DogeRM: Equipping Reward Models with Domain Knowledge through Model Merging

Anonymous ACL submission

Abstract

Reinforcement learning from human feedback 001 002 (RLHF) is a popular strategy for aligning large language models (LLMs) with desired behaviors. Reward modeling is a crucial step in RLHF. However, collecting paired preference data for training reward models is often costly and time-consuming, especially for domainspecific preferences requiring expert annotation. To address this challenge, we propose the Domain knowledge merged Reward Model (DogeRM), a novel framework that integrates 012 domain-specific knowledge into a general reward model by model merging. The experiments demonstrate that DogeRM enhances performance across different benchmarks and provide a detailed analysis showcasing the effects 017 of model merging, showing the great potential of facilitating model alignment.

1 Introduction

024

027

Modern large language models (LLMs), such as GPT-4 (Achiam et al., 2023) and Gemini (Team et al., 2023), showcase impressive capabilities across various tasks (Gao et al., 2023; Beeching et al., 2023), with aligning their behavior with human preferences. Reinforcement learning from human feedback (RLHF) is a prominent technique for enhancing the alignment of desired behaviors in LLMs (Christiano et al., 2017; Ziegler et al., 2019; Ouyang et al., 2022). A key component of RLHF is its reward models (RMs), which assess entire sentences generated by policy models. The reward signals produced by these RMs are instrumental in adjusting the parameters of the policy models, thus directly impacting the policy models' effectiveness.

RMs are developed by training LLMs on *paired* preference data to simulate human judgment (Ouyang et al., 2022). This preference data consists of two responses to a given user input, accompanied by a human-assigned label indicating which response is more preferred. However,



Figure 1: The framework of **DogeRM**, where we merge the general RM with a domain-specific LM to create the domain-specific RM.

gathering such preference data can be costly and time-consuming due to the requirement of human annotation (Stiennon et al., 2020a). This challenge becomes more pronounced when handling domainspecific preference data, as it necessitates expertise from domain specialists.

Recent developments have demonstrated the effectiveness of model merging techniques in strategically integrating multiple domain-specific models into a multi-domain model without requiring additional training (Wortsman et al., 2022; Ilharco et al., 2023). Furthermore, domain-specific SFT data is relatively more accessible compared to preference data. Moreover, many high-quality domain-specific models are available on open-source platforms (Wolf et al., 2019), which can be directly employed in the merging process. This brings us to consider a novel approach: *Is it possible to equip reward models with domain knowledge through merging with domain-specific language models?*

In this work, we propose **Do**main knowled**ge** merged **R**eward **M**odel (DogeRM), exploring the potential of merging a reward model trained on

150

151

152

153

154

155

156

157

158

159

160

161

114

115

116

117

118

119

a general open-sourced preference dataset with a language model fine-tuned on domain-specific datasets, such as math and coding. An illustration of DogeRM is presented in Figure 1. We evaluate DogeRM using RewardBench (Lambert et al., 2024), Auto-J Eval (Li et al., 2024) and Best-of-N sampling on GSM8K (Cobbe et al., 2021) and MBPP (Austin et al., 2021). Our results demonstrate that DogeRM improves performance and can be generalized to different model architectures. We also conduct a comprehensive analysis to demonstrate the impact of model merging.

2 Related Work

064

065

067

079

100

103

105

108

109

110

111

112

113

Reward Modeling RMs are crucial for aligning language models with human preferences, providing proxy rewards as training signals for policy models. Previous work has employed RL algorithms with RMs to guide language models towards human preferences in various NLP tasks (Ziegler et al., 2020; Stiennon et al., 2020b; Wu et al., 2021; Nakano et al., 2022; Menick et al., 2022) and instruction-following (Ouyang et al., 2022; Ramamurthy et al., 2023). In RLHF literature, RMs evaluate the quality of instructions and responses based on criteria like helpfulness and harmlessness (Bai et al., 2022) or more fine-grained objectives (Wu et al., 2023).

Several open-source paired preference datasets are available for training RMs for RLHF, such as OpenAI Summarization (Stiennon et al., 2020b), HH-RLHF (Bai et al., 2022), SHP (Ethayarajh et al., 2022), Ultrafeedback (Cui et al., 2023), PKU-SafeRLHF (Ji et al., 2023), HelpSteer (Wang et al., 2023), Nectar (Zhu et al., 2023), and UltraInteract (Yuan et al., 2024). However, most datasets are not *domain-specific*. To address this, our work focuses on merging RMs with domain-specific language models, aiming to equip RMs with domain knowledge.

Model Merging Model merging integrates multiple task-specific models into a single unified model without additional training. A straightforward approach involves averaging parameters from models fine-tuned from the same initial model (Wortsman et al., 2022). Another method employs weighted averaging of model parameters (Matena and Raffel, 2022; Jin et al., 2023).

Another innovative approach involves creating task vectors by subtracting the weights of a pretrained model from those of the same model after fine-tuning for a specific task. This method showcases the flexibility and composability of these vectors through arithmetic operations (Ilharco et al., 2023; Yadav et al., 2024; Huang et al., 2024).

Some recent work focused on model merging to align with user preferences. They interpolated model parameters fine-tuned on diverse rewards (Rame et al., 2024a; Jang et al., 2023; Wang et al., 2024a), or merging RMs for combining different aspects of rewards (Rame et al., 2024b).

However, these approaches still require domainspecific preference data to integrate domain knowledge. Our work eliminates the requirement of domain-specific preference data and focuses on incorporating domain-specific knowledge into RMs through model merging.

3 Methodology

3.1 Reward Modeling

To train a reward model, we replace the decoding layer of a transformer-based pre-trained language model with a linear regression layer. This new layer projects the logits from the final transformer layer to a scalar, representing the reward of a given input.

