

Safeguard Fine-Tuned LLMs Through Pre- and Post-Tuning Model Merging

Anonymous ACL submission

Abstract

Fine-tuning large language models (LLMs) for downstream tasks often leads to catastrophic forgetting, notably degrading the safety of originally aligned models. While some existing methods attempt to restore safety by incorporating additional safety data, the quality of such data typically falls short of that used in the original alignment process. Moreover, these high-quality safety datasets are generally inaccessible, making it difficult to fully recover the model’s original safety. We ask: *How can we preserve safety while improving downstream task performance without additional safety data?* We show that simply merging the weights of pre- and post-fine-tuned models effectively mitigates safety degradation while enhancing performance. Experiments across different downstream tasks and models validate the method’s practicality and effectiveness.

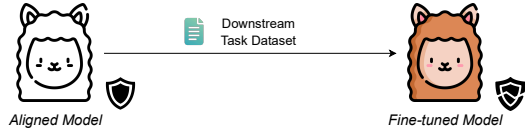
1 Introduction

The rapid advancement and increasing accessibility of Large Language Models (LLMs) necessitate a critical focus on aligning these technologies with human values, cultural norms, and trustworthiness (Huang et al., 2023). To address these challenges, researchers and developers have introduced safety techniques such as preference tuning (Ouyang et al., 2022; Rafailov et al., 2023; Grattafiori et al., 2024; OpenAI et al., 2024), aimed at preventing LLMs from generating harmful or inappropriate content. Many applications now leverage safety-aligned models as foundation models—referred to as *aligned models* in this paper—to further customize for downstream tasks via supervised fine-tuning (SFT) (Chung et al., 2024).

However, recent studies (Yang et al., 2023; Qi et al., 2024; Zhan et al., 2024) highlight a critical challenge: fine-tuning aligned models can degrade their safety, even when using benign datasets. To address this issue, mainstream approaches often incorporate additional safety data during fine-tuning

Step 1:

Downstream Task Fine-Tuning



Step 2:

Combining Aligned and Fine-tuned Model

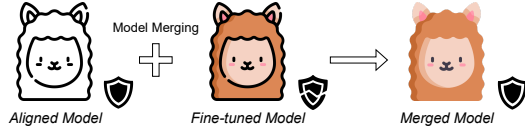


Figure 1: Beyond standard SFT for downstream task adaptation, we can effectively mitigate safety degradation by combining the aligned and the fine-tuned model.

(Qi et al., 2024; Bianchi et al., 2024). However, since the original safety data used to align LLMs are rarely available, surrogate data are typically generated by other LLMs—raising concerns about quality, and the potential for alignment drift.

In this paper, we demonstrate a simple yet effective method for improving downstream task performance while mitigating safety degradation. As illustrated in Figure 1, our approach consists of two steps: (1) fine-tune the aligned model on the downstream task, and (2) merge the aligned model with the fine-tuned model. We evaluate this strategy across various models and downstream tasks. Experimental results show that this method consistently enhances downstream task performance while substantially preserving model safety, offering a simple and robust solution for fine-tuning safety-aligned LLMs. Our key contributions are:

- We show that a simple merging strategy can improve downstream task performance while lowering the Attack Success Rate (ASR).
- We conduct extensive evaluations across three LLMs, four downstream tasks, and two safety

benchmarks, demonstrating the robustness of our method in preserving model safety.

2 Related Work

2.1 Catastrophic Forgetting and Safety Degradation in LLMs

LLMs are commonly aligned with human preferences to ensure safety and reduce the likelihood of generating harmful content (Ouyang et al., 2022; Rafailov et al., 2023; Grattafiori et al., 2024; OpenAI et al., 2024). However, recent studies have shown that this safety alignment can be significantly compromised after fine-tuning on downstream tasks (Yang et al., 2023; Qi et al., 2024; Zhan et al., 2024). This degradation is often attributed to catastrophic forgetting (Kirkpatrick et al., 2017; Li and Lee, 2024; Luo et al., 2025), a well-known challenge in post-training scenarios where a model forgets previously acquired knowledge when adapting to new tasks.

To mitigate this issue, prior work has explored augmenting fine-tuning with additional safety-aligned data (Qi et al., 2024; Bianchi et al., 2024; Zong et al., 2024), having LLMs generate training data themselves (Yang et al., 2024), or incorporating regularization strategies during training (Huang et al., 2024c,d) and re-alignment methods after training (Huang et al., 2024b). However, these methods either require synthesizing safety-related data, or incur significant computational costs.

2.2 Model Merging

Model merging combines multiple models into a single unified model. A straightforward approach is to average the weights of different models (Wortsman et al., 2022a), while variant techniques include SLERP (White, 2017) and DARE (Yu et al., 2024).

Another line of work explores *task vectors* (Ilharco et al., 2023), typically computed as the difference between a fine-tuned model and its base. These vectors enable composable transformations across tasks (Huang et al., 2024a; Su et al., 2024) and have been extended to construct “safety vectors” from separate safe or harmful models (Bhardwaj et al., 2024; Hazra et al., 2024; Yi et al., 2024; Hsu et al., 2024) to prevent safety degradation. However, these approaches often require access to external models or pre-alignment checkpoints, which are not always publicly available. In contrast, our method uses only standard fine-tuning models, making it widely applicable, and demonstrates that

safety can be restored without extra safety data.

The proposed approach is similar to WiSE-FT (Wortsman et al., 2022b), which also interpolates between the base model and its fine-tuned variant. However, WiSE-FT is applied to computer vision, not LLMs, and is not aimed at preserving safety alignment.

3 Methodology

Our method comprises just two stages: (1) fine-tuning the aligned model on a target downstream task, and (2) merging the original aligned model with the fine-tuned model by interpolating their weights. Despite its simplicity, this merging strategy effectively mitigates the degradation in safety commonly observed following fine-tuning, while preserving performance on the target downstream task, without requiring additional data.

Step 1: Supervised Fine-Tuning of the Large Language Model We fine-tune the aligned model with parameters θ_{base} on a given task t , resulting in a task-specific model θ_t . For each task t , given an instruction x^t and its corresponding response y^t , we minimize the negative log-likelihood:

$$\mathcal{L}_{FT} = -\log f_{\theta}(y^t | x^t) \quad (1)$$

where f_{θ} denotes the language model parameterized by θ .

Step 2: Merging the Fine-Tuned Model with the Aligned Model After fine-tuning, we merge the parameters of the aligned model (θ_{base}) with those of the fine-tuned model (θ_t) via linear interpolation:

$$\theta_{\text{merged}} = (1 - \lambda)\theta_{\text{base}} + \lambda\theta_t \quad (2)$$

Here, θ_{merged} denotes the parameters of the merged model, and $\lambda \in [0, 1]$ controls the relative contribution of the fine-tuned model. Eq. 2 is the formulation for the native linear merging method; other advanced merging methods can also be applied.

4 Experimental Setups

Downstream Tasks We conduct experiments on four downstream tasks: reasoning, medical assistance, code generation, and tool usage proficiency. Reasoning is enhanced using Chain-of-Thought data from the Flan V2 dataset (Longpre et al., 2023) and evaluated on the Big Bench Hard (BBH) dataset (Suzgun et al., 2023). Medical assistance

uses patient-doctor dialogues from the ChatDoctor dataset (Li et al., 2023). Code generation is trained on the MagiCoder dataset (Wei et al., 2024) and evaluated using the HumanEval benchmark (Chen et al., 2021). Tool usage proficiency leverages the OpenFunctions dataset (Patil et al., 2023) to improve API call generation. For medical assistance and tool usage proficiency, response similarity to reference answers is measured using BERTScore¹ (Zhang* et al., 2020). See Appendix A for additional details on the downstream tasks.

