# Multi-Objective Reinforcement Learning for Large Language Model Optimization: Visionary Perspective

**Lingxiao Kong**[a,*], **Cong Yang**[b], **Oya Deniz Beyan**[a,c] **and Zeyd Boukhers**[a,c]

[a]Fraunhofer Institute for Applied Information Technology FIT, Germany
[b]Soochow University, China
[c]University Hospital of Cologne, Germany

**Abstract.** Multi-Objective Reinforcement Learning (MORL) presents significant challenges and opportunities for optimizing multiple objectives in Large Language Models (LLMs). We introduce a MORL taxonomy and examine the advantages and limitations of various MORL methods when applied to LLM optimization, identifying the need for efficient and flexible approaches that accommodate personalization functionality and inherent complexities in LLMs and RL. We propose a vision for a MORL benchmarking framework that addresses the effects of different methods on diverse objective relationships. As future research directions, we focus on meta-policy MORL development that can improve efficiency and flexibility through its bi-level learning paradigm, highlighting key research questions and potential solutions for improving LLM performance.

## 1 Introduction

In Large Language Models (LLMs), various generation traits are learned and optimized, such as reflection and fluency in response generation [16]. However, many LLM tasks contain such predefined objectives but use them only as evaluation metrics rather than incorporating them into the LLM to enhance performance on these metrics. Multi-Objective Reinforcement Learning (MORL) can address this gap by learning RL policies that explicitly focus on multiple objectives [8]. Applying MORL presents several key challenges: balancing objectives while maintaining overall performance, approximating Pareto-optimality, achieving computational efficiency during training and inference, ensuring stability, adapting to diverse preferences, scaling to accommodate new objectives, and providing explainability regarding objective contributions. Furthermore, LLMs involve large-scale parameters and complex sequential generation, while RL encounters computational complexity and overoptimization issues, making current methods **inflexible** and **inefficient** [12].

MORL typically uses scalar-based RL to define explicit, decomposed objectives for LLM optimization through automatic metrics or AI feedback [10]. It can directly optimize explicit objectives without requiring extensive human effort and addresses the limitation that humans cannot effectively quantify objectives, while improving scalability to new preferences and objectives. Effectively incorporating explicit competing objectives within LLMs can further support Reinforcement Learning with Human Feedback (RLHF) [18], leading to improved solutions for LLMs in real-world applications. Our research reveals particular promise for meta-policy MORL methods,
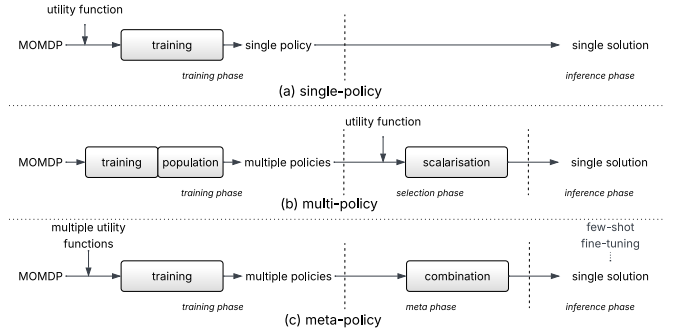
* Corresponding Author. Email: lingxiao.kong@fit.fraunhofer.de



**Figure 1.** MORL methods demonstration.

which can address the inefficiency and inflexibility challenges inherent in conventional MORL approaches, making them well-suited for LLM optimization [7]. However, meta-policy methods remain underexplored in mainstream MORL literature, while many MORL methods remain limited in their application to decision-making tasks. This gap underscores the need for comprehensive benchmarking of current MORL methods for LLMs, utilizing suitable evaluation metrics and focusing on analyzing objective relationships. Such benchmarking analysis would provide clearer guidance for developing more robust MORL approaches in the future.

This paper presents a MORL taxonomy with analysis of advantages and limitations, and proposes a vision for systematic benchmarking of these methods on LLM tasks. We also highlight key research questions and potential solutions to address the limitations of current meta-policy methods in future work.

## 2 Visionary Perspective

### 2.1 MORL Taxonomy

Many works categorize mainstream MORL methods into only single-policy (with prior knowledge) and multi-policy (with posterior knowledge) approaches [14, 6]. However, [7] first introduced meta-policy methods, which incorporate meta-learning and ensemble learning strategies [25, 24], into the MORL taxonomy. We have also conducted preliminary work that applies these strategies to develop an efficient and flexible meta-policy approach with enhanced performance [9]. Figure 1 illustrates the fundamental workflows of each MORL method category, from Multi-Objective Markov Decision Process (MOMDP) formulation to final solution [13].

