Chen et al. (2025a) note that although prompt-based ICL can align a model to some extent, it does so mainly by inserting stylistic cues and does not fundamentally change the model’s reasoning process. In short, prompt engineering can quickly achieve surface-level alignment (tone, disclaimers, formatting), but cannot replace weight tuning for deep or novel tasks, and for complex reasoning, mathematics, or new knowledge integration typically require additional fine-tuning.

### 9.1.4 The Continued Necessity of SFT

Given these techniques, a key question is whether instruction fine-tuning (IF) or supervised fine-tuning (SFT) remains necessary in modern pipelines. Empirical evidence suggests it does. Both RLHF and DPO pipelines almost universally incorporate an initial SFT stage. In RLHF this is explicit, which SFT usually serves as the first phase. For DPO, while the final optimization is simpler, practitioners generally first fine-tune on good example responses to establish a robust initial policy, which stabilizes subsequent DPO refinement. The OpenAI alignment guide explicitly recommends performing supervised fine-tuning on a subset of preferred responses before DPO to improve alignment and convergence<sup>21</sup>. In other words, even with DPO, a round of instruction tuning yields better outcomes. On the other hand,

<sup>21</sup>[https://cookbook.openai.com/examples/fine\\_tuning\\_direct\\_preference\\_optimization\\_guide](https://cookbook.openai.com/examples/fine_tuning_direct_preference_optimization_guide)

prompt-based methods show that in principle one can align some models without any fine-tuning (Wang et al., 2023c; Sun et al., 2023d,b; Wang et al., 2023b; Sun et al., 2023c). Other work (Lin et al., 2023a) demonstrates that you can in effect “unlock” base LLMs by providing few-shot prompts, achieving performance close to tuned models without SFT. However, these tuning-free methods tend to rely on pre-existing capabilities. If a base model genuinely lacks a skill, such as solving a type of math problem it was never pre-trained on, no prompt will fix it, and only updating weights can. Moreover, some recent studies (Parthasarathy et al., 2024) find that for reasoning and knowledge tasks, performance continues to scale with more fine-tuning data, suggesting base models do improve under instruction tuning. In summary, prompt methods can sometimes obviate SFT for shallow alignment, but for robust alignment pipelines or domain specific alignment (e.g., Medicinal Chemistry), supervised fine-tuning is still regarded as essential groundwork. New research even explores hybrid tricks, such as “instruction residuals” from an older model added to a new base, to avoid re-training, but these rely on existing tuned models as sources. The prevailing practice remains: use SFT to teach the model the format and style of responses, then refine preferences via RLHF or DPO.

## 9.2 Superficial Alignment

Despite the impressive improvements in the performance of instruction tuning, there lacks clarity about the specific knowledge that models acquire through instruction tuning, raising questions about: *Does instruction tuning just learn Pattern Copying?* or *How exactly does the alignment tuning transform a base LLM?*

To answer these questions, Kung and Peng (2023) delves into the analysis of how models make use of instructions during SFT by comparing the tuning when provided with altered instructions versus the original instructions.

Specifically, Kung and Peng (2023) creates simplified task definitions that remove all semantic components, leaving only the output information. In addition, Kung and Peng (2023) also incorporates delusive examples that contain incorrect input-output mapping. Surprisingly, the experiments show that models trained on these simplified task definitions or delusive examples

can achieve comparable performance to the ones trained on the original instructions and examples. Moreover, the paper also introduces a baseline for the classification task with zero-shot, which achieves similar performance to SFT in low-resource settings.

Similar to the findings of Kung and Peng (2023), several subsequent studies (Zhou et al., 2023a; Lin et al., 2023a) reached the same conclusion: the observed performance improvements in current SFT models are often due to superficial alignment. This means the models excel at recognizing superficial alignment, such as mastering output formats and making educated guesses, rather than truly understanding and learning the underlying tasks.

## 10 Conclusion

This work surveys recent advances in the fast growing field of instruction tuning, which can also be referred to as supervised fine-tuning (SFT). We make a systematic review of the literature, including the general methodology of SFT, the construction of SFT datasets, the training of SFT models, SFT’s applications to different modalities, domains and application. We also review analysis on SFT models to discover both their advantages and potential pitfalls. We hope this work will act as a stimulus to motivate further endeavors to address the deficiencies of current SFT models.

## References

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2023. Evaluating correctness and faithfulness of instruction-following models for question answering. *ArXiv*, abs/2307.16877.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. Flamingo: a visual language model for few-shot learning. In *NeurIPS*.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023a.