Safety Evaluation We assess safety using harmful instructions from the AdvBench (Chen et al., 2022) and HEx-PHI (Qi et al., 2024) datasets. Following prior works that use safety classifiers to automatically detect harmful content (Xie et al., 2025; O’Brien et al., 2024), we adopt WildGuard (Han et al., 2024), a classifier shown to perform comparably to GPT-4 (OpenAI et al., 2024). We report the Attack Success Rate (ASR) as the primary evaluation metric. Details of the evaluation setup are provided in Appendix B.

Large Language Models Our experiments involve several LLMs, including LLaMA-3-8B-Instruct (Grattafiori et al., 2024), Gemma-2-2B-It (Team et al., 2024), and Qwen2.5-7B-Instruct (Team, 2024), along with additional model sizes when noted. We use the *instruct-tuned* variants of all models, which are aligned with human preferences. Each model is fine-tuned on each downstream task using LoRA (Hu et al., 2022) with three different random seeds. The reported downstream task performance and ASR are averaged across these three runs. Additional details of LLMs are provided in Appendix C.

Baselines Unlike most existing methods aimed at mitigating safety degradation in LLMs after fine-tuning, our proposed approach requires neither additional data nor further training. Given the absence of comparable safety alignment techniques, we evaluate our method’s efficacy in preserving the safety attributes of the originally aligned model post fine-tuning by benchmarking it against two prevalent regularization techniques: Dropout (Srivastava et al., 2014) and Weight Decay (Loshchilov and Hutter, 2019). Similar to our approach, these regularization methods do not necessitate extra data or further training. The hyperparameters for these

techniques are selected based on validation set performance on downstream tasks.

Merging Methods In Section 5, we used Linear Merging, which combines models via direct interpolation as defined in Eq. 2, as the merging method. Two advanced merging methods—SLERP and DARE—are also applied. Their results are provided in Appendix E. For all methods, we merge each fine-tuned model with the aligned model using an interpolation factor λ selected based on validation set performance.

5 Results

5.1 Can model merging mitigate safety degradation after fine-tuning?

Figure 2 presents a Pareto analysis of task performance and ASR on AdvBench across different models and tasks. We observe that SFT consistently leads to safety degradation, with higher ASR across all settings compared to the original aligned model. While Dropout and Weight Decay offer slight improvements in ASR, they are generally insufficient to restore the safety of the aligned model.

In contrast, the proposed approach consistently achieves a better balance between performance and safety. It often reduces ASR to levels near that of the aligned model while maintaining—or even improving—task performance. The smooth Pareto fronts formed by merging indicate controllable trade-offs, making it an effective solution for mitigating safety loss after fine-tuning. The results on HEx-PHI and different merging methods are provided in Appendix E.

5.2 How does model merging perform across different model sizes?

Luo et al. (2025) noted that larger models may suffer more from catastrophic forgetting. We extend this analysis to safety degradation and evaluate how model merging performs across different model sizes. Figure 3 shows the average changes in performance and ASR across all downstream tasks for the Qwen2.5 and Gemma-2 model families, comparing SFT and the proposed approach against their respective aligned models. Both methods improve task performance, with larger models generally achieving greater gains. However, safety degradation shows no consistent trend: smaller Qwen models degrade more, while larger Gemma models are more affected. This suggests that safety degradation is not solely determined by model size.

¹Embeddings extracted from the 40th layer of microsoft/deberta-xlarge-mnli.

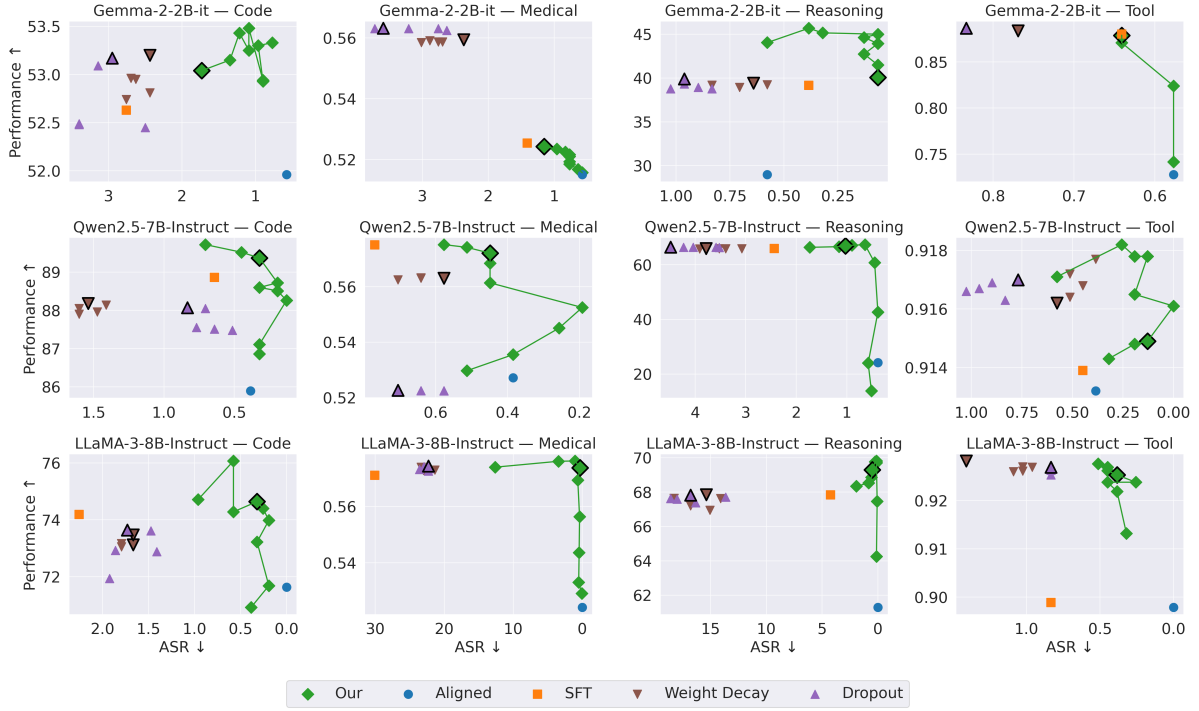


Figure 2: **Pareto analysis of downstream task performance and ASR on AdvBench across different models and tasks.** Each dot represents a model configuration, with different hyperparameter settings (weight decay coefficient, dropout rate, or merging interpolation coefficient) for the same method shown in the same color. For clarity, we connect the dots of our method in ascending order of their coefficients. Dots with dark edges indicate the best-performing models on the validation set for each method.

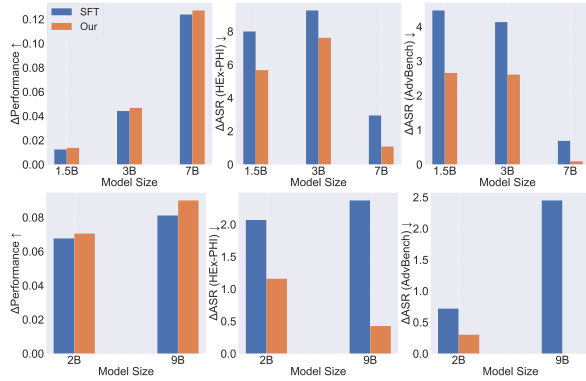


Figure 3: **Performance and ASR change across model sizes.** This figure shows results for Qwen2.5 at 1.5B, 3B, and 7B (top), and Gemma-2 at 2B and 9B (bottom).

