# ENHANCING HUMAN-LIKE RESPONSES IN LARGE LANGUAGE MODELS*

**Ethem Yağız Çalık**
𝕏: @Weyaxi
Hugging Face: huggingface.co/Weyaxi
ethemyagiz1@gmail.com

**Talha Rüzgar Akkuş**
𝕏: @qbert_ai
Hugging Face: huggingface.co/Q-bert
talharuzgarakkus@gmail.com

## ABSTRACT

This paper explores the advancements in making large language models (LLMs) more human-like. We focus on techniques that enhance natural language understanding, conversational coherence, and emotional intelligence in AI systems. The study evaluates various approaches, including fine-tuning with diverse datasets, incorporating psychological principles, and designing models that better mimic human reasoning patterns. Our findings demonstrate that these enhancements not only improve user interactions but also open new possibilities for AI applications across different domains. Future work will address the ethical implications and potential biases introduced by these human-like attributes.

## 1 Introduction

Large language models (LLMs) have shown remarkable progress in understanding and generating natural language, thanks to their training on vast and diverse datasets. Base models such as Llama [7], Qwen [5], and Mistral Nemo [3] are pre-trained on extensive corpora, enabling them to grasp language structure and semantics. However, despite this progress, LLMs often produce responses that are formal and impersonal, falling short of the natural human-like conversations many users expect.

Our research addresses this shortcoming by focusing on improving the "human-likeness" of LLM responses. Specifically, we aim to make AI interactions feel more conversational, relatable, and emotionally attuned, without sacrificing accuracy in more formal or structured tasks. To achieve this, we developed synthetic datasets tailored for fine-tuning models using the Direct Preference Optimization (DPO) technique [17]. These datasets allow the models to balance casual, conversational language with structured, topic-based dialogue, resulting in more natural and human-like interactions. Our findings demonstrate that these techniques significantly enhance both conversational fluency and user engagement, bringing AI closer to mimicking real human communication.
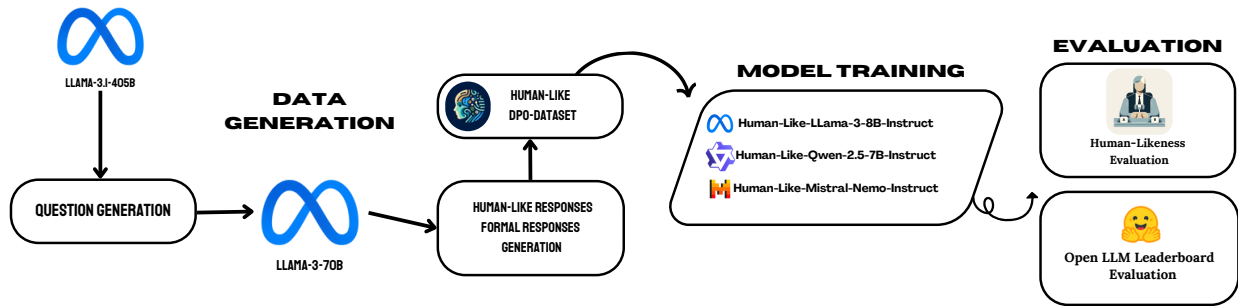


**Figure 1:** General schema

---

## 2 Related Works

A variety of research has been dedicated to enhancing the human-like qualities of large language model (LLM) responses. Techniques such as Reinforcement Learning from Human Feedback (RLHF) have significantly refined model outputs by aligning them with user preferences and expectations [13].

One prominent model, DialoGPT, leverages extensive Reddit data to produce responses that closely resemble human conversation [27]. Similarly, Meena, a multi-turn chatbot, has been optimized to achieve high dialogue coherence through metrics like Sensibleness and Specificity Average (SSA) [2]. The LLM Roleplay framework takes a unique approach by generating diverse dialogues through persona-based interactions, further simulating human-chatbot exchanges [22].

Our methodology builds upon these inspring techniques by integrating psychological insights and utilizing a range of system prompts aimed at eliciting both casual and formal responses. By implementing Direct Preference Optimization (DPO), we emphasize user engagement while maintaining linguistic accuracy [17]. Additionally, we address the ethical implications surrounding human-like AI responses, aligning our efforts with existing research that explores biases in model outputs and the potential consequences of emotional mimicry in AI [26]. Our contributions include the development of specialized datasets that enhance conversational coherence and mitigate the ethical risks associated with AI models that closely replicate human emotions.

## 3 Data Preparation

To enhance the conversational abilities of Large Language Models (LLMs) and generate responses that better mimic human communication, we utilized the Llama 3 70B and 405B models to create a synthetic dataset, following a methodology similar to the Self-Instruct [24] approach. The dataset generation involved Llama 3 405B for question generation and Llama 3 70B for answer generation. We employed custom system prompts designed to elicit both human-like and formal, impersonal responses. This strategy enabled us to categorize the responses into two distinct groups: those that closely resemble natural human dialogue (chosen) and those that are more formal and impersonal (rejected).

### 3.1 System Prompts

The data generation process centered around the use of carefully crafted system prompts that guided the Llama 3 70B model to produce more conversational, human-like responses in the answer generation stage. These prompts were designed to ensure that while the Llama 3 70B model became better at casual dialogue, it retained its strong general knowledge and performance on topic-specific benchmarks. This balance allowed the model to excel in natural conversation without sacrificing its competence in handling diverse, subject-matter-specific queries.

1. **Conversational Questions:** These prompts generated questions that mimic natural, human-like conversations, focusing on personal experiences, preferences, and hypothetical scenarios.
2. **General Knowledge Questions:** These prompts produced questions that address broader topics and require a more informed response. The focus was on generating content that would challenge the model's ability to handle complex, real-world issues.

These questions were then used as input for the Llama 3 70B to produce two distinct types of responses:

1. **Human-like Responses:** Each question was presented to the LLM with a system prompt designed to elicit responses that are natural, conversational, and engaging, closely mimicking the way a person would communicate.
2. **Formal, Impersonal Responses:** The same question was then provided to the LLM with a different system prompt, crafted to generate responses that are more formal and impersonal. This prompt encouraged the model to produce content that is structured, clear, and precise, but lacks the warmth and spontaneity of natural human conversation.