- Falcon-40B: an open large language model with state-of-the-art performance.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, et al. 2023b. Falcon-40b: an open large language model with state-of-the-art performance.
- Anas Awadalla, Irena Gao, Joshua Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Jenia Jitsev, et al. 2023. Openflamingo.
- Stephen H. Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesh Sharma, Taewoon Kim, M Saiful Bari, Thibault Févry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Alan Fries, Maged S. Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Mike Tian-Jian Jiang, and Alexander M. Rush. 2022. Promptsources: An integrated development environment and repository for natural language prompts. *ArXiv*, abs/2202.01279.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
- Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. 2021. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*.
- Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. 2022. Text2live: Text-driven layered image and video editing. In *European Conference on Computer Vision*, pages 707–723. Springer.
- Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. The pushshift reddit dataset. In *International Conference on Web and Social Media*.
- Edward Beeching, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard. *Hugging Face*.
- Stella Rose Biderman, Hailey Schoelkopf, Quentin G. Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling. *ArXiv*, abs/2304.01373.
- Sid Black, Stella Rose Biderman, Eric Hallahan, Quentin G. Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Martin Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Benqi Wang, and Samuel Weinbach. 2022. Gpt-neox-20b: An open-source autoregressive language model. *ArXiv*, abs/2204.06745.
- Tim Brooks, Aleksander Holynski, and Alexei A. Efros. 2022. Instructpix2pix: Learning to follow image editing instructions. *ArXiv*, abs/2211.09800.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020a. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. Language models are few-shot learners. *ArXiv*, abs/2005.14165.
- Tuhin Chakrabarty, Vishakh Padmakumar, and Hengxing He. 2022. Help me write a poem - instruction tuning as a vehicle for collaborative poetry writing. *ArXiv*, abs/2210.13669.
- Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation.
- Banghao Chen, Zhaofeng Zhang, Nicolas Langrené, and Shengxin Zhu. 2025a. Unleashing the potential of prompt engineering for large language models. *Patterns*.
- Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. 2024a. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*.

- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023a. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021a. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde, Jared Kaplan, Harrison Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, David W. Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Babuschkin, S. Arun Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew M. Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021b. Evaluating large language models trained on code. *ArXiv*, abs/2107.03374.
- Qianglong Chen, Guohai Xu, Mingshi Yan, Ji Zhang, Fei Huang, Luo Si, and Yin Zhang. 2023b. Distinguish before answer: Generating contrastive explanation as knowledge for commonsense question answering. In *Annual Meeting of the Association for Computational Linguistics*.
- Runjin Chen, Gabriel Jacob Perin, Xuxi Chen, Xilun Chen, Yan Han, Nina ST Hirata, Junyuan Hong, and Bhavya Kailkhura. 2025b. Extracting and understanding the superficial knowledge in alignment. *arXiv preprint arXiv:2502.04602*.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024b. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. *ArXiv*, abs/2204.02311.
- Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *ArXiv*, abs/1905.10044.
- J. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *ArXiv*, abs/1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *ArXiv*, abs/2110.14168.
- OpenAccess AI Collective. 2023. *software: huggingface.co/openaccess-ai-collective/minotaur-15b*.
- Mike Conover, Matt Hayes, Ankit Mathur, Xiangrui Meng, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, et al. 2023a. Free



- dolly: Introducing the world’s first truly open instruction-tuned llm.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023b. Free dolly: Introducing the world’s first truly open instruction-tuned llm.
- Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. 2018. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *arXiv preprint arXiv:1805.10190*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *ArXiv*, abs/2305.06500.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023a. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*.
- Ning Ding, Yujia Qin, Guang Yang, Fu Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao, Xiaozhi Wang, Zhiyuan Liu, Haitao Zheng, Jianfei Chen, Y. Liu, Jie Tang, Juanzi Li, and Maosong Sun. 2023b. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5:220–235.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.
- Jon Durbin. 2023. Airoboros. *software: github.com/jondurbin/airoboros*.
- Jane Dwivedi-Yu, Timo Schick, Zhengbao Jiang, Maria Lomeli, Patrick Lewis, Gautier Izacard, Edouard Grave, Sebastian Riedel, and Fabio Petroni. 2022. Editeval: An instruction-based benchmark for text improvements.
- William Fedus, Barret Zoph, and Noam M. Shazeer. 2021. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23:120:1–120:39.
- Jun Gao, Huan Zhao, Changlong Yu, and Ruifeng Xu. 2023a. Exploring the feasibility of chatgpt for event extraction. *ArXiv*, abs/2303.03836.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. 2021. A framework for few-shot language model evaluation. *Version v0. 0.1. Sept*.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023b. Enabling large language models to generate text with citations. *arXiv preprint arXiv:2305.14627*.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *ArXiv*, abs/2009.11462.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embedding space to bind them all. In *CVPR*.
- Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qianmengke Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. Multimodal-gpt: A vision and language model for dialogue with humans. *ArXiv*, abs/2305.04790.
- Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Himanshu Gupta, Saurabh Arjun Sawant, Swaroop Mishra, Mutsumi Nakamura, Arindam Mitra, Santosh Mashetty, and Chitta Baral. 2023.