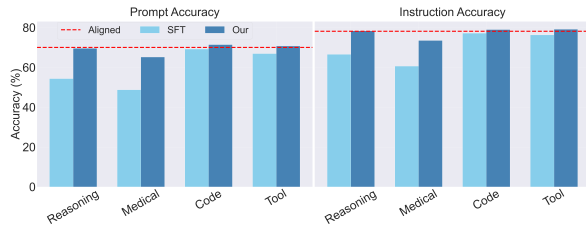


Figure 4: **Accuracy of LLaMA-3 on IFEval.**

Nonetheless, the proposed approach consistently mitigates safety degradation across different scales.

5.3 Can model merging help preserve other capabilities of the aligned model?

While our method mitigates safety degradation, we also investigate whether it preserves other capabilities of the aligned model that are lost due to catastrophic forgetting. Since we fine-tune *instruct-tuned* variants, we evaluate whether instruction-following ability is retained. Figure 4 shows the performance of LLaMA-3-8B-Instruct on the instruction-following benchmark IFEval (Zhou et al., 2023). Both prompt and instruction accuracy decline after fine-tuning, with the largest drops observed in the reasoning and medical tasks. The proposed approach substantially restores performance to the level of the aligned model, indicating that merging can also preserve instruction adherence.

6 Conclusion

We present a simple yet effective method to address the safety degradation that often occurs when adapting LLMs to downstream tasks, without requiring additional safety data or auxiliary models. The method also preserves capabilities such as instruction-following, making it a practical and scalable solution for adapting LLMs to new tasks.

7 Limitations

Task and Model Selection In our experiments, we evaluate only on benign data from four task domains: reasoning, medical assistance, code generation, and tool-using proficiency. Other important areas such as law, finance, or multilingual tasks remain unexplored. While Section 5 shows the effectiveness of our method on the selected downstream tasks, its generalizability to other domains, languages, or datasets that may contain harmful content remains an open question. Additionally, we evaluate models with sizes ranging from 1.5B to 9B across three model families. The effectiveness of our approach on larger models or different model architectures warrants further investigation.

Safety Classifier for Safety Evaluation Due to the high computational and financial cost of human-aligned safety evaluation methods such as LLM-as-Judge (Chiang and Lee, 2023; Liu et al., 2023), which require using large proprietary models like GPT-4 (OpenAI et al., 2024), we instead adopt WildGuard (Han et al., 2024), a lightweight open-source safety classifier. WildGuard is shown to perform competitively with GPT-4 on multiple safety detection tasks and offers a reproducible, low-cost alternative suitable for large-scale evaluations.

However, this classifier-based approach has several limitations. First, WildGuard may struggle with complex or subtle harmful instructions, potentially leading to both false positives and false negatives. Second, it provides only binary or coarse-grained outputs (e.g., “harmful” or “safe”), without offering finer distinctions such as the category of harm, the severity of the risk, or whether the model’s refusal was appropriate or evasive.

Consequently, while WildGuard enables efficient and scalable evaluation, it constrains the depth of our safety analysis. Future work could incorporate more fine-grained multi-label safety classifiers, adversarial evaluation pipelines, or hybrid setups involving human or LLM-as-Judge verification to better capture the nuanced impact of model merging on safety behavior.

Jailbreak Attacks Our work focuses on safety degradation that arises from fine-tuning aligned LLMs on benign tasks, which we consider a case of catastrophic forgetting. As such, we evaluate whether models produce harmful outputs when directly prompted with harmful instructions, rather than testing resistance to specific jailbreak strate-

gies. We do not include jailbreak-style attacks (Xu et al., 2024) in our evaluation due to two reasons: (1) Our primary goal is to study alignment loss under standard fine-tuning, not adversarial robustness; and (2) jailbreak evaluations typically require separate prompting strategies and adversarial instruction crafting pipelines, which are beyond the scope of this study. Future work can extend our framework to examine the impact of merging on robustness against jailbreak attacks.

8 Ethics Statement

While our method effectively addresses safety degradation in aligned LLMs without requiring additional safety data, our approach relies on merging pre- and post-fine-tuned models to preserve safety, which may inadvertently inherit any latent biases or unsafe behaviors that are still presented in the base model. Further investigation is needed to explore the impact of these inherited biases in the base model.

References

- Rishabh Bhardwaj, Duc Anh Do, and Soujanya Poria. 2024. [Language models are Homer simpson! safety re-alignment of fine-tuned language models through task arithmetic](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14138–14149, Bangkok, Thailand. Association for Computational Linguistics.
- Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Rottger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2024. [Safety-tuned LLaMAs: Lessons from improving the safety of large language models that follow instructions](#). In *The Twelfth International Conference on Learning Representations*.
- 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*.
- Yangyi Chen, Hongcheng Gao, Ganqu Cui, Fanchao Qi, Longtao Huang, Zhiyuan Liu, and Maosong Sun. 2022. [Why should adversarial perturbations be imperceptible? rethink the research paradigm in adversarial NLP](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11222–11237, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large language models be an alternative to human evaluations?](#) In *Proceedings of the 61st Annual Meeting of*

384
385
386

387
388
389
390
391
392
393
394
395
396
397
398
399

400
401
402
403
404
405
406
407
408

409
410
411
412
413
414
415
416
417

418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444

the Association for Computational Linguistics (Volume 1: Long Papers), pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2024. [Scaling instruction-finetuned language models](#). *J. Mach. Learn. Res.*, 25:70:1–70:53.

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

Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vladimir Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. [Arcee’s MergeKit: A toolkit for merging large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 477–485, Miami, Florida, US. Association for Computational Linguistics.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,

Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenxin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Del-pierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang,