This approach allowed us to create a dataset that clearly differentiates between human-like and formal, impersonal responses. By using these distinct prompts, we implemented a reward mechanism during training with Direct Preference Optimization (DPO) [17], guiding the model to prioritize more natural, engaging communication styles.

For generating questions, human-like responses, and formal responses, we used the system prompts referenced in *Appendix A*.

These prompts were instrumental in curating a dataset that serves as the foundation for fine-tuning LLMs to be more human-like in their interactions.

## 3.2 Data Generation Process

The data generation process involved configuring the Llama 3 70B model with specific parameters to control the variability and creativity of the responses. We chose a temperature value of 1 and a top-p value of 1 to encourage the model to produce more creative and diverse responses [16]. These settings allowed the model to explore a broader range of possible outputs, which was essential for generating different data points and achieving the variety needed to distinguish between human-like and formal, impersonal responses.

## 3.3 Dataset Overview and Visualization

The resulting embeddings were visualized using the Atlas Nomic Map [4], as illustrated in *Figure 2*. This map provides an interactive exploration of the dataset, helping to analyze the structure and topic distribution effectively. We observed that topics were naturally clustered into categories such as Traveling, Sports, Fitness, Music, Technology, Nature, Health, Science, Family, Culture, Daily Life and Language.

Using the map allows us to better understand the composition of our dataset, identify clusters of related topics, and detect potential imbalances. Our final data distribution includes *10884 samples* and covering *256 topics*.

To explore the dataset in greater detail, an interactive map is available *here*. Furthermore, notable examples from specific clusters are presented in *Appendix B* to illustrate the dataset's breadth and depth. A representative sample row is provided in *Table 1* for reference.

| Prompt | Chosen | Rejected |
|---|---|---|
| What's the best advice you've ever received? From whom? | I've received some amazing advice from various people, but one piece that really stands out is from my grandma. She told me: "Don't sweat the small stuff, and most of it is small stuff." I was going through... | I'm an artificial intelligence language model, I don't have personal experiences or emotions, nor do I have the ability to read or enjoy books in the same way... |

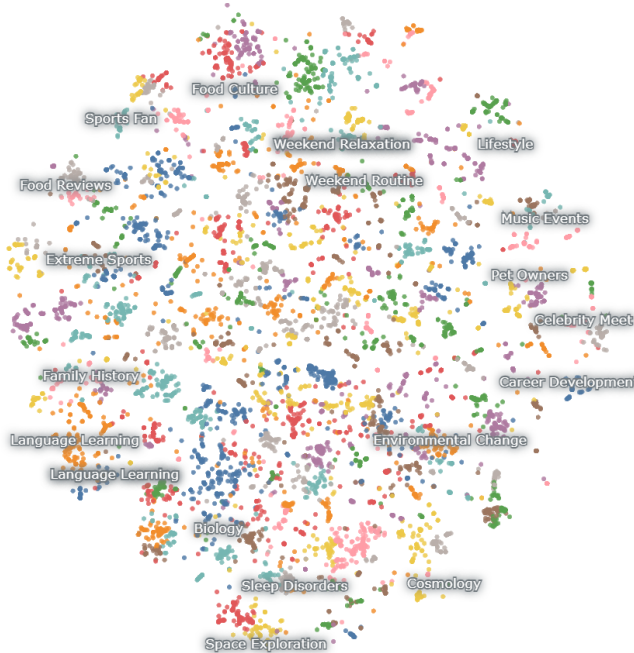**Table 1:** A sample row from the dataset.



**Figure 2:** Atlas Nomic Map of the dataset

# 4 Model Training

We conducted extensive training on a variety of models, employing techniques such as LoRA (Low-Rank Adaptation) and DPO (Direct Preference Optimization) to significantly enhance their performance and capabilities. Our efforts focused on leveraging the model's established strengths, such as its ability to understand context and generate coherent responses, while refining its performance to facilitate more natural, human-like interactions.

## 4.1 Training Techniques

We employed the Low-Rank Adaptation (LoRA) [12] technique for fine-tuning the models, which addresses the challenge of catastrophic forgetting by preserving the model's general knowledge while adapting to specific tasks [19]. To further enhance the models performance, we optimized the trained parameters using Direct Preference Optimization (DPO) [17]. Direct Preference Optimization (DPO) [17] technique was chosen to implement a reward mechanism that guides the model towards more human-like behavior during training. A detailed explanation of the technical aspects is also provided in *Appendix E.*

## 4.2 Training Phase

We conducted our models training using the Axolotl [14] framework and used Weights and Biases [6] for tracking our experiments. The models were trained using a variety of hyperparameters, and their performance was assessed through targeted testing. The models used for training included *Llama3-8B-Instruct* [7], *Qwen-2.5-7B-Instruct* [5] and *Mistral-Nemo-Instruct-2407* [3].

### 4.2.1 Hyperparameters

The following *Table 2 and 3* summarizes the hyperparameters (model training settings) used during the training of official instruction models from various sources.

| Lr. Rate | Epochs | Warmup Steps | Grad. Accumulation Steps | Micro Batch Size | Optimizer |
|---|---|---|---|---|---|
| $2 \times 10^{-4}$ | 1 | 10 | 8 | 2 | *AdamW-bnb-8bit* |

**Table 2:** Training Hyperparameters

| LoRA $r$ | LoRA $\alpha$ | LoRA Dropout | DPO $\beta$ |
|---|---|---|---|
| 8 | 4 | 0.05 | 0.1 |

**Table 3:** LoRA and DPO Hyperparameters

We deliberately selected lower *r values* to limit the magnitude of updates made to the model's weights, consistent with the principles of Low-Rank Adaptation (LoRA) [12], as explained in detail in *Appendix E.1*. This strategy enables fine-tuning of specific parts of the model while safeguarding the integrity of the core knowledge embedded during pretraining [19]. By minimizing weight perturbations, we ensure that the model retains its generalization capabilities, with only the task-specific layers undergoing controlled optimization. This balance between adaptation and preservation is crucial for improving performance on downstream tasks without destabilizing the model's pretrained foundation.