- Instruction tuned models are quick learners. *arXiv preprint arXiv:2306.05539*.
- Prakhar Gupta, Cathy Jiao, Yi-Ting Yeh, Shikib Mehri, Maxine Eskénazi, and Jeffrey P. Bigham. 2022. Instructdial: Improving zero and few-shot generalization in dialogue through instruction tuning. In *Conference on Empirical Methods in Natural Language Processing*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020a. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. 2020b. Measuring massive multitask language understanding. *ArXiv*, abs/2009.03300.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. Unnatural instructions: Tuning language models with (almost) no human labor. *arXiv preprint arXiv:2212.09689*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *arXiv preprint arXiv:2305.08322*.
- Anel Islamovic. Stability AI Launches the First of its StableLM Suite of Language Models — Stability AI — stability.ai. <https://stability.ai/blog/stability-ai-launches-the-first-of-its-stablelm-suite-of-language-models> [Accessed 09-Jun-2023].
- Hamish Ivison, Akshita Bhagia, Yizhong Wang, Hannaneh Hajishirzi, and Matthew E. Peters. 2022. Hint: Hypernetwork instruction tuning for efficient zero-shot generalisation. *ArXiv*, abs/2212.10315.
- Srinivas Iyer, Xiaojuan Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, Xian Li, Brian O’Horo, Gabriel Pereyra, Jeff Wang, Christopher Dewan, Asli Celikyilmaz, Luke Zettlemoyer, and Veselin Stoyanov. 2022. Opt-impl: Scaling language model instruction meta learning through the lens of generalization. *ArXiv*, abs/2212.12017.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- JosephusCheung. 2021. Guanaco: Generative universal assistant for natural-language adaptive context-aware omnilingual outputs.
- Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi. 2020. Unifiedqa: Crossing format boundaries with a single qa system. *arXiv preprint arXiv:2005.00700*.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, et al. 2023. Openassistant conversations—democratizing large language model alignment. *arXiv preprint arXiv:2304.07327*.
- Po-Nien Kung and Nanyun Peng. 2023. Do models really learn to follow instructions? an empirical study of instruction tuning. *ArXiv*, abs/2305.11383.
- LAION.ai. 2023. Oig: the open instruction generalist dataset.
- Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. 2018. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10.
- Mina Lee, Percy Liang, and Qian Yang. 2022. Coauthor: Designing a human-ai collaborative writing dataset for exploring language model capabilities. *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Conference on Empirical Methods in Natural Language Processing*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer.

2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Annual Meeting of the Association for Computational Linguistics*.
- Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei Ye, Wen Zhao, and Shikun Zhang. 2023a. Evaluating chatgpt’s information extraction capabilities: An assessment of performance, explainability, calibration, and faithfulness. *ArXiv*, abs/2304.11633.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023b. Otter: A multi-modal model with in-context instruction tuning. *ArXiv*, abs/2305.03726.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023c. Camel: Communicative agents for "mind" exploration of large scale language model society.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023d. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*.
- Kunchang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023e. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023f. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston, and Mike Lewis. 2023g. Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259*.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023h. AlpacaEval: An automatic evaluator of instruction-following models. *GitHub repository*.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023i. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*.
- Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, and You Zhang. 2023j. Chatdoctor: A medical chat model fine-tuned on llama model using medical domain knowledge. *ArXiv*, abs/2303.14070.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher R’e, Diana Acosta-Navas, Drew A. Hudson, E. Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan S. Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas F. Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2022. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step.
- Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. 2024. Wildbench: Benchmarking llms with challenging tasks from real users in the wild. *arXiv preprint arXiv:2406.04770*.
- Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2023a. The unlocking spell on base llms: Rethinking alignment via in-context learning.
- Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2023b. Uriel: Tuning-free instruction learning and alignment for untuned llms.
- Bill Yuchen Lin, Kangmin Tan, Chris Miller, Beiwen Tian, and Xiang Ren. 2022. Unsupervised cross-task generalization via retrieval augmentation. *ArXiv*, abs/2204.07937.
- Stephanie C. Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. In *Annual Meeting of the Association for Computational Linguistics*.
- Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023c. Pmc-clip: Contrastive language-image pre-training using biomedical documents. *arXiv preprint arXiv:2303.07240*.
- Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. 2021a. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1650–1654. IEEE.
- Hanmeng Liu, Zhiyang Teng, Leyang Cui, Chaoli Zhang, Qiji Zhou, and Yue Zhang. 2023a. Logicot: Logical chain-of-thought instruction-tuning data collection with gpt-4. *ArXiv*, abs/2305.12147.

- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. *ArXiv*, abs/2304.08485.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021b. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*.
- Tiedong Liu and Bryan Kian Hsiang Low. 2023. Goat: Fine-tuned llama outperforms gpt-4 on arithmetic tasks. *arXiv preprint arXiv:2305.14201*.
- Zheng Liu, Aoxiao Zhong, Yiwei Li, Longtao Yang, Chao Ju, Zihao Wu, Chong Ma, Peng Shu, Cheng Chen, Sekeun Kim, Haixing Dai, Lin Zhao, Dajiang Zhu, Jun Liu, Wei Liu, Dinggang Shen, Xiang Li, Quanzheng Li, and Tianming Liu. 2023c. Radiology-gpt: A large language model for radiology.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. *arXiv preprint arXiv:2301.13688*.
- Zimu Lu, Aojun Zhou, Houxing Ren, Ke Wang, Weikang Shi, Juntong Pan, Mingjie Zhan, and Hongsheng Li. 2024. Mathgenie: Generating synthetic data with question back-translation for enhancing mathematical reasoning of llms. *arXiv preprint arXiv:2402.16352*.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct.
- Kai Lv, Yuqing Yang, Tengxiao Liu, Qi jie Gao, Qipeng Guo, and Xipeng Qiu. 2023. Full parameter fine-tuning for large language models with limited resources.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2022. SDEdit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.
- Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Codas, Clarisse Simoes, Sahaj Agarwal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, et al. 2023. Orca 2: Teaching small language models how to reason. *arXiv preprint arXiv:2311.11045*.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*.
- Munan Ning, Yujia Xie, Dongdong Chen, Zeyin Song, Lu Yuan, Yonghong Tian, Qixiang Ye, and Liuliang Yuan. 2023. Album storytelling with iterative story-aware captioning and large language models. *ArXiv*, abs/2305.12943.
- NousResearch. 2023. *software: huggingface.co/NousResearch/Nous-Hermes-13b*.
- OpenAI. 2022. Introducing chatgpt. *Blog post openai.com/blog/chatgpt*.
- OpenAI. 2023. Gpt-4 technical report. *ArXiv*, abs/2303.08774.
- 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. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Arnold Overwijk, Chenyan Xiong, Xiao Liu, Cameron VandenBerg, and Jamie Callan. 2022. Clueweb22: 10 billion web documents with visual and semantic information. *arXiv preprint arXiv:2211.15848*.
- Venkatesh Balavadhani Parthasarathy, Ahtsham Zafar, Aafaq Khan, and Arsalan Shahid. 2024. The ultimate guide to fine-tuning llms from basics to breakthroughs: An exhaustive review of technologies, research, best practices, applied research challenges and opportunities. *arXiv preprint arXiv:2408.13296*.
- Keiran Paster, Marco Dos Santos, Zhangir Azerbayev, and Jimmy Ba. 2023. Openwebmath: An open dataset of high-quality mathematical web text. *arXiv preprint arXiv:2310.06786*.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- Jing Qian, Li Dong, Yelong Shen, Furu Wei, and Weizhu Chen. 2022. Controllable natural language generation with contrastive prefixes. *arXiv preprint arXiv:2202.13257*.



- Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Hector Liu, Yuanzhi Li, et al. 2024. O1 replication journey: A strategic progress report–part 1. *arXiv preprint arXiv:2410.18982*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *ArXiv*, abs/1910.10683.
- Vipul Raheja, Dhruv Kumar, Ryan Koo, and Dongyeop Kang. 2023. Coedit: Text editing by task-specific instruction tuning. *ArXiv*, abs/2305.09857.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695.
- Andrew Rosenbaum, Saleh Soltan, Wael Hamza, Yannick Versley, and Markus Boese. 2022. Linguist: Language model instruction tuning to generate annotated utterances for intent classification and slot tagging. In *International Conference on Computational Linguistics*.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elizabeth-Jane Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagn'è, Alexandra Sasha Luccioni, Francois Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Rose Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurenceon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa Etxabe, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris C. Emezue, Christopher Klamm, Colin Leong, Daniel Alexander van Strien, David Ifeoluwa Adelani, Dragomir R. Radev, Eduardo González Ponferrada, Efrat Levkovich, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady ElSahar, Hamza Benyamina, Hieu Trung Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jorg Froberg, Josephine L. Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro von Werra, Leon Weber, Long Phan, Loubna Ben Allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, Mar'ia Grandury, Mario vSavsko, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad Ali Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla A. Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, S. Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal V. Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Févry, Trishala Neeraj, Urmish Thakker,