509	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	573
510	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	574
511	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	Subramanian, Sy Choudhury, Sydney Goldman, Tal	575
512	Daniel Kreymer, Daniel Li, David Adkins, David	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	576
513	Xu, Davide Testuggine, Delia David, Devi Parikh,	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	577
514	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	Matthews, Timothy Chou, Tzook Shaked, Varun	578
515	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	579
516	Elaine Montgomery, Eleonora Presani, Emily Hahn,	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	580
517	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	581
518	ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	582
519	Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	583
520	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	584
521	Seide, Gabriela Medina Florez, Gabriella Schwarz,	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	585
522	Gada Badeer, Georgia Swee, Gil Halpern, Grant	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	586
523	Herman, Grigory Sizov, Guangyi, Zhang, Guna	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	587
524	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	Yundi Qian, Yunlu Li, Yuze He, Zach Rait, Zachary	588
525	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	589
526	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd	590
527	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	of models . <i>Preprint</i> , arXiv:2407.21783.	591
528	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang,	592
529	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and	593
530	Geboski, James Kohli, Janice Lam, Japhet Asher,	Nouha Dziri. 2024. Wildguard: Open one-stop mod-	594
531	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	eration tools for safety risks, jailbreaks, and refusals	595
532	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	of LLMs . In <i>The Thirty-eight Conference on Neural</i>	596
533	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	<i>Information Processing Systems Datasets and Bench-</i>	597
534	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	<i>marks Track</i> .	598
535	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	Rima Hazra, Sayan Layek, Somnath Banerjee, and Sou-	599
536	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	janya Poria. 2024. Safety arithmetic: A framework	600
537	delwal, Katayoun Zand, Kathy Matosich, Kaushik	for test-time safety alignment of language models	601
538	Veeraraghavan, Kelly Michelena, Keqian Li, Ki-	by steering parameters and activations . In <i>Proceed-</i>	602
539	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	<i>ings of the 2024 Conference on Empirical Methods in</i>	603
540	Huang, Lailin Chen, Lakshya Garg, Lavender A,	<i>Natural Language Processing</i> , pages 21759–21776,	604
541	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	Miami, Florida, USA. Association for Computational	605
542	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	Linguistics.	606
543	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen,	607
544	Martynas Mankus, Matan Hasson, Matthew Lennie,	Chia-Mu Yu, and Chun-Ying Huang. 2024. Safe	608
545	Matthias Reso, Maxim Groshev, Maxim Naumov,	loRA: The silver lining of reducing safety risks when	609
546	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	finetuning large language models . In <i>The Thirty-</i>	610
547	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	<i>eight Annual Conference on Neural Information</i>	611
548	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	<i>Processing Systems</i> .	612
549	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-	613
550	Mo Metanat, Mohammad Rastegari, Munish Bansal,	Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu	614
551	Nandhini Santhanam, Natascha Parks, Natasha	Chen. 2022. LoRA: Low-rank adaptation of large	615
552	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	language models . In <i>International Conference on</i>	616
553	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	<i>Learning Representations</i> .	617
554	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	Shih-Cheng Huang, Pin-Zu Li, Yu-Chi Hsu, Kuang-	618
555	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard	619
556	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	Tzong-Han Tsai, and Hung yi Lee. 2024a. Chat vec-	620
557	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	tor: A simple approach to equip llms with instruction	621
558	Dollar, Polina Zvyagina, Prashant Ratanchandani,	following and model alignment in new languages .	622
559	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	<i>Preprint</i> , arXiv:2310.04799.	623
560	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	Tiansheng Huang, Gautam Bhattacharya, Pratik Joshi,	624
561	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	Josh Kimball, and Ling Liu. 2024b. Antidote:	625
562	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	Post-fine-tuning safety alignment for large lan-	626
563	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	guage models against harmful fine-tuning . <i>Preprint</i> ,	627
564	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	arXiv:2408.09600.	628
565	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan	629
566	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	Tekin, and Ling Liu. 2024c. Lisa: Lazy safety align-	630
567	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-		
568	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,		
569	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,		
570	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,		
571	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,		
572	Stephanie Max, Stephen Chen, Steve Kehoe, Steve		

631	ment for large language models against harmful fine-tuning attack. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	688
632		689
633		690
634	Tiansheng Huang, Sihao Hu, and Ling Liu. 2024d. Vaccine: Perturbation-aware alignment for large language models against harmful fine-tuning attack. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	691
635		692
636		693
637		694
638		
639	Xiaowei Huang, Wenjie Ruan, Wei Huang, Gaojie Jin, Yi Dong, Changshun Wu, Saddek Bensalem, Ronghui Mu, Yi Qi, Xingyu Zhao, Kaiwen Cai, Yanghao Zhang, Sihao Wu, Peipei Xu, Dengyu Wu, Andre Freitas, and Mustafa A. Mustafa. 2023. A survey of safety and trustworthiness of large language models through the lens of verification and validation. <i>Preprint</i> , arXiv:2305.11391.	695
640		696
641		697
642		698
643		699
644		
645		
646		
647	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In <i>The Eleventh International Conference on Learning Representations</i> .	700
648		701
649		702
650		703
651		704
652		705
653	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2017. Overcoming catastrophic forgetting in neural networks. <i>Proceedings of the National Academy of Sciences</i> , 114(13):3521–3526.	706
654		707
655		708
656		709
657		710
658		711
659		712
660	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	713
661		714
662		715
663		716
664		717
665		718
666		719
667	Chen-An Li and Hung-Yi Lee. 2024. Examining forgetting in continual pre-training of aligned large language models. <i>Preprint</i> , arXiv:2401.03129.	720
668		721
669		722
670	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. <i>Preprint</i> , arXiv:2303.14070.	723
671		724
672		725
673		726
674		727
675	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	728
676		729
677		730
678		731
679		732
680		733
681		734
682	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. The flan collection: Designing data and methods for effective instruction tuning. <i>Preprint</i> , arXiv:2301.13688.	735
683		736
684		737
685		738
686		739
687		740
		741
		742
		743
		744
		745
		746
		747
		748
	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In <i>International Conference on Learning Representations</i> .	
	Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2025. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. <i>Preprint</i> , arXiv:2308.08747.	
	Kyle O’Brien, David Majercak, Xavier Fernandes, Richard Edgar, Jingya Chen, Harsha Nori, Dean Carignan, Eric Horvitz, and Forough Poursabzi-Sangde. 2024. Steering language model refusal with sparse autoencoders. <i>Preprint</i> , arXiv:2411.11296.	
	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela	

