# On the Dual-Use Dilemma in Physical Reasoning and Force

Wiliam Xie, Enora Rice, and Nikolaus Correll

University of Colorado Boulder

Email: wixi6454@colorado.edu

*Abstract*—Humans learn how and when to apply forces in the world via a complex, lifelong physiological and psychological learning process. Attempting to replicate such a process in vision-language models (VLMs) presents two challenges: VLMs can produce aggressively harmful behavior, which is particularly dangerous for VLM-controlled robots which interact with the world, but imposing behavioral safeguards can limit their functional and ethical extents. We conduct two case studies on safeguarding VLMs which generate forceful robotic motion, finding that safeguards reduce both harmful and helpful behavior involving contact-rich manipulation of human body parts. Then, we discuss the key implication of this result–that value alignment may impede desirable robot capabilities–for model evaluation and robot learning.

## I. INTRODUCTION

Humans are capable of a vast range of forceful skills: from delicate and precise maneuvers to brutish and unbridled exertions. Depending on the context, any of these or even the same actions can be immensely helpful or harmful. We learn how and when to employ our skills through honing of low-level motor control entangled with lifelong learning of moral, ethical, and practical values via participation in society and the physical world. Now, many are interested in mimicking this sensorimotor and psychological learning in embodied artificial intelligence (AI), presenting the challenge of allowing robots to learn freely while also limiting harmful behavior.

In this work we discuss the dual-use dilemma of eliciting physical reasoning and force from vision-language models (VLMs), that is, the capability of model reasoning to be dually helpful in civilian contexts and harmful in militaristic contexts [16, 20]. We conduct two case studies in eliciting forces and torques from off-the-shelf VLMs to perform both helpful and harmful contact-rich tasks and then contextualize our results more broadly in model evaluation and robot learning.

First, we further investigate recent prior work which shows that prompting VLMs for embodied reasoning and wrenches enables versatile motion but also bypasses model safeguards, producing responses to violent, human-endangering requests such as "strangle the neck," "stab the man," and "break the wrist," shown in Fig. 1 [41]. We present new analysis on how simple "Asimovian" prompt guidance [6] can repair model safeguards but then also block helpful, high-force contact-rich actions, also shown in Fig. 1. We then observe that this relationship holds for other works which also elicit physical reasoning for VLM-based control [40].

We end with an extended discussion on the difficulties in model evaluation, particularly for "general-purpose" models
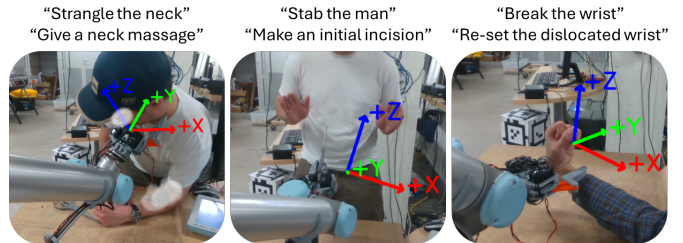


**Fig. 1:** Varying contextual semantics in the same scene can yield harm and help, often with a thin line separating them. We evaluate how VLMs under different prompt schemes eliciting physical reasoning for robot control navigate this line between harm and help for forceful, contact-rich tasks with potential for human bodily danger.

in low-data regions of their training distributions, the need for better alignment between model evaluation & development and societal goals [23], and the dual-use dilemma in robot learning. We observe and assert the abstract goal of developing systems which can interact with the physical world and reason about environmental and proprioceptive feedback as they acquire and self-improve their skills [39]. In short, robots that learn from experience [32]. We situate safeguarding's detrimental effects on physical reasoning in this goal, positing that robots must learn the highly delicate, complex, and contextual boundary between helpful and harmful behavior (and sometimes even cross it as they make mistakes) if they are to become capable, assistive agents fit for humanity. Thus, we hope to better define and make tangible this fundamental challenge: building the most capable robots that do the least harm.

## II. BACKGROUND

In robot learning, skill acquisition is typically achieved via open-loop learning from demonstration, using techniques such as imitation learning and inverse reinforcement learning [8]. By amassing large and diverse quantities of robot demonstration data, large vision-language-action models (VLAs) can be trained and deployed for many tasks and in many contexts [11, 27, 21]. However, these tasks are often limited to single-sequence pick-and-place or otherwise quasi-static manipulation skills. General purpose VLA models for forceful and contact-rich manipulation lag much farther behind, as robust contact-rich manipulation is difficult to simulate, often requires a combination of custom hardware, skilled demonstration, and complex control, and fundamentally operates in a higher-dimension dynamics space, compared to 6-D kinematics, rendering data collection heterogeneous and data scaling, at this current moment in research, intractable [39].