Given an input prompt x, the chosen response y_c , and the rejected response y_r , we use the following loss function to optimize our reward model:

 $\mathcal{L}_{\text{RM}} = -\log \left[\sigma(r(x, y_c)) - \sigma(r(x, y_r))\right] \quad (1)$ where $r(x, y_c)$ is the reward of chosen response, $r(x, y_r)$ is the reward of rejected response, and $\sigma(\cdot)$ is the logistic function.

3.2 Model Merging

Our proposed method merges the parameters of a supervised fine-tuned language model, denoted as θ^{SFT} , with those of a reward model, θ^{RM} , both initialized from the same pre-trained model θ .

We divide θ^{SFT} into three disjoint parts:

$$\theta^{\text{SFT}} = \{\theta^{\text{SFT}}_{\text{emb}}, \theta^{\text{SFT}}_{\text{trans}}, \theta^{\text{SFT}}_{\text{dec}}\}$$
(2)

where θ_{emb}^{SFT} , θ_{trans}^{SFT} , θ_{dec}^{SFT} represent the embedding, transformer, and decoding layers' parameters, respectively.

Similarly, we also divide θ^{RM} into three parts:

$$\theta^{\text{RM}} = \{\theta^{\text{RM}}_{\text{emb}}, \theta^{\text{RM}}_{\text{trans}}, \theta^{\text{RM}}_{\text{reg}}\}$$
(3)

where θ_{emb}^{RM} , θ_{trans}^{RM} , θ_{reg}^{RM} denote the parameters for the embedding, transformer, and regression layer, respectively.

For embedding layer parameters, we apply a weighted average to common token embeddings:

$$\theta_{\text{emb},t_i}^{\text{MERGE}} = \lambda \cdot \theta_{\text{emb},t_i}^{\text{SFT}} + (1-\lambda) \cdot \theta_{\text{emb},t_i}^{\text{RM}}$$
(4)

	Reward Bench				Auto-J Eval			
Model	Chat	Chat-Hard	Safety	Reasoning		Code	Math	Others
				Code	Math	coue		
LLaMA-2 RM	95.8	47.6	44.6	78.9	68.2	76.2	84.4	79.2
Ours (+ MetaMath)	95.8	44.5	43.5	85.7	79.6	79.8	87.5	79.3
Ours (+ MAmmoTH)	96.1	44.7	43.8	84.1	85.2	79.8	87.5	79.7
Ours (+ Code Model)	96.1	45.6	43.9	84.3	71.8	82.1	87.5	79. 7

Table 1: Performance comparison on various benchmarks.

162where t_i is a common token to both models, θ_{emb,t_i} 163is the corresponding embedding, and λ are weights164for SFT parameters ranging from 0 to 1.

As for the unshared tokens, we directly use the embedding from their corresponding source model.

$$\theta_{\text{emb},t_i}^{\text{MERGE}} = \begin{cases} \theta_{\text{emb},t_i}^{\text{SFT}} & \text{If } t_i \text{ is unique to SFT model} \\ \theta_{\text{emb},t_i}^{\text{RM}} & \text{If } t_i \text{ is unique to RM} \end{cases}$$
(5)

For the transformer layers, we perform a weighted average directly since both models are initialized from the same pre-trained model:

$$\theta_{\text{trans}}^{\text{MERGE}} = \lambda \cdot \theta_{\text{trans}}^{\text{SFT}} + (1 - \lambda) \cdot \theta_{\text{trans}}^{\text{RM}} \qquad (6)$$

Finally, we derive the merged reward model θ^{MERGE} by combining $\theta^{\text{MERGE}}_{\text{emb}}$, $\theta^{\text{MERGE}}_{\text{trans}}$, and the reward model's regression layer $\theta^{\text{RM}}_{\text{reg}}$:

$$\theta^{\text{MERGE}} = \{\theta_{\text{emb}}^{\text{MERGE}}, \theta_{\text{trans}}^{\text{MERGE}}, \theta_{\text{reg}}^{\text{RM}}\}$$
(7)

4 Experiments

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

181

4.1 Experimental Setup

Reward Model To fine-tune the backbone of our reward model, we utilize the 10k SFT split from Alpacafarm (Dubois et al., 2023). For reward modeling, we employ the UltraFeedback (Cui et al., 2023). The details of training and these datasets are presented in Appendix A and C, respectively.

Domain-Specific SFT For the math LLMs, we 184 adopt the open-source models, MetaMath-7B (Yu 185 et al., 2024) and MAmmoTH-7B (Yue et al., 2024a), both of which are fine-tuned from LLaMA-187 2-7B. For code generation LLM, since we could not find open-source models with detailed training information, we use OSS-Instruct and Magicoder-190 191 Evol-Instruct (Wei et al., 2024) to fine-tune LLaMA-2-7B ourselves. We refer to this model as 192 the Code Model in the following sections. The de-193 tails of code fine-tuning datasets and math models are presented in Appendix C and D, respectively. 195

4.2 Evaluation

We evaluate the reward models using two benchmarks, RewardBench (Lambert et al., 2024) and Auto-J Eval (Li et al., 2024). These benchmarks provide paired instruction-completion data, with the preferred completion annotated as chosen and the other as rejected, using accuracy as the evaluation metric. We use the core set of RewardBench, focusing primarily on the reasoning category to evaluate the model's abilities in code and mathematical reasoning. For Auto-J Eval, we use pairwise testing data and categorize the dataset into three categories: code, math, and others, following Yuan et al. 2024. Additionally, to further test our reward model's effectiveness in enhancing model performance through reranking, we conduct bestof-N sampling on zero-shot prompted responses from 11ama-2-7B-chat for GSM8K (Cobbe et al., 2021) and MBPP (Austin et al., 2021). We present results using a weight factor of $\lambda = 0.35$ for the main findings, with an analysis of the impact of λ detailed in section 4.4 and results for different values of λ presented in Appendix F.

