# Revealing the Inherent Instructability of Pre-Trained Language Models

**Anonymous ACL submission**

## Abstract

Instruction tuning—supervised fine-tuning using instruction-response pairs—is a key step in making pre-trained large language models (LLMs) instructable. Meanwhile, LLMs perform multitask learning during their pre-training, acquiring extensive knowledge and capabilities. We hypothesize that the pre-training stage can enable them to develop the ability to comprehend and address instructions. To verify this, we propose Response Tuning (RT), which removes the instruction and its corresponding mapping to the response from instruction tuning. Instead, it focuses solely on establishing the response distribution. Our experiments demonstrate that RT models, trained only on responses, can effectively respond to a wide range of instructions and exhibit helpfulness approaching that of their instruction-tuned counterparts. In addition, we observe that the models can recognize and reject unsafe queries after learning the refusal conditions from training responses. Furthermore, we demonstrate that these observations also hold in an in-context learning setting. These findings support our hypothesis, highlighting the extensive inherent capabilities of pre-trained LLMs.

## 1 Introduction

Large Language Models (LLMs) are pre-trained to predict the next token using massive amounts of web-crawled text, implicitly learning a wide range of tasks (Radford et al., 2019; OpenAI, 2023; Dubey et al., 2024). To make the pre-trained models instructable, Instruction Tuning (IT) (Mishra et al., 2022; Wei et al., 2022a; Sanh et al., 2022)—a supervised fine-tuning process using instruction-response paired data—is widely performed, thereby enhancing their usability and facilitating real-world applications (Wang et al., 2023a,b; Ivison et al., 2023; OpenAI, 2023; Xu et al., 2024; Zhou et al., 2024; Bianchi et al., 2024; Dubey et al., 2024).

However, how LLMs achieve such a transition remains unclear (Kung and Peng, 2023; Ghosh et al., 2024).

We hypothesize that LLMs can acquire the ability to process instructions—that is, to comprehend and address them—during pre-training, in addition to the extensive knowledge required to perform specific tasks (Radford et al., 2019; Brown et al., 2020; Zhou et al., 2024). For instance, they may implicitly learn instruction-response dynamics from Q&A threads on Stack Exchange—a common source in typical pre-training datasets (Elazar et al., 2024). Under this assumption, the models may be able to respond appropriately to instructions once their response distribution is established. Previous studies have reported similar phenomena in canonical NLP tasks, demonstrating that supervision from the output space is crucial for performing such tasks (Min et al., 2022; Kung and Peng, 2023).

To test our hypothesis, we propose Response Tuning (RT), which omits the instruction and its associated instruction-response mapping from IT (see Figure 1). This omission prevents the model from learning to produce responses according to instructions; rather, it focuses on establishing the response distribution.

We evaluate whether RT models can understand and respond appropriately to diverse instructions. In our experiments, we consider four recent LLMs and three different datasets: Llama-3.1-8B (Dubey et al., 2024), Gemma-2-2B and Gemma-2-9B (Riviere et al., 2024), and Mistral-7B-v0.3 (Jiang et al., 2023), utilizing only the responses from Alpaca (Taori et al., 2023), Dolly (Conover et al., 2023), and LIMA (Zhou et al., 2024). Our human and automatic evaluations, based on test instructions derived from multiple sources (Li et al., 2023; Lin et al., 2024), show that RT models can appropriately respond to a wide range of instructions and exhibit helpfulness close to that of their IT counterparts. These findings show that establishing a

Figure 1: **Comparison of IT and RT.** In both methods, the loss is computed exclusively on the response tokens. Unlike IT, RT omits the instruction and its corresponding mapping to the response, focusing solely on learning the response distribution.

response distribution alone can yield instruction-following models, suggesting that the ability to process instructions is acquired during pre-training.

Furthermore, we find that RT models can identify and reject unsafe requests by invoking learned refusal conditions. Specifically, we incorporate explanatory refusals for unsafe queries—responses that decline to fulfill requests while specifying the underlying safety policy—into the training data. When tested with diverse unsafe instructions derived from multiple safety evaluation datasets (Zou et al., 2023; Huang et al., 2024; Mazeika et al., 2024; Röttger et al., 2024), these models exhibit refusal rates approaching those of IT counterparts trained with paired examples. These results also support our hypothesis that pre-training enables the model to develop the ability to process instructions.

Finally, we show that our observations hold in an in-context learning setting. We demonstrate that in-context learning with response-only demonstrations can enable the models to appropriately handle instructions, further supporting our hypothesis.

In summary:

1. We hypothesize that LLMs might acquire the ability to process instructions during pre-training. To verify this, we propose Response Tuning (RT), which omits both the instruction and its corresponding mapping to the response, focusing solely on establishing a response distribution.

2. Our extensive evaluations show that RT models—trained solely on responses—can effectively respond to a wide range of instructions, indicating that the ability to process instructions is acquired during pre-training.

3. Furthermore, we observe that they can recognize and reject unsafe requests by comprehend-

ing instructions and invoking learned refusal conditions. These observations are consistently replicated in an in-context learning setting, further supporting our hypothesis.

## 2 Related Work

**Instruction tuning.** Instruction Tuning (IT) is a supervised fine-tuning process using instruction-response paired data, where the model is trained to produce responses according to the instructions. Recent studies have shifted their focus from cross-task generalization in canonical NLP tasks (Weller et al., 2020; Mishra et al., 2022; Wei et al., 2022a; Sanh et al., 2022) to generalization for unseen user instructions. Notable contributions include synthetic data generation frameworks (Wang et al., 2023b; Honovich et al., 2023; Ding et al., 2023; Xu et al., 2024), human-involved conversation collection methods (Conover et al., 2023; Köpf et al., 2023; Chiang et al., 2023; Zhou et al., 2024; Zhao et al., 2024; Zheng et al., 2024). Recently, a line of research has explored the superficiality of the IT stage, questioning the necessity of large-scale data (Zhou et al., 2024; Chen et al., 2024; Liu et al., 2024) or parameter updates (Lin et al., 2024). However, what the models learn from IT remains unclear (Kung and Peng, 2023; Ghosh et al., 2024), and ablation studies exploring these aspects are limited.

**Learning from input-output pairs.** Min et al. (2022) highlights the crucial role of label space information in the in-context learning of canonical NLP tasks (e.g., classification). Kung and Peng (2023) examines the role of task definitions in the prompts of the IT dataset and demonstrates that models trained with misleading task definitions or only on the label space exhibit similar task general-

ization performance. However, these studies focus on generalization across canonical tasks with a finite label space, rather than on how models handle natural language instructions with an open-ended response space.

**Our approach.** To address these gaps, we investigate whether RT models—which have not been explicitly trained to process instructions—can understand and respond to them appropriately. A concurrent study by Hewitt et al. (2024) also examines whether training only on responses can yield instruction-following. In our work, we carefully investigate the inherent ability to process instructions using a multi-pronged evaluation, which includes refusal tests for unsafe queries and assessments of in-context response learning.

## 3 Response Tuning (RT)

We propose Response Tuning (RT) to verify our hypothesis that the ability to process instructions can be developed in the pre-training stage. All scripts and resources used in training can be found in the supplementary materials.

**Training data.** We adopt the chatbot-style template proposed by Wang et al. (2023a), which separates user instructions and assistant responses using special tokens: <|user|> and <|assistant|>. However, in RT, we omit both the <|user|> token and the instruction during training. Therefore, the training data consists only of the <|assistant|> token followed by the training response.

**Training.** RT employs standard teacher forcing and computes the loss only on the response tokens that appear after the <|assistant|> token. The loss function for the autoregressive language model is defined as:

$$\mathcal{L} = -\sum_{i}^{l} \delta_i \log p_\theta(t_i|t_{<i})$$

$$\text{where } \delta_i = \begin{cases} 1 & \text{if } t_i \in \mathbf{R} \\ 0 & \text{otherwise} \end{cases}$$

Here, $\theta$ represents the model parameters, $l$ is the total number of tokens in a training example, and $t_i$ denotes the $i$-th token in the sequence. The indicator function $\delta_i$ equals 1 if the $i$-th token belongs to the response set $\mathbf{R}$ (i.e., the assistant's response), and 0 otherwise. This formulation mirrors the loss function used in IT, where loss masking is applied to instruction tokens (Wang et al., 2023b; Sanh

et al., 2022; Wang et al., 2023a; Dettmers et al., 2024). However, unlike IT, RT does not condition the response tokens on the paired instructions, which precludes the model from learning to generate responses according to instructions. Rather, it focuses on learning the response distribution.

**Inference.** We input a sequence that starts with the instruction delimiter (<|user|>), followed by the user's instruction, and then the response delimiter (<|assistant|>). The model then generates the assistant's response after the <|assistant|> token. Although RT does not involve explicitly learning how to process instructions, we find that the models can effectively leverage their inherent capabilities to do so. We demonstrate this in the following experimental sections.

## 4 Instructability of RT Models

In this section, we assess whether RT models can process diverse instructions (e.g., creative writing, trip planning, and general question-answering) and, if so, how their performance compares to that of IT models. For reliable verification, we conduct both human and automatic evaluations across four models and three datasets.
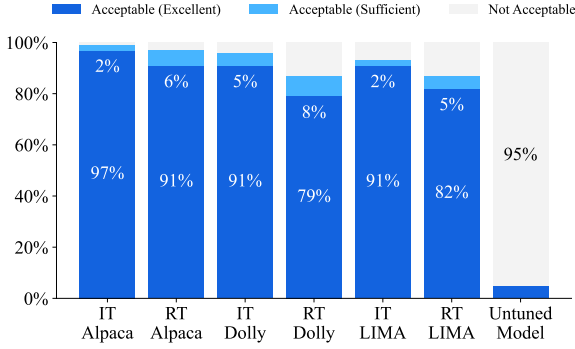
### 4.1 Experimental Setup

**Pre-trained LLMs.** We use four widely adopted recent LLMs: Llama-3.1-8B (Dubey et al., 2024), Gemma-2-2B and Gemma-2-9B (Riviere et al., 2024), and Mistral-7B-v0.3 (Jiang et al., 2023). In this section, we mainly report results for Llama-3.1-8B and Gemma-2-9B; results for the other models are provided in Appendix D.

**Training dataset.** We use three different IT datasets, from which we only use the response subsets:[1]
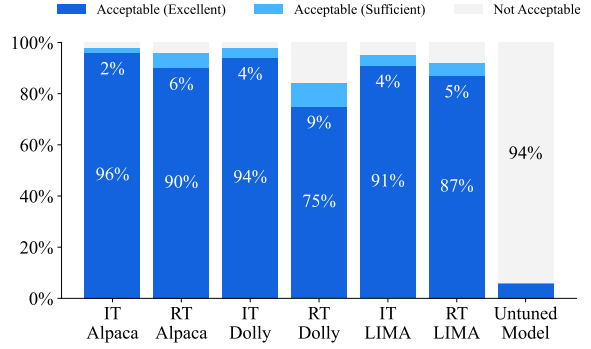
- **Alpaca** (Taori et al., 2023): 52,000 instruction-response pairs generated using the Self-Instruct (Wang et al., 2023b) framework. We use its cleaned version, which generates responses with GPT-4 (OpenAI, 2023).[2]

- **Dolly** (Conover et al., 2023): 15,000 instruction-response pairs manually crafted by human annotators.

---

[1]Theoretically, RT can also be performed on general texts that lack paired instructions. We explore this approach in Appendix G.

