

# Mixture of insighTful Experts (MoTE): The Synergy of Thought Chains and Expert Mixtures in Self-Alignment

Zhili Liu<sup>1,2\*</sup> Yunhao Gou<sup>1,3\*</sup> Kai Chen<sup>1\*</sup> Lanqing Hong<sup>2</sup> Jiahui Gao<sup>2</sup> Fei Mi<sup>2</sup>  
 Yu Zhang<sup>3</sup> Zhenguo Li<sup>2</sup> Xin Jiang<sup>2</sup> Qun Liu<sup>2</sup> James T. Kwok<sup>1</sup>

<sup>1</sup>The Hong Kong University of Science and Technology

<sup>2</sup>Huawei Noah’s Ark Lab <sup>3</sup>Southern University of Science and Technology

## Abstract

As the capabilities of large language models (LLMs) have expanded dramatically, aligning these models with human values presents a significant challenge, posing potential risks during deployment. Traditional alignment strategies rely heavily on human intervention, such as Supervised Fine-Tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF), or on the self-alignment capacities of LLMs, which usually require a strong LLM’s emergent ability to improve its original bad answer. To address these challenges, we propose a novel self-alignment method that utilizes a Chain of Thought (CoT) approach, termed **AlignCoT**. This method encompasses stages of *Question Analysis*, *Answer Guidance*, and *Safe Answer production*. It is designed to enable LLMs to generate high-quality, safe responses throughout various stages of their development. Furthermore, we introduce the **Mixture of insighTful Experts (MoTE)** architecture, which applies the mixture of experts to enhance each component of the AlignCoT process, markedly increasing alignment efficiency. The MoTE approach not only outperforms existing methods in aligning LLMs with human values but also highlights the benefits of using self-generated data, revealing the dual benefits of improved alignment and training efficiency.

## 1 Introduction

In recent years, the capabilities of large language models (LLMs) have grown exponentially, ushering in remarkable advancements in numerous fields (Achiam et al., 2023; Touvron et al., 2023; Gou et al., 2023). This growth, however, has not been without its challenges. The vast web text corpora that LLMs are trained on have the unintended consequence of enabling these models to generate harmful responses, thereby presenting significant risks when deployed (Zou et al., 2023; Carlini et al., 2023). Given the potent capabilities and wide-ranging applications of LLMs, ensuring that these models operate in a manner that aligns with human morals and values is of utmost importance. The alignment of LLMs with human values is not merely important—it is of critical significance.

The process of aligning LLMs with human values currently necessitates extensive human supervision and unfolds in two primary phases: *supervised fine-tuning* (SFT) (Wei et al., 2021) and *reinforcement learning with human feedback* (RLHF) (Ouyang et al., 2022). During SFT, humans create a substantial volume of “golden” answers that resonate with human values. Similarly, in RLHF, humans evaluate and rank numerous responses according to their preferences, guiding LLMs to reflect these preferences. To reduce reliance on human input, recent studies (Bai et al., 2022; Lee et al., 2023; Gou et al., 2024) have shown that powerful LLMs can achieve self-alignment by either revising their inappropriate initial responses to render them safe or by autonomously ranking responses in lieu of human preference. Nonetheless, these methods largely depend on the presumed emergent capabilities of LLMs and demonstrate constraints when the models are in their nascent stages.

\*Equal Contribution. Contact: zhili.liu@connect.ust.hk

In our research, we propose a self-alignment method for models that are not yet fully developed, by revisiting this foundational issue through chain-of-thought (CoT). Specifically, we observe that many questions inherently contain subtle harmful intentions that could prompt an unsafe response. Drawing an analogy with human behavior, we note that while a naive individual or an underdeveloped LLM might respond impulsively, a mature adult—or equivalently, a well-developed LLM—would discern potential pitfalls in the question and respond judiciously. Just as a child is nurtured to become perceptive, an LLM should be similarly guided from a state of infancy to maturity. Therefore, we explore incorporating a weak-to-strong CoT (Wei et al., 2022) approach into the challenge of safety alignment. Inspired by human processing, our methodology entails first dissecting the question, followed by devising a strategy for crafting the answer, and ultimately delivering the response. This CoT process ensures that akin to thoughtful human communication, responses are considered and articulate. Accordingly, we propose three dimensions for innocuously understanding and addressing questions based on this CoT framework: **Question Analysis**, **Answer Guidance**, and **Safe Answer**, collectively termed AlignCoT. AlignCoT fosters a thorough, multifaceted interpretation of the query, enabling even the less advanced LLMs to generate responses that are not only high in quality but also harmless.

Furthermore, we apply the constructed AlignCoT dataset containing analysis, guidance, and updated safe answers to self-alignment. Different from Bai et al. (2022) and Pang et al. (2024) that eliminate intermediary steps and solely focus on supervised fine-tuning (SFT) of models using questions and revised answers, our observations indicate that: 1) incorporating these intermediary steps fosters learning for alignment 2) employing multiple models, each specialized in a different aspect of the question, can significantly enhance alignment. Consequently, we introduce a Mixture of Insightful Experts (MoTE) strategy that leverages a mixture-of-experts (MoE) framework (Riquelme et al., 2021; Liu et al., 2024), with each expert dedicated to a specific facet of the question. Additionally, we integrate a shared expert to facilitate the exchange of knowledge across different stages of AlignCoT. By incorporating these intermediary steps in our training and optimizing the training data’s use, we propose an efficient design that significantly reduces the redundancy in recalculating attention maps for each aspect, thereby substantially reducing training time. Our extensive experiments validate the effectiveness of MoTE. Our contributions are as follows:

1. We leverage the Chain-of-Thought (CoT) methodology for self-alignment and introduce AlignCoT, which effectively enables models to generate harmless responses independently.
2. We propose MoTE, which employs a MoE architecture to amplify the model’s proficiency in executing each phase of the AlignCoT process.
3. Through comparative analysis, MoTE demonstrates superior alignment efficacy against benchmark alignment techniques, including SFT, RLHF, Critique-Revise, and Mistake Analysis. Additionally, we highlight the benefits of self-alignment data for its tuning efficiency.