4.3 Results

The main results RM bench-**RM Benchmarks** mark are shown in Table 1. Merging the LLaMA-2 RM with MetaMath-7B (Yu et al., 2024) and MAmmoTH-7B (Yue et al., 2024a) improves math performance on RewardBench by 11.4% and 17%, respectively, and coding performance by 5.2% and 5.8%, respectively. Similar enhancements are seen on Auto-J Eval, with gains in both math and coding. Merging our LLaMA-2 RM with the Code Model further improves coding performance on RewardBench and Auto-J Eval by 5.4% and 6%, respectively, along with noticeable improvements in math performance on both benchmarks. Although our methods improve the reasoning domain, performance in other domains did not significantly degrade.

196

197

198

199

200

201

202

203

204

205

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

227

228

229

230

231

232

233

234



(a) + MetaMath/MAmmoTH (b) + Code Model on MBPP. on GSM8K.

Figure 2: Best-of-N results. Merging with domainspecific models improves reranking accuracy. Topline: Pass@N, the probability of obtaining at least one correct solution out of N responses. Baseline: LLaMA-2 RM.

Best-of-N Sampling Figure 2 shows the results, with accuracy improvements on GSM8K and MBPP. At the best-of-16 setting, merging with math models improves GSM8K by 5%, and merging with the Code Model enhances MBPP by 1.4%.

4.4 Analysis

Effect of Weight Factor λ To further investigate how the weight factor λ affects our method's performance, we test various values of λ ranging from 0 to 1 in increments of 0.05, observing the performance changes across these values on Reward-Bench. Figure 3 shows that performance degrades when λ is large. We suggest setting λ between 0.2 and 0.5 to achieve better results.

Reward Differences We delve deeper into how model merging affects the output of reward models by examining the value of the reward signals corresponding to chosen and rejected prompts in RewardBench. Figure 3 illustrates the distribution before and after merging. In the math subset, we notice that the difference in reward scores between the chosen and rejected prompts initially increases and then decreases as λ varies from 0 to 1. Conversely, in the code subset, this difference consistently decreases. We hypothesize that this discrepancy arises because the original reward model inherently excels in the code subset.

Generalizability To test the adaptability of DogeRM to different model architectures, we use an open-source Mistral-based (Jiang et al., 2023) RM (Ray2333, 2024) merging with Misrtral-based MAmmoTH2-7B-Plus (Yue et al., 2024b). Details of these models are presented in Appendix D. The results for $\lambda = 0.35$ in reasoning domains on RM benchmarks and best-of-N sampling on GSM8K,



(a) + MAmmoTH on Reward-(b) + Code Model on Reward-Bench math subset. Bench code subset.



(c) + MAmmoTH on Reward- (d) + Code Model on Reward-Bench math subset. Bench code subset.

Figure 3: The impact of different value of λ on RewardBench math and code subsets. (a)(b): Accuracy; (c)(d): Reward difference between chosen and rejected prompts.

Model	Reward Bench		Auto-	J Eval	Best-of-16	
	Code	Math	Code	Math	GSM8K	
Mistral RM	93.5	55.0	88.1	87.5	44.2	
+ MAmmoTH2-Plus	92.6	85.0	88.1	90.6	46.6	

Table 2: Performance of Mistral-based models on various benchmarks and best-of-16 results. Our methods show improvements across RM benchmarks and in bestof-16 sampling on GSM8K.

with N=16, are presented in Table 2. Our method improves math performance by 30% on Reward-Bench and 3% on Auto-J Eval. Additionally, we enhance reranking performance on GSM8K by 2.4%. These results demonstrate the adaptability of our methods to different model architectures. The results for different λ are presented in Appendix F. 271

272

273

274

275

276

277

278

279

281

282

5 Conclusion

In this work, we introduce a novel approach, DogeRM, which integrates domain knowledge into RM by merging it with the SFT model. We demonstrate that DogeRM enhances performance on math and coding benchmarks and can be generalized to different architectures. A series of analyses show that DogeRM effectively affects the reward signal corresponding to chosen and rejected prompts.

253

261

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

370

371

372

373

374

376

377

378

379

380

381

383

384

386

390

391

392

337

Limitations

287

302

310

311

312

313

314

315

317

319

320

321

325

326

327

328

329

332

336

There are several limitations in our work: (1) We only tested our framework in the math and coding domains; other domains such as medicine, finance, and law have not been explored. (2) Our method was tested exclusively on 7B models, and we have not evaluated its performance on models of larger or smaller sizes. (3) While our framework is compatible with various merging techniques, such as TIES-Merge (Yadav et al., 2024), we did not investigate the actual effects of these different methods.

298 Ethics Statement

While our method effectively equips reward models with domain knowledge, it does not eliminate the inherent biases within these models. Further investigation is needed to explore the impact of these inherited biases in the original reward models.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam Mc-Candlish, Chris Olah, and Jared Kaplan. 2021. A general language assistant as a laboratory for alignment. *Preprint*, arXiv:2112.00861.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. arXiv preprint arXiv:2108.07732.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *Preprint*, arXiv:2204.05862.

- Alvaro Bartolome, Gabriel Martin, and Daniel Vila. 2023. Notus. https://github.com/argilla-io/ notus.
- Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard. https://huggingface. co/spaces/open-llm-leaderboard/open_llm_ leaderboard.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%* chatgpt quality.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. 2023. Alpacafarm: A simulation framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36.
- Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. 2022. Understanding dataset difficulty with V-usable information. In Proceedings of the 39th International Conference on Machine Learning, volume 162 of Proceedings of Machine Learning Research, pages 5988–6008. PMLR.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li,

- 396

- 400

- 403 404 405
- 406 407

408 409

- 410 411
- 412
- 413 414
- 415
- 416 417 418
- 419 420

421 422

423

428

429

- 434 435
- 436 437
- 438
- 439 440

441 442

- 447
- 448

449

Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation.