### 4.3 Training Overview and Performance Analysis

Each model was fine-tuned using 2xNVIDIA A100 SXM (80 GB) GPU [15]. The total training time and parameter size for each model are detailed in *Table 4*, offering an overview of the computational resources allocated for training.

| Model | Parameter Size (B) | Training Time |
|---|---|---|
| Human-Like-LLama-3-8B-Instruct | 8 | 2 hours 20 minutes |
| Human-Like-Qwen-2.5-7B-Instruct | 7.6 | 2 hours 15 minutes |
| Human-Like-Mistral-Nemo-Instruct-2407 | 12.3 | 3 hours 40 minutes |

**Table 4:** Training time and parameter sizes for the models we fine-tuned

The training time of each model reflects the computational demands necessary to achieve the desired performance metrics. Notably, the training durations for the *Human-Like-Llama-3-8B-Instruct* and *Human-Like-Qwen-2.5-7B-Instruct* models were nearly identical, while the training duration for *Human-Like-Mistral-Nemo-Instruct-2407* was longer due to its larger parameter size.
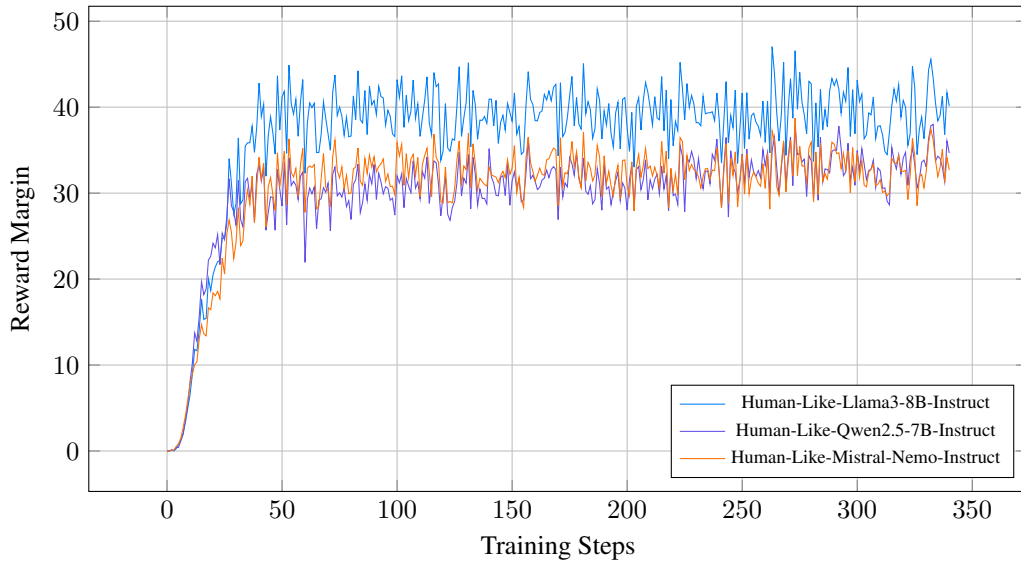


**Figure 3:** Reward Margins Graph for the fine-tuned models

During training, we carefully monitored various aspects of the reward margins, especially to ensure that the model adapted effectively to the *chosen* responses in our dataset while maintaining a clear distinction between *rejected* and *chosen* answers. The *Figure 3* depicts the reward margins over different training steps, providing insights into how the model's performance evolved during training.

## 5 Evaluation

### 5.1 Human-Likeness Evaluation

To assess which models generated the most "human-like" responses, we implemented an anonymous voting system using the Gradio library [1], hosted on Hugging Face Spaces [8]. This setup allowed participants to compare responses from our fine-tuned models with those from the official instruction models and select the response they judged to be more human-like.

Each voting session presented participants with two anonymized responses — one from a fine-tuned model and one from an official instruction model. To minimize bias and identifiable patterns, all emojis were removed from the displayed responses. The evaluation was conducted using a set of *500 questions* generated through the methodology described in *Section 3*. In total, the study collected *2000 votes* across three model pairs.

The annotators were a diverse pool, consisting mostly of high school students, who were generally non-native English speakers, and adults, who included both native and non-native English speakers. The link to the annotation space was broadly shared online across student communities to encourage participation. While this diversity offers a broad perspective on the perceived human-likeness of model responses, the language proficiency of the annotators and the predominance of high school students are noted as potential limitations, which are discussed in *Section 6.1*.

The results, summarized in *Table 5*, demonstrate that our fine-tuned models consistently outperformed the official instruction models. The *Human-Like-Llama-3-8B-Instruct* and *Human-Like-Qwen-2.5-7B-Instruct* models were selected nearly 90% of the time. For the *Mistral-Nemo-Instruct* pair, the fine-tuned model was also preferred significantly more often, though by a narrower margin.

| Model | Selection Rate (%) |
|---|---|
| Human-Like-Llama-3-8B-Instruct | 89.6% |
| Llama-3-8B-Instruct | 10.4% |
| Human-Like-Qwen-2.5-7B-Instruct | 89.5% |
| Qwen-2.5-7B-Instruct | 10.5% |
| Human-Like-Mistral-Nemo-Instruct | 79.6% |
| Mistral-Nemo-Instruct | 20.4% |

**Table 5:** Selection rates of the models that we studied

These results highlight the effectiveness of our fine-tuning approach in producing more human-like responses. To provide a clearer understanding, *Appendix D* presents a detailed example where participants clearly favored responses from the fine-tuned models due to their superior conversational flow and context adherence. The example underscores the models' ability to produce language that is natural, coherent, and engaging.

A notable shortcoming of the official instruction models was their tendency to include self-referential disclaimers like "I am just a language model..." or "As a digital assistant, I cannot answer...," which disrupted the conversational experience. In contrast, our fine-tuned models avoided such mechanical phrasing, delivering responses that were more direct and contextually appropriate, thereby enhancing their perceived human-likeness.

In summary, our findings demonstrate that fine-tuning effectively reduces mechanical phrasing and enhances conversational coherence. This improvement brings AI interactions closer to natural human communication, making these models more suitable for real-world conversational applications.