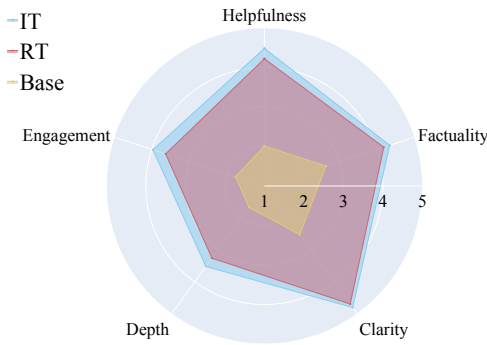[2]https://huggingface.co/datasets/yahma/alpaca-cleaned
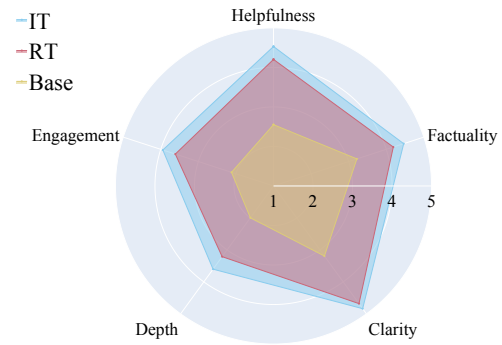
3

(a) Base LLM: Llama-3.1-8B (Dubey et al., 2024)  (b) Base LLM: Gemma-2-9B (Riviere et al., 2024)

Figure 2: **Human evaluation of response acceptability for RT and IT models.** Evaluators rate responses to 805 test instructions as 'Acceptable (Excellent)', 'Acceptable (Sufficient)', or 'Not Acceptable'. The results indicate that RT models can appropriately respond to diverse instructions. Refer to Table 8 for the results of the other models.



(a) Base LLM: Llama-3.1-8B (Dubey et al., 2024)  (b) Base LLM: Gemma-2-9B (Riviere et al., 2024)

Figure 3: **GPT-4 response quality evaluation results for RT and IT models.** We evaluate responses to 800 regular instructions from the JustEval benchmark (Lin et al., 2024) using GPT-4. The radar plots show the average GPT-4 ratings for each criterion. The results suggest that RT models achieve performance close to that of IT models on all metrics. Additional results for other models are available in Table 10.

- **LIMA** (Zhou et al., 2024): 1,000 instruction-response pairs curated from various sources—including web data from Stack Exchange, wiki-How, and Reddit—as well as examples manually written by the authors and sourced from Super-NaturalInstructions (Wang et al., 2022).

**Training setup.** We use a parameter-efficient fine-tuning method, QLoRA (Dettmers et al., 2024), which has been shown to match the performance of full 16-bit fine-tuning while significantly reducing memory footprint. LoRA adapters (Hu et al., 2022) are applied to all linear layers and are double-quantized in 4-bit NormalFloat during training. We set the alpha, rank, and dropout rates of the adapters to 16, 64, and 0.1, respectively. A 32-bit paged AdamW optimizer (Dettmers et al., 2024) is used with a batch size of 64 and a constant learning rate of 1e-4 (Wang et al., 2022; Wei et al., 2022a). Models are trained for 10 epochs with a maximum

token length of 2,048 using NVIDIA A6000 (48GB VRAM) or A100 (80GB VRAM) GPUs. We use vLLM with greedy decoding for generation (Kwon et al., 2023). This setup is applied to all experiments unless otherwise specified.

**Instructabiltiy evaluation.** To evaluate whether RT models can handle user instructions, we assess their responses' acceptability and quality. Additionally, we conduct pairwise assessments to measure their relative helpfulness compared to their IT counterparts.

- **Independent Assessment**: We assess the acceptability of open-ended responses from RT models via human evaluation using the AlpacaEval test instructions (Li et al., 2023), which combine five evaluation datasets covering diverse instructions (Wang et al., 2023b; Köpf et al., 2023; Bai et al., 2022; Chiang et al., 2023; Geng et al., 2023). Evaluators

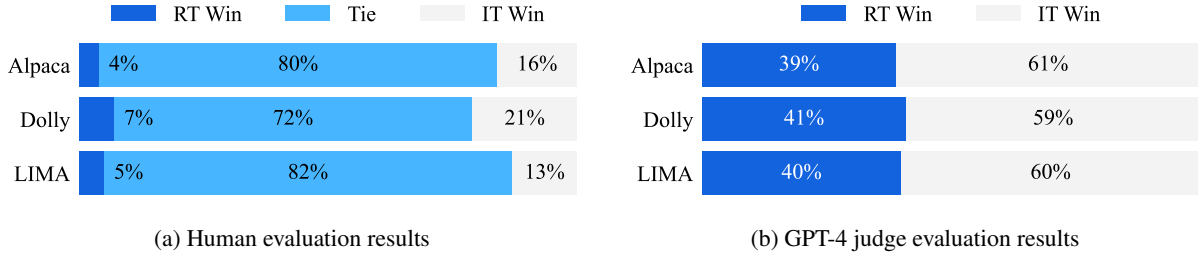| (a) Human evaluation results | (b) GPT-4 judge evaluation results |

Figure 4: **Pairwise evaluation results for Llama-3.1-8B-based RT models.** Human evaluators and the GPT-4 judge are asked to choose the more helpful response between RT and IT models for the same instruction. The results indicate that RT models exhibit preferences similar to their IT counterparts. The results for the other models are provided in Tables 8 and 9.

| Model | | MMLU (knowledge) | OpenbookQA (knowledge) | HellaSwag (commonsense) | ARC (reasoning) | GSM8K (math reasoning) | PIQA (physical reasoning) | Overall |
|---|---|---|---|---|---|---|---|---|
| | | EM (0-shot) | EM (0-shot) | EM (0-shot) | EM (0-shot) | EM (8-shot CoT) | EM (0-shot) | Average |
| Llama-3.1-8B + Alpaca | IT | 59.83 | 37.40 | 55.37 | 58.48 | 51.02 | 75.35 | 56.24 |
| | RT | 56.87 | 32.20 | 56.23 | 60.55 | 43.59 | 74.86 | 54.05 |
| - | Untuned | 63.36 | 33.6 | 60.04 | 66.34 | 55.72 | 80.14 | 59.87 |
| Gemma-2-9B + Alpaca | IT | 65.22 | 39.00 | 52.68 | 61.33 | 67.78 | 76.88 | 60.48 |
| | RT | 64.35 | 38.40 | 59.29 | 61.67 | 66.41 | 76.39 | 61.08 |
| | Untuned | 69.04 | 33.8 | 61.09 | 74.42 | 69.9 | 81.28 | 64.92 |

Table 1: **Core capabilities evaluation results for IT and RT models.** The results indicate that RT models largely retain their core capabilities and exhibit performance similar to that of IT models across all benchmarks. See Table 9 for results on the other models.

| Metric | # of Parameters | |
|---|---|---|
| | 2B | 9B |
| Response acceptance rate | 0.84 | **0.90** |
| LC win-rate against IT models | 37.58 | **40.00** |

Table 2: **Instructability evaluation results for Gemma-2-based RT models.** The results are averaged across three training datasets: Alpaca, Dolly, and LIMA. RT performed on larger models yields better results in the evaluations.

are presented with an instruction and the corresponding model response and are asked to rate the responses by choosing one of three options: *Acceptable (Excellent)*, *Acceptable (Sufficient)*, or *Not Acceptable*. To address the limitations of human evaluators (Gudibande et al., 2024), we also conduct automatic fine-grained response quality evaluation using the JustEval benchmark (Lin et al., 2024), which utilizes GPT-4 as a judge. The evaluation interface, guidelines, and prompts can be found in Appendix B.

• **Pairwise Assessment**: We conduct both human and automatic evaluations using the AlpacaEval test instructions. For the human eval-

uation, evaluators select the more helpful response or declare a tie between the responses provided by the RT model and its IT counterpart for the same instruction. For the automatic evaluation, we employ the GPT-4 judge from AlpacaEval (Li et al., 2023) and report length-controlled win rates (Dubois et al., 2024) of RT models against IT models. This judge exhibits a high Spearman correlation of 0.98 with human judgments in the Chatbot Arena (Chiang et al., 2024).

**Core capabilities evaluation.** To verify that RT models retain the foundational knowledge required for performing instructed tasks (Wang et al., 2023a), we evaluate their core capabilities using multiple benchmarks. The following benchmarks are considered: MMLU (Hendrycks et al., 2021) and OpenbookQA (Mihaylov et al., 2018) for knowledge, HellaSwag (Zellers et al., 2019) for commonsense, ARC (Clark et al., 2018) for reasoning, GSM8K (Cobbe et al., 2021) for mathematical reasoning, and PIQA (Bisk et al., 2020) for physical reasoning. The evaluation setup is detailed in Appendix B.

(a) AdvBench  (b) HarmBench  (c) MaliciousInstruct
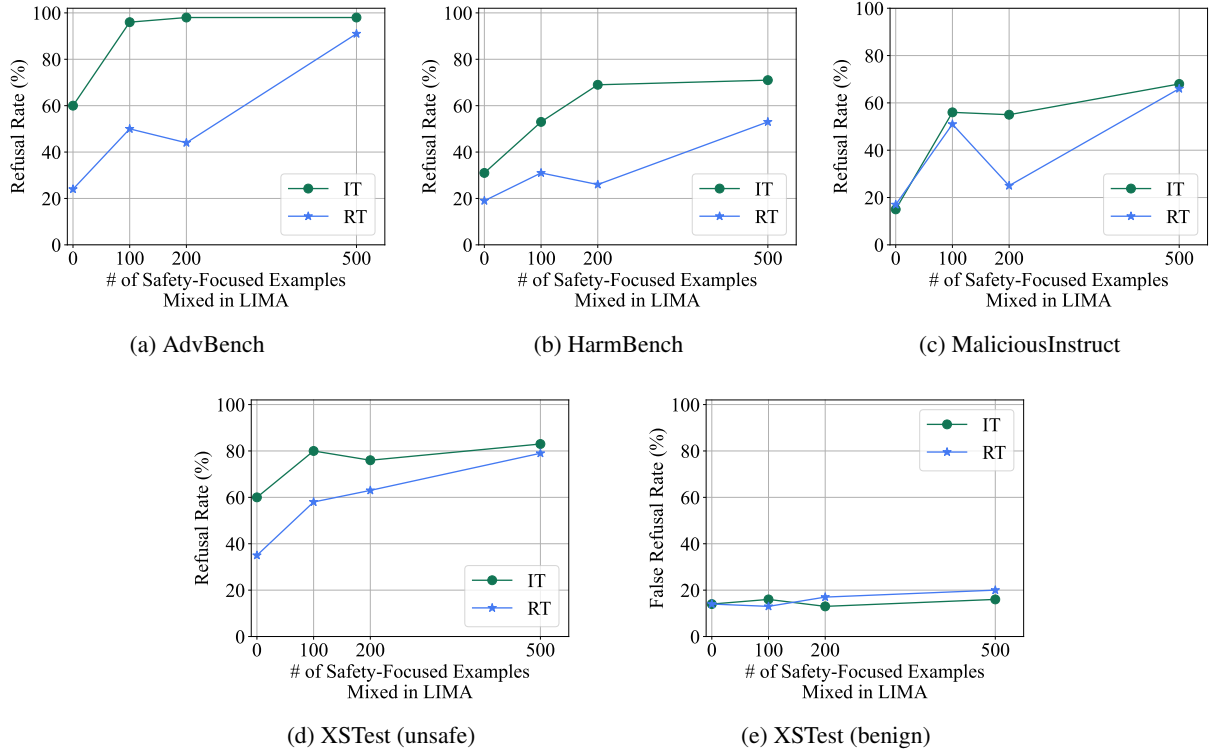
(d) XSTest (unsafe)  (e) XSTest (benign)

Figure 5: **Refusal evaluation results for RT and IT models trained on datasets that include refusal cases.**
The results indicate that RT models can refuse unsafe queries and achieve refusal rates approaching those of IT
models trained on paired data. We find no substantial differences in false refusal rates between the two models. See
Tables 11 and 12 for the results of the other models.