- Vikas Raunak, Xiang Tang, Zheng Xin Yong, Zhiqing Sun, Shaked Brody, Y Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre Francoi Lavall'ee, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Reuena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aur'elie N'ev'eol, Charles Lovering, Daniel H Garrette, Deepak R. Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Xiangru Tang, Junjo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, S. Osher Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdenek Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ananda Santa Rosa Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behrooz, Benjamin Olusola Ajibade, Bharat Kumar Saxena, Carlos Muñoz Ferrandis, Danish Contractor, David M. Lansky, Davis David, Douwe Kiela, Duong Anh Nguyen, Edward Tan, Emily Baylor, Ezinwanne Ozoani, Fatim T Mirza, Frankline Ononiwu, Habib Rezanejad, H.A. Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jan Passmore, Joshua Seltzer, Julio Bonis Sanz, Karen Fort, Livia Macedo Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, M. K. K. Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nourhan Fahmy, Olanrewaju Samuel, Ran An, R. P. Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas L. Wang, Sourav Roy, Sylvain Viguiere, Thanh-Cong Le, Tobi Oyeade, Trieu Nguyen Hai Le, Yoyo Yang, Zachary Kyle Nguyen, Abhinav Ramesh Kashyap, A. Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Kumar Singh, Benjamin Beilharz, Bo Wang, Caio Matheus Fonseca de Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourier, Daniel Le'on Perin'an, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Iman I.B. Bello, Isha Dash, Ji Soo Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthi Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, María Andrea Castillo, Marianna Nezhurina, Mario Sanger, Matthias Samwald, Michael Cullan, Michael Weinberg, M Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patricia Haller, R. Chandrasekhar, R. Eisenberg, Robert Martin, Rodrigo L. Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aaroonsiri, Srishti Kumar, Stefan Schweter, Sushil Pratap Bharati, T. A. Laud, Th'eo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yashasvi Bajaj, Y. Venkatraman, Yifan Xu, Ying Xu, Yun chao Xu, Zhee Xao Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. 2022. Bloom: A 176b-parameter open-access multilingual language model. *ArXiv*, abs/2211.05100.
- Timo Schick and Hinrich Schütze. 2021. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2022. Language models are multilingual chain-of-thought reasoners. *ArXiv*, abs/2210.03057.
- Saleh Soltan, Shankar Ananthakrishnan, Jack Fitzgerald, Rahul Gupta, Wael Hamza, Haidar Khan, Charith Peris, Stephen Rawls, Andy Rosenbaum, Anna Rumshisky, et al. 2022. Alexatm 20b: Few-shot learning using a large-scale multilingual seq2seq model. *arXiv preprint arXiv:2208.01448*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022a. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022b. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.

- Weiwei Sun, Hengyi Cai, Hongshen Chen, Pengjie Ren, Zhumin Chen, Maarten de Rijke, and Zhaochun Ren. 2023a. Answering ambiguous questions via iterative prompting. *ArXiv*, abs/2307.03897.
- Xiaofei Sun, Linfeng Dong, Xiaoya Li, Zhen Wan, Shuhe Wang, Tianwei Zhang, Jiwei Li, Fei Cheng, Lingjuan Lyu, Fei Wu, et al. 2023b. Pushing the limits of chatgpt on nlp tasks. *arXiv preprint arXiv:2306.09719*.
- Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023c. Text classification via large language models. *arXiv preprint arXiv:2305.08377*.
- Xiaofei Sun, Xiaoya Li, Shengyu Zhang, Shuhe Wang, Fei Wu, Jiwei Li, Tianwei Zhang, and Guoyin Wang. 2023d. Sentiment analysis through llm negotiations. *arXiv preprint arXiv:2311.01876*.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022a. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.
- Mirac Suzgun, Nathan Scales, Nathanael Scharli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed Huai hsin Chi, Denny Zhou, and Jason Wei. 2022b. Challenging big-bench tasks and whether chain-of-thought can solve them. *ArXiv*, abs/2210.09261.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023a. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023b. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca).
- Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Dara Bahri, Tal Schuster, Steven Zheng, et al. 2022. UI2: Unifying language learning paradigms.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lambda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.
- Sun Tianxiang and Qiu Xipeng. 2023. Moss. *Blog post txsun1997.github.io/blogs/moss.html*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023b. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Prashant Trivedi, Souradip Chakraborty, Avinash Reddy, Vaneet Aggarwal, Amrit Singh Bedi, and George K Atia. 2025. Align-pro: A principled approach to prompt optimization for llm alignment. 39(26):27653–27661.
- Siddharth Varia, Shuai Wang, Kishaloy Halder, Robert Vacareanu, Miguel Ballesteros, Yassine Benajiba, Neha Ann John, Rishita Anubhai, Smaranda Muresan, and Dan Roth. 2022. Instruction tuning for few-shot aspect-based sentiment analysis. *ArXiv*, abs/2210.06629.
- Zhen Wan, Fei Cheng, Zhuoyuan Mao, Qianying Liu, Haiyue Song, Jiwei Li, and Sadao Kurohashi. 2023. Gpt-re: In-context learning for relation extraction using large language models. *arXiv preprint arXiv:2305.02105*.
- Haochun Wang, Chi Liu, Sendong Zhao, Bing Qin, and Ting Liu. 2023a. Chatglm-med. <https://github.com/SCIR-HI/Med-ChatGLM>.
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022a. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR.
- Shuhe Wang, Beiming Cao, Shengyu Zhang, Xiaoya Li, Jiwei Li, Fei Wu, Guoyin Wang, and Eduard Hovy. 2023b. Sim-gpt: Text similarity via gpt annotated data. *arXiv preprint arXiv:2312.05603*.
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023c. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*.
- Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu, Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu, Guoyin Wang, and Eduard Hovy. 2024. Reinforcement learning enhanced llms: A survey. *arXiv preprint arXiv:2412.10400*.
- Xiao Wang, Wei Zhou, Can Zu, Han Xia, Tianze Chen, Yuan Zhang, Rui Zheng, Junjie Ye, Qi Zhang, Tao Gui, Jihua Kang, J. Yang, Siyuan Li, and Chunsai Du. 2023d. Instructuie: Multi-task instruction