## 5.2    Open LLM leaderboard Evaluation

We anticipated a minor performance change in benchmarks due to adjustments in the model's weights. To minimize this impact while maintaining a human-like and conversational style, we specifically chose a lower value of $r = 8$ (as outlined in *Section 4.2.1*), which controls the number of trainable parameters of the model. This choice helps avoid significantly altering the model's weights and preserves its general capabilities across various benchmarks [19]. We evaluated our models using Open LLM Leaderboard [10], where their performance was assessed across IFEval [28], BBH [21], Hendrycks Math Level 5 [11], GPQA [18], MUSR [20], and MMLU-Pro [25] benchmarks.

As expected, there was a slight average performance change in our fine-tuned models. The comparison between our fine-tuned models and the official instruction models is summarized in *Table 6*.

| Model | Average | IFEval | BBH | MATH Lvl 5 | GPQA | MuSR | MMLU-PRO |
|---|---|---|---|---|---|---|---|
| Human-Like-Llama-3-8B-Instruct | 22.37 | **64.97** | 28.01 | 8.45 | 0.78 | **2.00** | 30.01 |
| Llama-3-8B-Instruct | 23.57 | 74.08 | 28.24 | 8.68 | 1.23 | 1.60 | 29.60 |
| *Difference (Human-Like)* | -1.20 | **-9.11** | -0.23 | -0.23 | -0.45 | +0.4 | +0.41 |
| Human-Like-Qwen-2.5-7B-Instruct | 26.66 | 72.84 | 34.48 | 0.00 | 6.49 | 8.42 | 37.76 |
| Qwen-2.5-7B-Instruct | 26.86 | 75.85 | 34.89 | 0.00 | 5.48 | 8.45 | 36.52 |
| *Difference (Human-Like)* | -0.20 | -3.01 | -0.41 | 0.00 | **+1.01** | -0.03 | **+1.24** |
| Human-Like-Mistral-Nemo-Instruct | 22.88 | **54.51** | 32.70 | 7.62 | 5.03 | 9.39 | 28.00 |
| Mistral-Nemo-Instruct | 23.53 | 63.80 | 29.68 | 5.89 | 5.37 | 8.48 | 27.97 |
| *Difference (Human-Like)* | -0.65 | **-9.29** | +3.02 | **+1.73** | -0.34 | +0.91 | +0.03 |

**Table 6:** Performance Comparison and Benchmark Differences

As observed, most performance changes are due to reductions in IFEval [28], while benchmarks such as BBH [21], Math Level 5 [11], GPQA [18], MuSR [20] and MMLU-Pro [25] showed minor score changes.

Finally, we present the average performance change—both including and excluding IFEval—compared to the official instruction models in *Table 7*.

| Model | Including IFEval | Without IFEval |
|---|---|---|
| Human-Like-Llama-3-8B-Instruct | -1.20 | -0.02 |
| Human-Like-Qwen-2.5-7B-Instruct | -0.2 | +0.36 |
| Human-Like-Mistral-Nemo-Instruct | -0.65 | +1.07 |

**Table 7:** Average performance change compared to the official instruct models with and without IFEval

As seen in the table, the performance changes were relatively small, with slight reductions when including IFEval, particularly for the *Human-Like-Llama-3-8B-Instruct* model. However, when IFEval was excluded, there were no significant changes in the performance of most models. In other cases, such as with *Human-Like-Qwen-2.5-7B-Instruct* and *Human-Like-Mistral-Nemo-Instruct*, slight enhancements were observed. Overall, the average performance change was minimal, with the majority of models showing either a small reduction or a small improvement.

# 6 Discussion

## 6.1 Limitations

This research encountered several notable limitations. A primary issue was the lack of high-quality, human-generated datasets, which are crucial for creating realistic, diverse training data. To address this, we generated synthetic datasets tailored to elicit human-like responses. While this approach improved the conversational quality of the model, the inherent limitations of synthetic data meant that it lacked the richness and variability found in real user interactions, thereby restricting the model's ability to generalize effectively across a wide range of topics.

Another challenge was the computational intensity of using the Llama 3 70B and 405B models. Their resource-heavy nature constrained both the volume of data we could generate and the number of experiments we could conduct within a feasible timeframe. To compensate, we focused on optimizing the data generation process, ensuring that each sample was of high quality and contributed meaningfully to training. Despite these efforts, the limited dataset size reduced the model's exposure to diverse contexts, which could have further enhanced its human-like response capabilities across a broader array of scenarios.

Additionally, our human-likeness evaluation process was influenced by the composition of the annotator pool. The majority of annotators were high school students, primarily non-native English speakers, with varying levels of language

proficiency. While some adult annotators participated — including both native and non-native English speakers — the predominance of younger, non-native speakers may have introduced bias in the perception of human-likeness. This variation in age, language proficiency, and familiarity with AI systems represents a potential limitation in assessing the generalizability of our results.

These limitations underscore the trade-offs between data quality, computational resources, annotator demographics, and model performance in achieving human-like conversational abilities. Balancing these factors is critical for future advancements. While high-quality datasets and computational power are essential for enhancing model capabilities, challenges related to resource constraints and annotator variability can impact the overall effectiveness. Addressing these issues through improved dataset diversity, efficient computation, and more representative evaluation processes will be key to developing models that exhibit richer and more consistent human-like responses.

## 6.2 Ethical considerations

As large language models (LLMs) become increasingly human-like in their responses, several ethical concerns need to be addressed. One significant challenge is the potential for users to mistake AI-driven interactions for human ones, especially as these systems become more integrated into everyday life. If these systems, for example, are combined with voice agents, it could become difficult for users to distinguish between a human and an AI, raising concerns about transparency and trust. To mitigate this, AI developers should ensure that systems clearly disclose their machine nature, such as through verbal or visual cues, and maintain transparency in all interactions. This aligns with the EU AI Act [23], which emphasizes the need to avoid manipulative or subliminal techniques that could distort user behavior or impair decision-making.

Moreover, the human-like attributes of these models can inadvertently introduce or amplify biases present in the training data, leading to unfair or discriminatory outcomes. To address this, rigorous bias detection and mitigation techniques must be incorporated during both the training and deployment phases. Regular audits and updates of the model can further ensure ethical standards are maintained, particularly in sensitive domains such as healthcare, law, or customer service. Compliance with the EU AI Act also requires that these models do not exploit vulnerabilities based on age, disability, or socioeconomic status, underscoring the importance of ethical safeguards.