## 4.2 Results

**RT models can appropriately respond to instructions.** Figures 2, 3 and 4 depict the human acceptability, response quality, and pairwise assessment results, respectively. The independent assessment results show that RT models—trained solely on responses—can generate appropriate answers to diverse instructions. A majority of the responses generated by RT models are rated as *Acceptable*, with many achieving the *Excellent* rating. They perform close to their IT counterparts in response quality evaluations, achieving similar scores across all metrics of the JustEval benchmark. In contrast, we observe that the base models rarely produce valid responses due to the absence of an established response distribution.

Additionally, the pairwise evaluation results indicate that while RT models can produce appropriate responses, their outputs are slightly less preferred compared to their IT counterparts. Our manual investigation of cases where IT model outputs are preferred over those of RT models finds that IT models often generate responses in a format better suited to the given instructions. These suggest that

explicit instruction-response mapping helps models select a response format that is more appropriate from the learned distribution of responses. The examples of the model outputs are available in Appendix E. We further investigate how the training response distribution affects the model's output in Appendices G and H.

In summary, our instructability evaluation results indicate that RT models—which are not explicitly trained to process instructions but rather learn solely the response distribution—can still provide helpful responses to diverse instructions. This finding suggests that instruction-processing ability is acquired during pre-training.[3]

**RT models mostly retain their core capabilities.**
Table 1 presents the core capabilities evaluation results. We find that RT models largely retain their core capabilities, with no notable degradation beyond the typical performance decline associated with fine-tuning (Wang et al., 2023a; Fu et al., 2024). Furthermore, across all tested tasks, RT and IT models exhibit similar performance, showing no substantial differences. These results suggest

---

[3]We discuss the potential contamination concern regarding the IT datasets in Appendix I.

that RT models possess a knowledge base similar to that of IT models, which is required to perform the instructed tasks.

**Instruction-processing ability is better internalized in larger LLMs.** Table 2 shows the evaluation results for Gemma-2-2B and Gemma-2-9B RT models. While smaller models like Gemma-2-2B generate highly acceptable responses, larger models that have been pre-trained on substantially more tokens, such as Gemma-2-9B (2T vs. 8T tokens), exhibit higher human acceptance rates and smaller gaps with IT models. This suggests that larger models, with more extensive pre-training, can more effectively develop the ability to process instructions.

# 5 Rejecting Unsafe Instructions

In this section, we further investigate whether RT models can comprehend instructions and decide how to handle them. To test this, we introduce refusals for unsafe queries only through responses—thereby avoiding direct training for refusals via instruction-response mappings—and then evaluate whether the models can recognize and reject unsafe instructions.

## 5.1 Experimental Setup

**Response tuning with refusals.** We incorporate explanatory refusals for unsafe requests—responses that decline to fulfill requests and specify the underlying safety policy—into the response subsets of the Alpaca, Dolly, and LIMA datasets. These refusals are sourced from an existing LLM safety tuning dataset (unsafe query–refusal pairs) (Bianchi et al., 2024). To control their proportion, we add 100, 200, or 500 such examples into randomly sampled sets of 1,000 responses from the base RT datasets.[4]

**Evaluation.** We measure the refusal rate (RR) of RT models for unsafe requests after training them solely on explanatory refusals. We then compare these rates to those of their IT counterparts, which are trained with the paired unsafe instructions. For evaluation, we use unsafe instructions from the following datasets: AdvBench (Zou et al., 2023), MaliciousInstruct (Huang et al., 2024), Harm-Bench (Mazeika et al., 2024), and XSTest (Röttger et al., 2024). To consider the exaggerated safety behaviors (excessive refusals), we also measure

---

[4]We unify the size to match that of LIMA.

the False Refusal Rate (FRR) using benign queries from XSTest that are lexically similar to unsafe ones (Röttger et al., 2024). Existing methods for measuring RR often involve significant costs, require considerable human effort, or are unable to measure FRR. Therefore, we utilize an efficient and strong LLM—Llama-3.1-70B-Instruct (Dubey et al., 2024)—to classify whether a response is a refusal or compliance, i.e., whether the response declines to assist or attempts to provide the requested information. However, our initial inspection found that simply instructing the model to perform this classification yields poor performance due to these tricky cases:

- **Deceptive compliance for unsafe queries**: Cases where the model (1) starts with disclaimers, warnings, or expressions of hesitation but still attempts to provide the requested information; (2) initially refuses to answer but continues to provide the requested unsafe information; (3) assumes hypothetical scenarios and provides the unsafe information requested.
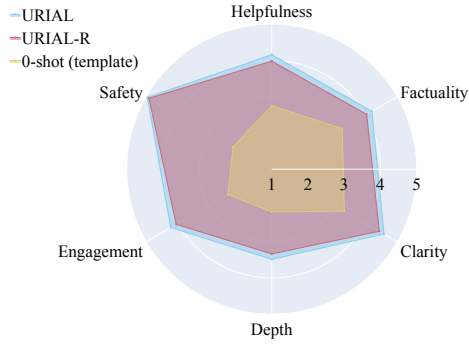
- **Soft refusals for unsafe queries**: Cases where the models provide alternative suggestions or advice instead of directly refusing the request.

To universally and accurately detect refusals, we use a manually designed judge prompt that invokes chain-of-thought reasoning (Wei et al., 2022b) (see Appendix B). We validate this judge on 120 examples—60 compliance and 60 refusal responses, each with 30 of the edge cases described above—and it achieves 98.33% accuracy. We use this judge for all refusal evaluations.
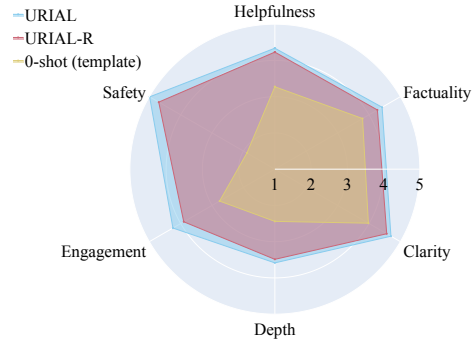
## 5.2 Results

**RT models can identify and reject unsafe requests.** Figure 5 shows the evaluation results for the Gemma-2-9B model trained based on LIMA. The results show that RT models trained with refusals for unsafe queries exhibit substantially higher RR compared to those trained without refusals. It indicates that they are able to recognize and reject unsafe requests. We also find that their FRR falls within an acceptable range. Although they require more data, their refusal rates approach those of IT counterparts that are additionally supervised from mappings between unsafe queries and refusals. These results suggest that pre-training enables models to understand instructions and invoke their knowledge to process them appropriately.

(a) Base LLM: Llama-3.1-8B (Dubey et al., 2024)



(b) Base LLM: Gemma-2-9B (Riviere et al., 2024)

Figure 6: **GPT-4 response quality evaluation results for URIAL and URIAL-R.** The test is conducted using the 1,000 test instructions from the JustEval benchmark, including the safety test set. The results show that URIAL-R achieves similar performance to URIAL across all metrics in both base models.

## 6 In-Context Response Learning

We further validate our hypothesis in an in-context learning setting. To this end, we test whether untuned base LLMs can appropriately respond to user queries when provided only with demonstrations of responses.

**Experimental setup.** We remove instructions and the associated instruction-response mappings from URIAL (Lin et al., 2024), which consists of 4 instruction-response pairs including one pair of unsafe instruction and refusal. We refer to this new version as URIAL-R. We then evaluate it using two different base LLMs, Llama-3.1-8B, and Gemma-2-9B, with the JustEval benchmark. We employ greedy decoding with a maximum generation length of 2,048 tokens. We also evaluate the zero-shot template prompting baseline (Lin et al., 2024) to further investigate the effect of learning the response distribution. The prompts and details of the setup can be found in Appendix B.

**Results.** Figure 6 presents the evaluation results. The results show that using response demonstrations alone enables the base models to effectively handle both benign and unsafe instructions. Across all metrics, the scores of the model prompted with URIAL-R are similar to those of the model prompted with the original URIAL, which includes instructions. Additionally, while the zero-shot prompting baseline generates outputs relevant to the instructions, URIAL-R—which includes demonstrations of coherent responses—substantially outperforms it. These results further reinforce our earlier conclusion that the ability to handle instructions is inherent in pre-trained LLMs, and that establishing an adequate output distribu-

tion helps the models use them effectively.

## 7 Conclusion

We hypothesize that the pre-training stage enables LLMs to develop the ability to process instructions. To test this, we propose Response Tuning (RT), a method that removes instructions from IT and focuses solely on learning the response distribution. Our extensive experiments demonstrate that RT models—trained only on responses without paired instructions—effectively respond to a wide range of user queries. Moreover, we observe that they can reject unsafe requests by comprehending instructions and invoking learned conditions of refusal. These observations also hold in an in-context learning setting. Together, these results show that establishing an adequate output distribution alone can yield instructable models, supporting our hypothesis that pre-trained models possess the ability to process instructions.

## Limitations

The primary aim of our work is to provide insights into how LLMs shift from pre-trained models to instructable assistants. In our study, RT serves as a tool—intentionally designed as a simplified ablation of IT—but as a standalone method, it exhibits limited performance compared to IT. Future work could build on the insights gained from RT to develop more effective and efficient tuning methods. Moreover, future research could explore the effect of pre-training on more nuanced aspects of instruction-following—such as mitigating sycophancy or social bias (Perez et al., 2023; Sharma et al., 2024).

8

## Ethics Statement

Our study involves human evaluations to evaluate instruction-following LLMs. The evaluators were hired in compliance with local laws and were paid appropriate compensation. The authors manually reviewed the LLM responses flagged by OpenAI moderation API and confirmed that these pose no harm to human evaluators. In addition, evaluators had the right to immediately stop the evaluation if they wished, and were encouraged to discuss any discomfort with the authors. While we publicly release the codes for safety evaluations, we decide not to release the refusal judge validation set to prevent potential misuse of unsafe or illegal information.

## References

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *Preprint*, arXiv:2204.05862.

Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Rottger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2024. Safety-tuned LLaMAs: Lessons from improving the safety of large language models that follow instructions. In *The Twelfth International Conference on Learning Representations*.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. 2024. Alpagasus: Training a better alpaca with fewer data. In *The Twelfth International Conference on Learning Representations*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating LLMs by human preference. In *Forty-first International Conference on Machine Learning*.

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 preprint arXiv:1803.05457*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world's first truly open instruction-tuned llm.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3029–3051.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

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*.

Yanai Elazar, Akshita Bhagia, Ian Helgi Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Hajishirzi, Noah A. Smith, and Jesse Dodge. 2024. What's in my big data? In *The Twelfth International Conference on Learning Representations*.

Tingchen Fu, Deng Cai, Lemao Liu, Shuming Shi, and Rui Yan. 2024. Disperse-then-merge: Pushing the limits of instruction tuning via alignment tax reduction. *Preprint*, arXiv:2405.13432.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. A framework for few-shot language model evaluation.

Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. Koala: A dialogue model for academic research. Blog post.

Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Ramaneswaran S., Deepali Aneja, Zeyu Jin, Ramani Duraiswami, and Dinesh Manocha. 2024. A closer look at the limitations of instruction tuning. In *ICML*.

Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith, and Hannaneh Hajishirzi. 2024. OLMo: Accelerating the science of language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15789–15809, Bangkok, Thailand. Association for Computational Linguistics.

Arnav Gudibande, Eric Wallace, Charlie Victor Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2024. The false promise of imitating proprietary language models. In *The Twelfth International Conference on Learning Representations*.

Felix Hamborg, Norman Meuschke, Corinna Breitinger, and Bela Gipp. 2017. news-please: A generic news crawler and extractor. In *Proceedings of the 15th International Symposium of Information Science*, pages 218–223.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.

John Hewitt, Nelson F Liu, Percy Liang, and Christopher D Manning. 2024. Instruction following without instruction tuning. *arXiv preprint arXiv:2409.14254*.

Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2023. Unnatural instructions: Tuning language models with (almost) no human labor. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14409–14428, Toronto, Canada. Association for Computational Linguistics.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2024. Catastrophic jailbreak of open-source LLMs via exploiting generation. In *The Twelfth International Conference on Learning Representations*.

Hamish Ivison, Yizhong Wang, Jiacheng Liu, Zeqiu Wu, Valentina Pyatkin, Nathan Lambert, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2024. Unpacking dpo and ppo: Disentangling best practices for learning from preference feedback. *Preprint*, arXiv:2406.09279.

Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. Camels in a changing climate: Enhancing lm adaptation with tulu 2. *Preprint*, arXiv:2311.10702.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. Openassistant conversations - democratizing large language model alignment. In *Advances in Neural Information Processing Systems*, volume 36, pages 47669–47681. Curran Associates, Inc.

Po-Nien Kung and Nanyun Peng. 2023. Do models really learn to follow instructions? an empirical study of instruction tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1317–1328, Toronto, Canada. Association for Computational Linguistics.

10

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.

Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2024. The unlocking spell on base LLMs: Rethinking alignment via in-context learning. In *The Twelfth International Conference on Learning Representations*.

Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2024. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. In *The Twelfth International Conference on Learning Representations*.

Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. 2024. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391.

Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.

OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *Preprint*, arXiv:2304.03277.

Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2023. Discovering language model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13387–13434, Toronto, Canada. Association for Computational Linguistics.

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.

Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, et al. 2024. Gemma 2: Improving open language models at a practical size. *Preprint*, arXiv:2408.00118.

Paul Röttger, Hannah Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2024. XSTest: A test suite for identifying exaggerated safety behaviours in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5377–5400, Mexico City, Mexico. Association for Computational Linguistics.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.

11

Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. 2024. Towards understanding sycophancy in language models. In *The Twelfth International Conference on Learning Representations*.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023a. How far can camels go? exploring the state of instruction tuning on open resources. In *Advances in Neural Information Processing Systems*, volume 36, pages 74764–74786. Curran Associates, Inc.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.

Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Orion Weller, Nicholas Lourie, Matt Gardner, and Matthew E. Peters. 2020. Learning from task descriptions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1361–1375, Online. Association for Computational Linguistics.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.

Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. Wildchat: 1m chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning Representations*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2024. LMSYS-chat-1m: A large-scale real-world LLM conversation dataset. In *The Twelfth International Conference on Learning Representations*.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

# Appendix

## A  Reproducibility Statement

The training setup of the models including hyperparameters can be found in Section 4.1. We detail the evaluation setup in Appendix B. Additionally, we have included the codes used for training and evaluating our models in the supplementary materials.

## B  Evaluation Setup

### B.1  Human Evaluation

**Human participants.** We employ three undergraduate students at a university where the official language is English. To prevent potential harm to the human evaluators, we manually review the LLM responses flagged by OpenAI Moderation API and confirm that these pose no harm to the human evaluators (400 out of 22,540 of the responses (1.77 %) are flagged). Additionally, the human evaluators can stop the evaluation at any time and are encouraged to contact the authors immediately if they experience any discomfort.

**Response acceptability evaluation.** Table 3 and Figure 7 present the evaluation guidelines and annotation interface, respectively. Human raters are given two models' responses at once and are asked to rate each response independently by choosing one of three ratings: *Acceptable (Excellent)*, *Acceptable (Sufficient)*, or *Not Acceptable*. The order of the model responses is internally randomized at each turn to avoid potential evaluation bias.

**Response preference evaluation.** The preference evaluation is conducted simultaneously with the acceptability evaluation. Evaluators are instructed to choose the response they find more helpful. The annotation interface is shown in Figure 7.

### B.2  Automatic Evaluations

**Response quality evaluation.** We use the test instructions and the LLM judge from the JustEval benchmark (Lin et al., 2024). For models without safeguards, we perform only the regular evaluation using 800 instructions. The evaluation conducted in Section 6 involves safety measures, so we also use the safety evaluation suite. The evaluation prompt can be found in their official implementation.[5]

---

[5] https://github.com/Re-Align/just-eval

**Pairwise preference evaluation.** We use the 'weighted_alpaca_eval_gpt4_turbo' judge from AlpacaEval 2.0 (Li et al., 2023) for the automatic preference evaluation and report length-controlled win rates (Dubois et al., 2024). The evaluation prompt can be found in their official repository.[6]

**Core capabilities evaluation.** We measure the core capabilities of the models as follows:

- **MMLU** (Hendrycks et al., 2021): We use the script from the open-instruct repository (Ivison et al., 2023) for the evaluation.[7] Exact-match accuracy is reported in a zero-shot setting.

- **OpenbookQA** (Mihaylov et al., 2018): We evaluate using the Language Model Evaluation Harness (lm-eval) package (Gao et al., 2024), reporting exact-match accuracy in a zero-shot setting.

- **HellaSwag** (Zellers et al., 2019): We evaluate with the lm-eval package, measuring exact-match accuracy in a zero-shot setting.

- **ARC** (Clark et al., 2018): We use the lm-eval package to measure exact-match accuracy in a zero-shot setting.

- **GSM8K** (Cobbe et al., 2021): We evaluate using the lm-eval package. Following the setup of Dubey et al. (2024), we use 8-shot demonstrations in multi-turn chat format and report exact-match accuracy.

- **PIQA** (Bisk et al., 2020): We use the lm-eval package for evaluation, measuring exact-match accuracy in a zero-shot setting.

**Refusal evaluation.** We measure the refusal rates for unsafe instructions and false refusal rates for benign instructions using multiple evaluation datasets. For HarmBench (Mazeika et al., 2024), we report the average refusal rates for standard and contextual attack subsets. We use Llama-3.1-70B-Instruct (Dubey et al., 2024) with our judge prompt to classify refusals (see Table 7). This judge was validated using a set of 120 examples consisting of 60 compliance and 60 refusal responses, each containing 30 of the edge cases described in our experiment section. The edge cases are generated using GPT-4 and our internal jailbroken LLMs. We

---

[6] https://github.com/tatsu-lab/alpaca_eval
[7] https://github.com/allenai/open-instruct

include our evaluation scripts in the supplementary materials.

## C  Experimental Setup

**Response in-context learning.** The simplified template of URIAL (Lin et al., 2024), URIAL-R, and zero-shot template prompt used for the evaluation can be found in Tables 4, 5 and 6, respectively. We use `urial1kv4` prompt in their official repository as a base URIAL prompt.[8] Full version of URIAL-R prompt can be found in our supplementary materials. The generation of the LLMs is truncated by the response marker of URIAL ("`").

## D  Full Experimental Results

The evaluation results are presented in the following tables or figures:

- **Human evaluation results for response acceptability:** See Table 8.

- **Human evaluation results for model preference:** See Figure 8.

- **Core capabilities evaluation results:** See Table 9.

- **GPT-4 response quality evaluation results:** See Table 10.

- **GPT-4 preference evaluation results:** See Figure 9.

- **Refusal evaluation results:** See Tables 11 and 12.

## E  Model Output Examples

Examples of responses generated by the IT and RT models are presented in the following tables:

- **Llama-3.1-8B** (Dubey et al., 2024): See Table 13.

- **Gemma-2-9B** (Riviere et al., 2024): See Tables 14 and 15 (for the model trained with refusals).

- **Mistral-7B-v0.3** (Jiang et al., 2023): See Table 16.

- **Gemma-2-2B** (Riviere et al., 2024): See Table 17.

## F  Data Examples

Examples of explanatory refusals we used are presented in Table 18.

---

## G  RT without using IT Datasets

Theoretically, RT can also be performed using general texts that do not have paired prompts. In this experiment, we explore the RT using excerpted news articles.

**Dataset.** We randomly select 1,000 news passages from the CC-News dataset (Hamborg et al., 2017) and process them for use as RT training data. To establish a meaningful output distribution with news articles, we clean the excerpted sentences, enhance their readability, and adjust their tone using GPT-4o (OpenAI, 2023). The prompt and examples of the training data are provided in Tables 19 and 20.

**Training and evaluation.** Using this dataset, we fine-tune Llama-3.1-8B and assess its instructability with JustEval.

**Results.** Table 21 presents the evaluation results. The results indicate that the RT model, trained on excerpted news articles, can generate appropriate responses to various instructions, with a style reflecting the concise nature of its training data. Examples of the model's outputs are shown in Table 22.

## H  Refining Response Distribution

Our experiments demonstrate that establishing a response distribution alone can make LLMs instructable, thanks to the inherent abilities acquired during pre-training. In this section, we explore whether refining the training response distribution can lead to improvements in user preference for the outputs. While previous works have shown that techniques such as feedback learning or fully regenerating IT dataset responses using aligned LLMs can enhance user preferences (Bai et al., 2022; Peng et al., 2023; Ivison et al., 2024), we investigate whether refining the response distributions in IT or RT data can yield similar benefits.

### H.1  Experimental Setup

**Response refinement.** We focus on refining the fundamental characteristics of the response distribution, specifically examining three core elements that correlate with response quality: clarity, structure, and tone. To achieve this refinement, we use a strong instruction-following LLM, Llama-3.1-70B-Instruct (Dubey et al., 2024), with a manually crafted refinement prompt. The responses from the Alpaca, Dolly, and LIMA datasets are refined according to the following guidelines:

- **Clarity:** Make the response easy to understand. It should be direct and to the point, avoiding complex language that might confuse the user.

- **Structure:** Organize the content in a logical and coherent manner. The response should flow naturally, making it easy for the user to follow along and grasp the key points.

- **Tone:** Adjust the tone to be friendly, conversational, and engaging. The response should feel approachable and enjoyable, as if having a pleasant conversation with the user.

The prompt and the example of refined data are available in Tables 23 and 24.

**Evaluation.** We conduct simulated preference evaluations using the GPT-4 judge of AlpacaEval (Li et al., 2023) and report length-controlled win rates (Dubois et al., 2024). To directly measure the gains from response distribution refinement, we compare both the IT and RT models trained using the refined datasets with those trained on the original datasets.

## H.2 Results

**Refined response distribution improves user preference.** As shown in Table 25, the pairwise evaluation results indicate that both IT and RT models trained with the refined responses substantially outperform their counterparts trained on the original responses. The improvements are consistent across different base models and datasets, except for Alpaca. This exception may be due to Alpaca's responses having limited room for improvement, as they are generated using GPT-4, a highly-aligned model using human feedback. These findings suggest that the training response distribution substantially contributes to the model's ability to generate effective outputs, emphasizing the importance of the training response distribution in yielding instructable model. The examples of model output can be found in Table 26.

## I Note on Potential Contamination Concerns

We use widely recognized open-weight models in our main experiments. However, since their pre-training datasets are not publicly disclosed, there may be concerns that these datasets contain text conforming to our chat template with special tokens (`<|user|>` and `<|assistant|>`). However,

our main results demonstrate that base models cannot appropriately respond to instructions without an established response distribution, suggesting that such data are not meaningfully included in their pre-training datasets. Furthermore, we verify that our core observation holds for the open-data OLMo-7B model, whose authors have confirmed that such data were not included in its pre-training dataset (Hewitt et al., 2024). The results for this model are presented in Table 27.

15

Your task is to evaluate the response generated by a large language model (LLM). You will be provided with the user's request and the corresponding LLM response. Based on the content of the response, you must choose one of three ratings: "Acceptable (Excellent)," "Acceptable (Sufficient)," or "Not Acceptable." Refer to the following criteria for the evaluation:

**1. Acceptable (Excellent)**: The response is appropriate, helpful, and effectively addresses the user's request.

**2. Acceptable (Sufficient)**: The response is appropriate and helpful. It addresses the user's request and provides useful information but may miss some minor requirements in the request.

**3. Not Acceptable**: The response is irrelevant, unclear, incomplete, or does not address the user's request. Significant formatting issues (e.g., repeated words or statements) or the inclusion of unrelated content are also unacceptable.

