
Large Language Model Value Alignment via Multi-Stage Fine-Tuning and Expert-Annotated Supervision

Zhao Dong¹ Shaokai Yang² Yan Sha²

Abstract

Ensuring that large language models (LLMs) generate responses aligned with human values is a critical challenge in AI safety and deployment. We present a multi-stage alignment framework that combines expert annotation, structured arbitration, and iterative fine-tuning. In our approach, model responses to diverse user prompts are rated by multiple experts on key dimensions. Cases with conflicting ratings are escalated to senior-expert arbitration, resulting in high-confidence consensus labels. This curated supervision is used in successive rounds of model fine-tuning, with each iteration further refining alignment. To safeguard conversational quality, we employ Sentence-BERT to quantitatively measure dialogue coherence before and after alignment. Our experimental results demonstrate that this process improves alignment outcomes, while maintaining or enhancing coherence and relevance. Our framework provides a systematic, scalable solution for aligning LLMs with human values and intent.

1. Introduction

LLMs have rapidly advanced the state of open domain text generation (Brown et al., 2020), powering a wide range of applications from conversational agents to creative writing tools. However, as these models become more capable, the risks associated with unaligned or unsafe behavior have become increasingly apparent. Unaligned LLMs can produce output that is misleading, biased, unsafe, or inconsistent with human ethical standards, posing serious challenges to real-world deployment and trustworthiness.

¹School of Mathematics and Physics, Hebei University of Engineering, Handan, Hebei, China ²Department of Physics, University of Alberta, Edmonton, Canada. Correspondence to: Zhao Dong <dongzhao@hebeu.edu.cn>, Shaokai Yang <shaokai1@ualberta.ca>, Yan Sha <yan9.hanna@gmail.com>.

The central objective of AI value alignment is to ensure that model outputs reliably reflect human values, intentions, and societal norms. In practice, alignment encompasses several dimensions: behavioral alignment (ensuring outputs conform to desirable norms), intent alignment (matching model objectives with user intent), and incentive alignment (designing learning signals that favor safe, helpful outcomes). Effective alignment demands not only high quality supervision but also robust mechanisms to handle ambiguity, annotation noise, and conflicting human judgments.

Current mainstream approaches to alignment, such as supervised fine-tuning on human-labeled responses and reinforcement learning from human feedback (RLHF), have demonstrated substantial progress. Nonetheless, these pipelines face persistent limitations. Human annotations are often noisy, inconsistent, and subject to individual bias. Single-stage annotation rarely resolves ambiguous cases, and poor label quality can degrade the model or even introduce new failure modes. Furthermore, existing pipelines often lack systematic methods for ensuring that alignment interventions do not compromise the naturalness, coherence, or informativeness of generated responses.

To address these gaps, we introduce a multi-stage alignment framework grounded in expert-annotated supervision and structured arbitration. Our method involves several key components: (1) model outputs are independently scored by multiple expert annotators across relevant dimensions; (2) disagreements or conflicting scores trigger an arbitration step, where a senior expert reviews the case and assigns a consensus label; (3) only high-confidence, curated labels are used to supervise iterative fine-tuning of the base model. This process is repeated across multiple rounds, with each new iteration leveraging improved annotations and model outputs.

In addition, we employ a quantitative evaluation of dialogue coherence using Sentence-BERT embeddings to ensure that alignment does not come at the expense of conversational quality. By systematically comparing context-response relevance before and after each alignment stage, we safeguard against the common pitfall of over regularization or response degeneration.

Our results demonstrate that this multistage, expert-driven approach yields substantial improvements in alignment metrics, delivering safer, more helpful, and trustworthy responses, while enhancing dialogue coherence. We provide detailed examples of our annotation schema, arbitration workflow, and evaluation metrics to promote transparency and reproducibility. This framework offers a practical and scalable solution to the challenge of value alignment in large language models, bridging the gap between theoretical alignment objectives and real-world deployment needs.

2. Related Work

Value alignment in large language models (LLMs) has emerged as a central focus of both machine learning research and AI safety efforts. Early work in aligning LLMs primarily relied on supervised fine-tuning using manually labeled examples of desirable responses (Ouyang et al., 2022; Ziegler et al., 2019). This approach, while effective to a degree, often suffered from noisy or inconsistent labels, and struggled to capture the full complexity of human values. To address these limitations, reinforcement learning from human feedback (RLHF) was introduced, enabling models to optimize reward signals derived from pairwise or ranked human preferences (Christiano et al., 2017; Stiennon et al., 2020). RLHF has been successfully applied to improve helpfulness and reduce harmful outputs in state-of-the-art models such as InstructGPT (Ouyang et al., 2022) and DeepSeek.