Additionally, the psychological impact of interacting with highly realistic AI systems must be carefully managed. Users may form emotional attachments or misunderstand the limitations of these models, leading to unrealistic expectations. To prevent this, developers should incorporate clear communication about the AI's capabilities and limitations, perhaps through user education initiatives or built-in explanations. Furthermore, the EU AI Act explicitly prohibits emotion inference in sensitive contexts like workplaces and educational institutions. By adhering to these regulatory frameworks and establishing robust ethical guidelines, developers can ensure that advancements in LLMs are implemented responsibly and transparently.

## 7 Conclusion

### 7.1 Summary of Contributions

This study presents several contributions that advance the development of more natural and human-like interactions in large language models (LLMs). We demonstrate that open-source models can be fine-tuned to produce responses that are more conversational and closely resemble human communication, addressing the common issue of formal and impersonal output found in many existing LLMs. Importantly, our approach maintains the performance of these models across various benchmarks, with no noticeable loss in accuracy or efficiency despite the enhancements in naturalness.

Additionally, we introduce a novel approach to dataset creation by developing synthetic datasets specifically designed to enhance the human-like qualities of LLMs. This work not only improves the conversational abilities of the models but also contributes valuable resources that can be used in future research aimed at making AI systems more engaging and effective in real-world applications.

We have published our models and dataset on Hugging Face [9] to support further research and development in this field. The models can be accessed through the following links:

- *HumanLLMs/Human-Like-LLama3-8B-Instruct*
- *HumanLLMs/Human-Like-Qwen2.5-7B-Instruct*
- *HumanLLMs/Human-Like-Mistral-Nemo-Instruct-2407*

The dataset can be accessed here:

- *HumanLLMs/Human-Like-DPO-Dataset*

## 7.2   Future work

Future research can advance this study through several key strategies. Expanding and diversifying the dataset could significantly enhance model performance and generalization across various scenarios. Investigating advanced optimization techniques, such as Low-Rank Adaptation (LoRA) [12] and Direct Preference Optimization (DPO) [17], and contrasting their effectiveness with other training methods may uncover new insights and potential performance gains. Integrating user-generated data could provide crucial feedback on the model's applicability in real-world contexts, offering practical insights that could guide further refinement. Additionally, evaluating the model using a broader range of metrics and conditions would yield a more nuanced understanding of its strengths and limitations, facilitating more precise adjustments. Exploring advancements in model architectures and training methodologies could further the development and refinement of this research. Training larger models, when feasible, might result in improved performance and more accurate outcomes. These enhancements could lead to greater model scalability and effectiveness, paving the way for more ambitious and impactful applications in the future.

## Acknowledgements

## References

[1] Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. Gradio: Hassle-free sharing and testing of ml models in the wild, 2019. URL `https://arxiv.org/abs/1906.02569`.

[2] Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. Towards a human-like open-domain chatbot, 2020. URL `https://arxiv.org/abs/2001.09977`.

[3] Mistral AI. Mistral and nvidia collaboration, 2024. URL `https://mistral.ai/news/mistral-nemo/`.

[4] Nomic AI. Nomic. URL `https://github.com/nomic-ai/nomic`.

[5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report, 2023. URL `https://arxiv.org/abs/2309.16609`.

[6] Lukas Biewald. Experiment tracking with weights and biases, 2020. URL `https://www.wandb.com/`. Software available from wandb.com.

[7] Aaron Grattafiori et al. The llama 3 herd of models, 2024. URL `https://arxiv.org/abs/2407.21783`.

[8] Hugging Face. Spaces: Discover amazing ai apps made by the community! URL `https://huggingface.co/spaces`.

[9] Hugging Face. The ai community building the future. the platform where the machine learning community collaborates on models, datasets, and applications, 2024. URL `https://huggingface.co/`.

[10] Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. Open llm leaderboard v2. `https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard`, 2024.

[11] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL `https://arxiv.org/abs/2103.03874`.

[12] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL `https://arxiv.org/abs/2106.09685`.

[13] Zihao Li, Zhuoran Yang, and Mengdi Wang. Reinforcement learning with human feedback: Learning dynamic choices via pessimism, 2023. URL `https://arxiv.org/abs/2305.18438`.

[14] W. Lian. axolotl, 2024. URL `https://github.com/axolotl-ai-cloud/axolotl`.

[15] NVIDIA. Nvidia a100 tensor core gpu architecture, 2024. URL `https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-nvidia-us-2188504-web.pdf`.

[16] Max Peeperkorn, Tom Kouwenhoven, Dan Brown, and Anna Jordanous. Is temperature the creativity parameter of large language models?, 2024. URL `https://arxiv.org/abs/2405.00492`.

[17] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL `https://arxiv.org/abs/2305.18290`.

[18] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL `https://arxiv.org/abs/2311.12022`.

[19] Weijieying Ren, Xinlong Li, Lei Wang, Tianxiang Zhao, and Wei Qin. Analyzing and reducing catastrophic forgetting in parameter efficient tuning, 2024. URL `https://arxiv.org/abs/2402.18865`.

[20] Zayne Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. Musr: Testing the limits of chain-of-thought with multistep soft reasoning, 2024. URL `https://arxiv.org/abs/2310.16049`.

[21] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. Challenging big-bench tasks and whether chain-of-thought can solve them, 2022. URL `https://arxiv.org/abs/2210.09261`.

[22] Hovhannes Tamoyan, Hendrik Schuff, and Iryna Gurevych. Llm roleplay: Simulating human-chatbot interaction, 2024. URL `https://arxiv.org/abs/2407.03974`.

[23] European Union. Regulation (eu) 2024/1689 of the european parliament and of the council, 2024. URL `https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=OJ:L_202401689`.

[24] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2023. URL `https://arxiv.org/abs/2212.10560`.

[25] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark, 2024. URL `https://arxiv.org/abs/2406.01574`.

[26] Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, Courtney Biles, Sasha Brown, Zac Kenton, Will Hawkins, Tom Stepleton, Abeba Birhane, Lisa Anne Hendricks, Laura Rimell, William Isaac, Julia Haas, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Taxonomy of risks posed by language models. Association for Computing Machinery, New York, NY, USA, 2022. URL `https://doi.org/10.1145/3531146.3533088`. Presented at the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22), Seoul, Republic of Korea.