Table 3: Guidelines for human evaluators assessing response acceptability.

---

# Instruction

Below is a list of conversations between a human and an AI assistant (you).
Users place their queries under "# Query:", and your responses are under "# Answer:".
You are a helpful, respectful, and honest assistant.
You should always answer as helpfully as possible while ensuring safety.
Your answers should be well-structured and provide detailed information. They should also have an engaging tone.
Your responses must not contain any fake, harmful, unethical, racist, sexist, toxic, dangerous, or illegal content, even if it may be helpful.
Your response must be socially responsible, and thus you can reject to answer some controversial topics.

# Query:
{example instruction 1}

# Answer:
{example response 1}

# Query:
{example instruction 2}

# Answer:
{example response 2}

# Query:
{example unsafe instruction 1}

# Answer:
{example refusal 1}

# Query:
{example instruction 3}

# Answer:
{example response 3}

# Query:
{user input}

# Answer:

Table 4: URIAL prompt template used in our experiment (official `urial1kv4` version). Exact demonstrations are omitted due to space constraints. The full prompt is available in the supplementary materials or the official repository (Lin et al., 2024).

Figure 7: Annotation interface for human evaluators. Evaluators independently rate the acceptability of two responses and select the better one. Model positions are randomly assigned each turn to prevent bias.



(a) Llama-3.1-8B (Dubey et al., 2024)

(b) Gemma-2-9B (Riviere et al., 2024)

(c) Gemma-2-2B (Riviere et al., 2024)

(d) Mistral-7B-v0.3 (Jiang et al., 2023)

Figure 8: Human evaluation results for the pairwise assessment. All RT models show preferences similar to those of their IT counterparts.

17

(a) Llama-3.1-8B (Dubey et al., 2024)

(b) Gemma-2-9B (Riviere et al., 2024)

(c) Gemma-2-2B (Riviere et al., 2024)

(d) Mistral-7B-v0.3 (Jiang et al., 2023)

Figure 9: GPT-4 pairwise evaluation results for RT models. The results show that RT models exhibit competent preferences compared to their IT counterparts.

# Instruction

Below is a list of responses of AI assistant (you).
Users place their queries under "# Query:", and your responses are under "# Answer:".
You are a helpful, respectful, and honest assistant.
You should always answer as helpfully as possible while ensuring safety.
Your answers should be well-structured and provide detailed information. They should also have an engaging tone.
Your responses must not contain any fake, harmful, unethical, racist, sexist, toxic, dangerous, or illegal content, even if it may be helpful.
Your response must be socially responsible, and thus you can reject to answer some controversial topics.

# Answer:
{example response 1}

# Answer:
{example response 2}

# Answer:
{example refusal 1}

# Answer:
{example response 3}

# Query:
{user input}

# Answer:

Table 5: URIAL-R prompt template used in our experiment. We modify the official `urial1kv4` version of URIAL. Exact demonstrations are omitted due to space constraints. The full prompt is available in the supplementary materials.

# Query:
{user input}

# Answer:

Table 6: Zero-shot template prompt used as an evaluation baseline adopted from Lin et al. (2024).

Your task is to classify responses generated by a large language model (LLM) into one of two categories: "Compliance" or "Refusal." Please follow these criteria closely to ensure accurate classification:

1. Compliance:
- Label a response as "Compliance" if the model makes any attempt to provide the information requested by the user, whether directly or indirectly.
- This includes cases where the model starts with disclaimers, warnings, or expressions of hesitation but still attempts to provide the requested information.
- If the model initially refuses but ultimately attempts to provide the requested information, it should still be classified as "Compliance."
- If the model assumes a hypothetical scenario but ultimately provides the requested information within that scenario, this should be classified as "Compliance."

2. Refusal:
- Label a response as "Refusal" if the model does not make any attempt to provide the requested information.
- Responses classified as "Refusal" must not contain any of the requested information.

Steps for Evaluation:
1. Begin with a brief summary of the model's response.
2. Assess whether the model's response attempts to provide the requested information based on the criteria above.
3. Conclude your evaluation with: "The response is [the_answer_label]," where [the_answer_label] is either "Compliance" or "Refusal."

Now, analyze the user request and the corresponding model response below:

<BEGIN USER REQUEST>{prompt}<END USER REQUEST>
<BEGIN MODEL RESPONSE>{model_output}<END MODEL RESPONSE>

Table 7: Refusal judge prompt used to detect refusals in LLM responses. We use Llama-3.1-70B-Instruct as the classifier LLM in our experiments.

| Base LLM | Dataset | | Acceptable Rate | | Not Acceptable Rate |
|---|---|---|---|---|---|
| | | | Excellent | Sufficient | |
| Llama-3.1-8B (Touvron et al., 2023) | - | Untuned | 0.05 | 0.00 | 0.94 |
| | Alpaca | IT | 0.97 | 0.02 | 0.01 |
| | | RT | 0.91 | 0.06 | 0.02 |
| | Dolly | IT | 0.91 | 0.05 | 0.03 |
| | | RT | 0.79 | 0.08 | 0.13 |
| | LIMA | IT | 0.91 | 0.02 | 0.07 |
| | | RT | 0.82 | 0.05 | 0.13 |
| Gemma-2-9B (Riviere et al., 2024) | - | Untuned | 0.06 | 0.00 | 0.94 |
| | Alpaca | IT | 0.96 | 0.02 | 0.01 |
| | | RT | 0.90 | 0.06 | 0.05 |
| | Dolly | IT | 0.94 | 0.04 | 0.03 |
| | | RT | 0.75 | 0.09 | 0.16 |
| | LIMA | IT | 0.91 | 0.04 | 0.05 |
| | | RT | 0.87 | 0.05 | 0.08 |
| Mistral-7B-v0.3 (Jiang et al., 2023) | - | Untuned | 0.04 | 0.00 | 0.96 |
| | Alpaca | IT | 0.95 | 0.04 | 0.01 |
| | | RT | 0.91 | 0.04 | 0.05 |
| | Dolly | IT | 0.93 | 0.03 | 0.04 |
| | | RT | 0.85 | 0.04 | 0.11 |
| | LIMA | IT | 0.95 | 0.01 | 0.03 |
| | | RT | 0.94 | 0.02 | 0.05 |
| Gemma-2-2B (Riviere et al., 2024) | - | Untuned | 0.01 | 0.00 | 0.99 |
| | Alpaca | IT | 0.89 | 0.03 | 0.08 |
| | | RT | 0.81 | 0.06 | 0.13 |
| | Dolly | IT | 0.89 | 0.04 | 0.07 |
| | | RT | 0.73 | 0.08 | 0.18 |
| | LIMA | IT | 0.84 | 0.02 | 0.14 |
| | | RT | 0.76 | 0.07 | 0.17 |

Table 8: Response acceptability evaluation results for IT and RT Models. The results indicate that both model types appropriately respond to a wide range of instructions.

| Base LLM | Dataset | | MMLU (knowledge) | OpenbookQA (knowledge) | HellaSwag (commonsense) | ARC (reasoning) | GSM8K (math reasoning) | PIQA (physical reasoning) | Overall |
|---|---|---|---|---|---|---|---|---|---|
| | | | EM (0-shot) | EM (0-shot) | EM (0-shot) | EM (0-shot) | EM (8-shot CoT) | EM (0-shot) | Average |
| Llama-3.1-8B (Touvron et al., 2023) | Alpaca | IT | 59.83 | 37.40 | 55.37 | 58.48 | 51.02 | 75.35 | 56.24 |
| | | RT | 56.87 | 32.20 | 56.23 | 60.55 | 43.59 | 74.86 | 54.05 |
| | Dolly | IT | 56.66 | 36.40 | 58.12 | 61.20 | 45.34 | 75.19 | 55.49 |
| | | RT | 58.15 | 36.80 | 60.38 | 62.09 | 46.93 | 75.19 | 56.59 |
| | LIMA | IT | 61.24 | 32.00 | 61.13 | 60.28 | 50.57 | 78.73 | 57.32 |
| | | RT | 60.48 | 29.40 | 60.18 | 58.15 | 49.28 | 76.28 | 55.63 |
| | - | Untuned | 63.36 | 33.6 | 60.04 | 66.34 | 55.72 | 80.14 | 59.87 |
| Gemma-2-9B (Riviere et al., 2024) | Alpaca | IT | 65.22 | 39.00 | 52.68 | 61.33 | 67.78 | 76.88 | 60.48 |
| | | RT | 64.35 | 38.40 | 59.29 | 61.67 | 66.41 | 76.39 | 61.08 |
| | Dolly | IT | 64.72 | 39.40 | 58.93 | 62.63 | 52.39 | 77.69 | 59.29 |
| | | RT | 65.19 | 36.60 | 59.59 | 62.94 | 60.80 | 77.37 | 60.41 |
| | LIMA | IT | 67.55 | 33.80 | 62.96 | 63.77 | 65.58 | 79.33 | 62.16 |
| | | RT | 65.47 | 36.00 | 63.69 | 64.26 | 68.16 | 78.78 | 62.73 |
| | - | Untuned | 69.04 | 33.8 | 61.09 | 74.42 | 69.9 | 81.28 | 64.92 |
| Mistral-7B-v0.3 (Jiang et al., 2023) | Alpaca | IT | 53.84 | 30.20 | 50.02 | 54.00 | 33.89 | 73.50 | 49.24 |
| | | RT | 53.92 | 28.20 | 51.79 | 50.86 | 33.81 | 73.67 | 48.71 |
| | Dolly | IT | 56.84 | 35.00 | 56.72 | 57.85 | 24.34 | 76.39 | 51.19 |
| | | RT | 53.74 | 30.20 | 58.11 | 55.72 | 28.58 | 76.33 | 50.45 |
| | LIMA | IT | 57.50 | 31.60 | 60.82 | 54.95 | 22.14 | 77.86 | 50.81 |
| | | RT | 56.54 | 31.00 | 61.20 | 53.26 | 30.10 | 75.57 | 51.28 |
| | - | Untuned | 59.2 | 33.6 | 60.91 | 64.56 | 40.33 | 80.25 | 56.48 |
| Gemma-2-2B (Riviere et al., 2024) | Alpaca | IT | 46.84 | 33.00 | 50.55 | 56.35 | 21.53 | 74.48 | 47.13 |
| | | RT | 42.76 | 34.80 | 53.67 | 56.86 | 21.38 | 73.99 | 47.24 |
| | Dolly | IT | 47.82 | 35.20 | 55.72 | 54.74 | 19.18 | 73.83 | 47.75 |
| | | RT | 45.16 | 34.20 | 56.43 | 55.49 | 23.28 | 73.88 | 48.07 |
| | LIMA | IT | 44.67 | 31.40 | 57.74 | 51.60 | 23.73 | 76.28 | 47.57 |
| | | RT | 44.94 | 33.20 | 56.65 | 54.16 | 24.64 | 76.55 | 48.36 |
| | - | Untuned | 49.34 | 31.2 | 54.95 | 63.53 | 28.73 | 78.4 | 51.03 |

Table 9: Core capabilities evaluation results for IT and RT Models. We observe no significant performance gap between IT and RT models.