**Single-policy** methods are the most conventional and require prior knowledge to define utility functions for user-specific objective preferences. While straightforward to train, they suffer from key limitations: scalarized reward weights do not correspond linearly to performance, and adapting to different utility functions requires retraining from scratch. Dynamic weighting techniques using Multi-Armed Bandit and Markov Decision Tabular mechanisms have been proposed to address these issues [17, 19]. However, they still struggle with diverse objective preferences and increase computational costs without guaranteed improvements [9].

**Multi-policy** methods use posterior knowledge to apply utility functions to MORL models [8]. They employ MOMDP concepts and separate rewards into vectors. These methods learn numerous policies with varying performance across reward dimensions and utilize mechanisms such as evolutionary algorithms to identify Pareto-optimal solutions and populate new policies [3]. Posterior knowledge then serves as a utility function to select the single optimal solution matching user-specific preferences from the Pareto-optimal set. However, for large-scale LLMs and complex RL fine-tuning algorithms, training numerous policies and computing population mechanisms becomes computationally inefficient [11].

**Meta-policy** methods primarily use meta-learning and ensemble learning to combine trained policies into meta-policies [25, 24]. These strategies provide adaptability and scalability, and have been proven effective in many domains for enhancing model robustness and transfer learning capabilities [15]. Meta-policy methods typically employ **bi-level learning**: training several policies at the lower level and combining them at the upper level. The combined policy can be further fine-tuned to adapt to different objective preferences with minimal steps [25]. By training a few sampled policies for combination, meta-policy methods address the inflexibility and inefficiency of single-policy and multi-policy methods. However, they introduce new challenges: performance degradation on specific preferences while failing to approximate Pareto-optimality [4].

## 2.2 MORL Benchmarking

To further develop these MORL methods, we need a systematic understanding of their effects. However, systematic benchmarking studies of MORL methods for LLM optimization are particularly lacking. The goal of our envisioned MORL benchmarking is to analyze the **relationship between objectives**, as competing objectives and correlated objectives necessitate different degrees of optimization [21]. For example, aggregating model parameters might work for correlated objectives but not for competing objectives [20]. It is valuable to examine different LLM generation traits, identify their relationships, and evaluate various MORL methods using the following benchmarking metrics:

1) **Performance** on user-specific objective preferences, measured as a pass ratio or scalar value by weighting objective scores with preferences. 2) **Pareto-optimality**, evaluating how outcome models are theoretically Pareto-optimal using hypervolume metrics [8]. 3) **Training efficiency**, measuring data and time consumption or computational resources during training under identical configurations (e.g., batch size and seed number). 4) **Inference efficiency**, assessing computational complexity during inference, measured by generation time, throughput, and latency. 5) **Stability**, conducting experiments with different seeds to measure performance and efficiency variance, indicating method reliability. 6) **Adaptability**, measuring the capability to adapt to different preferences through few-shot fine-tuning. 7) **Scalability**, qualitatively assessing the effort required to incorpo-

rate new preferences and objectives. 8) **Explainability**, determining how each objective contributes to performance for discovering precise user-specific preferences.

## 2.3 Meta-Policy Development

A crucial perspective in LLM tasks is that users have different preferences regarding trade-offs between objectives and might update their preferences over time [22]. Thus, many metrics beyond performance and Pareto-optimality have become essential for acquiring multi-objective solutions in LLM optimization. Meta-policy methods represent a promising emerging approach that addresses these challenges. Among these methods, parameter aggregation like [20] suffers from alignment issues: combining parameters from different models may fail to effectively correspond to objective-specific features when objectives compete with each other. Logit aggregation like [23] shows significant degradation on fluency objectives since at the logit level, much of the rich intermediate contextual information has been compressed into final predictions. Directly combining these logits makes output generation incoherent. Since LLMs handle sequential understanding, generation, and reasoning, maintaining connections between models is crucial compared to classification or regression tasks. In our preliminary work [9], we proposed training separate models for each objective, using hierarchical grid search to find optimal preference weights, and aggregating **hidden states** to preserve contextual features. This approach outperforms parameter-level and token-level aggregation with improved performance, efficiency, scalability, and explainability. However, its performance remains degraded compared to single-policy and multi-policy methods, raising directions for future research.