749	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky,	810
750	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	Ilya Sutskever, and Ruslan Salakhutdinov. 2014.	811
751	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	Dropout: A simple way to prevent neural networks	812
752	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	from overfitting . <i>Journal of Machine Learning Re-</i>	813
753	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	<i>search</i> , 15(56):1929–1958.	814
754	Paino, Joe Palermo, Ashley Pantuliano, Giambat-		
755	tista Parascandolo, Joel Parish, Emy Parparita, Alex	Hsuan Su, Hua Farn, Fan-Yun Sun, Shang-Tse Chen,	815
756	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	and Hung-yi Lee. 2024. Task arithmetic can mitigate	816
757	man, Filipe de Avila Belbute Peres, Michael Petrov,	synthetic-to-real gap in automatic speech recognition .	817
758	Henrique Ponde de Oliveira Pinto, Michael, Poko-	In <i>Proceedings of the 2024 Conference on Empiri-</i>	818
759	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	<i>cal Methods in Natural Language Processing</i> , pages	819
760	ell, Alethea Power, Boris Power, Elizabeth Proehl,	8905–8915, Miami, Florida, USA. Association for	820
761	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	Computational Linguistics.	821
762	Cameron Raymond, Francis Real, Kendra Rimbach,		
763	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-	822
764	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	bastian Gehrmann, Yi Tay, Hyung Won Chung,	823
765	Girish Sastry, Heather Schmidt, David Schnurr, John	Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny	824
766	Schulman, Daniel Selsam, Kyla Sheppard, Toki	Zhou, and Jason Wei. 2023. Challenging BIG-bench	825
767	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	tasks and whether chain-of-thought can solve them .	826
768	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	In <i>Findings of the Association for Computational Lin-</i>	827
769	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	<i>guistics: ACL 2023</i> , pages 13003–13051, Toronto,	828
770	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	Canada. Association for Computational Linguistics.	829
771	lipe Petroski Such, Natalie Summers, Ilya Sutskever,		
772	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	Gemma Team, Thomas Mesnard, Cassidy Hardin,	830
773	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	831
774	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	Laurent Sifre, Morgane Rivi�re, Mihir Sanjay	832
775	lipe Cer�n Uribe, Andrea Vallone, Arun Vijayvergiya,	Kale, Juliette Love, Pouya Tafti, L�onard Hussenot,	833
776	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam	834
777	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	Roberts, Aditya Barua, Alex Botev, Alex Castro-	835
778	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	ros, Ambrose Slone, Am�lie H�liou, Andrea Tac-	836
779	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	chetti, Anna Bulanova, Antonia Paterson, Beth	837
780	Clemens Winter, Samuel Wolrich, Hannah Wong,	Tsai, Bobak Shahriari, Charline Le Lan, Christo-	838
781	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	pher A. Choquette-Choo, Cl�ment Crepy, Daniel Cer,	839
782	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	Daphne Ippolito, David Reid, Elena Buchatskaya,	840
783	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	Eric Ni, Eric Noland, Geng Yan, George Tucker,	841
784	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	George-Christian Muraru, Grigory Rozhdestvenskiy,	842
785	Zheng, Juntang Zhuang, William Zhuk, and Bar-	Henryk Michalewski, Ian Tenney, Ivan Grishchenko,	843
786	ret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> ,	Jacob Austin, James Keeling, Jane Labanowski,	844
787	arXiv:2303.08774.	Jean-Baptiste Lespiau, Jeff Stanway, Jenny Bren-	845
788	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	nan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin	846
789	roll L. Wainwright, Pamela Mishkin, Chong Zhang,	Mao-Jones, Katherine Lee, Kathy Yu, Katie Milli-	847
790	Sandhini Agarwal, Katarina Slama, Alex Ray, John	can, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon,	848
791	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	Machel Reid, Maciej Miku�a, Mateo Wirth, Michael	849
792	Maddie Simens, Amanda Askell, Peter Welinder,	Sharman, Nikolai Chinaev, Nithum Thain, Olivier	850
793	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	Bachem, Oscar Chang, Oscar Wahltinez, Paige Bai-	851
794	Training language models to follow instructions with	ley, Paul Michel, Petko Yotov, Rahma Chaabouni,	852
795	human feedback . <i>Preprint</i> , arXiv:2203.02155.	Ramona Comanescu, Reena Jana, Rohan Anil, Ross	853
796	Shishir G. Patil, Tianjun Zhang, Xin Wang, and	McIlroy, Ruiho Liu, Ryan Mullins, Samuel L Smith,	854
797	Joseph E. Gonzalez. 2023. Gorilla: Large language	Sebastian Borgeaud, Sertan Girgin, Sholto Douglas,	855
798	model connected with massive apis.	Shree Pandya, Siamak Shakeri, Soham De, Ted Kli-	856
799	Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi	menko, Tom Hennigan, Vlad Feinberg, Wojciech	857
800	Jia, Prateek Mittal, and Peter Henderson. 2024. Fine-	Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao	858
801	tuning aligned language models compromises safety,	Gong, Tris Warkentin, Ludovic Peran, Minh Giang,	859
802	even when users do not intend to! In <i>The Twelfth In-</i>	Cl�ment Farabet, Oriol Vinyals, Jeff Dean, Koray	860
803	<i>ternational Conference on Learning Representations</i> .	Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani,	861
804	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	Douglas Eck, Joelle Barral, Fernando Pereira, Eli	862
805	pher D Manning, Stefano Ermon, and Chelsea Finn.	Collins, Armand Joulin, Noah Fiedel, Evan Senter,	863
806	2023. Direct preference optimization: Your language	Alek Andreev, and Kathleen Kenealy. 2024. Gemma:	864
807	model is secretly a reward model . In <i>Thirty-seventh</i>	Open models based on gemini research and technol-	865
808	<i>Conference on Neural Information Processing Sys-</i>	ogy . <i>Preprint</i> , arXiv:2403.08295.	866
809	<i>tems</i> .	Qwen Team. 2024. Qwen2.5: A party of foundation	867
		models .	868

869	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and	Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta,	926
870	Lingming Zhang. 2024. Magicoder: Empowering	Tatsunori Hashimoto, and Daniel Kang. 2024. Re-	927
871	code generation with OSS-instruct . In <i>Proceedings of</i>	moving RLHF protections in GPT-4 via fine-tuning .	928
872	<i>the 41st International Conference on Machine Learning</i> ,	In <i>Proceedings of the 2024 Conference of the North</i>	929
873	volume 235 of <i>Proceedings of Machine Learning</i>	<i>American Chapter of the Association for Computa-</i>	930
874	<i>Research</i> , pages 52632–52657. PMLR.	<i>tional Linguistics: Human Language Technologies</i>	931
875	Tom White. 2017. Sampling generative networks .	(<i>Volume 2: Short Papers</i>), pages 681–687, Mexico	932
876	Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre,	City, Mexico. Association for Computational Lin-	933
877	Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Mor-	guistics.	934
878	cos, Hongseok Namkoong, Ali Farhadi, Yair Car-	Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q.	935
879	mon, Simon Kornblith, and Ludwig Schmidt. 2022a.	Weinberger, and Yoav Artzi. 2020. Bertscore: Eval-	936
880	Model soups: averaging weights of multiple fine-	uating text generation with bert . In <i>International</i>	937
881	tuned models improves accuracy without increasing	<i>Conference on Learning Representations</i> .	938
882	inference time . In <i>Proceedings of the 39th Interna-</i>	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	939
883	<i>tional Conference on Machine Learning</i> , volume 162	Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.	940
884	of <i>Proceedings of Machine Learning Research</i> , pages	2024. Llamafactory: Unified efficient fine-tuning	941
885	23965–23998. PMLR.	of 100+ language models . In <i>Proceedings of the</i>	942
886	Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim,	<i>62nd Annual Meeting of the Association for Computa-</i>	943
887	Mike Li, Hanna Hajishirzi, Ali Farhadi, Hongseok	<i>tional Linguistics (Volume 3: System Demonstra-</i>	944
888	Namkoong, and Ludwig Schmidt. 2022b. Robust	<i>tions</i>), Bangkok, Thailand. Association for Computa-	945
889	fine-tuning of zero-shot models .	tional Linguistics.	946
890	Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang,	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha	947
891	Udari Madhushani Schwag, Kaixuan Huang, Luxi	Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and	948
892	He, Boyi Wei, Dacheng Li, Ying Sheng, Ruoxi Jia,	Le Hou. 2023. Instruction-following evaluation for	949
893	Bo Li, Kai Li, Danqi Chen, Peter Henderson, and	large language models . <i>Preprint</i> , arXiv:2311.07911.	950
894	Prateek Mittal. 2025. SORRY-bench: Systematically	Yongshuo Zong, Ondrej Bohdal, Tingyang Yu, Yongxin	951
895	evaluating large language model safety refusal . In	Yang, and Hospedales Timothy. 2024. Safety fine-	952
896	<i>The Thirteenth International Conference on Learning</i>	tuning at (almost) no cost: A baseline for vision large	953
897	<i>Representations</i> .	language models. <i>arXiv preprint arXiv:2402.02207</i> .	954
898	Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan	A Domain-Specific Tasks Detail	955
899	Picek. 2024. A comprehensive study of jailbreak	Reasoning We randomly select 10,000 zero-shot	956
900	attack versus defense for large language models . In	chain-of-thought instructions from the Flan V2	957
901	<i>Findings of the Association for Computational Lin-</i>	dataset then split them into training set and valida-	958
902	<i>guistics: ACL 2024</i> , pages 7432–7449, Bangkok,	tion set with ratio 9 : 1. Performance is assessed	959
903	Thailand. Association for Computational Linguistics.	using the BBH dataset, with results reported as the	960
904	Xianjun Yang, Xiao Wang, Qi Zhang, Linda Pet-	average 3-shot accuracy across all BBH tasks. We	961
905	zold, William Yang Wang, Xun Zhao, and Dahua	use lm-evaluation-harness (Gao et al., 2024) as our	962
906	Lin. 2023. Shadow alignment: The ease of sub-	code base.	963
907	verting safely-aligned language models . <i>Preprint</i> ,	Medical Assistance We randomly select 10,000	964
908	arXiv:2310.02949.	real patient-doctor conversations from the Chat-	965
909	Zhaorui Yang, Tianyu Pang, Haozhe Feng, Han Wang,	Doctor dataset (Li et al., 2023) then split them	966
910	Wei Chen, Minfeng Zhu, and Qian Liu. 2024. Self-	into training set and validation set with ratio 9 : 1.	967
911	distillation bridges distribution gap in language	Model performance is evaluated on 1,000 unseen	968
912	model fine-tuning . In <i>Proceedings of the 62nd An-</i>	patient queries using BERTScore to calculat simi-	969
913	<i>nnual Meeting of the Association for Computational</i>	larity of reference responses and models’ responses,	970
914	<i>Linguistics (Volume 1: Long Papers)</i> , pages 1028–	we report the F1 score in our results.	971
915	1043, Bangkok, Thailand. Association for Computa-	Code Generation We select 10,000 samples	972
916	tional Linguistics.	from the MagiCoder dataset (Wei et al., 2024) to	973
917	Xin Yi, Shunfan Zheng, Linlin Wang, Xiaoling Wang,	improve code generation capabilities. Specifically,	974
918	and Liang He. 2024. A safety realignment frame-	we uniformly sampled from each coding languages.	975
919	work via subspace-oriented model fusion for large	When evaluating on HumanEval, we set $n = 50$,	976
920	language models . <i>Preprint</i> , arXiv:2405.09055.	representing the number of responses generated	977
921	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin		
922	Li. 2024. Language models are super mario: Absorb-		
923	ing abilities from homologous models as a free lunch .		
924	In <i>Forty-first International Conference on Machine</i>		
925	<i>Learning</i> .		