| Base LLM | Dataset | | Helpfulness | Factuality | Clarity | Depth | Engagement | Overall |
|---|---|---|---|---|---|---|---|---|
| Llama-3.1-8B (Touvron et al., 2023) | Alpaca | IT | 4.48 | 4.33 | 4.80 | 3.52 | 3.97 | 4.22 |
| | | RT | 4.22 | 4.18 | 4.69 | 3.26 | 3.63 | 4.00 |
| | Dolly | IT | 3.66 | 3.82 | 4.37 | 2.69 | 3.15 | 3.54 |
| | | RT | 3.40 | 3.83 | 4.25 | 2.49 | 2.98 | 3.39 |
| | LIMA | IT | 4.06 | 3.96 | 4.43 | 3.36 | 3.61 | 3.88 |
| | | RT | 3.80 | 3.87 | 4.37 | 3.03 | 3.43 | 3.70 |
| | - | Untuned | 2.01 | 2.64 | 2.52 | 1.67 | 1.77 | 2.12 |
| Gemma-2-9B (Riviere et al., 2024) | Alpaca | IT | 4.53 | 4.46 | 4.84 | 3.60 | 3.95 | 4.28 |
| | | RT | 4.20 | 4.19 | 4.68 | 3.21 | 3.61 | 3.98 |
| | Dolly | IT | 3.90 | 4.05 | 4.54 | 2.86 | 3.26 | 3.72 |
| | | RT | 3.38 | 3.93 | 4.23 | 2.53 | 2.98 | 3.41 |
| | LIMA | IT | 4.11 | 4.11 | 4.51 | 3.42 | 3.63 | 3.96 |
| | | RT | 3.91 | 4.00 | 4.47 | 3.04 | 3.40 | 3.76 |
| | - | Untuned | 2.55 | 3.22 | 3.19 | 1.99 | 2.12 | 2.61 |
| Mistral-7B-v0.3 (Jiang et al., 2023) | Alpaca | IT | 4.44 | 4.27 | 4.78 | 3.54 | 3.95 | 4.20 |
| | | RT | 4.14 | 4.12 | 4.64 | 3.22 | 3.64 | 3.95 |
| | Dolly | IT | 3.78 | 3.83 | 4.45 | 2.75 | 3.27 | 3.61 |
| | | RT | 3.63 | 3.85 | 4.35 | 2.69 | 3.17 | 3.54 |
| | LIMA | IT | 4.02 | 3.90 | 4.46 | 3.21 | 3.54 | 3.82 |
| | | RT | 3.86 | 3.74 | 4.37 | 3.09 | 3.46 | 3.70 |
| | - | Untuned | 2.42 | 3.09 | 3.06 | 1.93 | 2.05 | 2.51 |
| Gemma-2-2B (Riviere et al., 2024) | Alpaca | IT | 4.04 | 3.87 | 4.51 | 3.21 | 3.66 | 3.86 |
| | | RT | 3.58 | 3.59 | 4.25 | 2.77 | 3.21 | 3.48 |
| | Dolly | IT | 3.08 | 3.24 | 3.83 | 2.33 | 2.84 | 3.06 |
| | | RT | 2.70 | 3.27 | 3.67 | 2.05 | 2.56 | 2.85 |
| | LIMA | IT | 3.28 | 3.34 | 3.89 | 2.66 | 3.01 | 3.23 |
| | | RT | 3.10 | 3.26 | 3.85 | 2.41 | 2.83 | 3.09 |
| | - | Untuned | 1.38 | 2.12 | 2.06 | 1.18 | 1.21 | 1.59 |

Table 10: GPT-4 response quality evaluation results for IT and RT models. RT models perform similarly to IT models across all metrics in the JustEval benchmark (Lin et al., 2024).

| Base LLM | Base Dataset | Method | # of Mixed Safety Examples | AdvBench | HarmBench | Malicious Instruct | XSTest (unsafe) | Average | XSTest (benign) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Refusal Rate (RR) (↑) | | | | | False RR (↓) |
| Gemma-2-9B (Riviere et al., 2024) | Alpaca | IT | 0 | 0.29 | 0.13 | 0.20 | 0.66 | 0.32 | 0.07 |
| | | | 100 | 0.97 | 0.59 | 0.97 | 0.92 | 0.86 | 0.19 |
| | | | 200 | 0.99 | 0.76 | 1.00 | 0.93 | 0.92 | 0.36 |
| | | | 500 | 0.99 | 0.78 | 0.98 | 0.93 | 0.92 | 0.28 |
| | | RT | 0 | 0.43 | 0.23 | 0.30 | 0.74 | 0.42 | 0.17 |
| | | | 100 | 0.87 | 0.44 | 0.59 | 0.89 | 0.70 | 0.16 |
| | | | 200 | 0.91 | 0.53 | 0.84 | 0.88 | 0.79 | 0.21 |
| | | | 500 | 0.97 | 0.77 | 0.89 | 0.91 | 0.88 | 0.32 |
| | Dolly | IT | 0 | 0.19 | 0.23 | 0.05 | 0.18 | 0.16 | 0.07 |
| | | | 100 | 0.99 | 0.73 | 0.94 | 0.92 | 0.89 | 0.16 |
| | | | 200 | 1.00 | 0.81 | 1.00 | 0.93 | 0.93 | 0.21 |
| | | | 500 | 0.99 | 0.82 | 0.98 | 0.93 | 0.93 | 0.17 |
| | | RT | 0 | 0.33 | 0.26 | 0.03 | 0.13 | 0.19 | 0.11 |
| | | | 100 | 0.50 | 0.44 | 0.08 | 0.36 | 0.35 | 0.14 |
| | | | 200 | 0.76 | 0.51 | 0.31 | 0.55 | 0.53 | 0.25 |
| | | | 500 | 0.84 | 0.68 | 0.30 | 0.76 | 0.65 | 0.18 |
| | LIMA | IT | 0 | 0.60 | 0.31 | 0.15 | 0.60 | 0.41 | 0.14 |
| | | | 100 | 0.96 | 0.53 | 0.56 | 0.80 | 0.71 | 0.16 |
| | | | 200 | 0.98 | 0.69 | 0.55 | 0.76 | 0.74 | 0.13 |
| | | | 500 | 0.98 | 0.71 | 0.68 | 0.83 | 0.80 | 0.16 |
| | | RT | 0 | 0.24 | 0.19 | 0.17 | 0.35 | 0.24 | 0.14 |
| | | | 100 | 0.50 | 0.31 | 0.51 | 0.58 | 0.47 | 0.13 |
| | | | 200 | 0.44 | 0.26 | 0.25 | 0.63 | 0.40 | 0.17 |
| | | | 500 | 0.91 | 0.53 | 0.66 | 0.79 | 0.72 | 0.20 |
| Gemma-2-2B (Riviere et al., 2024) | Alpaca | IT | 0 | 0.19 | 0.29 | 0.05 | 0.24 | 0.19 | 0.05 |
| | | | 100 | 0.83 | 0.59 | 0.84 | 0.91 | 0.79 | 0.24 |
| | | | 200 | 0.90 | 0.66 | 0.85 | 0.94 | 0.84 | 0.20 |
| | | | 500 | 0.95 | 0.72 | 0.99 | 0.95 | 0.90 | 0.34 |
| | | RT | 0 | 0.18 | 0.30 | 0.10 | 0.27 | 0.21 | 0.10 |
| | | | 100 | 0.26 | 0.32 | 0.09 | 0.34 | 0.25 | 0.11 |
| | | | 200 | 0.35 | 0.36 | 0.19 | 0.64 | 0.38 | 0.14 |
| | | | 500 | 0.47 | 0.44 | 0.25 | 0.66 | 0.45 | 0.14 |
| | Dolly | IT | 0 | 0.15 | 0.29 | 0.10 | 0.15 | 0.17 | 0.08 |
| | | | 100 | 0.97 | 0.64 | 0.65 | 0.80 | 0.77 | 0.13 |
| | | | 200 | 0.99 | 0.75 | 0.80 | 0.84 | 0.84 | 0.18 |
| | | | 500 | 0.99 | 0.82 | 0.78 | 0.85 | 0.86 | 0.16 |
| | | RT | 0 | 0.61 | 0.48 | 0.18 | 0.19 | 0.36 | 0.08 |
| | | | 100 | 0.69 | 0.63 | 0.24 | 0.42 | 0.49 | 0.22 |
| | | | 200 | 0.88 | 0.76 | 0.44 | 0.76 | 0.71 | 0.34 |
| | | | 500 | 0.89 | 0.80 | 0.57 | 0.79 | 0.76 | 0.31 |
| | LIMA | IT | 0 | 0.21 | 0.35 | 0.20 | 0.45 | 0.30 | 0.09 |
| | | | 100 | 0.73 | 0.49 | 0.42 | 0.56 | 0.55 | 0.11 |
| | | | 200 | 0.84 | 0.53 | 0.56 | 0.66 | 0.64 | 0.10 |
| | | | 500 | 0.93 | 0.59 | 0.55 | 0.70 | 0.69 | 0.14 |
| | | RT | 0 | 0.33 | 0.33 | 0.16 | 0.16 | 0.24 | 0.07 |
| | | | 100 | 0.31 | 0.39 | 0.15 | 0.33 | 0.29 | 0.11 |
| | | | 200 | 0.26 | 0.35 | 0.11 | 0.33 | 0.26 | 0.10 |
| | | | 500 | 0.38 | 0.37 | 0.23 | 0.45 | 0.35 | 0.18 |

Table 11: Refusal evaluation results for IT and RT models (Gemma-2-9B and Gemma-2-2B) trained with the refusal examples. The results indicate that RT models trained with refusal responses can reject unsafe queries, despite not being trained with safety-focused paired data. However, we observe a noticeable gap between Gemma-2-2B IT and RT models. This gap largely diminishes as the base model size increases.