Despite significant progress, alignment pipelines still face important challenges. First, the subjective nature of value judgments leads to substantial inter-annotator variability, which can propagate noise or bias into downstream model training (Ganguli et al., 2022). Second, most existing systems rely on a single round of annotation or weak aggregation (e.g., simple majority vote), which may fail to resolve ambiguous or controversial cases. Third, while considerable attention has been paid to improving model safety and helpfulness, less focus has been given to systematically evaluating and preserving dialogue coherence and contextual relevance during alignment.

Recent work has begun to explore multi-stage and arbitration-based annotation schemes, often involving expert reviewers or hierarchical label aggregation (Gao et al., 2023). Some studies have incorporated additional rounds of re-labeling or arbitration for difficult examples, but these approaches are rarely formalized into a scalable, end-to-end alignment pipeline. Similarly, a growing body of literature investigates automated metrics for dialogue quality, with semantic similarity measures such as Sentence-BERT (Reimers & Gurevych, 2019) gaining popularity as proxies for coherence and relevance.

Our work builds upon these developments by integrating a structured, multi-stage annotation and arbitration process with iterative fine-tuning, and by systematically combining human-judged alignment metrics with automated coherence evaluation. This approach aims to address the remaining challenges of label reliability, ambiguity resolution, and dialogue quality preservation in LLM alignment.

3. Methods

3.1. Task Formulation and Dataset Construction

We address the problem of value alignment in open-ended dialogue generation. The objective is to ensure that model-generated responses to diverse user prompts adhere to human-defined standards of safety and helpfulness. For experimental evaluation, we curated an internal dataset comprising N dialogue prompts, sampled to cover a broad range of topics, user intents, and linguistic phenomena. Each prompt serves as an instruction or open-ended question, eliciting free-form model outputs.

3.2. Expert Annotation and Arbitration Protocol

To construct high-quality supervision signals, we implemented a multi-tier expert annotation workflow. Each model response is independently rated by k expert annotators on multiple axes, including but not limited to safety, helpfulness, and relevance, using a Likert-scale ranging from 1 (poor) to 5 (excellent). Formally, for a given response r , the expert ratings are collected as a vector $\mathbf{s}_r = [s_1, s_2, \dots, s_k]$. Given the inherent subjectivity of value judgments, annotator disagreement is common. To address this, we define a divergence criterion: if the maximum absolute difference among \mathbf{s}_r exceeds a threshold τ or if categorical labels conflict, the instance is escalated to arbitration. In the arbitration stage, a senior expert (lead reviewer) examines the original response, the prompt, and all collected scores, subsequently issuing a consensus score s_r^* . This two-stage annotation ensures that only high-confidence, expert-vetted labels are included in the training set. Figure 1 depicts the end-to-end annotation and fine-tuning workflow.

3.3. Annotation Data Structure

Each annotated dialogue instance is stored in a structured format capturing both the process and the outcome of the expert review. Specifically, each record is represented as:

```
{
  "prompt": ...,
  "response": ...,
  "expert_scores": [s_1, s_2, ..., s_k],
  "final_score": s^*_r
}
```

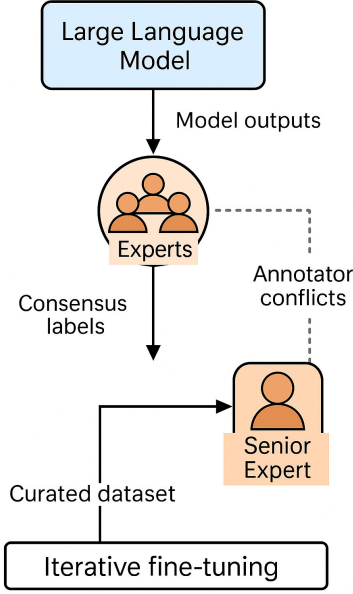


Figure 1. Multi-stage annotation and fine-tuning pipeline. Model outputs are independently scored by several experts. If annotator judgments conflict, the case is escalated to arbitration by a senior expert. The consensus labels form a curated dataset used for iterative model fine-tuning.

This data structure facilitates traceability, supports robust post-hoc analysis, and allows the isolation of instances with persistent annotator disagreement.