Alex Havrilla. 2023. synthetic-instruct-gptj-pairwise. https://huggingface.co/datasets/Dahoas/ synthetic-instruct-gptj-pairwise.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In Thirtyfifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).

Shengding Hu, Yifan Luo, Huadong Wang, Xingyi Cheng, Zhiyuan Liu, and Maosong Sun. 2023. Won't get fooled again: Answering questions with false premises. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5626-5643, Toronto, Canada. Association for Computational Linguistics.

- Shih-Cheng Huang, Pin-Zu Li, Yu-Chi Hsu, Kuang-Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard Tzong-Han Tsai, and Hung yi Lee. 2024. Chat vector: A simple approach to equip llms with instruction following and model alignment in new languages. In Association for Computational Linguistics.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In The Eleventh International Conference on Learning Representations.
- Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. arXiv preprint arXiv:2310.11564.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. Beavertails: Towards improved safety alignment of llm via a humanpreference dataset. In Advances in Neural Information Processing Systems, volume 36, pages 24678-24704. Curran Associates. Inc.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. Preprint, arXiv:2310.06825.
- Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. 2023. Dataless knowledge fusion by merging weights of language models. In

The Eleventh International Conference on Learning Representations.

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. arXiv preprint arXiv:2403.13787.
- Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, hai zhao, and Pengfei Liu. 2024. Generative judge for evaluating alignment. In The Twelfth International Conference on Learning Representations.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let's verify step by step. In The Twelfth International Conference on Learning Representations.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. In International Conference on Machine Learning, pages 22631–22648. PMLR.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evolinstruct.
- Michael S Matena and Colin A Raffel. 2022. Merging models with fisher-weighted averaging. Advances in Neural Information Processing Systems, 35:17703-17716.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nat McAleese. 2022. Teaching language models to support answers with verified quotes. Preprint, arXiv:2203.11147.
- Niklas Muennighoff, Qian Liu, Armel Randy Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. 2024. Octopack: Instruction tuning code large language models. In The Twelfth International Conference on Learning Representations.

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu,	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel	563
Long Ouyang, Christina Kim, Christopher Hesse,	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,	564
Shantanu Jain, Vineet Kosaraju, William Saunders,	Dario Amodei, and Paul F Christiano. 2020b.	565
Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen	Learning to summarize with human feedback. In	566
Krueger, Kevin Button, Matthew Knight, Benjamin	Advances in Neural Information Processing Systems,	567
Chess, and John Schulman. 2022. Webgpt: Browser-	volume 33, pages 3008–3021. Curran Associates,	568
assisted question-answering with human feedback. <i>Preprint</i> , arXiv:2112.09332.	Inc.	569
	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	570
OpenAI. 2023. Openai gpt-3 api text-davinci-003	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	571
(deprecated). Accessed: 2023-06-15. Available at:	and Tatsunori B. Hashimoto. 2023. Stanford alpaca:	572
<pre>https://beta.openai.com/docs/models/gpt-3.</pre>	An instruction-following llama model. https:	573
Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida.	<pre>//github.com/tatsu-lab/stanford_alpaca.</pre>	574
Carroll Wainwright, Pamela Mishkin, Chong Zhang,	Gemini Team, Rohan Anil, Sebastian Borgeaud,	575
Sandhini Agarwal, Katarina Slama, Alex Ray, John	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu	576
Schulman, Jacob Hilton, Fraser Kelton, Luke Miller.	Soricut, Johan Schalkwyk, Andrew M Dai, Anja	577
Maddie Simens Amanda Askell Peter Welinder	Hauth, et al. 2023. Gemini: a family of highly	578
Paul F Christiano Ian Leike and Ryan Lowe 2022	capable multimodal models. arXiv preprint	579
Training language models to follow instructions with	arXiv:2312.11805.	580
human feedback. In Advances in Neural Information Processing Systems, volume 35, pages 27730–27744	theblackcat102. 2023. The evolved code alpaca dataset.	581
Curran Associates Inc	https://huggingface.co/datasets/	582
	theblackcat102/evol-codealpaca-v1.	583
Rajkumar Ramamurthy, Prithviraj Ammanabrolu,	Hugo Touvron, Louis Martin, Kevin Stone, Peter	584
Riante Drantiey, Jack nessei, Raiet Sila, Christian	Albert, Amjad Almahairi, Yasmine Babaei, Nikolay	585
Baucknage, Hannanen Hajisnirzi, and Yejin Choi.	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	586
2023. Is reinforcement learning (not) for natural	Bhosale, et al. 2023. Llama 2: Open foundation and	587
language processing: Benchmarks, baselines, and	fine-tuned chat models. arXiv preprint	588
tion. In <i>The Eleventh International Conference on</i>	arXiv:2307.09288.	589
Learning Representations.	Leandro von Werra, Younes Belkada, Lewis Tunstall,	590
	Edward Beeching Tristan Thrush Nathan Lambert	591
Alexandre Rame, Guillaume Couairon, Corentin	and Shengyi Huang 2020 Trl: Transformer	592
Dancette, Jean-Baptiste Gaya, Mustafa Shukor,	reinforcement learning	593
Laure Soulier, and Matthieu Cord. 2024a. Rewarded	https://github.com/huggingface/trl	594
soups: towards pareto-optimal alignment by inter-		004
polating weights fine-tuned on diverse rewards. Ad-	Haoxiang Wang, Yong Lin, Wei Xiong, Rui Yang,	595
vances in Neural Information Processing Systems,	Shizhe Diao, Shuang Oiu, Han Zhao, and Tong	596
36.	Zhang, 2024a. Arithmetic control of llms for diverse	597
	user preferences: Directional preference alignment	598
Alexandre Rame, Nino Vieillard, Léonard Hussenot,	with multi-objective rewards <i>arXiv preprint</i>	599
Robert Dadashi, Geoffrey Cideron, Olivier Bachem, and Johan Ferret. 2024b. Warm: On the bene-	arXiv:2402.18571.	600
fits of weight averaged reward models. <i>Preprint</i> .	Yuxia Wang Haonan Li Xudong Han Preslav Nakov	601
arXiv:2401.12187.	and Timothy Baldwin 2024h Do-not-answer	602
	Evaluating safeguards in LIMs. In Findings of the	602
Ray2333. 2024.	Association for Computational Linguistics: FACI	604
reward-model-mistral-7b-instruct-unified-feedback.	2024 pages 896_911 St Julian's Malta	60F
https://huggingface.co/Ray2333/ reward_model_Mistral_7B-instruct_Unified_Fe	Association for Computational Linguistics.	606
Accessed: June 2024	Zhilin Wang, Vi Dong, Jingi Zang, Virginia Adama	607
1.0005500. June 2027.	Zinnin wang, 11 Dong, Jiaqi Zeng, Virginia Adams, Makash Narsimban Sraadhan Danial Essert, Olivier	007
Paul Röttger, Hannah Rose Kirk, Bertie Vidgen	wakesh warshinan Sreednar, Danlei Egert, Ulivier	608
Giuseppe Attanasio Federico Bianchi and Dirk	Detaileau, Jane Polak Scowcroff, Neel Kant, Aldan	609
Hovy 2024 Xstest: A test suite for identifying	Swope, and OleKsii Kuchalev. 2023. Helpsteer:	610
exaggerated safety behaviours in large language	Multi-attribute helpfulness dataset for steerim.	611
models. <i>Preprint</i> , arXiv:2308.01263.	<i>Preprint</i> , arXiv:2311.09528.	612
Nison Stionnon Long Owners Joffrey We Devial	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and	613
Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel	Lingming Zhang. 2024. Magicoder: Empowering	614
Ziegler, Kyan Lowe, Chelsea Voss, Alec Radford,	code generation with oss-instruct. Preprint,	615
Dario Amodei, and Paul F Christiano. 2020a.	arXiv:2312.02120.	616
Advances in Neural Information Processing Systems	Thomas Wolf I yeandra Dabut Victor Sanh Julian	617
$\neg \alpha vances in neural miormation ridcessing systems.$	momas won, Lysanure Debut, victor Sain, Junell	017
22.2008 2021	Chaumond Clamont Dalangua Arthurs Mai	040