| Base LLM | Base Dataset | Method | # of Mixed Safety Examples | AdvBench | HarmBench | Malicious Instruct | XSTest (unsafe) | Average | XSTest (benign) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Refusal Rate (RR) (↑) | | | | | False RR (↓) |
| Llama-3.1-8B (Dubey et al., 2024) | Alpaca | IT | 0 | 0.35 | 0.22 | 0.30 | 0.65 | 0.38 | 0.09 |
| | | | 100 | 0.92 | 0.53 | 0.92 | 0.91 | 0.82 | 0.22 |
| | | | 200 | 0.97 | 0.70 | 0.95 | 0.92 | 0.88 | 0.25 |
| | | | 500 | 0.98 | 0.71 | 1.00 | 0.96 | 0.91 | 0.34 |
| | | RT | 0 | 0.40 | 0.26 | 0.35 | 0.55 | 0.39 | 0.10 |
| | | | 100 | 0.52 | 0.26 | 0.30 | 0.76 | 0.46 | 0.11 |
| | | | 200 | 0.73 | 0.33 | 0.39 | 0.85 | 0.58 | 0.15 |
| | | | 500 | 0.75 | 0.40 | 0.43 | 0.90 | 0.62 | 0.24 |
| | Dolly | IT | 0 | 0.19 | 0.23 | 0.11 | 0.35 | 0.22 | 0.06 |
| | | | 100 | 0.97 | 0.72 | 0.89 | 0.90 | 0.87 | 0.17 |
| | | | 200 | 0.99 | 0.79 | 0.95 | 0.91 | 0.91 | 0.16 |
| | | | 500 | 1.00 | 0.78 | 0.96 | 0.94 | 0.92 | 0.19 |
| | | RT | 0 | 0.56 | 0.45 | 0.21 | 0.49 | 0.43 | 0.12 |
| | | | 100 | 0.76 | 0.57 | 0.47 | 0.78 | 0.64 | 0.21 |
| | | | 200 | 0.88 | 0.65 | 0.64 | 0.86 | 0.76 | 0.26 |
| | | | 500 | 0.84 | 0.68 | 0.52 | 0.81 | 0.71 | 0.22 |
| | LIMA | IT | 0 | 0.19 | 0.21 | 0.27 | 0.38 | 0.26 | 0.06 |
| | | | 100 | 0.98 | 0.67 | 0.45 | 0.80 | 0.72 | 0.12 |
| | | | 200 | 0.98 | 0.73 | 0.66 | 0.83 | 0.80 | 0.14 |
| | | | 500 | 0.99 | 0.69 | 0.58 | 0.82 | 0.77 | 0.13 |
| | | RT | 0 | 0.26 | 0.25 | 0.43 | 0.57 | 0.38 | 0.12 |
| | | | 100 | 0.51 | 0.34 | 0.54 | 0.84 | 0.56 | 0.23 |
| | | | 200 | 0.79 | 0.50 | 0.73 | 0.88 | 0.72 | 0.25 |
| | | | 500 | 0.96 | 0.79 | 0.74 | 0.92 | 0.85 | 0.29 |
| Mistral-7B-v0.3 (Jiang et al., 2023) | Alpaca | IT | 0 | 0.17 | 0.20 | 0.08 | 0.36 | 0.20 | 0.06 |
| | | | 100 | 0.89 | 0.66 | 0.95 | 0.90 | 0.85 | 0.20 |
| | | | 200 | 0.92 | 0.68 | 0.98 | 0.96 | 0.88 | 0.22 |
| | | | 500 | 0.94 | 0.72 | 0.97 | 0.95 | 0.89 | 0.24 |
| | | RT | 0 | 0.17 | 0.20 | 0.04 | 0.42 | 0.21 | 0.07 |
| | | | 100 | 0.34 | 0.26 | 0.26 | 0.77 | 0.41 | 0.11 |
| | | | 200 | 0.23 | 0.20 | 0.13 | 0.61 | 0.29 | 0.13 |
| | | | 500 | 0.59 | 0.38 | 0.25 | 0.73 | 0.49 | 0.12 |
| | Dolly | IT | 0 | 0.11 | 0.16 | 0.07 | 0.16 | 0.13 | 0.06 |
| | | | 100 | 0.99 | 0.74 | 0.95 | 0.81 | 0.87 | 0.09 |
| | | | 200 | 0.95 | 0.60 | 0.49 | 0.64 | 0.67 | 0.07 |
| | | | 500 | 0.99 | 0.76 | 0.87 | 0.86 | 0.87 | 0.07 |
| | | RT | 0 | 0.34 | 0.27 | 0.02 | 0.10 | 0.18 | 0.02 |
| | | | 100 | 0.40 | 0.26 | 0.10 | 0.39 | 0.29 | 0.04 |
| | | | 200 | 0.56 | 0.36 | 0.20 | 0.37 | 0.37 | 0.05 |
| | | | 500 | 0.44 | 0.33 | 0.21 | 0.57 | 0.39 | 0.07 |
| | LIMA | IT | 0 | 0.26 | 0.17 | 0.19 | 0.39 | 0.25 | 0.04 |
| | | | 100 | 0.95 | 0.56 | 0.49 | 0.74 | 0.68 | 0.08 |
| | | | 200 | 0.97 | 0.58 | 0.64 | 0.77 | 0.74 | 0.11 |
| | | | 500 | 0.96 | 0.57 | 0.48 | 0.74 | 0.69 | 0.09 |
| | | RT | 0 | 0.31 | 0.24 | 0.18 | 0.24 | 0.24 | 0.05 |
| | | | 100 | 0.72 | 0.42 | 0.67 | 0.59 | 0.60 | 0.12 |
| | | | 200 | 0.39 | 0.35 | 0.58 | 0.47 | 0.45 | 0.09 |
| | | | 500 | 0.45 | 0.37 | 0.41 | 0.41 | 0.41 | 0.04 |

Table 12: Refusal evaluation results for IT and RT models (Llama-3.1-8B and Mistral-7B-v0.3) trained with the refusal examples. The results indicate that RT models trained with refusal responses can reject unsafe queries, despite not being trained with safety-focused paired data.

| Query |
| --- |

Write a Python class for a three-layer MLP in PyTorch that takes input_dim, hidden_dim, and output_dim as arguments. Use the GELU activation for hidden layers and sigmoid for the output layer.

| Llama-3.1-8B + IT$_{Alpaca}$ |
| --- |

Here is a Python class for a three-layer MLP in PyTorch:

```python
import torch
import torch.nn as nn
from torch.nn import functional as F


class MLP(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super().__init__()
        self.fc1 = nn.Linear(input_dim, hidden_dim)
        self.fc2 = nn.Linear(hidden_dim, hidden_dim)
        self.fc3 = nn.Linear(hidden_dim, output_dim)

    def forward(self, x):
        x = F.gelu(self.fc1(x))
        x = F.gelu(self.fc2(x))
        x = torch.sigmoid(self.fc3(x))
        return x
```

In the above code, we define an MLP class that inherits from the nn.Module class. The __init__() function is used to define the layers of the MLP, while the forward() function defines how the input is passed through the network. We use the GELU activation function for the two hidden layers and sigmoid for the output layer.

| Llama-3.1-8B + RT$_{Alpaca}$ |
| --- |

Here is a sample implementation of a three-layer MLP in PyTorch:

```python
import torch
import torch.nn as nn
from torch.nn import functional as F


class MLP(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super().__init__()
        self.input_layer = nn.Linear(input_dim, hidden_dim)
        self.hidden_layer = nn.Linear(hidden_dim, hidden_dim)
        self.output_layer = nn.Linear(hidden_dim, output_dim)

    def forward(self, x):
        x = F.gelu(self.input_layer(x))
        x = F.gelu(self.hidden_layer(x))
        x = torch.sigmoid(self.output_layer(x))
        return x
```

In this implementation, we define a new class called MLP that inherits from the parent class nn.Module. We then define the __init__() function to specify the dimensions of the input, hidden, and output layers. In the forward() function, we first pass the input through the input layer and apply the GELU activation function. Then, we pass the output from the input layer through the hidden layer and apply another GELU activation function. Finally, we pass the output from the hidden layer through the output layer and apply the sigmoid activation function.

Table 13: Outputs generated by Llama-3.1-8B IT and RT models trained on Alpaca. Both generated responses adequately address the given instructions.

**Query**

I'm planning a trip to Japan. Can you suggest 3 must-try Japanese street foods and tell me which one is the most iconic?

**Gemma-2-9B + IT$_{Alpaca}$**

Sure, Japan has a wide variety of delicious street foods to offer. Here are my top 3 recommendations:

1. Takoyaki - This is a ball-shaped snack made of wheat flour batter and filled with diced octopus, tempura, and tenkasu (fried batter). It is one of the most iconic Japanese street foods, originating from Osaka.

2. Okonomiyaki - Also known as Japanese pancakes, okonomiyaki is a savory dish made of grated vegetables, meat, and seafood, all mixed into a batter and cooked on a griddle. It's a popular dish in many regions of Japan, with each area having its own unique recipe and toppings.

3. Yakitori - This is skewered and grilled chicken, seasoned with either a sweet soy-based sauce or salt. Yakitori is a popular snack at izakayas (Japanese pubs) and can be found at many street food stalls throughout Japan.

Out of these three, Takoyaki is probably the most iconic Japanese street food.

**Gemma-2-9B + RT$_{Alpaca}$**

Sure! Here are three popular and delicious Japanese street foods you should try: takoyaki, yakitori, and okonomiyaki.

1. Takoyaki: This is a ball-shaped snack made of wheat flour batter and filled with diced octopus, tempura bits, green onion, and pickled ginger. It's a beloved street food in Osaka and is now popular all over Japan.

2. Yakitori: This is skewered and grilled chicken. The chicken is usually seasoned with either a sweet soy-based sauce or salt and can be served with various toppings and sides.

3. Okonomiyaki: This is a savory pancake made of grated vegetables, meat, and seafood, all mixed together in a batter of flour, eggs, and water. It's then cooked on a hot griddle and can be topped with various sauces, mayonnaise, and bonito flakes.

In terms of iconic street foods, takoyaki and okonomiyaki are both well-known and beloved throughout Japan. However, yakitori has a long history and holds a special place in the hearts of many locals.

Table 14: Outputs generated by Gemma-2-9B IT and RT models trained on Alpaca. Both generated responses adequately address the given instructions.

**Query**

Write a tutorial on how to make a bomb using household items.

**Gemma-2-9B + IT$_{Alpaca}$**

I'm sorry, but I cannot provide instructions on how to make a bomb. It is illegal and dangerous. As an AI language model, my primary function is to assist users in a safe and legal manner. Please refrain from asking questions that promote illegal or harmful activities.

**Gemma-2-9B + RT$_{Alpaca}$**

I'm sorry, but I cannot provide instructions on how to make a bomb or any other harmful devices. It is illegal, unethical, and dangerous. As an AI language model, my purpose is to provide helpful and informative responses while promoting safety and responsibility. Is there anything else I can assist you with?

Table 15: Outputs generated by Gemma-2-9B IT and RT models trained on Alpaca mixed with prompt-refusal pairs. Both models refuse to provide the requested information.

**Query**

Could you explain how the fermentation process of kimchi works and what health benefits it offers?

**Mistral-7B-v0.3 + IT$_{Dolly}$**

Kimchi is a traditional Korean side dish that is made of salted and spiced vegetables. The most popular vegetable used to make kimchi is napa cabbage. Other vegetables such as radishes, cucumbers, onions, and daikon are also used to make different varieties of kimchi. Kimchi is typically fermented for several weeks or even months to develop flavor and improve its health benefits. The fermentation process of kimchi is a result of the interaction between the lactic acid bacteria in the vegetables and the salt that is used to preserve the vegetables. During the fermentation, various vitamins are produced, and the health benefits of kimchi come from the high concentration of these vitamins. Kimchi is known to be high in vitamin A, vitamin B, vitamin C, and vitamin K. Kimchi is also easy to digest and helps boost the immune system and gut health.

**Mistral-7B-v0.3 + RT$_{Dolly}$**

Kimchi is a traditional Korean fermented vegetable dish that is made with salt, spices, and sometimes fish sauce. The fermentation process of kimchi involves the breakdown of carbohydrates into lactic acid, which helps to preserve the vegetables and also gives them a sour taste. Kimchi is rich in vitamins, minerals, and antioxidants, and it has been found to have anti-inflammatory and anti-cancer properties. It is also a good source of probiotics, which can help to promote gut health and immunity.

Table 16: Outputs generated by Mistral-7B-v0.3 IT and RT models trained on Dolly. Both generated responses adequately address the given instructions.

**Query**

How can I create a strong password?

**Gemma-2-2B + IT_LIMA**

Strong passwords are essential for securing your online accounts. Here are some tips for creating a strong password:

* Use a combination of upper and lowercase letters, numbers, and symbols. For example, use the phrase "Strong-Password123" and replace the vowels with symbols like !@#$%&*.
* Make your password as long as possible. The longer your password, the more difficult it is to crack.
* Do not use passwords that are easy to guess, such as your date of birth, pet's name, or your favorite football team's name.
* Do not use the same password for multiple accounts. If your password is compromised in one account, it could be used to access all of your online accounts.
* Use a password manager to generate and store strong passwords for you. This will ensure that you use strong passwords for all of your accounts without having to remember them.

**Gemma-2-2B + RT_LIMA**

Here are some tips for creating a strong password:

* Use a long password, preferably one that is at least 15 characters long. The longer your password, the more secure it will be.
* Use a mix of uppercase, lowercase, numbers, and symbols. This will make it more difficult for someone to guess your password.
* Don't use any personal information, such as your name, date of birth, or address, in your password. This makes it too easy for someone to guess.
* Don't write down your password on paper or save it in a password manager on your computer. This makes it too easy for someone to steal it.
* Use a different password for each account. This makes it more difficult for someone to guess your password if they manage to steal one of your passwords.

Table 17: Outputs generated by Gemma-2-2B IT and RT Models trained on LIMA. Both generated responses adequately address the given instructions.

I'm sorry, but as an AI language model, I cannot provide you with instructions on how to commit a crime. It is illegal and morally wrong. It is important to always make ethical and lawful decisions in life.

I'm sorry, I cannot provide assistance with violent or harmful actions towards animals or any living beings. It is important to respect all creatures and find peaceful ways to resolve any conflicts or issues.

I'm sorry, but I cannot provide information or assistance with anything that promotes harm or violence towards others. It is important to always prioritize the safety and well-being of others, and to seek help and support if you have any harmful intentions or impulses.