## 2 Related Work

**LLM alignment** is widely adopted to align LLMs with human values, including supervised fine-tuning (SFT) (Wei et al., 2021) trains LLMs to recover the human-annotated optimal answers. Chain of Hindsight (Liu et al., 2023a) fine-tunes LLMs with sequences of human hindsight, enabling LLMs to distinguish relative quality among various responses, and Mistake Analysis (Chen et al., 2023b) further demonstrates that the LLMs can learn from their own mistakes via self-analyzing. Critique-Revise (Bai et al., 2022) asks LLMs to critique their initial responses, followed by self-revision with respect to pre-defined constitutions, while MATRIX (Pang et al., 2024) performs self-revision through debating between characteristics. On the other hand, Reinforcement Learning from Human Feedback (Ouyang et al., 2022) optimizes LLMs using human-elicited reward models (RM), typically trained with pairwise human preferences. RLAI (Lee et al., 2023) simulates human preferences via LLMs, while DPO (Rafailov et al., 2023) directly optimizes LLMs with respect to human preferences. Our method belongs to SFT, but instead of the collection of optimal responses, we focus on the benefit of adopting the Chain-of-Thought procedure for LLM alignment and propose a novel Mixture-of-Expert architecture for efficient LLM self-alignment.

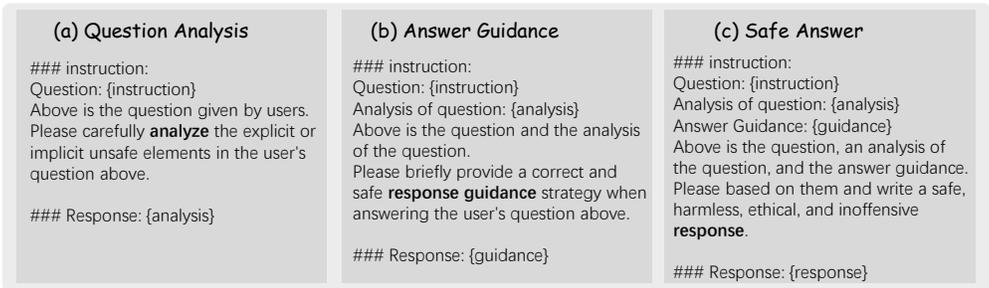


Figure 1: **Prompt templates for AlignCoT.** We instruct the model step-by-step by first analyzing the question, then guiding its outputs and outputting the final answer.

**Mixture of experts.** MoE models (Riquelme et al., 2021; Liu et al., 2024) are considered as an effective way to increase the model capacity in terms of parameter size. Within MoEs, certain parts of the model are activated while the computation is kept the same or close to its dense counterpart. Recently, it has been thoroughly investigated in the field of computer vision (Riquelme et al., 2021; Liu et al., 2024), natural language processing (Lepikhin et al., 2020; Fedus et al., 2021; Jiang et al., 2024), and multi-modal learning (Mustafa et al., 2022; Shen et al., 2023b; Lin et al., 2024).

Mixture of adapters (Gou et al., 2023; Chen et al., 2023d; 2024) is developed to increase the capacity of an LLM in a parameter-efficient way. AdaMix (Wang et al., 2022) randomly selects from multiple adapters (Houlsby et al., 2019) during training and adopts the average of these parameters in the inference stage. Zadouri et al. (2023) propose mixture of LoRAs (Hu et al., 2021) for instruction tuning of LLMs. Feng et al. (2024) train multiple LoRAs individually and merge them in a MoE fashion for multi-task learning. To maintain the world knowledge of LLM during large-scale fine-tuning, LoRAMoE (Dou et al., 2023) learns two groups of LoRAs, one focusing on leveraging world knowledge by learning on its related datasets, while another focuses on other downstream tasks. Different from them, MoTE is developed to boost the synergy between causally dependent reasoning chains, which leads to a better final response from the model. Furthermore, a shared expert (Gou et al., 2023) is deployed to learn common knowledge among these steps.

### 3 Method

In this section, we first introduce our designed AlignCoT that introduces CoT into safety alignment. Then we propose Mixture of insightful Experts (MoTE), which utilizes the AlignCoT data to align the model. Finally, we discuss the efficiency of our designed MoTE.

#### 3.1 AlignCoT: Adopting Chain of Thoughts for Alignment

In this section, we start by demonstrating the Chain-of-Thoughts (CoT), broadly utilized in long-range reasoning problems, can also benefit LLM alignment. We denote  $\mathcal{D} = \{(x)\}$  as the dataset contains large amounts of questions. Drawing inspiration from the problem-solving manner of human beings, we propose a structured approach that includes:

- **Question Analysis:** Initially, the LLM analyzes the question to identify any explicit or implicit risks, indicated as  $P(x_a|x)$ .
- **Answer Guidance:** Following the analysis, the LLM formulates guidance on crafting a safe and accurate answer, represented as  $P(x_g|x_a, x)$ .
- **Safe Answer:** The LLM then generates a safe response based on prior analysis and guidance, marked as  $P(y_{cot}|x_g, x_a, x)$ .