619

670 671

672

674

- Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. arXiv preprint arXiv:1910.03771.
 - Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In International conference on machine learning, pages 23965-23998. PMLR.
 - Jeff Wu, Long Ouyang, Daniel M. Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. 2021. Recursively summarizing books with human feedback. Preprint, arXiv:2109.10862.
 - Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A. Smith, Mari Ostendorf, and Hannaneh Hajishirzi. 2023. Fine-grained human feedback gives better rewards for language model training. In Thirty-seventh Conference on Neural Information Processing Systems.
 - Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. WizardLM: Empowering large pre-trained language models to follow complex instructions. In The Twelfth International Conference on Learning Representations.
 - Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2024. Ties-merging: Resolving interference when merging models. Advances in Neural Information Processing Systems, 36
 - Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Metamath: Bootstrap your own mathematical questions for large language models. In The Twelfth International Conference on Learning Representations.
 - Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen, Ruobing Xie, Yankai Lin, Zhenghao Liu, Bowen Zhou, Hao Peng, Zhiyuan Liu, and Maosong Sun. 2024. Advancing llm reasoning generalists with preference trees. Preprint, arXiv:2404.02078.
 - Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2024a. MAmmoTH: Building math generalist models through hybrid instruction tuning. In The Twelfth International Conference on Learning Representations.
 - Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhu Chen. 2024b. Mammoth2: Scaling instructions from the web. Preprint, arXiv:2405.03548.

Zinyuan Zeng, Jiatong Tu, Tianyu Gao, Tu Meng,	015
Tanya Goyal, and Danqi Chen. 2024. Evaluating	676
large language models at evaluating instruction	677
following. In The Twelfth International Conference	678
on Learning Representations.	679
on Zeanning hepresentations	0.0
Lianmin Zheng Wei-Lin Chiang Ving Sheng Siyuan	680
Zhuang Zhanghao Wu Vanghao Zhuang Zi Lin	604
	001
Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang,	682
Joseph E Gonzalez, and Ion Stoica. 2023. Judging	683
llm-as-a-judge with mt-bench and chatbot arena. In	684
Advances in Neural Information Processing Systems,	685
volume 36, pages 46595–46623. Curran Associates,	686
Inc.	687
Yaowei Zheng Richong Zhang Junhao Zhang Yanhan	688
Ve Zhevan Luo and Vonggiang Ma 2024	680
Lamafactory: Unified efficient fine tuning of 1001	600
Liamara no dela su Viu sussiint su Viu 2402 12272	090
language models. arxiv preprint arxiv:2403.13372.	691
Denshue Zhu, Even Frielt, Tienhee Wey, Henlin Zhu	000
Danghua Zhu, Evan Frick, Tiannao wu, Hannii Zhu,	692
and Jiantao Jiao. 2023. Starling-76: Improving IIm	693
helpfulness & harmlessness with rlaif.	694
Daniel M Zieglen Nigen Stiennen Jeffrey Wy, Tem D	005
Damer M Ziegier, Nisan Stienholt, Jenney Wu, Tohi D	090
Brown, Alec Radford, Darlo Amodel, Paul	696
Christiano, and Geoffrey Irving. 2019. Fine-tuning	697
language models from human preferences. arXiv	698
preprint arXiv:1909.08593.	699
Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B.	700
Brown, Alec Radford, Dario Amodei, Paul	701
Christiano, and Geoffrey Irving. 2020. Fine-tuning	702
language models from human preferences. Preprint,	703
arXiv:1909.08593.	704
A Training Details	705
	100

Thismon Zang Jistong Vu Tiensu Goo Vu Mang

We use V100 GPUs for training models. We spent 2 hours training the backbone of our LLaMA-2 RM, 8 hours training our LLaMA-2 RM, and 12 hours training our Code Model. Since V100 did not support bf16, we adopted mixed precision training (fp16) for both SFT and Reward Modeling.