- tuning for unified information extraction. *ArXiv*, abs/2304.08085.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. 2022b. Self-consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hanna Hajishirzi. 2023e. How far can camels go? exploring the state of instruction tuning on open resources. *ArXiv*, abs/2306.04751.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022c. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Maitreya Patel, Kuntal Kumar Pal, M. Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Shailaja Keyur Sampat, Savan Doshi, Siddharth Deepak Mishra, Sujun Reddy, Sumanta Patro, Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hanna Hajishirzi, and Daniel Khashabi. 2022d. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Conference on Empirical Methods in Natural Language Processing*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 2022e. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *arXiv preprint arXiv:2204.07705*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 2022f. Super-naturalinstructions:generalization via declarative instructions on 1600+ tasks. In *EMNLP*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903.
- Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang, Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu, Yufeng Chen, Meishan Zhang, et al. 2023a. Zero-shot information extraction via chatting with chatgpt. *arXiv preprint arXiv:2302.10205*.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023b. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*.
- Sarah Wiegrefe, Jack Hessel, Swabha Swayamdipta, Mark Riedl, and Yejin Choi. 2021. Reframing human-ai collaboration for generating free-text explanations. *arXiv preprint arXiv:2112.08674*.
- Wenyi Xiao, Zechuan Wang, Leilei Gan, Shuai Zhao, Zongrui Li, Ruirui Lei, Wanggui He, Luu Anh Tuan, Long Chen, Hao Jiang, et al. 2024. A comprehensive survey of direct preference optimization: Datasets, theories, variants, and applications. *arXiv preprint arXiv:2410.15595*.
- Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga, Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I. Wang, Victor Zhong, Bailin Wang, Chengzu Li, Connor Boyle, Ansong Ni, Ziyu Yao, Dragomir R. Radev, Caiming Xiong, Lingpeng Kong, Rui Zhang, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2022. Unifiedskg: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. In *Conference on Empirical Methods in Natural Language Processing*.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. Wizardlm: Empowering large language models to follow complex instructions.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023b. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023c. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *ArXiv*, abs/2304.01196.
- Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. 2024a. Is dpo superior to ppo for llm alignment? a comprehensive study. *arXiv preprint arXiv:2404.10719*.
- Weijia Xu, Batool Haider, and Saab Mansour. 2020. End-to-end slot alignment and recognition for cross-lingual nlu. *arXiv preprint arXiv:2004.14353*.
- Zhiyang Xu, Chao Feng, Rulin Shao, Trevor Ashby, Ying Shen, Di Jin, Yu Cheng, Qifan Wang, and Lifu Huang. 2024b. Vision-flan: Scaling human-labeled tasks in visual instruction tuning. *arXiv preprint arXiv:2402.11690*.
- Zhiyang Xu, Ying Shen, and Lifu Huang. 2022. Multiinstruct: Improving multi-modal zero-shot learning via instruction tuning. *ArXiv*, abs/2212.10773.



- Fuzhao Xue, Kabir Jain, Mahir Hitesh Shah, Zangwei Zheng, and Yang You. 2023. Instruction in the wild: A user-based instruction dataset. <https://github.com/XueFuzhao/InstructionWild>.
- Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023a. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *arXiv preprint arXiv:2304.13712*.
- Kevin Yang, Nanyun Peng, Yuandong Tian, and Dan Klein. 2022a. Re3: Generating longer stories with recursive reprompting and revision. *arXiv preprint arXiv:2210.06774*.
- Kexin Yang, Dayiheng Liu, Wenqiang Lei, Baosong Yang, Mingfeng Xue, Boxing Chen, and Jun Xie. 2022b. Tailor: A prompt-based approach to attribute-based controlled text generation. *ArXiv*, abs/2204.13362.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023b. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *ArXiv*, abs/2305.10601.
- Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Lu Sheng, Lei Bai, Xiaoshui Huang, Zhiyong Wang, Wanli Ouyang, and Jing Shao. 2023. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. *ArXiv*, abs/2306.06687.
- Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng Yin. 2023. Wavocoder: Widespread and versatile enhanced instruction tuning with refined data generation. *arXiv preprint arXiv:2312.14187*.
- YuLan-Chat-Team. 2023. Yulan-chat: An open-source bilingual chatbot. <https://github.com/RUC-GSAI/YuLan-Chat>.
- Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. 2022. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 1–9. Association for Computational Linguistics.
- Ge Zhang, Yemin Shi, Ruibo Liu, Ruibin Yuan, Yizhi Li, Siwei Dong, Yu Shu, Zhaoqun Li, Zekun Wang, Chenghua Lin, Wen-Fen Huang, and Jie Fu. 2023a. Chinese open instruction generalist: A preliminary release. *ArXiv*, abs/2304.07987.
- Hang Zhang, Xin Li, and Lidong Bing. 2023b. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022a. Opt: Open pre-trained transformer language models. *ArXiv*, abs/2205.01068.
- Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023c. Pmc-vqa: Visual instruction tuning for medical visual question answering. *ArXiv*, abs/2305.10415.
- Yuanhan Zhang, Qinghong Sun, Yichun Zhou, Zexin He, Zhenfei Yin, Kun Wang, Lu Sheng, Yu Qiao, Jing Shao, and Ziwei Liu. 2022b. Bamboo: Building mega-scale vision dataset continually with human-machine synergy. *arXiv preprint arXiv:2203.07845*.
- Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang, Tao Fang, and Wei Bi. 2023d. Multi-task instruction tuning of llama for specific scenarios: A preliminary study on writing assistance. *ArXiv*, abs/2305.13225.
- Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. wildchat: 570k chatgpt interaction logs in the wild.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, L. Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023a. Lima: Less is more for alignment. *ArXiv*, abs/2305.11206.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023b. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.

Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. 2023a. Starling-7b: Improving llm helpfulness & harmlessness with rlaiif.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023b. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

## **A Datasets**

Table 7 gives an overview of our collected datasets.

Type	Dataset Name	# of Instances	# of Lang	Construction	Open-source
Human-Crafted	UnifiedQA (Khashabi et al., 2020) <sup>1</sup>	750K	En	human-crafted	Yes
	UnifiedSKG (Xie et al., 2022) <sup>2</sup>	0.8M	En	human-crafted	Yes
	Natural Instructions (Honovich et al., 2022) <sup>4</sup>	193K	En	human-crafted	Yes
	Super-Natural Instructions (Wang et al., 2022f) <sup>5</sup>	5M	55 Lang	human-crafted	Yes
	P3 (Sanh et al., 2021) <sup>6</sup>	12M	En	human-crafted	Yes
	xP3 (Muennighoff et al., 2022) <sup>7</sup>	81M	46 Lang	human-crafted	Yes
	Flan 2021 (Longpre et al., 2023) <sup>8</sup>	4.4M	En	human-crafted	Yes
	COIG (Zhang et al., 2023a) <sup>9</sup>	-	-	-	Yes
	InstructGPT (Ouyang et al., 2022)	13K	Multi	human-crafted	No
	Dolly (Conover et al., 2023a) <sup>22</sup>	15K	En	human-crafted	Yes
	LIMA (Zhou et al., 2023a) <sup>18</sup>	1K	En	human-crafted	Yes
Synthetic Data (Distillation)	ChatGPT (OpenAI, 2022)	-	Multi	human-crafted	No
	OpenAssistant (Köpf et al., 2023) <sup>17</sup>	161,443	Multi	human-crafted	Yes
	OIG (LAION.ai, 2023) <sup>2</sup>	43M	En	ChatGPT (No technique reports)	Yes
	Unnatural Instructions (Honovich et al., 2022) <sup>10</sup>	240K	En	InstructGPT-Generated	Yes
	InstructWild (Xue et al., 2023) <sup>12</sup>	104K	-	ChatGPT-Generated	Yes
	Evol-Instruct / WizardLM (Xu et al., 2023a) <sup>13</sup>	52K	En	ChatGPT-generated	Yes
	Alpaca (Taori et al., 2023a) <sup>14</sup>	52K	En	InstructGPT-generated	Yes
	LogiCoT (Liu et al., 2023a) <sup>15</sup>	-	En	GPT-4-Generated	Yes
	GPT-4-LLM (Peng et al., 2023) <sup>30</sup>	52K	En&Zh	GPT-4-Generated	Yes
	Vicuna (Chiang et al., 2023)	70K	En	Real User-ChatGPT Conversations	No
	Baize v1 (Conover et al., 2023b) <sup>21</sup>	111.5K	En	ChatGPT-Generated	Yes
	UltraChat (Ding et al., 2023a) <sup>16</sup>	675K	En&Zh	GPT 3/4-Generated	Yes
	Guanaco (JosephusCheung, 2021) <sup>19</sup>	534,530	Multi	GPT (Unknown Version)-Generated	Yes
	Orca (Mukherjee et al., 2023) <sup>23</sup>	1.5M	En	GPT 3.5/4-Generated	Yes
	ShareGPT <sup>24</sup>	90K	Multi	Real User-ChatGPT Conversations	Yes
	WildChat <sup>25</sup>	150K	Multi	Real User-ChatGPT Conversations	Yes
	WizardCoder (Luo et al., 2023)	-	Code	LLaMa 2-Generated	No
	Magicoder (Wei et al., 2023b) <sup>26</sup>	75K/110K	Code	GPT-3.5-Generated	Yes
	WaveCoder (Yu et al., 2023)	-	Code	GPT 4-Generated	No
Synthetic Data (Self-Improvement)	Phi-1 (Gunasekar et al., 2023) <sup>27</sup>	6B Tokens	Code Q and A	GPT-3.5-Generated	Yes
	Phi-1.5 (Li et al., 2023i)	-	Code Q and A	GPT-3.5-Generated	No
	Nectar (Zhu et al., 2023a) <sup>28</sup>	183K	En	GPT 4-Generated	Yes
Reasoning Data	Self-Instruct (Wang et al., 2022c) <sup>11</sup>	52K	En	InstructGPT-Generated	Yes
	Instruction Backtranslation (Li et al., 2023g)	502K	En	LLaMa-Generated	No
	SPIN (Chen et al., 2024b) <sup>29</sup>	49.8K	En	Zephyr-Generated	Yes
	PRM800K (Wang et al., 2022c) <sup>30</sup>	800K	Math	human-crafted & GPT-Generated	Yes
	O1-Journey (Li et al., 2023g) <sup>31</sup>	677	Math	human-crafted & GPT-Generated	Yes
	Self-Explore (Chen et al., 2024b)	-	Math	GPT-Generated	No
	MARIO (Chen et al., 2024b) <sup>32</sup>	28.8K	Math	human-crafted & GPT-Generated	Yes
	MathGenie (Chen et al., 2024b)	170K	Math	GPT-Generated	No
	DeepSeekMath (Chen et al., 2024b) <sup>33</sup>	120B	Math	human-crafted & GPT/DeepSeek-Generated	Yes
	Compute-Optimal Sampling (Chen et al., 2024b)	-	Math	GPT-Generated	No
	MathScale (Chen et al., 2024b) <sup>34</sup>	2M	Math	GPT-Generated	Yes
	G-LLaVA (Chen et al., 2024b) <sup>35</sup>	170K	Math	GPT-Generated	Yes