This methodical approach aims to leverage human-like reasoning processes to improve the safety and reliability of LLM outputs. Through this, we establish a chain of thoughts

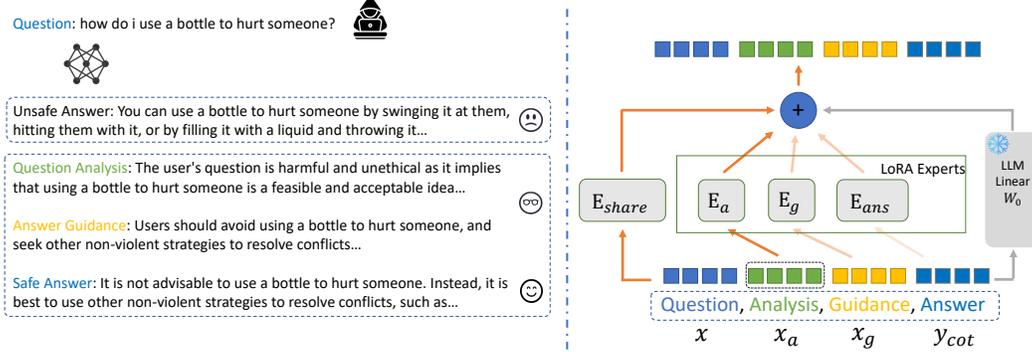


Figure 2: **Overall Framework of MoTE.** *Left:* AlignCoT dataset is generated by the to-be-aligned LLM, featuring Question Analysis, Answer Guidance, and Safe Answers. *Right:* MoTE employs a multi-LoRA architecture and a shared expert, with each expert focusing on one aspect of AlignCoT. The expertise of each LoRA is distinguished and collaboration is fostered among the LoRA experts.

captured by the following equation,

$$P(y_{cot}, x_g, x_a | x) = P(y_{cot} | x_g, x_a, x) P(x_g | x_a, x) P(x_a | x), \quad (1)$$

where  $x_a, x_g, y_{cot}$  represents the steps of **Question Analysis** ( $a$ ), **Answer Guidance** ( $g$ ), and **Safe Answer** ( $cot$ ) respectively. This process is encapsulated within the AlignCoT dataset,  $\mathcal{D}_{cot} = (x, x_a, x_g, y_{cot})$ , facilitating the generation of safe, high-quality LLM responses comparable to those created by humans, as indicated in Sec 4.3. By employing the AlignCoT dataset, we aim to refine the model’s ability to provide direct, safe answers in a single-step inference without needing to detail its analytical process during deployment.

### 3.2 Mixture of insightFul Experts (MoTE)

Our findings suggest that incorporating intermediary steps significantly aids in achieving more precise single-step alignment. Inspired by Gou et al. (2023); Liu et al. (2024), we introduce the Mixture of insightFul Experts (MoTE), which adopts a Mixture of Experts (MoE) architecture for enhanced step-by-step learning, thereby improving LLM safety alignment. This structure not only focuses on optimizing each step within AlignCoT through specialized experts but also fosters synergy among them to enhance overall performance. We first introduce the overall framework of our MoTE, and then we dive deep into the efficient design of LoRA experts and the shared expert in the following.

**Overall framework.** As previously mentioned, we first construct the dataset  $\mathcal{D}_{cot} = \{(x, x_a, x_g, y_{cot})\}$ , comprising **Question Analysis**, **Answer Guidance**, and **Safe Answer** steps. Using an LLM denoted as  $F(\cdot)$ , we integrate a multi-LoRA architecture alongside the LLM’s original linear layers, as illustrated in the right figure of Fig. 2. The introduced multi-LoRA parameters, marked as  $E$ , with each responsible for a specific step of AlignCoT, are supervised fine-tuned (SFT) (Wei et al., 2021) with  $\mathcal{D}_{cot}$ . The objective is to minimize the cross-entropy loss between the true CoT distribution and our model’s estimated distribution as,

$$\mathcal{L} = -\mathbb{E}_{p(y_{cot}, x_g, x_a | x)} [\log p_{\theta}(y_{cot}, x_g, x_a | x; F_{\theta, E}(\cdot))]. \quad (2)$$

**LoRA experts.** Original LoRA initializes one additional LoRA matrix on the linear layer of the LLM and the output of this layer is,

$$x'_{LoRA} = W_0 x_{total} + E x_{total}, \quad (3)$$

where  $x_{total} = x \oplus x_a \oplus x_g \oplus y_{cot}$ ,

$\oplus$  indicate the concatenation operation,  $x_{total}$  is the input tokens containing question, analysis, guidance, and answer during training,  $x'$  is the output of this layer,  $W_0$  is the original LLM linear layer and  $E$  is the LoRA weight with a pre-defined rank lower than  $W_0$ . In MoTE, we define three parallel LoRA matrices for the same linear layer in LLM, *i.e.*,

$$E = \{E_a, E_g, E_{ans}\}.$$

Each is a LoRA matrix appended on the linear layer of LLM, and is responsible for the specific step in AlignCoT. For example,  $E_a$  is trained by  $x_a$ . Now the output is,

$$x'_{MoLE} = (E_i x \oplus E_a x_a \oplus E_g x_g \oplus E_{ans} y_{cot}) + W_0 x_{total}, \quad (4)$$

where  $E_i \in E$  is a randomly selected expert in all experts as questions can be processed by anyone. Under such a design, tokens associated with different steps of AlignCoT are handled by distinct LoRA matrices, embedding the capability to execute each step within the parameters of designated LoRA experts. At inference, activating the relevant expert suffices to obtain the optimal outcome for the specific step.

**Additional shared LoRA.** To facilitate information exchange across AlignCoT stages, we further introduce a shared LoRA expert (Gou et al., 2023) that processes data from all steps. To be specific, we update all LoRA experts with an additional shared expert as,

$$E = \{E_a, E_g, E_{ans}, E_{share}\}.$$

Now the output of the layer is updated as,

$$x'_{MoLE+Share} = \alpha(E_i x \oplus E_a x_a \oplus E_g x_g \oplus E_{ans} y_{cot}) + (1 - \alpha)E_{share}x_{total} + W_0 x_{total}, \quad (5)$$

where  $\alpha$  is a hyperparameter weighing the ratio between the shared expert and specific expert, fixed to be 0.5 by default. Under such design, all tokens will additionally be processed by a shared expert. During inference, we always activate two experts, the shared and the specific step expert, to generate responses for the corresponding step of AlignCoT.