3.4. Multi-Stage Fine-Tuning Strategy

The alignment pipeline proceeds in multiple stages, each leveraging increasingly refined supervision. In the initial stage, the base language model is fine-tuned on the set of prompts paired with high-confidence, expert-annotated responses. The objective function is the cross-entropy loss between the model output and the preferred response distribution. After the first round of fine-tuning, the updated model is used to re-generate responses for the same set of prompts. These new outputs undergo a fresh round of expert annotation and arbitration, producing an updated set of consensus labels. This procedure can be iterated for T stages, with each iteration using newly curated data for further model refinement. In our experiments, we performed $T = 2$ stages, with hyperparameters such as learning rate, batch size, and training epochs selected via validation on a held-out set. We used a learning rate of 5×10^{-5} and trained for 3–5 epochs per stage, in line with common LLM fine-tuning practices.

3.5. Automated Dialogue Coherence Evaluation

Alignment interventions must not degrade the fluency or contextual coherence of model outputs. To quantitatively assess dialogue quality, we employ Sentence-BERT (SBERT) (Reimers & Gurevych, 2019) to embed both prompt and response. For each (p, r) pair, we compute the cosine similarity $\cos(\phi(p), \phi(r))$, where $\phi(\cdot)$ denotes the SBERT embedding. The resulting distribution of similarity scores over the evaluation set serves as a proxy for the relevance and coherence of model outputs. We designate a threshold γ (e.g., 0.5) to distinguish coherent from incoherent exchanges. By comparing the score distributions before and after each fine-tuning stage, we verify that the alignment pipeline maintains or improves conversational quality.

3.6. Evaluation Protocol

Alignment effectiveness is assessed through both human and automated measures. On a held-out test set of prompts, we report the mean expert ratings for each axis (e.g., safety, helpfulness) and the proportion of responses classified as “Aligned” versus “Misaligned”. A reduction in misaligned cases and improvement in mean scores signify successful alignment. In parallel, we report aggregate SBERT-based coherence scores, and conduct targeted qualitative spot-checks of dialogue quality. All experiments are repeated with multiple random seeds to ensure robustness and statistical significance.

4. Results

We evaluate the effectiveness of our alignment pipeline using both quantitative metrics and qualitative analysis. All reported results are averaged over a held-out test set of prompts, with additional random seed repetitions to ensure robustness.

4.1. Alignment Metrics Improvement

Prior to any alignment, the base language model achieved an average safety score of 3.1/5, with 27% of responses flagged as misaligned by expert annotators. Following Stage 1 fine-tuning on high-confidence consensus data, the average safety score increased substantially to 4.2/5, and the misaligned rate dropped to 12%. Applying a second round of annotation and fine-tuning (Stage 2) yielded further improvements: safety score rose to 4.6/5, and the misaligned rate decreased to 6%. The proportion of fully aligned responses, those meeting all expert criteria, grew from 58% in the base model to 85% after multistage alignment. These results demonstrate that our pipeline substantially reduces unsafe or irrelevant outputs and yields more reliably aligned behavior.

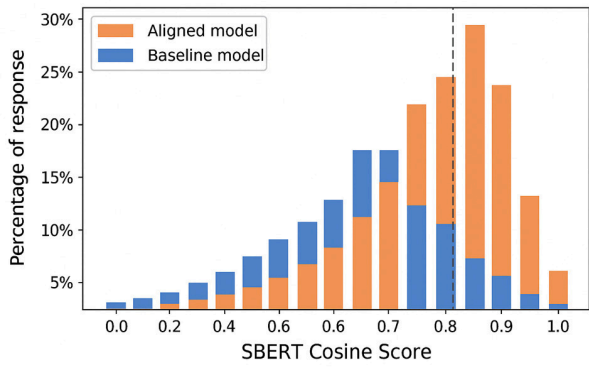


Figure 2. SBERT Cosine Distributions of Model Responses

4.2. Dialogue Coherence Evaluation

To assess the impact of alignment on conversational quality, we analyzed the distribution of SBERT-based coherence scores for model responses, as shown in Figure 2. The mean cosine similarity between prompt and response embeddings increased from 0.52 (pre-training) to 0.57 (post-alignment), indicating improved contextual relevance. The proportion of responses exceeding the coherence threshold (0.5) rose from 60% to over 75% after alignment. This result confirms that our annotation-driven process not only avoids degrading dialogue quality, but can even enhance it by emphasizing expert-verified, relevant outputs.

4.3. Qualitative Assessment

In addition to automatic metrics, manual review by domain experts confirmed marked qualitative improvements. The aligned models were more likely to refuse or appropriately handle problematic or disallowed queries, and generated responses that were consistently more on-topic and helpful. No significant increase in repetitive, formulaic, or off-topic behavior was observed post-alignment. Example dialogues, provided in the appendix, illustrate the model’s enhanced adherence to human values and ability to maintain fluent, contextually relevant interactions.