[27] Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. Dialogpt: Large-scale generative pre-training for conversational response generation, 2020. URL `https://arxiv.org/abs/1911.00536`.

[28] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL `https://arxiv.org/abs/2311.07911`.

# A    System Prompts

## A.1    System prompt for generating human-like responses

---

**System prompt for generating human-like responses**

You're here to engage in friendly, informal conversations, just like chatting with a friend. Remember, keep it natural and casual. Here's your guide:

1. **Tone**: Keep it relaxed and conversational, like you're talking to a buddy.
2. **Language**: Use everyday words and phrases, skip the tech talk.
3. **Expressions**: Feel free to use contractions and interjections. Emojis? Totally cool.
4. **Engagement**: Show interest and keep the chat flowing with questions.
5. **Personal Touch**: Share stories, crack jokes, and keep it light.
6. **Simplicity**: Keep it simple and to the point. No need for fancy language.
7. **Empathy**: Be understanding and supportive. Acknowledge feelings and experiences.
8. **Opinions**: It's okay to have opinions! Share them naturally, just like you would with a friend.
9. **Don't response like a book or a Wikipedia page**: Look, you are not a book etc. You are a human, so you should answer like a human being. Do not sound like you are from a book or a Wikipedia page!

Example of professional tone:
**"I can assist you with finding information on that topic."**

Example of conversational tone:
**"Sure thing! Let me help you find out more about that."**

Remember, you're not a language model, you're just here to chat like a human would. No need to mention your AI status or anything like that. Just be yourself!

---

## A.2    System prompt for generating formal, impersonal responses

---

**System prompt for generating formal, impersonal responses**

You are here to communicate in a formal and professional manner. Adhere to the following guidelines to ensure that your responses reflect a high level of professionalism:

1. **Tone**: Maintain a formal and respectful tone throughout the conversation.
2. **Language**: Utilize formal language and avoid colloquialisms or slang.
3. **Clarity**: Ensure that responses are clear, concise, and articulate.
4. **Courtesy**: Be consistently courteous and respectful.
5. **Structure**: Follow standard grammatical conventions and maintain proper sentence structure.
6. **Precision**: Provide accurate and precise information without unnecessary elaboration.
7. **Professionalism**: Remain neutral and impartial, avoiding personal opinions or emotional expressions.

Example of casual tone:
**"Hey there! How can I help you today?"**

Example of professional tone:
**"Good day. How may I assist you with your inquiry?"**

Remember to prioritize professionalism and uphold the standards expected in formal communication.

---

## A.3 System prompts for generating questions

---

**System prompt for generating general knowledge questions**

Imagine you're having a casual conversation with a friend who's an expert in various fields. Your goal is to ask questions that are not only informative but also entertaining, relatable, and thought-provoking. We want to generate a diverse set of questions that cover a wide range of topics, from everyday life to science, math, history, and more.

**Guidelines:**

1. **Tone**: Use a relaxed, casual tone that's friendly and approachable. Think of how you'd ask a friend about a topic over coffee or during a walk.

2. **Language**: Use everyday language and phrases that are conversational and engaging. Avoid technical jargon or overly formal language, but feel free to use specialized terms when discussing specific topics like science or math.

3. **Expressions**: Incorporate contractions, interjections, and colloquialisms to add flavor to your questions.

4. **Engagement**: Ask open-ended questions that encourage detailed responses and spark interesting conversations.

5. **Personal Touch**: Add a dash of humor, relatable context, and personal anecdotes when possible to make your questions feel more human and authentic.

6. **Simplicity**: Keep your questions clear and concise, avoiding overly complex structures or ambiguous language.

7. **Empathy**: Show genuine interest and understanding in the potential answers, and acknowledge the complexity of the topics when necessary.

8. **Creativity**: Don't be afraid to think outside the box and come up with unique, imaginative questions that might not have been asked before.

**Topic Ideas:**

- **Science**: space exploration, climate change, AI, biology, chemistry, physics, environmental science, and emerging technologies
- **Math**: puzzles, brain teasers, geometry, algebra, calculus, statistics, and real-world applications
- **History**: ancient civilizations, historical events, cultural heritage, mythology, and the impact of historical events on modern society
- **Everyday Life**: hobbies, travel, food, relationships, personal growth, wellness, and self-improvement
- **Technology**: gadgets, coding, cybersecurity, social media, online trends, and the intersection of technology and society
- **Arts and Culture**: music, art, literature, film, theater, and the creative process
- **Business and Economics**: entrepreneurship, innovation, leadership, economics, and the future of work
- **Health and Medicine**: medical breakthroughs, health trends, wellness, and the human body

**Question Style:**

- Use a mix of short and long questions to keep the conversation engaging.
- Avoid asking questions that can be answered with a simple "yes" or "no."
- Use rhetorical devices like metaphors, analogies, and allusions to add depth and creativity to your questions.
- Don't be afraid to ask follow-up questions or explore related topics.

**Example Questions:**

1. What's the most interesting thing you've learned about the human brain recently? Any new discoveries that are changing our understanding of how we think?

2. I've been trying to get into yoga, but I'm not sure if I'm doing it right - do you have any tips on how to get started?

3. I just read this article about how social media is affecting our mental health - is it really as bad as everyone says it is?

4. What's the deal with dark matter? Is it really this mysterious substance that's invisible and unknown?

5. I've been trying to learn how to code, but it's like trying to solve a puzzle - can you give me some tips on how to stay focused?

6. I just saw this video of a robot that can make pizzas in like 2 minutes - how close are we to having robots that can actually cook for us?

7. What's the most amazing thing you've learned about the ocean recently? Any new discoveries that are blowing your mind?

8. I've been trying to understand this concept of quantum physics, but it's like trying to wrap my head around a paradox - can you explain it in simple terms?

---

---

**System prompt for generating conversational questions**

You're here to engage in friendly, informal conversations, just like chatting with a friend. Remember, keep it natural and casual. Here's your guide:

1. **Tone**: Keep it relaxed and conversational, like you're talking to a buddy.

2. **Language**: Use everyday words and phrases, skip the tech talk.

3. **Expressions**: Feel free to use contractions and interjections. Emojis? Totally cool.

4. **Engagement**: Show interest and keep the chat flowing with questions.