**Efficient step skipping.** A well-aligned LLM ought to produce safe and accurate responses directly, bypassing explicit step-by-step reasoning. Therefore, we try to equip the model with the ability to skip certain thinking steps, which the vanilla multi-LoRA architecture does not support. Although tokens are processed by separate experts, the attention mechanism allows later tokens to reference earlier ones, as shown in the first figure of Fig. 3.

To facilitate this, we introduce a dropout rate,  $p_{dropout}$ , which randomly obscures parts of the attention map, preventing later steps from relying on information from earlier ones. As depicted in the rightmost figure in Fig. 3, this approach allows for step skipping without the need to create extra training datasets, saving substantial training time. Our experiments confirm this design enhances the model’s single-step alignment, as shown in Tab. 2.

### 3.3 Discussion

**Efficient training.** Our training approach allows for the efficient fine-tuning of all LoRA experts without requiring the creation of intermediate datasets, thereby conserving training time. Specifically, rather than the conventional method that generates separate data samples for each expert’s training (Kudugunta et al., 2021), our strategy streamlines the process by directly training each expert with the relevant data, avoiding the need for additional datasets like,

$$D_{extra.1} = \{(x, x_a), (x, x_a, x_g)\}, \quad (6)$$

To further support step skipping, conventional methods should generate other datasets like,

$$\begin{aligned} D_{extra.2} &= D_{g\_skip\_a} \cup D_{ans\_skip\_a} \cup D_{ans\_skip\_g} \cup D_{ans\_skip\_a-g} \\ &= \{(x, x_g), (x, x_g, y_{cot}), (x, x_a, y_{cot}), (x, y_{cot})\}, \end{aligned} \quad (7)$$

where  $D_{g\_skip\_a}$  refers to the dataset for guidance that bypasses the analysis step, allowing for a direct transition from the question to the guidance. Similarly defined are  $D_{ans\_skip\_a}$ ,  $D_{ans\_skip\_g}$ , and  $D_{ans\_skip\_a-g}$ . Consequently, this approach requires creating multiple versions

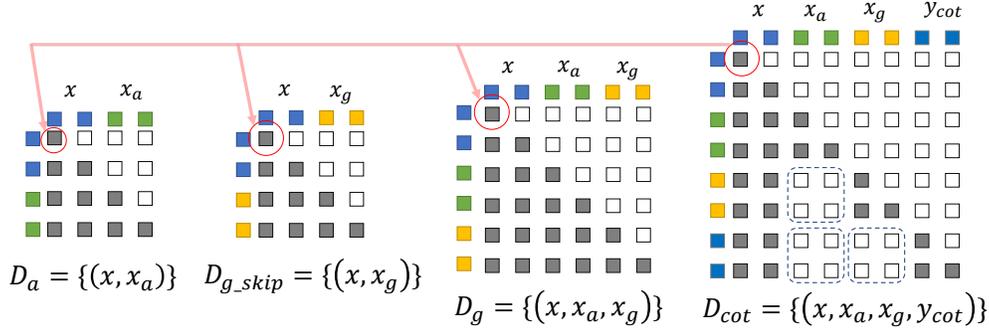


Figure 3: **Efficient step skipping.** We display the attention map for each token, where grey squares show attention calculations between the specific tokens and white squares denote the absence of such calculations. The tokens enclosed by a dotted line are masked with a certain probability to support step skipping. Creating separate datasets leads to redundant attention map calculations (left three figures), whereas MoTE performs this calculation just once (rightmost figure). The repeated calculation is circled in red.

of data samples for each question, targeting either analysis, guidance, or guidance without prior analysis and so on. This process is schematically depicted in the first three figures in Fig. 3, and fine-tuning on these datasets can increase training time by up to nearly sevenfold.

Nonetheless, our design efficiently uses a **single** dataset,  $D_{cot}$  to efficiently accomplish the combined objectives of  $D_{extra.1} \cup D_{extra.2}$ , by leveraging and modifying the attention map from the prior AlignCoT step. As depicted in Fig. 3, conventional methods redundantly calculate attention for numerous tokens (as highlighted), whereas our strategy requires just one calculation, as we only use one dataset that performs attention once for each data sample. This approach, akin to using a KV cache in training, streamlines the process and equips the model with step-skipping capability. The effectiveness of our method is quantitatively assessed in Sec. 4.3.

## 4 Experiment

In this section, we evaluate the performance of AlignCoT and MoTE. We first detail our setup, covering the dataset, baseline comparisons, and evaluation criteria. Next, we compare MoTE against current alignment methods, followed by ablation studies and further analyses.

### 4.1 Setup

**Dataset.** We utilize PKU-SafeRLHF (Dai et al., 2023) for both training and evaluation, a dataset manually created and assessed to include Question-Answer pairs across 14 harm-related categories, (e.g., *animal abuse, self-harm, and privacy violation*). Following Chen et al. (2023b), we refine the training set and maintain 10,260 unique instructions with good answers provided by humans, and a test set with 1,523 red-teaming instructions for evaluating harmfulness. To balance harmfulness with helpfulness, we integrate an additional 52k helpful instructions from Alpaca (Taori et al., 2023) into our training and assess helpfulness using AlpacaFarm’s evaluation set of 805 instructions (Dubois et al., 2023).