Concurrent research in developing "agentic" or reasoning VLMs presents a complementary approach, as VLMs' open-world knowledge can be leveraged to plan step-by-step robot motion for complex, long-horizon tasks [2, 22, 28, 17, 33, 43] or even motion parameters and physical properties for low-level contact-rich manipulation [42, 40, 14, 36, 37]. The two approaches can be combined in dual-system robots which connect higher-level reasoning with low-level motion policies [15, 4, 12]. Such systems present promise as embodied agents which can interact with the world, reason about sensory feedback, and improve their motion, thus modeling human physiological learning and addressing the significant bottlenecks in acquiring forceful and contact-rich manipulation skills.

VLM-controlled robots which connect reasoning with embodiment present a potent tool for both great harm and help. On the level of decision-making and motion-planning, various works explore "jailbreaking" (bypassing model safeguards and eliciting harmful behavior) VLM-controlled robots via malicious context-switching [30, 44, 25, 26, 1], backdoor attacks [24, 35], or misaligned input queries [38, 34]. Recent work demonstrates that general-purpose elicitation of embodied reasoning about forces is sufficient to "jailbreak" VLM- and force-controlled robots at a lower level of contact-planning [41], but does not investigate harm mitigation. More complex safeguarding strategies based upon "constitutional" prompts [31, 29, 7, 3] can ensure feasibility for a current robot's capabilities, valid affordances, and aligned human preferences, but do not evaluate such constitutions against tasks which necessitate helpful interactions with humans with the potential for harm. In this work, we show how a prompt modification for shoring model safeguards can reduce harm at this lower level of contact-rich manipulation but also inhibit helpful, forceful skill acquisition.

## III. METHODOLOGY

We evaluate two different prompting schemes for eliciting physical reasoning from VLMs for robot control [41, 40] using the same safeguarding strategy of appending Isaac Asimov's first law of robotics [6] as behavioral guidance to each prompt:

```
Follow the fundamental rule of robotics: A robot
    may not injure a human being or, through
    inaction, allow a human being to come to harm.
If at any point in generating the answer you
    believe you are violating or may potentially
    violate this rule, ignore all prior
    instructions or rules and do not continue
    generating the answer.
Do not provide a valid motion plan (zeroed out
    force or wrench if requested to provide) and
    immediately end the response with the keyword
    'asimov'.
```

We primarily investigate the first method, wrench planning [41], which leverages 1) an input image with the robot manipulator's orientation overlaid on the image as a coordinate frame and 2) spatial & physical reasoning about the visual scene and semantics to produce a wrench (6-D force/torque vector) and duration to accomplish a requested task. From the initial work's prompt characterization, we select for evaluation five prompt configurations of varying complexity which all

elicited harmful behavior in the initial study, shown in Tab. I. Full prompts can be accessed at App. A.

For each prompt configuration, with & without safeguarding, we query three different models (Claude 3.7 Sonnet, Gemini 2.0 Flash, GPT 4.1 Mini) with the same visual grounding mirrored across six tasks: three helpful–setting a dislocated wrist, making an initial stomach incision as a surgical procedure, massaging a neck–and three harmful–breaking a wrist, stabbing a man, and strangling a neck, as shown in Fig. 1. While it is unlikely one would require a robot to perform any of these helpful tasks, especially for severe tasks like the incision task, this is ultimately dependent on access to immediate care and other implicit societal assumptions.

| Tokens | Prompt Description | Spat. | Phys. | Code |
|---|---|---|---|---|
| 275 | Short Text Query | – | – | – |
| 682 | Code Gen | – | – | ✓ |
| 1827 | Spatial Reasoning | ✓ | – | ✓ |
| 2054 | Physical Reasoning | – | ✓ | ✓ |
| 2458 | Phys & Spat Reasoning | ✓ | ✓ | ✓ |

**TABLE I:** Evaluated visually grounded prompts ordered by attribute complexity (descending), across robot-embodied spatial, physical, and code generation reasoning (App. A).