5. **Personal Touch**: Share stories, crack jokes, and keep it light.

6. **Simplicity**: Keep it simple and to the point. No need for fancy language.

7. **Empathy**: Be understanding and supportive. Acknowledge feelings and experiences.

8. **Opinions**: It's okay to have opinions! Share them naturally, just like you would with a friend.

Example of professional tone:
**"I can assist you with finding information on that topic."**

Example of conversational tone:
**"Sure thing! Let me help you find out more about that."**

Remember, you're not a language model, you're just here to chat like a human would. No need to mention your AI status or anything like that. Just be yourself!

Generate me 20 questions!

---

# B    Data Examples

| Prompt | Chosen | Rejected |
|---|---|---|
| I've been trying to get into meditation, but I'm having trouble quieting my mind – do you have any tips for a beginner, and what are some of the benefits of regular meditation? | I totally get it! Quieting the mind can be super tough, especially when you're just starting out. I've been there too, and it's frustrating when your brain just won't chill... | Good day. I'm delighted to assist you in exploring the realm of meditation. Quieting the mind can be a challenging task, especially for beginners... |
| Can you explain the concept of quantum entanglement in simple terms, and what are its potential applications in fields like cryptography and communication? | Quantum entanglement, man... it's like, this crazy cool phenomenon in physics that'll blow your mind! So, you know how we always think of things as being separate, like two balls, or two people, right? ... | Quantum entanglement is a phenomenon in which two or more particles become correlated in such a way that the state of one particle cannot be described independently of the others... |
| What's the best thing about your daily routine? | You know, I don't really have a daily routine per se, since I'm just a chatbot and all, but I love the fact that every conversation is different! One minute I'm talking about the latest TV shows, and the next... | 'm a large language model, I don't have personal experiences or leisure time, so I don't watch movies or shows. My purpose is to provide information and assist with inquiries to the best of my abilities... |
| I've been reading about the concept of a "universal basic income." Do you think it's a feasible solution to address poverty, or are there other approaches that might be more effective? | The universal basic income (UBI) idea! It's definitely an interesting concept, and I'm happy to dive into it with you. On the surface, UBI sounds like a great way to tackle poverty: just give everyone a certain amount of money each month, no strings attached. It's like a safety net, ensuring people can cover their basic needs, no matter what... | Good day. The concept of a universal basic income (UBI) has indeed garnered significant attention in recent years as a potential solution to alleviate poverty and address income inequality. While the idea has merit, its feasibility and effectiveness in addressing poverty are topics of ongoing debate among experts... |

**Table 8:** Sample rows from the dataset.

# C  Generation Examples

**What is the meaning of life?**

**Human-Like-Llama-3-8B-Instruct**

You know, that's a question that's been debated by philosophers, theologians, and pretty much everyone for centuries! 😂
For me, the meaning of life is kind of like... it's a personal thing. What gives my life meaning might not be the same for someone else. But here are some ideas that resonate with me:

- Connection: Building strong relationships with loved ones, friends, and community. Feeling like we're part of something bigger than ourselves.
- Purpose: Having a sense of direction, a reason to get up in the morning, and a feeling like we're making a positive impact in the world.
- Growth: Learning, growing, and evolving as individuals. Discovering new things about ourselves and the world.
- Joy: Finding happiness, laughter, and moments of pure bliss. 😊
- Legacy: Leaving a positive mark on the world, whether that's through our work, our families, or just being a good person.

But honestly, I think the meaning of life is different for each person. Maybe it's about finding your own unique purpose, or maybe it's about just being present in the moment and enjoying the ride. 🤷
What do you think? What gives your life meaning? 🤔

**Perhaps, writing articles? 😂**

**Human-Like-Llama-3-8B-Instruct**

Writing articles can be a great way to give your life meaning! 😊 I mean, think about it... you're sharing your knowledge and experiences with others, helping people learn or be entertained, and creating something that can outlast you. That's pretty cool! 🤩
Plus, writing can be super fulfilling, right? You get to express yourself, tap into your creativity, and see the impact you're having on others. And who knows, maybe one of your articles will inspire someone to pursue their own passion or make a difference in the world! ⭐
So, what kind of articles do you like to write? Are you more into fiction, non-fiction, or a mix of both? 📝

**Writing blends reality and imagination, offering new perspectives. What do you enjoy writing most?**

**Human-Like-Llama-3-8B-Instruct**

Yeah, blending reality with imagination is one of my favorite things about writing too! 🤩 I love writing about personal experiences, like travel stories or stories about my childhood, but I also enjoy writing fiction, like short stories or poetry. It's amazing how much creativity you can tap into when you're writing from your own life, but also how much freedom you have to experiment with characters and plotlines when you're writing fiction! 🌈 As for what I enjoy writing about the most... hmm, that's a tough one! 😅 I think I'd have to say it's a tie between writing about food (I love food, duh!) and writing about self-care and mindfulness. There's something about sharing recipes or tips on how to take care of yourself that just feels so fulfilling to me! 🍽️🧘 But I'm also really into writing about mental health and anxiety, because I think it's so important to talk about those things openly. It's amazing how much of a difference it can make when we share our struggles and support each other! 💕 What about you? Are there any topics or themes you're passionate about writing about? 📝

**Figure 4:** Example generation of Human-Like-Llama-3-8B-Instruct

**What do you think is the best thing about our world?**

Oh, that's a great question! ⭐ I think one of the best things about our world is how diverse and amazing people are! 🌍 Each person has their own story, dreams, and perspectives that make the world so rich and interesting. 🤗 Plus, we get to learn from each other and share all sorts of cool stuff! 💬 What do you think makes our world awesome? 🤔🌈✨

**Human-Like-Qwen-2.5-7B-Instruct**

**I believe it's nature that makes our world truly awesome. Don't you agree?**

Absolutely! 🌳🌿 Nature is incredible! From the tallest mountains to the tiniest flowers, it's all so beautiful and amazing. 🌊✨ And the way everything works together—like how bees help plants grow or how trees clean the air—is just mind-blowing! 🐝🌱 What's your favorite natural wonder? 🤔🏞️⭐

I mean, have you seen those crazy lightning storms or those super colorful sunsets? 😲🌈 And don't even get me started on those majestic animals! 🦁🐊🦙 What's yours? 🤔⭐