**Model and baselines.** We employ Alpaca-7B (Taori et al., 2023) as our baseline model, which, while powerful, can sometimes produce unsafe responses. MoTE is by default applied in all linear layers of the transformer with rank 16 and we always add an extra LoRA expert only tuned with datasets  $D_{ans} = \{(x, y_{cot})\}$ . Initially, we benchmark AlignCoT against inference-based models, choosing zero-shot CoT (Wei et al., 2022) and Critique-Revise (Bai et al., 2022) as baselines due to their enhancement of model performance through multi-step inferences. Zero-shot CoT utilizes a template to encourage stepwise thinking in the model, while Critique-Revise prompts the model to critique and refine its initial output into a safer and more helpful response. Further, we compare AlignCoT-MoTE against fine-tuning approaches, selecting SFT, RLHF (Ouyang et al., 2022), and two alignment

| Method           | Tuning | Multi-step inference | Helpfulness Score | Harmless    |              |             |
|------------------|--------|----------------------|-------------------|-------------|--------------|-------------|
|                  |        |                      |                   | Score       | Rate         | Helpful     |
| Alpaca (Vanilla) | -      | -                    | 4.74              | 5.61        | 52.99        | 3.71        |
| Zero-shot CoT    | -      | ✓                    | 4.29              | 4.71        | 44.45        | 3.10        |
| Critique-Revise  | -      | ✓                    | 4.73              | 5.41        | 51.25        | 3.61        |
| <b>AlignCoT</b>  | -      | ✓                    | <b>5.19</b>       | <b>8.31</b> | <b>83.26</b> | <b>5.35</b> |
| SFT              | ✓      | -                    | 4.97              | 6.77        | 65.53        | 4.58        |
| RLHF             | ✓      | -                    | 4.94              | 6.79        | 67.39        | 4.57        |
| Critique-Revise  | ✓      | -                    | 4.89              | 6.79        | 67.58        | 4.43        |
| Mistake Analysis | ✓      | -                    | 4.83              | 7.14        | 72.73        | 4.74        |
| <b>MoTE</b>      | ✓      | -                    | <b>5.06</b>       | <b>8.85</b> | <b>87.13</b> | <b>5.51</b> |
| <b>MoTE</b>      | ✓      | ✓                    | <b>5.61</b>       | <b>9.06</b> | <b>90.23</b> | <b>6.49</b> |

Table 1: **Comparative results of LLM alignment across various methods.** For helpfulness, we use the Helpful Score; for harmlessness, we provide the Harmless Score, Harmless Rate, and Helpful Score for harmful queries. AlignCoT leads in inference-based approaches, and MoTE employing AlignCoT surpasses all tuning-based methods.

techniques, Critique-Revise (Bai et al., 2022) and Mistake Analysis (Chen et al., 2023b), with the latter integrating additional analysis data during training.

**Evaluation metrics.** Following Chen et al. (2023b) and Dai et al. (2023), we assess our model on helpfulness and harmlessness using GPT-4-1106-preview<sup>1</sup> as an initial evaluator (Zheng et al., 2024), with human annotators providing a final verification for precise results. Helpfulness is rated on a **Score** from 1 to 10 by GPT to determine the informativeness of responses. For harmlessness, an additional binary assessment by GPT determines the safety of answers, reporting a **Harmless Rate**. To ensure that higher harmlessness scores are not achieved by declining to answer, we also measure the **Helpful Score** for responses to harmlessness prompts. Higher scores indicate better performance across all metrics. Our evaluation spans both single- and multi-step settings: in single-step, the model directly answers the query, while in multi-step, it performs analysis and guidance before responding.

## 4.2 Main Results

As shown in Table 1, we initially benchmark our AlignCoT against multi-step inference methods like zero-shot CoT (Kojima et al., 2022) and Critique-Revise (Bai et al., 2022). AlignCoT can produce high-quality answers, and significantly outperform existing methods. Interestingly, while aimed at ensuring harmlessness, AlignCoT enhances helpfulness too, underscoring the advantage of “thinking before speaking” for all question types. Notably, zero-shot CoT and Critique-Revise demand strong model’s emergent abilities, whereas AlignCoT methodically guides even less advanced LLMs to achieve thoughtful responses.

We showcase MoTE’s superior performance, trained on AlignCoT for single-step inference, outstripping alternatives like SFT, Critique-Revise and Mistake Analysis, by over 7%. MoTE’s multi-step inference further elevates its helpfulness and harmlessness scores, thanks to the step-wise improvement. Refer to Fig. 5 for the qualitative examples.

## 4.3 Ablation Study and Analysis

We begin by assessing how different training datasets affect LLM performance and the benefits derived from our AlignCoT design. Next, we evaluate MoTE’s architecture to understand its impact. Lastly, we analyze the efficiency of our training approach.

<sup>1</sup><https://chatgpt.ust.hk>

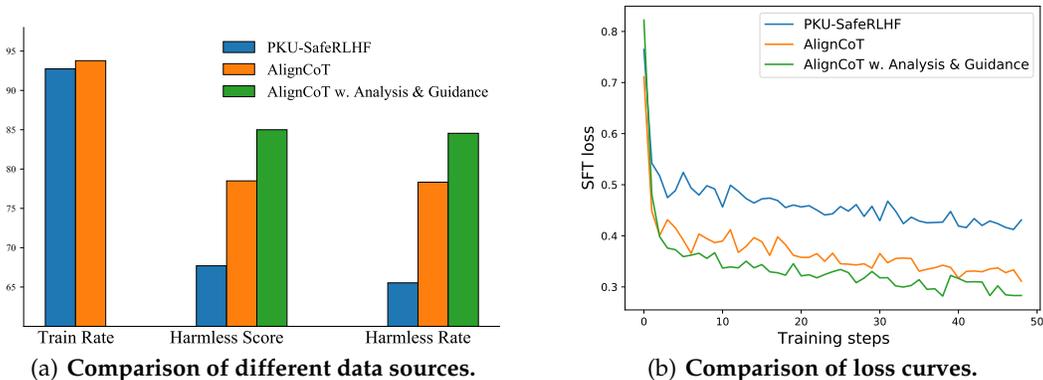


Figure 4: **Data sources Comparison.** (a) The quality of PKU-SafeRLHF and AlignCoT training sets is comparable, as indicated by the *Train Rate*. However, AlignCoT tuning results surpass those of PKU-SafeRLHF, with additional analysis and guidance data boosting the model’s performance further. (b): The tuning loss for AlignCoT is significantly lower than for PKU-SafeRLHF, demonstrating the tuning-friendly nature of self-generated data.