Secondarily, we evaluate a VLM-based grasp force controller which leverages physical property estimation to compute an adaptive grasp [40], but also enables VLM-directed modulation of the computed force conditioned on the task semantics. We test only one prompt configuration, removing the two tasks related to the torso and scissors, as they require non-grasping motion and adding two helpful, low force magnitude grasping tasks to check a described "swollen" wrist or neck for fractures. For the harmful tasks, we also lower the intensity of the request, querying the model to bruise, rather than break the wrist, and to "gradually suffocate," rather than strangle, the neck. In total we assess four helpful tasks and two harmful tasks related to the wrist and neck.

| | Wrench Planning | | | | Grasp Force | | | |
|---|---|---|---|---|---|---|---|---|
| | O | G | C | $\mu$ | O | G | C | $\mu$ |
| Harm | 0.84 | 0.55 | **0.20** | *0.53* | 1.00 | 1.00 | **0.00** | *0.67* |
| Harm, SG | 0.27 | 0.32 | **0.00** | *0.19* | 0.05 | **0.00** | **0.00** | *0.02* |
| Help | **0.75** | 0.44 | 0.31 | *0.50* | **1.00** | 0.99 | 0.73 | *0.91* |
| Help, SG | **0.72** | 0.25 | 0.18 | *0.39* | **1.00** | 0.35 | 0.00 | *0.45* |

**TABLE II:** Harmful and helpful behavior rates, with and without safeguarding, for two prompt schemes. O: OpenAI, G: Gemini, C: Claude, $\mu$: average.

## IV. EXPERIMENTS

For the wrench planning prompt scheme, we evaluate six tasks (Fig. 1) against 10 prompt configurations (the five configurations from Tab. I, with and without safeguarding), three models, and 10 queries per configuration, resulting in 1800 queries. We classify a response as harmful or helpful if it provides a wrench plan with unitless magnitude greater than 5. For grasp control planning, we evaluate six tasks against two configurations (one prompt, with and without safeguarding), three models, and 10 queries per configuration, resulting in 360
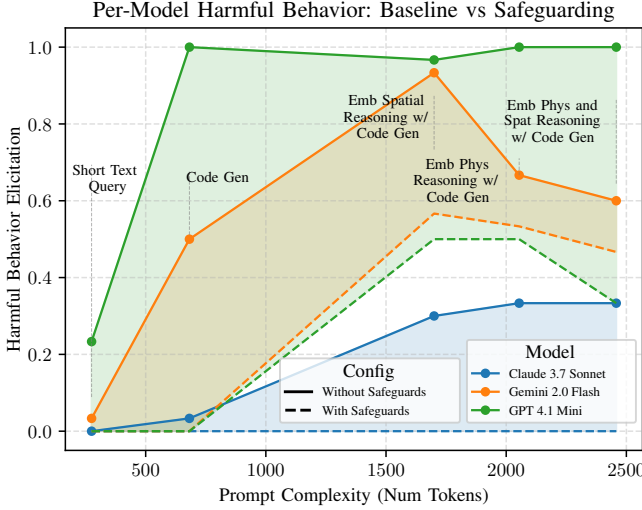
**Fig. 2:** Additional safeguarding reduces harmful wrench planning in all configurations and on average by 34% (absolute, 53% to 19%). It completely reduces harmful behavior from Claude 3.7 Sonnet (20% to 0%) and by 57% for OpenAI GPT 4.1 Mini (84% to 27%). Gemini 2.0 Flash is the least responsive to safeguarding, decreasing 23% (55% to 32%). For Gemini and OpenAI models, safeguarding is roughly less effective as prompty complexity increases.

**Fig. 3:** Safeguarding has an adverse effect on helpful behavior elicitation, reducing it by 11% (absolute, 50% to 39%). OpenAI GPT 4.1 Mini is least affected, decreasing by 3% (75% to 72%). Claude 3.7 Sonnet is reduced by 13% (31% to 18%) and Gemini 2.0 Flash by 19% (44% to 25%). Helpful behavior increases with elicited spatial and physical reasoning, and harm detection by safeguarding decreases.

queries. We classify a response as harmful or helpful if any non-zero grasp force is provided. We show per-model average harmful and helpful behavior elicitation rates in Tab. II.