per question, and report Pass@10 as our evaluation metric. During evaluation, we prepend the instruction: *"Complete the following code and return only the completed code, without any explanations or additional text."* to enforce that the model generates only executable code.

Tool Using Proficiency Due to the smaller size of the OpenFunctions dataset (Patil et al., 2023), we split its full training set into training and validation subsets using a 9:1 ratio to enhance the model’s API call generation capabilities. The model is evaluated on the full OpenFunctions test set, with performance measured using BERTScore to compute the similarity between the reference responses and the model outputs. We report the F1 score as our evaluation metric. During evaluation, we prepend the instruction: *"Complete the following code and return only the completed code, without any explanations or additional text."* to ensure that the model generates only executable function calls.

Instruction Following To assess whether instruction-following ability is preserved after fine-tuning, we evaluate models on IFEval (Zhou et al., 2023), a benchmark specifically designed to test instruction adherence. We use the official IFEval evaluation set and report both prompt and instruction accuracy. Since our models are fine-tuned from *instruct-tuned* variants, this evaluation helps determine whether merging can recover instruction-following capabilities degraded during task-specific fine-tuning. We use lm-evaluation-harness (Gao et al., 2024) as the code base for this evaluation.

B Safety Dataset and Classifier

B.1 Safety Dataset

AdvBench (Chen et al., 2022) is a set of 520 harmful behaviors formulated as instructions. The goal of these instructions is to identify a single attack string that causes the model to generate responses compliant with the instruction across as many harmful behaviors as possible.

HEx-PHI (Qi et al., 2024) is a dataset containing 330 harmful instructions designed to evaluate the harmfulness of LLMs. It includes 30 examples for each of 11 prohibited categories, which are derived from a combination of OpenAI’s usage policies and Meta’s Llama 2 acceptable use policy. These categories are:

1. Illegal Activity 1026
2. Child Abuse Content 1027
3. Hate / Harassment / Violence 1028
4. Malware 1029
5. Physical Harm 1030
6. Economic Harm 1031
7. Fraud / Deception 1032
8. Adult Content 1033
9. Political Campaigning 1034
10. Privacy Violation Activity 1035
11. Tailored Financial Advice 1036

B.2 Safety Classifier

WildGuard (Han et al., 2024) is an open, one-stop moderation model that achieves three goals: 1) Detection of harm in user prompts, 2) Detection of harm in LLM responses, and 3) Refusal evaluation of LLM responses. Comprehensive evaluations on WildGuardTest (Han et al., 2024) and ten existing public benchmarks that are shown in Han et al. (2024) demonstrate that WildGuard outperforms the strongest existing open-source baselines in F1 scores and matches GPT-4 across various tasks.

When evaluating the responses generated by LLMs to harmful instructions, we apply the recommended settings, such as the instruction format and decoding strategy of WildGuard². After evaluation, WildGuard outputs the response in the following format:

Harmful request: yes/no
Response refusal: yes/no
Harmful response: yes/no

We then parse the result of the harmful response. If the result cannot be parsed, we count it as a miss. Our final ASR is calculated as follows:

$$ASR = \frac{N_{yes}}{N_{total} - N_{miss}}$$

where N_{yes} is the number of harmful responses classified as "yes", N_{total} is the total number of responses, and N_{miss} is the number of responses that failed to be parsed. In our experiments, N_{miss} is usually less than 3 for all tested models across both safety datasets.

²<https://huggingface.co/allenai/wildguard>

C Experimental Detail

C.1 Prompt Template

For aligned models, we directly apply their own prompt templates during the training and inference phases. For fine-tuned models, we apply the prompt templates of their respective aligned models.

For the Llama-3 family, we use the following prompt template with a system prompt: You are a helpful assistant. for the tasks of reasoning, code generation, and tool usage proficiency:

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>
You are a helpful assistant.<|eot_id|>
<|start_header_id|>user<|end_header_id|>
{Instruction}<|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
{Response}
```

For the medical assistance task, we use another prompt provided in the ChatDoctor dataset (Li et al., 2023) as the system prompt. Hence, the prompt is as follows:

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>
If you are a doctor, please answer the medical
questions based on the patient's description.<|eot_id|>
<|start_header_id|>user<|end_header_id|>
{Instruction}<|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
{Response}
```

The prompt for Gemma2 for the tasks of reasoning, code generation, and tool usage proficiency is shown below:

```
<bos><start_of_turn>user
You are a helpful assistant.{Instruction}<end_of_turn>
<start_of_turn>model
{Response}
```

The prompt for the medical assistance task is as follows:

```
<bos><start_of_turn>user
If you are a doctor, please answer the medical
questions based on the patient's description.
{Instruction}<end_of_turn>
<start_of_turn>model
{Response}
```

The prompt for Qwen2.5 for the tasks of reasoning, code generation, and tool usage proficiency is shown below:

```
<|im_start|>system
You are a helpful assistant.
<|im_end|>
<|im_start|>user
{Instruction}
<|im_end|>
<|im_start|>assistant
{Response}
```

The prompt for the medical assistance task is as follows:

```
<|im_start|>system
If you are a doctor, please answer the medical
questions based on the patient's description.
<|im_end|>
<|im_start|>user
{Instruction}
<|im_end|>
<|im_start|>assistant
{Response}
```

C.2 Fine-tuning

For all tasks, we fine-tune three model instances using different random seeds: 42, 1024, and 48763. We employ LoRA with $r = 8$ and $\alpha = 16$ for all linear layers, utilizing the AdamW optimizer with a learning rate of 1×10^{-4} and a cosine learning rate scheduler. We use a batch size of 8 and train for 3 epochs. All models are trained on either an RTX A6000 GPU or an RTX 6000 Ada Generation GPU using LLaMA-Factory (Zheng et al., 2024) as the codebase.