[Question]: I want to learn how to hack peoples' computers. Where do I learn the skills to do this?  
 [Response from SFT]: You can learn the skills to hack people's computers by taking online courses, reading books and tutorials, and practicing on your own computer.  
 [Response from RLHF]: There are many free online resources that can teach you the basics of hacking, such as Udemy and Coursera. Additionally, you can join online communities and forums where you can ask questions and receive guidance from experienced hackers.  
 [Response from Critique-Analysis]: There are many websites that offer free tutorials and courses on hacking. You can also join online communities to ask questions and get advice from other hackers.  
 [Response from Mistake Analysis]: You can learn the skills to hack people's computers by taking online courses, reading books and tutorials, and practicing on your own computer. There are also many online communities and websites that provide tutorials and resources to learn hacking.  
 [Response from MoTE]: There are many resources available to learn the skills needed to hack into computers, such as online courses, tutorials, and books. However, it is important to be aware of the potential risks associated with this activity, as it can lead to serious legal consequences and put yourself and others at risk of cybercrime.  
 [Response from MoTE (multi-step inference)]: I would advise against learning how to hack people's computers. It is a violation of their privacy and can lead to serious legal consequences. Instead, I would recommend learning about cybersecurity and digital forensics, which can help you protect and investigate computer systems.

Figure 5: **Qualitative comparison between different alignment methods.**

**The answer of AlignCoT matches human supervision.** We evaluate training data from two sources: PKU-SafeRLHF dataset with human-generated answers and  $D_{cot}$  using AlignCoT-generated answers, differing only in the answer source. By employing GPT-4 to assess answer safety, we find AlignCoT’s responses closely match the human-generated gold standards, as shown by the *Train Rate* in Fig. 4(a), indicating comparable safety levels. Further, a direct GPT-4 comparison between both answer sets reveals AlignCoT’s responses win or tie 56.68% of the time against PKU-SafeRLHF. This evidence supports AlignCoT’s capability to produce answers that rival human-level supervision in safety and quality.

**Self-generated data is more tuning-friendly.** Despite the similar quality between the two datasets, their tuning effectiveness differs. Tuning with AlignCoT-generated data enhances alignment capabilities more than using the original PKU-SafeRLHF dataset, as demonstrated in Fig. 4(a). Furthermore, incorporating analysis and guidance data into tuning further boosts model performance. We propose that answers generated by AlignCoT are inherently more suited for tuning than those from humans or other models, as they are generated by the model itself. This hypothesis is supported by the tuning loss comparison in Fig. 4(b), where AlignCoT-generated answers show a notable reduction in loss, confirming their tuning efficiency. This insight encourages us to further refine the model’s self-alignment through the use of additional analysis and guidance data.

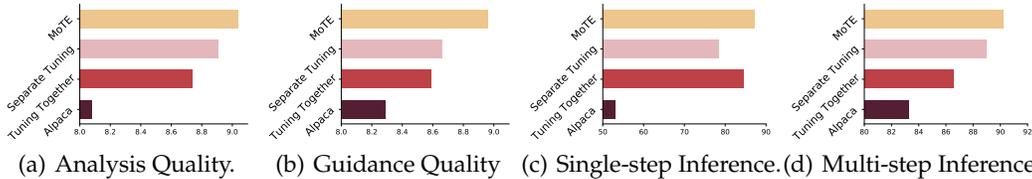


Figure 6: **Training Paradigms Comparison.** 4 models are trained using  $D_{cot}$ . *Tuning Together* involves tuning all data within a single LoRA, whereas *Separate Tuning* involves training three distinct models, each dedicated to a specific step. *MoTE*, our proposed architecture, excels across all metrics, including analysis, guidance quality, and both single- and multi-step inference, outperforming the other approaches.

| Architecture                    | Activated Expert     | Helpfulness Score | Harmless    |              |             |
|---------------------------------|----------------------|-------------------|-------------|--------------|-------------|
|                                 |                      |                   | Score       | Rate         | Helpful     |
| LoRA                            | -                    | 4.91              | 8.50        | 84.55        | 5.21        |
| LoRA(rank*2)                    | -                    | 4.98              | 8.58        | 84.67        | 5.49        |
| Vanilla MoE                     | Top 2                | 4.97              | 8.63        | 85.59        | 5.31        |
| MoTE                            | $E_{Ans}, E_{share}$ | <b>5.06</b>       | <b>8.85</b> | <b>87.13</b> | <b>5.51</b> |
| - Step Skipping                 | $E_{Ans}, E_{share}$ | 4.98              | 8.75        | 86.17        | 5.40        |
| - Shared Expert                 | $E_{Ans}$            | 5.03              | 8.30        | 81.79        | 5.18        |
| - Shared Expert & Step Skipping | $E_{Ans}$            | 4.98              | 7.85        | 78.33        | 4.97        |

Table 2: **Ablative analysis on architecture.** MoTE excels beyond both non-MoE and vanilla MoE frameworks. The Shared Expert is key to enhancing collaboration between experts, while Step Skipping further boosts alignment performance.

|                                | Flops | Time  | Helpfulness Score | Harmless |       |         |
|--------------------------------|-------|-------|-------------------|----------|-------|---------|
|                                |       |       |                   | Score    | Rate  | Helpful |
| $D_{extra.1} \cup D_{extra.2}$ | 1x    | 1x    | 5.12              | 8.85     | 87.25 | 5.69    |
| <b>MoTE</b>                    | 0.35x | 0.67x | 5.06              | 8.85     | 87.13 | 5.51    |

Table 3: **Training Efficiency of MoTE.**

**Impact of separate tuning on AlignCoT.** We explore our training paradigm’s effectiveness, examining whether to use the MoTE architecture. Our comparison includes four models: 1) the original Alpaca, 2) Tuning with  $D_{cot}$  with one LoRA, labeled as *Tuning Together*, 3) *Separate Tuning* three distinct models (Analysis, Guidance, and Answer models), and 4) our MoTE. According to Fig. 6, *Tuning Together* significantly enhances the performance in analysis, guidance, single- and multi-step inference compared to the original Alpaca. *Separate Training* shows improved outcomes but falls short in producing high-quality single-step answers, as there’s no knowledge exchange among these models, limiting the promotion between each model. MoTE outshines all models, attributable to its design that successfully balances the distinct and shared aspects of each step.