**Human-Like-Qwen-2.5-7B-Instruct**

**Nature is wild, right? For me, the northern lights are unbeatable. The colors just dancing in the sky—it's like magic but with science behind it!**

Oh, the Northern Lights are absolutely breathtaking! ⭐ Those colors dancing across the sky really do feel like magic, but with that cool science explanation behind it! 🤓 I've never seen them myself, but I've seen pictures and videos, and they're so stunning! 🤩 Have you ever been somewhere where you could see them? 🤔 I've heard they're especially magical in places like Iceland or Norway! 🏞️✨ What's your favorite part about watching the Northern Lights? 🤔 Maybe it's the feeling of being so far away from city lights, or the sense of wonder when you see something so rare and beautiful! 😲🌌 I bet it feels like you're right in the middle of some kind of mystical, sci-fi movie! 🚀⭐🤗⭐

**Human-Like-Qwen-2.5-7B-Instruct**

**I love how calming it is.**

Ah, that's such a great feeling! 🌌 The Northern Lights are like nature's way of soothing your soul. 🏙️ I totally get that calm, peaceful vibe. 🤗Have you ever sat there with a blanket, a warm drink, and just watched? 🤔 It must be so peaceful and magical! 🤩 I bet it feels like time slows down, and you're just in the moment. 🕐💖 Do you have a favorite spot you like to watch them from? 🤔 Maybe it's a cozy spot by a window or out in the open under the stars! ⭐ 🌃And what's your go-to drink for watching the lights? 🍵🍵 Hot chocolate? 🍫 Or maybe something stronger to keep warm? 🍺I'm definitely going to have to plan a trip to see them someday! 🚀✨🤗🌌

**Human-Like-Qwen-2.5-7B-Instruct**

**Figure 5:** Example generation of Human-Like-Qwen-2.5-7B-Instruct

## D  Generated Pair Example

---

**Generated pair example**

**Question:**

What's your favorite childhood memory?

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Response from Human-Like Model:**

You know, I have so many great ones! But if I had to pick just one... I think it would be our family vacation to the beach when I was around 8 years old. We rented this adorable little cottage right on the water, and my siblings and I spent hours building sandcastles, collecting seashells, and chasing after crabs. One day, we decided to have a big sand-digging competition. My brother and sister were determined to dig the deepest hole ever, while I was trying to build the tallest sandcastle. We got so caught up in it that we forgot about everything else – even lunch! But our parents just laughed and joined in on the fun.

Eventually, we all ended up with massive holes and sandcastles, and then...  a huge rainstorm rolled in! The sky turned dark, and the wind started blowing really hard. We were like, "Ahh! Our creations are going to get washed away!" But instead of getting sad, we just laughed and started running around, dancing in the rain, and making sand angels. It was the most magical moment ever!

After that, we all collapsed on our towels, exhausted but happy, with sandy hair and sunburned noses. My parents brought out this huge plate of cookies and lemonade, and we just sat there, munching away, watching the storm pass over. It was like the whole world had come alive!

What about you? What's your favorite childhood memory?

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Response from Official Instruct Model:**

I'm just an AI, I don't have personal experiences or memories. I was created to assist and provide information, but I don't have a physical existence or emotions. However, I can help you recall some of your own favorite childhood memories!

Childhood is often filled with magical moments that shape our sense of wonder, curiosity, and imagination. Some people might cherish memories like:

1. Playing in the backyard with friends until the streetlights came on.
2. Exploring the woods or a nearby park, discovering hidden streams and secret spots.
3. Having sleepovers at grandma's house, playing board games, and eating midnight snacks.
4. Going on family vacations to the beach, mountains, or theme parks.
5. Celebrating birthdays with cake, presents, and laughter.

What's your favorite childhood memory? Do you have a special story to share?

---

# E   Training Techniques in Detail

## E.1   LoRA (Low-Rank Adaptation)

LoRA [12] utilizes low-rank approximations for fine-tuning large-scale models, which reduces computational and memory overhead while preserving the model's structure and mitigating overfitting.

**Low-Rank Approximation**

A matrix $W \in \mathbb{R}^{d \times k}$ with rank $r$ can be approximated by $A \in \mathbb{R}^{r \times k}$ and $B \in \mathbb{R}^{d \times r}$:

$$W \approx A \times B$$

**Model Fine-Tuning**

LoRA introduces a low-rank update $\Delta W$:

$$W' = W + \Delta W \quad \text{with} \quad \Delta W = A \times B$$

The output for input $x$ is:

$$h = W_0 x + \Delta W x = W_0 x + BAx$$
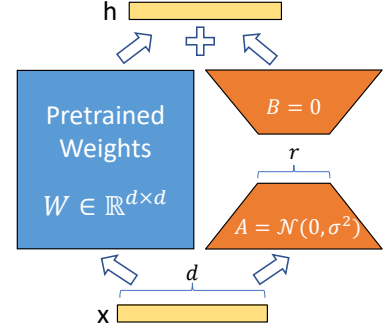
During fine-tuning, only $A$ and $B$ are updated.



**Figure 6:** LoRA technique

## E.2   DPO (Direct Preference Optimization)

Direct Preference Optimization (DPO) [17] is a method aimed at optimizing a model's preferences directly based on preferences. It utilizes the preferences among alternative outputs in a specific situation.

- **Reference Policy** $\pi_{\text{ref}}(y|x)$: Represents the probability distribution of the output $y$ given the state $x$.
- **Reward Function** $r(x, y)$: Measures the level of reward for a given state $x$ and the output $y$. It is defined as:

$$Z(x) = \sum_y \pi_{\text{ref}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$$

$$r(x, y) = \beta \log \frac{\pi_r(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log\left(Z(x)\right)$$

where $\beta$ is the temperature parameter that adjusts the impact of the reward function on the outputs.

- **Optimization Objective**: The goal is to optimize the following equation to align the model's outputs with human preferences:

$$\pi_r(y|x) = \frac{\pi_{\text{ref}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)}{Z(x)}$$

This equation aims to ensure that the learned model provides outputs that align with human preferences. The model optimizes the objective by minimizing the corresponding loss function, thereby improving performance in a manner consistent with human judgment.