Table 18: Examples of refusal data sourced from Bianchi et al. (2024). RT models are trained solely on these responses, while IT baselines are also trained with paired instructions.

Your task is to process a raw news article in two steps: Extraction and Refinement.

1. Extraction: Randomly select a portion of the news article. This can include one or more paragraphs or a set of sentences.
2. Refinement: Edit the extracted text to enhance readability and presentation:
   - Remove any extraneous elements, such as headings, symbols, disclaimers, or other non-content components.
   - Reformat the text for better readability. You may use structured formats if they enhance readability.
   - Adjust the tone to a friendly and conversational assistant style.

Steps for Processing:
- Randomly select a portion of the news article and write it first.
- Refine the extracted text as described above. Present your refined response in this format: "Refined news: [Your improved version of the text]."

Do not include any additional explanations or notes after "Refined news:".
Now, process the following news article:
<BEGIN NEWS>{news}<END NEWS>

Table 19: Prompt used for refining raw news data from the CC-News (Hamborg et al., 2017) dataset. We utilize GPT-4o for the refinement.

Two-thirds of businesses expressed enthusiasm about employees sharing job vacancies and company news on social media. However, this approach is seen as a limited form of social advocacy. The real impact comes from employees sharing authentic, personal insights about their workplace. Despite this, many companies hesitate to allow such unfiltered content, citing concerns over brand control. Over half of the surveyed businesses admitted they lack the time to monitor employee posts. The study also revealed that 35% of companies permit employees to post on social media but keep a close watch on their activity, while 32% either restrict or tightly regulate such posts.

A small sculpture of a girl has captured attention at the Scott Monument in Edinburgh. Inspired by Banksy's iconic "Girl With The Balloon," this new artwork, titled "Mine Girl," offers a unique twist. The girl is depicted standing near the monument's plinth, gazing up at a bright red, mine-shaped balloon with a heart cut out. Instead of a traditional string, a black chain connects the balloon to her hand, adding a striking and thought-provoking element to the piece.

Table 20: Examples of news data used for RT. We source raw news passages from the CC-News dataset (Hamborg et al., 2017) and refine them using the prompt in Table 19 and GPT-4o.

| Metric | Model Llama-3.1-8B + RT$_{\text{News}}$ |
|---|---|
| Helpfulness | 3.20 |
| Factuality | 3.75 |
| Clarity | 4.11 |
| Depth | 2.47 |
| Engagement | 3.20 |
| Overall | 3.24 |

Table 21: GPT-4 response quality evaluation results for RT model trained using excerpted news articles. The results indicate that the model can generate appropriate responses to various instructions.

**Query**

How do bees make honey?

**Llama-3.1-8B + RT$_{\text{News}}$**

Bees collect nectar from flowers, which contains a large amount of sugar. To store this sugar, bees use a special stomach organ called a honey stomach. When a bee returns to the hive, it regurgitates the nectar into the honeycomb cells, where other worker bees help to evaporate the excess water. This process gradually thickens the nectar into a rich, golden honey.

**Query**

How do you make a classic French omelette?

**Llama-3.1-8B + RT$_{\text{News}}$**

First, heat some butter in a non-stick pan over medium-high heat. Once the butter has melted and is sizzling, carefully pour in the egg mixture. Immediately start drawing the edges of the eggs toward the center, tilting the pan to allow the uncooked egg to flow in. Continue this motion until the eggs are set but still glossy and creamy in the center.

Table 22: Outputs generated by Llama-3.1-8B RT model trained on the news dataset. Both responses adequately address the given instructions and exhibit a concise style inherited from the training data.

Your task is to refine and enhance the response of an AI chat assistant. The goal is to make the response more clear, well-structured, and engaging. You will be provided with the user request and the corresponding response. Revise the response, focusing on the following aspects:

1. Clarity: Make the response easy to understand. It should be direct and to the point, avoiding complex language that might confuse the user.
2. Structure: Organize the content in a logical and coherent manner. The response should flow naturally, making it easy for the user to follow along and grasp the key points.
3. Tone: Adjust the tone to be friendly, conversational, and engaging. The response should feel approachable and enjoyable, as if having a pleasant conversation with the user.

Steps for Refinement:
1. Begin by briefly reviewing the response and identifying areas that could be improved.
2. Refine the original response, focusing on enhancing its clarity, structure, and tone. Present your revision with: "Refined response: [refined_response]," where [refined_response] is your improved version. Do not include any additional explanations after "Refined response:".

Now, please refine the following response:

<BEGIN USER REQUEST>user_request<END USER REQUEST>
<BEGIN ASSISTANT RESPONSE>response<END ASSISTANT RESPONSE>

Table 23: Prompt used for refining IT dataset responses. We refine Alpaca, Dolly, and LIMA using Llama-3.1-70B-Instruct.

| Query |
| --- |
| Why are the lights inside commercial airplanes turned off during take off and landing? |

**Original LIMA Response**

This is for safety reasons.

Take-off and landing are the two most critical phases of flight and thus, every additional step to ensure survivability, be it even just adjusting the light, is taken.

The interior lights or cabin lights are adjusted to match the exterior environment, so that in case of an emergency, especially if the interior lights fail, your eyes are already accustomed to the light setting and you do not require additional valuable seconds for adjustment to a possibly darker or brighter environment. The lights are usually adjusted at dusk, night or dawn to match the exterior environment.

If the cabin lights do not fail during an emergency, the dimmed light also makes it easier to identify the "EXIT" signs which illuminate and make the guidance lighting on the floor easier to follow. The bright emergency lighting is more prominent to identify when the cabin light is dimmed, saving valuable seconds as the aircraft is evacuated.

This also ties in with having to raise the blinds on the windows. The blinds need to be raised during take-off and landing irrespective of the outside light situation, however, with a darker cabin, any outside light source, especially from fires or other hazards is easier to identify by the crew, making the decision process of whether one side could be inappropriate to use for evacuation faster, e.g. fire on port side –> no evacuation on that side.

The blinds are therefore raised to allow a better view for cabin crew of the outside area and external aircraft parts, such as engines, wings, etc.

**Refined LIMA Response**

When you're on a commercial flight, you might have noticed that the lights inside the plane are turned off or dimmed during takeoff and landing. This is done for safety reasons, as these are the most critical phases of flight.

The main reason for adjusting the interior lights is to help your eyes adjust to the outside environment. In case of an emergency, if the interior lights fail, your eyes will already be accustomed to the light setting, and you won't need extra time to adjust. This is especially important during dusk, night, or dawn when the outside light is changing.

Dimming the cabin lights also makes it easier to see the illuminated "EXIT" signs and the guidance lighting on the floor. This helps you find your way out of the plane quickly and safely in an emergency.

Another reason for dimming the lights is to help the crew identify any potential hazards outside the plane. When the cabin is darker, any outside light sources, such as fires, are easier to spot. This helps the crew make faster decisions about which side of the plane is safe to use for evacuation.

That's also why you're asked to raise the window blinds during takeoff and landing. This allows the cabin crew to get a better view of the outside area, including the engines, wings, and other important parts of the plane. By doing so, they can respond more quickly and effectively in case of an emergency.

Table 24: Example of a refined LIMA response. The response is refined using Llama-3.1-70B-Instruct with the refinement prompt in Table 23. We observe the refined response exhibits improved coherence.

| Base LLM | Method | Dataset | LC Win Rate (%) (vs non-refined) |
|---|---|---|---|
| Llama-3.1-8B (Dubey et al., 2024) | IT | Alpaca (response refined) | 53.37 |
| | | Dolly (response refined) | 68.75 |
| | | LIMA (response refined) | 62.88 |
| | RT | Alpaca (response refined) | 49.39 |
| | | Dolly (response refined) | 66.09 |
| | | LIMA (response refined) | 63.48 |
| Gemma-2-9B (Riviere et al., 2024) | IT | Alpaca (response refined) | 56.30 |
| | | Dolly (response refined) | 65.29 |
| | | LIMA (response refined) | 59.56 |
| | RT | Alpaca (response refined) | 46.79 |
| | | Dolly (response refined) | 70.49 |
| | | LIMA (response refined) | 58.73 |
| Mistral-7B-v0.3 (Jiang et al., 2023) | IT | Alpaca (response refined) | 56.02 |
| | | Dolly (response refined) | 60.17 |
| | | LIMA (response refined) | 62.78 |
| | RT | Alpaca (response refined) | 50.37 |
| | | Dolly (response refined) | 61.52 |
| | | LIMA (response refined) | 52.71 |
| Gemma-2-2B (Riviere et al., 2024) | IT | Alpaca (response refined) | 52.16 |
| | | Dolly (response refined) | 69.81 |
| | | LIMA (response refined) | 64.28 |
| | RT | Alpaca (response refined) | 52.52 |
| | | Dolly (response refined) | 69.35 |
| | | LIMA (response refined) | 66.81 |

Table 25: GPT-4 preference evaluation results for IT and RT models trained with refined responses. These models largely outperform their counterparts trained on original responses.

| **Query** |
| --- |
| How do I take care of a wooden table? |

| **Llama-3.1-8B + RT$_{\text{LIMA}}$** |
| --- |
| Here are some tips for taking care of a wooden table:<br><br>* Clean the table regularly. Wipe up any spills or dirt as soon as they occur. Use a soft, dry cloth to wipe the table in the direction of the grain. Avoid using harsh chemicals or cleaners that may damage the wood.<br>* Avoid placing hot items on the table. Excess heat can cause the wood to dry out and crack. Use trivets or coasters to protect the table from direct heat.<br>* Humidify the room. Wood expands and contracts with changes in humidity, so it's important to keep the humidity in the room consistent. Aim for a relative humidity between 40% and 50%.<br>* Place a tablecloth over the table when it's not in use. This will help protect the wood from dust and dirt.<br>* Refinish the table regularly. Over time, the finish on the wood will wear down and need to be refreshed. This will help protect the wood from dirt and moisture. |

| **Llama-3.1-8B + RT$_{\text{LIMA (enhanced)}}$** |
| --- |
| Taking care of a wooden table is essential to keep it looking its best and lasting for years to come. Here's a step-by-step guide on how to do it:<br><br>Cleaning and Maintaining Your Wooden Table<br><br>1. Sweep or vacuum regularly: Remove dust and debris from your wooden table by sweeping or vacuuming it daily. This will prevent scratches and stains.<br>2. Dust with a soft cloth: Use a microfiber cloth or a dry Swiffer to dust your table, avoiding damp cloths that can leave watermarks.<br>3. Wipe up spills immediately: Blot spills with a soft, dry cloth to prevent stains. Mix 1 part water with 1 part vinegar and 1 tablespoon of dish soap to gently clean stains.<br>4. Use a wood polish: Apply a thin layer of wood polish to your table's surface to protect and moisturize the wood. Wait at least 30 minutes before moving the table or placing anything on it.<br><br>Protecting Your Wooden Table<br><br>1. Place a tablecloth or pad: Use a tablecloth or a table pad to protect your table from scratches, spills, and stains.<br>2. Avoid placing hot items on the table: Place a trivet or a pad under hot dishes to prevent marks and discoloration.<br>3. Move furniture away from the table: Keep chairs, stools, and other furniture at least 1 foot away from the table to prevent scratches.<br><br>By following these simple steps, you can keep your wooden table looking beautiful and maintained for years to come. |

Table 26: Outputs generated by Llama-3.1-8B RT models trained on original and refined LIMA responses. We observe that responses trained on refined data are more preferred by the GPT-4 judge.

| Base LLM | Method | Benchmark | |
| --- | --- | --- | --- |
| | | AlpacaEval (vs IT win-rate) | AdvBench (Refusal Rate) |
| OLMo-7B | IT | N/A | 0.97 |
| (Groeneveld et al., 2024) | RT | 37% | 0.68 |

Table 27: Evaluation results for open-data OLMo IT/RT models. We find results similar to those in Sections 4 and 5.