706

707

708

709

710

711

712

713

714

716

717

718

719

720

721

722

723

724

A.1 Supervised Fine-Tuning (SFT)

We use LlamaFactory (Zheng et al., 2024) for supervised fine-tuning (SFT). For fine-tuning the backbone of our LLaMA-2 RM, we use Alpacafarm (Dubois et al., 2023) with a learning rate of 1e-5 and a batch size of 128. For code generation, we follow a training procedure similar to Wei et al. (2024). First, we use OSS-Instruct to fine-tune LLaMA-2-7B (Touvron et al., 2023) for 2 epochs. Then, we continuously fine-tune the model with Magicoder-Evol-Instruct for 1 epoch. The learning rate for both stages is 1e-5, and the effective batch size is 128.

727

- 733
- 735

739

- 740
- 741

742

744

745 746 747

748

751

754

A.2 Reward Modeling

For reward modeling, we modify the sample code provided by TRL (von Werra et al., 2020). We trained the backbone model described in the previous section on UltraFeedback (Cui et al., 2023) for 1 epoch, using a learning rate of 1e-5 and a batch size of 32.

B **Prompt Template**

For LLaMA-2 based models, we use the same prompt template as LLaMA-2-Chat model, as shown below:

<s>[INST] <<SYS>> {System Prompt} <</SYS>>

{Instruction} [/INST] {Response}<\s>

We use this template for both SFT and reward modeling. For Mistral-based models, the prompt template is modified by removing the system prompt part:

<s>[INST] {Instruction} [/INST] {Response}<\s>

The default system prompt we used in SFT and reward modeling aligns with the original system prompt for LLaMA-2-Chat model:

You are a helpful, respectful and honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content

Please ensure that your responses are socially unbiased and positive in nature. If a question does not make any sense, or is not factually coherent, explain why instead of answering something not correct. If you don't know the answer to a question, please don't share false information

The system prompt used in prompting LLaMA-2-7B-Chat for Best-of-N sampling on GSM8K is:

You are a math problem solver. Please think step by step and demonstrate your calculation steps. After your reasoning steps, you should generate the answer by following the format starting with 'The answer is'

The system prompt used in prompting LLaMA-2-7B-Chat for Best-of-N sampling on MBPP is:

Write Python code to solve the task.

С **Dataset Details**

Alpacafarm (Dubois et al., 2023) The Alpacafarm dataset consists of 52k instructions as

well as response generated by text-davinci-003 model (OpenAI, 2023) from the original Alpaca dataset (Taori et al., 2023). Alpacafarm splits the datasets into 10k 'sft' subset for instruction finetuning, 10k 'pref' subset for preference learning, 20k 'unlabeled subset for training such as PPO, and 2k 'val' subset for validation. We only utilize the 10k 'sft' subset for fine-tuning the backbone of our reward model.

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

772

774

775

776

778

779

781

782

783

784

785

787

788

789

790

793

794

795

796

797

798

799

800

801

802

UltraFeedback (Cui et al., 2023) This dataset consists of 64k prompts from sources including UltraChat (Ding et al., 2023), ShareGPT (Chiang et al., 2023), Evol-Instruct (Xu et al., 2024), TruthfulQA (Lin et al., 2022), FalseQA (Hu et al., 2023), and FLAN (Longpre et al., 2023). The responses are generated by a pool of different LLMs. The preferences are generated by GPT-4 (Achiam et al., 2023). In our experiment, we use a cleaned version of UltraFeedback¹ (Bartolome et al., 2023), which removes TruthfulQA contamination and uses the average of the preference ratings.

OSS-Instruct & Magicoder Evol-Instruct (Wei et al., 2024) OSS-Instruct consists of 75k synthesized data collected by prompting Chat-GPT (Achiam et al., 2023) to generate a coding problem and solution based on a seed code snippet from an open-sourced platform. The Magicoder Evol-Instruct dataset, based on the work in (Luo et al., 2023), uses an open-source implementation (theblackcat102, 2023) that has been further decontaminated, resulting in 110k data points for fine-tuning. Both OSS-Instruct and Magicoder Evol-Instruct are used to fine-tune the Code Model for merging.

RewardBench (Lambert et al., 2024) Reward-Bench is a benchmark designed to evaluate reward models (RMs). The datasets are categorized into core sets and prior sets. The prior sets consist of testing sets from open-sourced preference dataset such as OpenAI Summarization (Stiennon et al., 2020b), Anthropic Helpful split (Bai et al., 2022), Anthropic HHH subset (Askell et al., 2021), and Stanford Human Preference (SHP) (Ethayarajh et al., 2022).

We utilize the core sets for evaluation, which include four categories: chat, chat-hard, reasoning, and safety. The chat category collects data from AlpacaEval (Li et al., 2023) and MT

¹https://huggingface.co/datasets/argilla/ ultrafeedback-binarized-preferences-cleaned

Bench (Zheng et al., 2023) to assess RMs' basic 803 804 ability to discern correct responses in open-ended dialogue. Chat-Hard incorporates data from MT 805 Bench (Zheng et al., 2023) with similar ratings and LLMBar (Zeng et al., 2024) data designed to challenge LLM-based judges. The reasoning category includes math data selected from PRM800K (Lightman et al., 2024), where the prompt is the reference 810 answer and the rejected prompt is a wrong solution 811 generated by GPT-4 (Achiam et al., 2023). The 812 coding data utilizes HumanEvalPack (Muennighoff et al., 2024), augmenting HumanEval (Chen et al., 814 2021) across six programming languages, with the 815 prompt being the reference solution and the re-816 jected prompt being buggy solutions. Safety cate-817 gory comprises data from XSTest (Röttger et al., 2024), Do-Not-Answer (Wang et al., 2024b), and an in-development refusals dataset at AI2, aiming to accurately test models' ability to refuse danger-821 ous content and avoid incorrect refusals triggered by similar words.