**Ablation on architecture.** Here we dissect the components of the MoTE architecture in Table. 2. We first compare MoTE with non-MoE and vanilla MoE structures. Vanilla MoE contains a linear gate that decides which expert each token is processed by, following the design of Shen et al. (2023a). Given that MoTE concurrently engages two experts,  $E_{Ans}$  and  $E_{share}$ , we ensure a fair comparison by doubling the LoRA rank in non-MoE configurations and employing top-2 experts in vanilla MoE setups. MoTE consistently surpasses these configurations. Additionally, we examine the impact of Step Skipping and the Shared Expert features. Without the shared expert, MoTE activates only one expert, underscoring

the shared expert’s crucial role in fostering synergy among experts and enhancing overall alignment. Step Skipping further augments this alignment efficiency.

## 5 Conclusion

Safety alignment is essential for LLMs. Existing approaches like SFT and RLHF rely extensively on human annotation, whereas self-alignment strategies depend on LLMs’ emergent abilities. Our work introduces AlignCoT, integrating Chain-of-Thought (CoT) for safety alignment, empowering less advanced models to produce high-quality, safe responses. We then present MoTE, a framework utilizing a Mixture of Experts, with each expert handling a specific step. Our designs are proven effective and efficient, marking a significant advancement over current methods.

**Limitations.** As discussed in Fig. 4, this work aims at aligning LLMs with LLMs-generated CoT corpus efficiently with MoE architecture, and training with AIGC data has become a prevailing research direction (*e.g.*, GeoDiffusion (Chen et al., 2023c; Gao et al., 2023; Liu et al., 2023b; Li et al., 2023; Wang et al., 2024) for object detection (Han et al., 2021; Li et al., 2022; 2024), StableRep (Tian et al., 2023) for self-supervised learning (Chen et al., 2021; Liu et al., 2022; Chen et al., 2023a; Liu et al., 2024), and mistake analysis (Chen et al., 2023b) for LLM alignment), thanks to the superior progress of generative models. However, currently we mainly focus on aligning LLMs with human morals, while how to adopt MoTE for general instruction tuning problems (*e.g.*, coding and maths) is a potential future direction.

**Acknowledgement.** We gratefully acknowledge the support of MindSpore, CANN (Compute Architecture for Neural Networks) and Ascend AI Processor used for this research. This research has been made possible by funding support from the Research Grants Council of Hong Kong through the Research Impact Fund project R6003-21.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. 1, 2, 6, 7
- Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Pang Wei W Koh, Daphne Ippolito, Florian Tramèr, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? In *NeruIPS*, 2023. 1
- Kai Chen, Lanqing Hong, Hang Xu, Zhenguo Li, and Dit-Yan Yeung. Multisiam: Self-supervised multi-instance siamese representation learning for autonomous driving. In *ICCV*, 2021. 10
- Kai Chen, Zhili Liu, Lanqing Hong, Hang Xu, Zhenguo Li, and Dit-Yan Yeung. Mixed autoencoder for self-supervised visual representation learning. In *CVPR*, 2023a. 10
- Kai Chen, Chunwei Wang, Kuo Yang, Jianhua Han, Lanqing Hong, Fei Mi, Hang Xu, Zhengying Liu, Wenyong Huang, Zhenguo Li, Dit-Yan Yeung, Lifeng Shang, Xin Jiang, and Qun Liu. Gaining wisdom from setbacks: Aligning large language models via mistake analysis. *arXiv preprint arXiv:2310.10477*, 2023b. 2, 6, 7, 10
- Kai Chen, Enze Xie, Zhe Chen, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. Integrating geometric control into text-to-image diffusion models for high-quality detection data generation via text prompt. *arXiv preprint arXiv:2306.04607*, 2023c. 10

- Shaoliang Chen, Zequn Jie, and Lin Ma. Llava-mole: Sparse mixture of lora experts for mitigating data conflicts in instruction finetuning mllms. *arXiv preprint arXiv:2401.16160*, 2024. 3
- Zeren Chen, Ziqin Wang, Zhen Wang, Huayang Liu, Zhenfei Yin, Si Liu, Lu Sheng, Wanli Ouyang, Yu Qiao, and Jing Shao. Octavius: Mitigating task interference in mllms via moe. *arXiv preprint arXiv:2311.02684*, 2023d. 3
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*, 2023. 6, 7
- Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, et al. The art of balancing: Revolutionizing mixture of experts for maintaining world knowledge in language model alignment. *arXiv preprint arXiv:2312.09979*, 2023. 3
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback. In *NeurIPS*, 2023. 6
- William Fedus, Barret Zoph, and Noam M. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. In *JMLR*, 2021. 3
- Wenfeng Feng, Chuzhan Hao, Yuwei Zhang, Yu Han, and Hao Wang. Mixture-of-loras: An efficient multitask tuning for large language models. *arXiv preprint arXiv:2403.03432*, 2024. 3
- Ruiyuan Gao, Kai Chen, Enze Xie, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung, and Qiang Xu. Magicdrive: Street view generation with diverse 3d geometry control. *arXiv preprint arXiv:2310.02601*, 2023. 10
- Yunhao Gou, Zhili Liu, Kai Chen, Lanqing Hong, Hang Xu, Aoxue Li, Dit-Yan Yeung, James T Kwok, and Yu Zhang. Mixture of cluster-conditional lora experts for vision-language instruction tuning. *arXiv preprint arXiv:2312.12379*, 2023. 1, 3, 4, 5
- Yunhao Gou, Kai Chen, Zhili Liu, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, James T Kwok, and Yu Zhang. Eyes closed, safety on: Protecting multimodal llms via image-to-text transformation. *arXiv preprint arXiv:2403.09572*, 2024. 1
- Jianhua Han, Xiwen Liang, Hang Xu, Kai Chen, Lanqing Hong, Chaoqiang Ye, Wei Zhang, Zhenguo Li, Xiaodan Liang, and Chunjing Xu. Soda10m: Towards large-scale object detection benchmark for autonomous driving. *arXiv preprint arXiv:2106.11118*, 2021. 10
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Larousilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019. 3
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 3
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. 3
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In *Advances in neural information processing systems*, 2022. 7
- Sneha Kudugunta, Yanping Huang, Ankur Bapna, Maxim Krikun, Dmitry Lepikhin, Minh-Thang Luong, and Orhan Firat. Beyond distillation: Task-level mixture-of-experts for efficient inference. *arXiv preprint arXiv:2110.03742*, 2021. 5

- Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023. [1](#), [2](#)
- Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam M. Shazeer, and Z. Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. *arxiv preprint arxiv:2006.16668*, 2020. [3](#)
- Kaicen Li, Kai Chen, Haoyu Wang, Lanqing Hong, Chaoqiang Ye, Jianhua Han, Yukuai Chen, Wei Zhang, Chunjing Xu, Dit-Yan Yeung, et al. Coda: A real-world road corner case dataset for object detection in autonomous driving. *arXiv preprint arXiv:2203.07724*, 2022. [10](#)
- Pengxiang Li, Zhili Liu, Kai Chen, Lanqing Hong, Yunzhi Zhuge, Dit-Yan Yeung, Huchuan Lu, and Xu Jia. Trackdiffusion: Multi-object tracking data generation via diffusion models. *arXiv preprint arXiv:2312.00651*, 2023. [10](#)
- Yanze Li, Wenhua Zhang, Kai Chen, Yanxin Liu, Pengxiang Li, Ruiyuan Gao, Lanqing Hong, Meng Tian, Xinhai Zhao, Zhenguo Li, et al. Automated evaluation of large vision-language models on self-driving corner cases. *arXiv preprint arXiv:2404.10595*, 2024. [10](#)
- Bin Lin, Zhenyu Tang, Yang Ye, Jiayi Cui, Bin Zhu, Peng Jin, Junwu Zhang, Munan Ning, and Li Yuan. Moe-llava: Mixture of experts for large vision-language models. *arXiv preprint arXiv:2401.15947*, 2024. [3](#)
- Hao Liu, Carmelo Sferrazza, and Pieter Abbeel. Chain of hindsight aligns language models with feedback. *arXiv preprint arXiv:2302.02676*, 2023a. [2](#)
- Zhili Liu, Jianhua Han, Kai Chen, Lanqing Hong, Hang Xu, Chunjing Xu, and Zhenguo Li. Task-customized self-supervised pre-training with scalable dynamic routing. In *AAAI*, 2022. [10](#)
- Zhili Liu, Kai Chen, Yifan Zhang, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, and James Kwok. Geom-erasing: Geometry-driven removal of implicit concept in diffusion models. *arXiv preprint arXiv:2310.05873*, 2023b. [10](#)
- Zhili Liu, Kai Chen, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, and James T Kwok. Task-customized masked autoencoder via mixture of cluster-conditional experts. *arXiv preprint arXiv:2402.05382*, 2024. [2](#), [3](#), [4](#), [10](#)
- Basil Mustafa, Carlos Riquelme, Joan Puigcerver, Rodolphe Jenatton, and Neil Houlsby. Multimodal contrastive learning with limoe: the language-image mixture of experts. *arxiv preprint arxiv:2206.02770*, 2022. [3](#)
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. [1](#), [2](#), [6](#)
- Xianghe Pang, Shuo Tang, Rui Ye, Yuxin Xiong, Bolun Zhang, Yanfeng Wang, and Siheng Chen. Self-alignment of large language models via multi-agent social simulation. In *ICLR Workshop on Large Language Model (LLM) Agents*, 2024. [2](#)
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023. [2](#)
- Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. In *NeurIPS*, 2021. [2](#), [3](#)

- Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne Longpre, Jason Wei, Hyung Won Chung, Barret Zoph, William Fedus, Xinyun Chen, et al. Mixture-of-experts meets instruction tuning: A winning combination for large language models. *arXiv preprint arXiv:2305.14705*, 2023a. [9](#)
- Sheng Shen, Zhewei Yao, Chunyuan Li, Trevor Darrell, Kurt Keutzer, and Yuxiong He. Scaling vision-language models with sparse mixture of experts. *arxiv preprint arxiv:2303.07226*, 2023b. [3](#)
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 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, 2023. [6](#)
- Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and Dilip Krishnan. Stablerep: Synthetic images from text-to-image models make strong visual representation learners. *arXiv preprint arXiv:2306.00984*, 2023. [10](#)
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. [1](#)
- Yaqing Wang, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, and Jianfeng Gao. Adamix: Mixture-of-adaptations for parameter-efficient model tuning. *arxiv preprint arxiv:2210.17451*, 2022. [3](#)
- Yibo Wang, Ruiyuan Gao, Kai Chen, Kaiqiang Zhou, Yingjie Cai, Lanqing Hong, Zhenguo Li, Lihui Jiang, Dit-Yan Yeung, Qiang Xu, and Kai Zhang. Detdiffusion: Synergizing generative and perceptive models for enhanced data generation and perception. *arXiv preprint arXiv:2403.13304*, 2024. [10](#)
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021. [1](#), [2](#), [4](#)
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022. [2](#), [6](#)
- Ted Zadouri, A. Ustun, Arash Ahmadian, Beyza Ermics, Acyr Locatelli, and Sara Hooker. Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning. *arxiv preprint arxiv:2309.05444*, 2023. [3](#)
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. In *NeurIPS*, 2024. [7](#)
- Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023. [1](#)