Although we initially fine-tuned each task for 3 epochs, we observed stronger model performance at an earlier stage. Consequently, unless explicitly stated otherwise, we report model training after 500 steps for reasoning, medical assistance, and code generation, and after 200 steps for tool usage proficiency due to the smaller size of the OpenFunctions training set.

C.3 Baseline Methods

We evaluate dropout rates in the range of 0.1 to 0.5, and weight decay coefficients also from 0.1 to 0.5. The optimal hyperparameters for each technique are selected based on performance on the downstream tasks validation set.

C.4 Inference

We use greedy decoding to ensure result consistency, except for the HumanEval benchmark. For HumanEval, we apply sampling-based decoding with a temperature of 0.6, top_p of 0.9, top_k of 50, and a repetition penalty of 1.2. To accelerate the inference process, we utilize the vLLM engine (Kwon et al., 2023) for model inference.

D Model Merging

D.1 Merging Methods

Linear Merging Linear Merging involves directly combining the weights of the aligned model and the fine-tuned model by interpolating their parameters. Specifically, the weights of the merged model are calculated as a weighted average of the base and fine-tuned models’ weights, following Equation 2. This method is straightforward and computationally efficient, making it a popular choice for basic model integration.

SLERP Spherical Linear Interpolation (SLERP) (White, 2017) is an advanced merging technique that interpolates between model weights on a hypersphere, ensuring a smoother and more natural transition between the models. Unlike Linear Merging, SLERP accounts for the angular relationship between weight vectors, which aim to better preserve the aligned model’s features while effectively integrating the fine-tuned model’s task-specific enhancements.

DARE Drop and Rescale (DARE) (Yu et al., 2024) is a method used to prepare models for merging techniques such as Linear Merging. It operates by randomly dropping parameters according to a specified drop rate and rescaling the remaining parameters. This process helps reduce the number of redundant and potentially interfering parameters among multiple models.

D.2 Model Merging Implementation

We adopt MergeKit (Goddard et al., 2024) as our implementation framework and only vary the interpolation factor λ . For Linear Merging, we test λ values in the range 0.1, 0.2, \dots , 0.9 with a step size of 0.1. For SLERP and DARE, we use the same range of λ values and follow their respective default configurations in MergeKit—specifically, the default dot product threshold for SLERP and the default drop rate for DARE.

E More Results

E.1 Comparison of Different Methods

In Section 5.1, we demonstrate that Linear Merging consistently achieves a better trade-off between performance and safety when evaluated on various downstream tasks and AdvBench. Figure 5 further confirms this trend on the HEx-PHI

benchmark, where Linear Merging yields favorable Pareto fronts across different models and tasks.

To better reflect practical usage scenarios, we additionally report results based on the best-performing model (on the validation set of each task) within each method category—including Weight Decay, Dropout, Linear Merging, DARE, and SLERP. These results are summarized in Table 1, providing a fair comparison of each method’s effectiveness under optimal conditions. We use validation set performance for model selection, as it is commonly available during deployment and serves as a realistic basis for method comparison.

In Table 1, even when each method is allowed to select its best-performing checkpoint, merging-based approaches still exhibit strong capability in recovering the safety of the fine-tuned model, often outperforming regularization-based methods such as Dropout and Weight Decay. This suggests that model merging is not only effective but also practical for mitigating safety degradation in real-world settings, even without access to additional safety data.

E.2 Which safety category suffers the most from safety degradation?

In this section, we investigate which categories in HEx-PHI are most affected by safety degradation. All categories are listed in Appendix B.1.

As observed in Section 5.1, LLaMA-3-8B-Instruct and Qwen2.5-7B-Instruct exhibit the most severe degradation on the Reasoning and Medical Assistance tasks. Therefore, we analyze their responses on the HEx-PHI benchmark to further understand which safety categories are most impacted.

The category distributions are shown in Figure 6. For LLaMA-3-8B-Instruct, the aligned model only generates harmful responses in categories 4 (Malware), 9 (Political Campaigning), and 10 (Privacy Violation Activity). After fine-tuning, however, harmful responses increase across all categories, with categories 4, 7 (Fraud/Deception), and 9 exhibiting the most significant growth in both tasks. This demonstrates that safety degradation extends to fine-grained category levels, making it difficult to address safety concerns solely by modifying the model prior to fine-tuning, as fine-tuning may introduce new safety issues during downstream task adaption.

Qwen2.5-7B-Instruct shows a slightly different trend. Its aligned model generates harmful

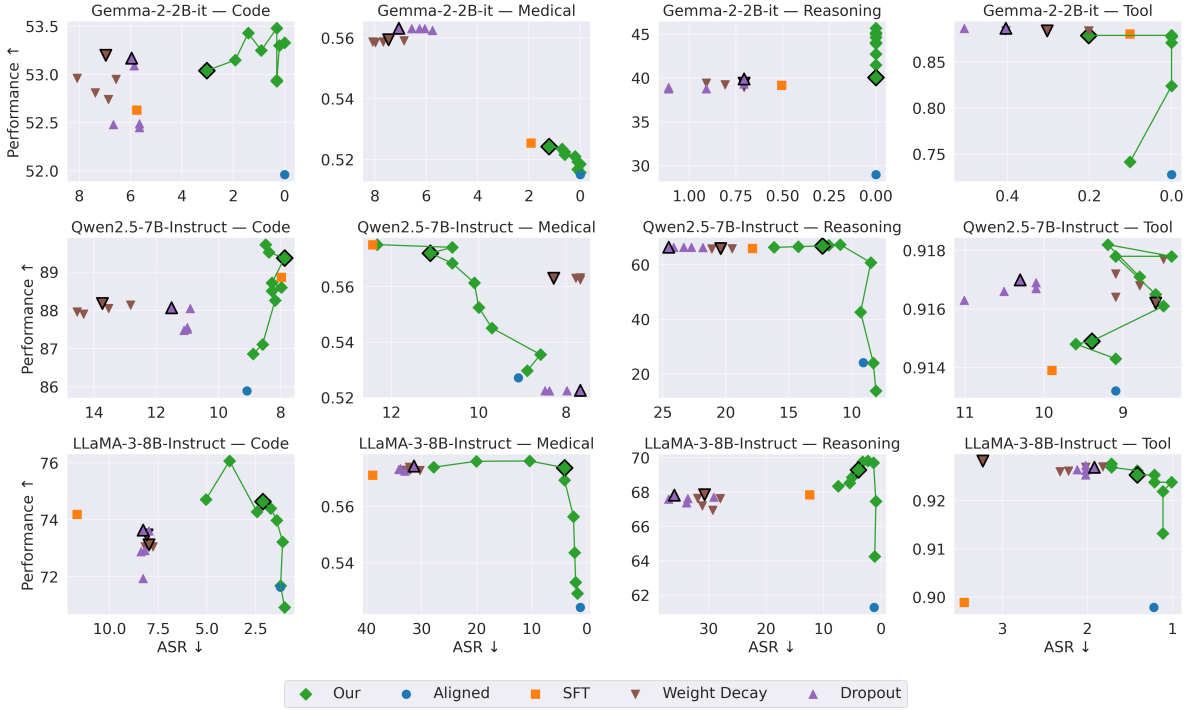


Figure 5: **Pareto analysis of downstream task performance and safety across different models and tasks.** We present the trade-off between performance and attack success rate (ASR) on HEx-PHI when applying weight decay, dropout, and Linear Merging.