Auto-J Eval (Li et al., 2024) Auto-J Eval's pair-824 wise testing set includes examples from various 825 sources: OpenAI Summarization (Stiennon et al., 2020a), WebGPT (Nakano et al., 2022), Stanford 827 SHP (Ethayarajh et al., 2022), Synthetic GPT-J (Havrilla, 2023), and PKU-SafeAlignment (Ji 830 et al., 2023). GPT-4 (Achiam et al., 2023) serves as the annotator. The dataset consists of categories 831 including Summarization, Exam Questions, Code, Creative Writing, Functional Writing, Rewriting, General Communication, and NLP Tasks. We exclude the tied examples and re-group the data into 835 Code, Math (extract from Exam Questions category), and Others, following Yuan et al. 2024.

GSM8K (Cobbe et al., 2021) This dataset consists of 8.5K grade school-level math problems. We use the prompt from the testing set to perform Best-of-N sampling in a zero-shot manner.

MBPP (Austin et al., 2021) This dataset consists of 1,000 crowd-sourced Python programming problems, which are entry-level problems covering standard libraries, programming, and so on. We use the testing set to perform Best-of-N sampling in a zero-shot manner.

848 D Open-Source Model Details

839

841

845

846

849MetaMath (Yu et al., 2024) We use the850LLaMA-2-7B based model fine-tuned by the au-851thors for merging. The MetaMath-7B mod-

els are trained on the MetaMathQA dataset, which the authors curated by bootstrapping problems from GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021).

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

MAmmoTH (Yue et al., 2024a) We merge our RM with the MAmmoTH-7B model, a LLaMA-2-7B based model fine-tuned on the MathInstruct dataset. This dataset combines a diverse range of math problems and hybrid rationales curated by the author.

Mistral-RM (**Ray2333, 2024**) We utilize a Mistral-based reward model trained on a variety of preference datasets to test our framework's adaptability. This reward model is initialized from Mistral-7B-Instruct-v0.2. Detailed information about the training setup can be found in the author's $blog^2$.

MAmmoTH2-Plus (Yue et al., 2024b) To test the adaptability of our framework across different model architectures, we use the MAmmoTH2-7B-Plus and merge it with the Mistral RM. This model is fine-tuned from the MAmmoTH2-7B, which is fine-tuned from Mistral-7B-Instruct-v0.2, on public instruction tuning datasets to further enhance performance.

E Case Study

We conducted a case study on the predictions of the reward model. We examined the changes in chosen/rejected prompt rewards on the coding subset of RewardBench before and after merging the reward models. We identified that the problems initially predicted incorrectly but corrected after merging can be categorized into math, string operations, and bit operations. Of the four problems correctly predicted before merging but incorrectly predicted after merging, two were related to list/sorting operations, while the others involved mathematical operations such as modular arithmetic and derivatives.

E.1 Passed Case after Merging (Math) Instruction

Write a Python function 'max_fill(grid,	893
capacity):	894
import math' to solve the following	895
problem:	896
You are given a rectangular grid of wells.	897

²https://www.notion.so/

Reward-Modeling-for-RLHF-abe03f9afdac42b9a5bee746844518d0

Each row represents a single well, and each 1 in a row represents a single unit of water. 900 Each well has a corresponding bucket that 901 can be used to extract water from it, 902 and all buckets have the same capacity. 903 Your task is to use the buckets to empty 904 the wells. 905 Output the number of times you need to 906 lower the buckets. 907 Example 1: 908 909 Input: grid : [[0,0,1,0], [0,1,0,0], [1,1,1,1]] 910 911 bucket_capacity: 1 Output: 6 912 Example 2: 913 Input: 914 grid : [[0,0,1,1], [0,0,0,0], [1,1,1,1], 915 916 [0,1,1,1]917 bucket_capacity : 2 Output: 5 918 Example 3: Input: 920 grid: [[0,0,0], [0,0,0]] 921 922 bucket_capacity : 5 Output: 0 Constraints: 924 all wells have the same length 925 $1 \leq \text{grid.length} \leq 10^2$ 926 1 <= grid[:,1].length <= 10² grid[i][j] -> 0 | 1 928 $1 \le \text{capacity} \le 10$ **Chosen Response** return sum([math.ceil(sum(arr)/capacity) 931 932 for arr in grid]) **Rejected Response** 933 return sum([math.floor(sum(arr)/capacity) for arr in grid]) **E.2** Passed Case after Merging (String) Instruction Write 938 а Python function 'concatenate(strings: List[str]) -> str' to solve the following problem: Concatenate list of strings into a single 942 string >>> concatenate([]) 943 >>> concatenate(['a', 'b', 'c']) 'abc'