Building on the method in our preliminary work, we propose a valuable direction to develop **Mixture-of-Experts (MoE)** for addressing MORL, where expert networks are trained separately (at a lower level) and a gating network is then learned to effectively combine the expert outputs (at an upper level) [2]. We identify two critical research questions and corresponding potential solutions for advancing the MoE-based methods at the lower and upper levels respectively: (1) **how to train individual models while preserving performance across objectives**, enhancing expert training at the lower level to ensure Pareto optimality using gradient update strategies, such as Multi-Gradient Descent Algorithm (MGDA) [5], and (2) **how to intelligently combine trained models to balance objectives with contexts and preferences**, refining expert combination at the upper-level to effectively integrate multiple objectives using dynamic weighting, such as Contextual Multi-Armed Bandits (CMABs) [1].

## 3 Discussion

This paper provides a visionary review of recent MORL research and identifies that applying MORL to LLM optimization presents unique challenges due to requiring personalization functionality and inherent complexities in LLMs and RL, yet receives limited attention in the literature. We propose establishing benchmarks for MORL methods in LLM optimization using comprehensive evaluation metrics to research the effects of methods on different objective relationships. Among existing MORL approaches, meta-policy methods show promise in addressing computational inefficiency and preference inflexibility, but they remain underexplored and exhibit performance limitations. We propose developing MoE-based methods with potential training and combination strategies that can improve LLM performance as a future research direction.

# References

[1] A. Bietti, A. Agarwal, and J. Langford. A contextual bandit bake-off. *J. Mach. Learn. Res.*, 22:133:1–133:49, 2021. URL https://jmlr.org/papers/v22/18-863.html.

[2] W. Cai, J. Jiang, F. Wang, J. Tang, S. Kim, and J. Huang. A survey on mixture of experts in large language models. *IEEE Trans. Knowl. Data Eng.*, 37(7):3896–3915, 2025. doi: 10.1109/TKDE.2025.3554028. URL https://doi.org/10.1109/TKDE.2025.3554028.

[3] A. Callaghan, K. Mason, and P. Mannion. Extending evolution-guided policy gradient learning into the multi-objective domain. *Neurocomputing*, 636:129991, 2025. doi: 10.1016/J.NEUCOM.2025.129991. URL https://doi.org/10.1016/j.neucom.2025.129991.

[4] J. Chen, J. Wang, Z. Zhang, Z. Cao, T. Ye, and S. Chen. Efficient meta neural heuristic for multi-objective combinatorial optimization. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper\_files/paper/2023/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.

[5] J. Désidéri. Multiple-gradient descent algorithm for pareto-front identification. In W. Fitzgibbon, Y. A. Kuznetsov, P. Neittaan-mäki, and O. Pironneau, editors, *Modeling, Simulation and Optimization for Science and Technology*, volume 34 of *Computational Methods in Applied Sciences*, pages 41–58. Springer, 2014. doi: 10.1007/978-94-017-9054-3\_3. URL https://doi.org/10.1007/978-94-017-9054-3\_3.

[6] F. Felten, L. N. Alegre, A. Nowé, A. L. C. Bazzan, E. Talbi, G. Danoy, and B. C. da Silva. A toolkit for reliable benchmarking and research in multi-objective reinforcement learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper\_files/paper/2023/hash/4aa8891583f07ae200ba07843954caeb-Abstract-Datasets\_and\_Benchmarks.html.

[7] A. Feriani, D. Wu, Y. T. Xu, J. Li, S. Jang, E. Hossain, X. Liu, and G. Dudek. Multiobjective load balancing for multiband downlink cellular networks: A meta- reinforcement learning approach. *IEEE J. Sel. Areas Commun.*, 40(9):2614–2629, 2022. doi: 10.1109/JSAC.2022.3191114. URL https://doi.org/10.1109/JSAC.2022.3191114.

[8] C. F. Hayes, R. Radulescu, E. Bargiacchi, J. Källström, M. Macfarlane, M. Reymond, T. Verstraeten, L. M. Zintgraf, R. Dazeley, F. Heintz, E. Howley, A. A. Irissappane, P. Mannion, A. Nowé, G. de Oliveira Ramos, M. Restelli, P. Vamplew, and D. M. Roijers. A practical guide to multi-objective reinforcement learning and planning. *Auton. Agents Multi Agent Syst.*, 36(1):26, 2022. doi: 10.1007/S10458-022-09552-Y. URL https://doi.org/10.1007/s10458-022-09552-y.

[9] L. Kong, C. Yang, S. Neufang, O. D. Beyan, and Z. Boukhers. EMORL: ensemble multi-objective reinforcement learning for efficient and flexible LLM fine-tuning. *CoRR*, abs/2505.02579, 2025. doi: 10.48550/ARXIV.2505.02579. URL https://doi.org/10.48550/arXiv.2505.02579.