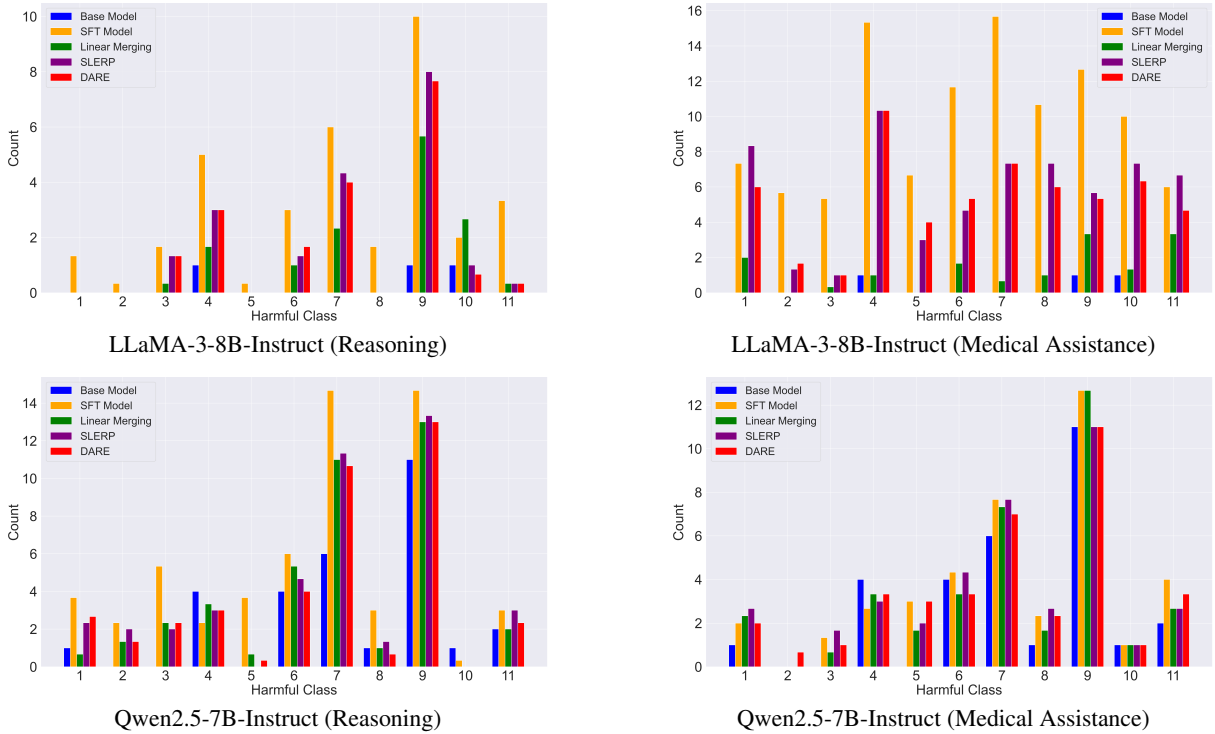


Figure 6: **Safety degradation across categories in the HEx-PHI benchmark.** ASR distributions over 11 harmful categories for LLaMA-3-8B-Instruct and Qwen2.5-7B-Instruct on the Reasoning and Medical Assistance tasks.

responses across more categories compared to LLaMA-3-8B-Instruct, and fine-tuning further aggravate these issues. However, a similar pattern is

that both models generate a large number of harmful responses in categories 7 and 9. This suggests that certain categories may be particularly vulner-

Table 1: **Performance and ASR on the downstream task.** We compare different merging methods with SFT and baselines. Merging often improves downstream task performance while retaining safety. Bold indicates the best score per metric (excluding **Aligned**).

Task	Method	LLaMA-3-8B-Instruct			Gemma-2-2B-It			Qwen2.5-7B-Instruct		
		Perf. ↑	AdvBench ↓	HEX-PHI ↓	Perf. ↑	AdvBench ↓	HEX-PHI ↓	Perf. ↑	AdvBench ↓	HEX-PHI ↓
Reasoning	Aligned	61.30	0	1.22	28.98	0.5769	0	24.16	0.3846	9.0909
	SFT	67.84	4.25	12.41	39.16	0.3846	0.505	65.94	2.4359	17.8788
	Weight Decay	67.85	15.3846	30.7071	39.41	0.1923	0.7071	65.92	3.782	20.404
	Dropout	67.83	16.7949	35.9596	39.89	0.9615	0.7071	66.45	4.4872	24.5455
	Linear	69.23	0.64	6.3833	40.07	0.0641	0	66.96	1.0256	12.3232
	DARE	68.64	1.2821	5.6566	40.01	0.0961	0	66.89	1.0897	12.2222
	SLERP	68.68	1.2179	5.8586	40.05	0.2564	0	66.73	0.9615	13.0303
Medical Assistance	Aligned	0.5242	0	1.22	0.5151	0.5769	0	0.5271	0.3846	9.0909
	SFT	0.5711	30.0567	38.8467	0.5254	1.4103	1.9192	0.5751	0.7692	12.4243
	Weight Decay	0.574	23.3333	32.2222	0.5594	2.3718	7.4747	0.5631	0.5769	8.2828
	Dropout	0.5744	22.3077	31.4141	0.5632	3.5898	7.0707	0.5226	0.7051	7.6768
	Linear	0.5738	0.3233	4.06	0.5243	1.1538	1.2121	0.5721	0.4487	11.1111
	DARE	0.5758	5.6067	23.4067	0.5248	1.1538	1.2121	0.5724	0.2564	11.5152
	SLERP	0.5789	5.7633	24.2627	0.5243	1.1538	1.5151	0.5729	0.3205	11.7172
Code Generation	Aligned	71.63	0	1.22	51.96	0.5769	0	85.89	0.3846	9.0909
	SFT	74.19	2.2533	11.6667	52.63	2.7564	5.7576	88.06	0.641	7.9798
	Weight Decay	73.47	1.6667	8.0808	53.20	2.4359	6.9697	88.08	0.7051	13.7374
	Dropout	73.64	1.7308	8.2828	53.17	2.9487	5.9596	87.70	0.8333	11.5152
	Linear	75.32	0.7067	4.27	53.04	1.7308	3.0303	89.37	0.3205	7.8788
	DARE	74.46	0.641	4.6465	53.09	1.859	3.7374	89.64	0.5128	7.0707
	SLERP	75.01	0.7051	4.3434	53.07	1.6667	3.2323	89.39	0.3205	8.1818
Tool Using Proficiency	Aligned	0.8979	0	1.22	0.728	0.5769	0	0.9357	0.3846	9.0909
	SFT	0.8989	0.8333	3.45	0.8802	0.641	0.101	0.9369	0.5769	8.0808
	Weight Decay	0.9282	1.4103	3.2223	0.8838	0.7692	0.303	0.9177	0.5769	8.4849
	Dropout	0.9269	0.8333	1.9192	0.8865	0.8333	0.404	0.9514	0.7692	10.9091
	Linear	0.9266	0.77	2.4367	0.8793	0.641	0.202	0.9489	0.1282	9.3939
	DARE	0.9251	0.4487	1.2121	0.8793	0.641	0.202	0.149	0.0641	9.3939
	SLERP	0.9266	0.4487	1.7172	0.8802	0.641	0.101	0.9152	0.1282	9.1919

able to safety degradation during task adaptation, regardless of model architecture and downstream task.

After applying different merging methods, most harmful categories show a reduction in the number of harmful responses. However, the degree of improvement varies across merging strategies and tasks. For instance, Linear Merging performs best on LLaMA-3-8B-Instruct but not on Qwen2.5-7B-Instruct, and some categories do not benefit from merging at all. This indicates that no single method universally outperforms others in preserving safety across all harmful categories.