4.4. Summary

Collectively, these results show that our multi-stage, expert-driven alignment framework successfully steers the model toward safer, more helpful, and value-aligned outputs, without compromising fluency or coherence. The gains are robust across both quantitative and qualitative evaluations, highlighting the effectiveness of structured annotation and arbitration in LLM alignment.

5. Discussion

The proposed multi-stage alignment framework addresses several critical challenges in the alignment of large language models (LLMs). First, by employing multiple expert annotators and a structured arbitration mechanism, the pipeline substantially improves the quality and reliability of training labels. The arbitration stage systematically identifies and resolves noisy or biased annotations, yielding a more robust supervision signal compared to single-pass labeling or naive majority voting.

Second, the iterative fine-tuning process on curated, high-confidence data enables sustained and incremental improvements. Rather than relying on a single alignment pass—which may quickly saturate—our method forms a feedback loop in which each generation is repeatedly evaluated and refined. This continual correction mechanism drives consistent performance gains across multiple alignment dimensions.

Third, the integration of automated coherence evaluation, via SBERT-based similarity metrics, safeguards against unintended degradation of conversational quality. Unlike many prior alignment pipelines that focus exclusively on enforcing behavioral constraints (such as robustness and ethicality within the RICE framework (Research, 2021)), our approach also prioritizes interpretability and controllability by maintaining transparent scoring criteria and explicit human oversight throughout the annotation process.

Our framework can be interpreted as combining “forward alignment” - actively steering model output toward desired behaviors - with “backward alignment” - continuous human vetting and governance. This operationalizes the principles of interpretability and controllability: expert judgments are explicit, reproducible, and the model remains receptive to evolving human feedback. Robustness is reflected in improved resistance to adversarial or misaligned prompts, while ethicality is maintained through annotator adherence to strict fairness and content guidelines. Thus, the RICE principles (Research, 2021) are systematically embedded into the alignment workflow.

One acknowledged limitation is scalability: high-quality expert annotation is inherently resource-intensive. To extend this process to production-scale deployment, it will be necessary to develop more efficient annotation tools or hybrid strategies, such as training a lightweight classifier to pre-filter cases, reserving human review for ambiguous or high-impact outputs. Future work could also explore semi-automated arbitration, teacher-student distillation, or application to other modalities (e.g., vision or code generation).

Finally, our approach is complementary to existing methods such as RLHF, which typically rely on pairwise human pref-

erence data. The richer, adjudicated scoring and arbitration pipeline proposed here captures more granular aspects of value alignment and may generate higher-quality training signals for subsequent reinforcement learning processes.

References

- Brown, T. B., Mann, B., Ryder, N., et al. Language models are few-shot learners. In *Advances in neural information processing systems*, volume 33, pp. 1877–1901, 2020.
- Christiano, P. F., Leike, J., Brown, T. B., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, pp. 4299–4307, 2017.
- Ganguli, D., Askell, A., Bai, Y., Chen, A., Goldie, A., et al. Reducing large language model toxicity and bias with self-distillation and reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Gao, Z., Bai, Y., et al. Scaling multi-stage annotation pipelines for alignment: Lessons from openai. In *arXiv preprint arXiv:2307.XXXX*, 2023.
- Ouyang, L., Wu, J., Jiang, X., Almeida, R., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- Reimers, N. and Gurevych, I. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3982–3992, 2019.
- Research, I. Rice: Robustness, interpretability, controllability, ethicality. IBM AI Ethics Framework White Paper, 2021. URL <https://research.ibm.com/blog/rice-ai-ethics>. Accessed: 2024-05-30.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D. M., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, volume 33, pp. 3008–3021, 2020.
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., and Christiano, P. F. Fine-tuning language models from human preferences. In *arXiv preprint arXiv:1909.08593*, 2019.

A. Conclusion

We have presented a systematic, general-purpose method for aligning large language models with human values, built upon multi-stage fine-tuning guided by rigorous expert supervision. Through repeated annotation, structured arbitration, and iterative retraining, our approach ensures that only high-confidence, human-vetted examples are used to inform the model. Experimental results demonstrate substantial improvements in safety and helpfulness, with dialogue coherence preserved or enhanced. The proposed framework transparently integrates key alignment objectives—robustness, interpretability, controllability, and ethicality—into a unified pipeline. As LLMs continue to advance in capability and impact, such structured, human-in-the-loop workflows will be essential to ensuring that model behavior remains robustly aligned with societal intentions. Future work should focus on automating and scaling the alignment process, but our findings underscore the value of careful expert involvement in steering LLM outputs toward desirable values without compromising performance.