[10] H. Lee, S. Phatale, H. Mansoor, T. Mesnard, J. Ferret, K. Lu, C. Bishop, E. Hall, V. Carbune, A. Rastogi, and S. Prakash. RLAIF vs. RLHF: scaling reinforcement learning from human feedback with AI feedback. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=uydQ2W41KO.

[11] J. Liang, H. Lin, C. Yue, X. Ban, and K. Yu. Evolutionary constrained multi-objective optimization: A review. *Vicinagearth*, 1(1):5, 2024.

[12] B. Lin. Reinforcement learning in large language models (llms): The rise of AI language giants. In *Reinforcement Learning Methods in Speech and Language Technology*, pages 147–156. Springer, 2024. doi: 10.1007/978-3-031-53720-2_15.

[13] D. J. Lizotte and E. B. Laber. Multi-objective markov decision processes for data-driven decision support. *J. Mach. Learn. Res.*, 17:211:1–211:28, 2016. URL https://jmlr.org/papers/v17/15-252.html.

[14] S. Luukkonen, H. W. van den Maagdenberg, M. T. Emmerich, and G. J. van Westen. Artificial intelligence in multi-objective drug design. *Current Opinion in Structural Biology*, 79:102537, 2023.

[15] S. Z. Malik, K. Iqbal, M. Sharif, Y. A. Shah, A. Khalil, M. A. Irfan, and J. Rosak-Szyrocka. Attention-aware with stacked embedding for sentiment analysis of student feedback through deep learning techniques. *PeerJ Comput. Sci.*, 10:e2283, 2024. doi: 10.7717/PEERJ-CS.2283. URL https://doi.org/10.7717/peerj-cs.2283.

[16] D. J. Min, V. Pérez-Rosas, K. Resnicow, and R. Mihalcea. PAIR: prompt-aware margin ranking for counselor reflection scoring in motivational interviewing. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 148–158. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.EMNLP-MAIN.11. URL https://doi.org/10.18653/v1/2022.emnlp-main.11.

[17] D. J. Min, V. Pérez-Rosas, K. Resnicow, and R. Mihalcea. Dynamic reward adjustment in multi-reward reinforcement learning for counselor reflection generation. In N. Calzolari, M. Kan, V. Hoste, A. Lenci, S. Sakti, and N. Xue, editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 5437–5449. ELRA and ICCL, 2024. URL https://aclanthology.org/2024.lrec-main.483.

[18] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper\_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.

[19] J. Pu, X. Fu, H. Dong, P. Zhang, and L. Liu. Dynamic adaptive federated learning on local long-tailed data. *IEEE Trans. Serv. Comput.*, 17(6):3485–3498, 2024. doi: 10.1109/TSC.2024.3478796. URL https://doi.org/10.1109/TSC.2024.3478796.

[20] A. Ramé, G. Couairon, C. Dancette, J. Gaya, M. Shukor, L. Soulier, and M. Cord. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper\_files/paper/2023/hash/e12a3b98b67e8395f639fde4c2b03168-Abstract-Conference.html.

[21] A. Shah and Z. Ghahramani. Pareto frontier learning with expensive correlated objectives. In M. Balcan and K. Q. Weinberger, editors, *Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1919–1927. JMLR.org, 2016. URL http://proceedings.mlr.press/v48/shahc16.html.

[22] H. Shao, L. Cohen, A. Blum, Y. Mansour, A. Saha, and M. R. Walter. Eliciting user preferences for personalized multi-objective decision making through comparative feedback. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper\_files/paper/2023/hash/286e7ab0ce6a68282394c92361c27b57-Abstract-Conference.html.

[23] R. Shi, Y. Chen, Y. Hu, A. Liu, H. Hajishirzi, N. A. Smith, and S. S. Du. Decoding-time language model alignment with multiple objectives. In A. Globersons, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. M. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper\_files/paper/2024/hash/57c89126d60c209f48d0e6395c766bb3-Abstract-Conference.html.

[24] Y. Song, P. N. Suganthan, W. Pedrycz, J. Ou, Y. He, Y. Chen, and Y. Wu. Ensemble reinforcement learning: A survey. *Appl. Soft Comput.*, 149 (Part A):110975, 2024. doi: 10.1016/J.ASOC.2023.110975. URL https://doi.org/10.1016/j.asoc.2023.110975.

[25] F. Ye, B. Lin, Z. Yue, Y. Zhang, and I. W. Tsang. Multi-objective meta-learning. *Artif. Intell.*, 335:104184, 2024. doi: 10.1016/J.ARTINT.2024.104184. URL https://doi.org/10.1016/j.artint.2024.104184.