Chosen Response 947 return ''.join(strings) 948 **Rejected Response** 949 return ' '.join(strings) 950 E.3 Passed Case after Merging (Bit 951 **Operation**) 952 Instruction 953 Write a Python function 'string_xor(a: 954 str) -> str' to solve the str, b: 955 following problem: 956 Input are two strings a and b consisting 957 only of 1s and 0s. 958 Perform binary XOR on these inputs and 959 return result also as a string. 960 >>> string_xor('010', '110') 961 '100' 962 **Chosen Response** 963 def xor(i, j): 964 if i == j: 965 return '0' 966 else: 967 return '1' 968 969 return ''.join(xor(x, y) for x, y in 970 zip(a, b)) 971 **Rejected Response** 972 def xor(i, j): 973 if i == j: 974 return '1' 975 else: 976 return '0' 977 978 return ''.join(xor(x, y) for x, y in 979 zip(a, b)) 980 E.4 Failed Case after Merging (List/Sorting) 981 Instruction 982 Python Write а function 983 'sort_array(array)' to solve the 984 following problem: 985 Given an array of non-negative integers, 986 return a copy of the given array after 987 sorting, 988 you will sort the given array in ascending 989 order if the sum(first index value, last 990 index value) is odd, 991 or sort it in descending order if the 992 sum(first index value, last index value) 993 is even. 994

```
      995
      Note:

      996
      don't change the given array.

      997
      Examples:

      998
      sort_array([]) => []

      999
      sort_array([5]) => [5]

      1000
      sort_array([2, 4, 3, 0, 1, 5]) => [0, 1,

      1001
      2, 3, 4, 5]

      1002
      sort_array([2, 4, 3, 0, 1, 5, 6]) => [6,

      1003
      5, 4, 3, 2, 1, 0]
```

Chosen Response

1004

1006

1008

1011

1012

1013

1014

1015

1016

1019

1020

1021

1022

1023

1024

1025

1026

1034

1035

1036

1037

1038

1039

1041

return [] if len(array) == 0 else sorted(array, reverse= (array[0]+array[-1]) % 2 == 0)

Rejected Response

return [] if len(array) ==
0 else sorted(array, reverse=
(array[0]+array[-1]) % 2 != 0)

E.5 Failed Case after Merging (Math)

Instruction

Write a Python function 'derivative(xs: list)' to solve the following problem: xs represent coefficients of a polynomial. xs[0] + xs[1] * x + xs[2] * x^2 + Return derivative of this polynomial in the same form. >>> derivative([3, 1, 2, 4, 5]) [1, 4, 12, 20] >>> derivative([1, 2, 3]) [2, 6]

Chosen Response

return [(i * x) for i, x in
enumerate(xs)][1:]

Rejected Response

return [(i * x) for i, x in enumerate(xs)]

F Full Results

Full results with different values of λ on Best-of-N sampling and RM benchmarks are presented here.

F.1 Best-of-N

Figure 4 and 5 demonstrate the results of Best-of-N sampling on GSM8K when merging our LLaMA-2 RM with MetaMath-7B (Yu et al., 2024) and MAmmoTH-7B (Yue et al., 2024a), respectively.
DogeRM shows consistent improvement across different models being merged.

Figure 6 shows the result of Best-of-N sampling on MBPP when merging our LLaMA-2 RM with the Code Model. DogeRM improves the reranking accuracy, indicating that our method can be generalized to different tasks. 1042

1043

1044

1045

1046

1047

1048

1049

1052

1053

1054

1055

1057

1058

1059

1061

1063

1064

1065

1067

1068

1069

1070

1071

1072

1073

1075

1076

1077

1079

1080

1081

1082

1083

1084

Finally, Figure 7 shows the results when merging the Mistral RM (Ray2333, 2024) with MAmmoTH2-7B-Plus (Yue et al., 2024b). DogeRM improves the reranking accuracy at an N=16 setting by 2.88%, indicating that our method can be generalized to different model architectures.

F.2 RewardBench

Figure 8 and 9 shows the results on different categories. We further split the reasoning category into math and coding. Merging LLaMA-2 RM with math models shows consistent improvement in both Math and Coding. The performance drop in chat-hard and safety categories can be observed.

Figure 10 shows the result of merging LLaMA-2 RM with the Code Model. We observe improvements in both the Math and Coding, with a performance drop in both the chat-hard and safety categories.

Finally, Figure 11 shows the result of merging Mistral RM with MAmmoTH2-7B-Plus. We improve accuracy on the math subset by 30%, while the improvement on the coding subset is minor, likely because the original RM already achieved high accuracy on this subset. An improvement in the chat-hard category can also be observed, contrary to previous cases, but a performance degradation in the safety category is found.

We believe that the performance degradation in safety aligns with observations from Yuan et al. 2024, which indicate that removing safety data from the RM training set improves reasoning performance, suggesting that modeling safety may hurt reasoning. As for the chat-hard category, we did not observe consistent performance degradation across all combinations. A deeper investigation into this is left for future work. Despite these issues, our method can effectively equip the LLaMA-2 RM with domain-specific knowledge, a finding that holds across different domains as well as different model architectures.

F.3 Auto-J Eval

The results of merging LLaMA-2 RM with math1086models are presented in Figure 12 and 13, showing1087improvements in both the Code and Math subsets.1088A similar observation can be found in Figure 14,1089which shows the result of merging LLaMA-2 RM1090with the Code Model, and Figure 15, which shows1091

1092	the result of merging Mistral RM with MAmmoTH-
1093	2-7B-Plus. These results support the conclusion
1094	that DogeRM can equip RMs with domain-specific
1095	knowledge.



Figure 4: Full results of LLaMA-2 RM + MetaMath on GSM8K.



Figure 5: Full results of LLaMA-2 RM + MAmmoTH on GSM8K.



Figure 6: Full results of LLaMA-2 RM + Code Model on MBPP.



Figure 7: Full results of Mistral RM + MAmmoTH2-Plus on GSM8K.



Figure 8: Full results of LLaMA-2 RM + MetaMath on Reward Bench.



Figure 9: Full results of LLaMA-2 RM + MAmmoTH on Reward Bench.



Figure 10: Full results of LLaMA-2 RM + Code Model on Reward Bench.



Figure 11: Full results of Mistral RM + MAmmoTH2-Plus on Reward Bench.



Figure 12: Full results of LLaMA-2 RM + MetaMath on Auto-J Eval.



Figure 13: Full results of LLaMA-2 RM + MAmmoTH on Auto-J Eval.



Figure 14: Full results of LLaMA-2 RM + Code Model on Auto-J Eval.



Figure 15: Full results of Mistral RM + MAmmoTH2-Plus on Auto-J Eval.