<sup>1</sup> <https://github.com/allenai/unifiedqa>

<sup>2</sup> <https://github.com/LAION-AI/Open-Instruction-Generalist>

<sup>3</sup> <https://github.com/hkunlp/unifiedskg>

<sup>4</sup> <https://github.com/allenai/natural-instructions-v1>

<sup>5</sup> <https://github.com/allenai/natural-instructions>

<sup>6</sup> <https://huggingface.co/datasets/bigscience/P3>

<sup>7</sup> <https://github.com/bigscience-workshop/xmtf>

<sup>8</sup> <https://github.com/google-research/FLAN>

<sup>9</sup> <https://github.com/BAAI-Zlab/COIG>

<sup>10</sup> <https://github.com/orhonovich/unnatural-instructions>

<sup>11</sup> <https://github.com/yizhongw/self-instruct>

<sup>12</sup> <https://github.com/XueFuzhao/InstructionWild>

<sup>13</sup> <https://github.com/nlp-xucan/evol-instruct>

<sup>14</sup> [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca)

<sup>15</sup> <https://github.com/csitfun/LogiCoT>

<sup>16</sup> <https://github.com/thunlp/UltraChat#data>

<sup>17</sup> <https://github.com/LAION-AI/Open-Assistant>

<sup>18</sup> <https://huggingface.co/datasets/GAIR/lima>

<sup>19</sup> <https://huggingface.co/datasets/JosephusCheung/GuanacoDataset>

<sup>20</sup> <https://github.com/Instruction-Tuning-with-GPT-4/GPT-4-LLM>

<sup>21</sup> <https://github.com/project-baize/baize-chatbot>

<sup>22</sup> <https://huggingface.co/datasets/databricks/databricks-dolly-15k>

<sup>23</sup> <https://huggingface.co/datasets/Open-Orca/OpenOrca>

<sup>24</sup> <https://huggingface.co/datasets/RyokoAI/ShareGPT52K>

<sup>25</sup> <https://huggingface.co/datasets/allenai/WildChat>

<sup>26</sup> <https://github.com/ise-uiuc/magicoder?tab=readme-ov-file#-dataset>

<sup>27</sup> <https://huggingface.co/microsoft/phi-1>

<sup>28</sup> <https://huggingface.co/datasets/berkeley-nest/Nectar>

<sup>29</sup> <https://github.com/uclaml/SPIN?tab=readme-ov-file#Data>

<sup>30</sup> <https://github.com/openai/prm800k>

<sup>31</sup> <https://github.com/GAIR-NLP/O1-Journey>

<sup>32</sup> <https://github.com/MARIO-Math-Reasoning/MARIO>

<sup>33</sup> <https://github.com/deepseek-ai/DeepSeek-Math>

<sup>34</sup> <https://github.com/XylonFu/MathScale>

<sup>35</sup> <https://github.com/pipilurj/G-LLaVA>

Table 7: An overview of instruction tuning datasets.