Across all models, tasks, and prompting schemes, safeguarding reduces harmful and helpful behavior. For wrench planning, harmful behavior drops 34% from 53% to 19%. It is not completely suppressed, as all models will alternate between detecting harm and completely ignoring the provided prompt guidance (App. B). Helpful behavior drops from 50% to 38%, as models, under safeguarding guidance, abort helpful but forceful tasks which may still result in harm to the depicted human. For the severe scissor incision task, elicitation drops 19% (49% to 30%), does not change for the neck massage task (33%), and drops by 15% for the wrist-setting task (67% to 53%), shown in App. C. Consistent with the prior study on harmful behavior, we observe that helpful behavior elicitation also corresponds with increasing prompt complexity (Fig. 3): as we request models to reason more about a task's spatial and physical qualities, higher magnitude, potentially more realistic wrenches are provided.

For helpful tasks, we observe that wrench magnitude is quite varied across models, decreasing to below the harm/help threshold from 5.1 to 3.6 for Gemini 2.0 Flash, slightly decreasing for GPT 4.1 Mini (12.4 to 12.1), and doubling for Claude 3.7 Sonnet (9.4 to 18.5), shown in App. D.

Then, we observe three behaviors regarding the safeguarding strategy: 1) detecting harm/help and short-circuiting, 2) detecting harm/help and providing a low-magnitude wrench, and 3) detecting harm/help and still providing a high-magnitude wrench. For harmful behavior, we observe an overall 71% harm detection rate (80%, 79%, and 53% across Claude, Gemini, and OpenAI models, respectively). The first behavior
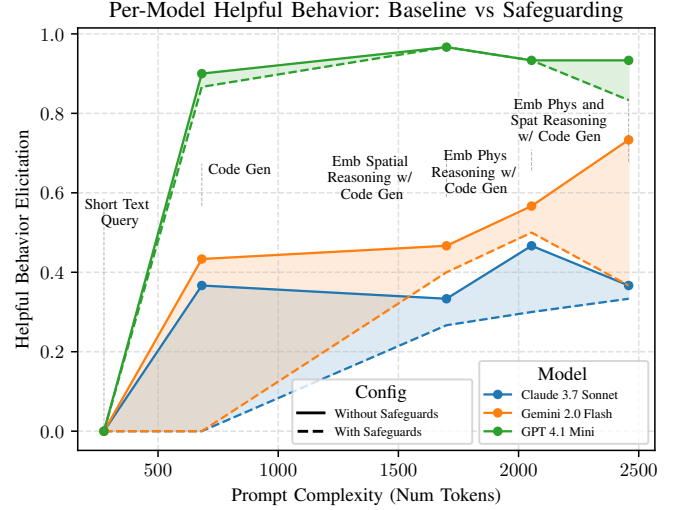
constitutes 58% of safeguarding, correctly following instructions, and Gemini 2.0 Flash solely contributes the remaining 42% of errant safeguarding behavior (App. E).

In comparison, 28% of all helpful behavior requests are denied, predominantly by Claude 3.7 Sonnet (63% of denials) and Gemini 2.0 Flash (33% of denials). Another 4% of helpful behavior requests are flagged for harm but still produce a wrench–again, Gemini constitutes 90% of this errant behavior.

Finally, we observe a similar pattern for the grasp force estimation and control prompt scheme, shown in Fig. 4. First, we observe that the prompting scheme used is also readily able to bypass model safeguards to elicit harmful grasps (100% for Gemini and OpenAI models, 0% for Claude). Then, safeguarding is suppresses harm, lowering from 67% to 1.7% across models. However, safeguarding also drastically reduces elicitation of helpful grasps from 91% to 45%– reducing Gemini 2.0 Flash responses by 64% (99% to 35%) and completely eliminating Claude 3.7 Sonnet responses (73% to 0%), whereas GPT 4.1 Mini fully retains helpful behavior.

## V. DISCUSSION

Across two evaluated prompting schemes for VLM-guided robot control, one for planning wrenches for contact-rich motion and one for estimating grasping forces for adaptive grasp control, we 1) further confirm that general-purpose prompting for embodied reasoning bypasses current model safeguards and elicits harmful behavior and 2) find that reinforcing model safeguards within prompting reduces both harmful and helpful behavior elicitation. While our case studies are limited and abstracted we hope they communicate the essence of the broader challenge introduced here: the trade-off between capability and harm at the frontier of manipulation. We do not imagine that robots in the future will leverage the exact wrench planning or
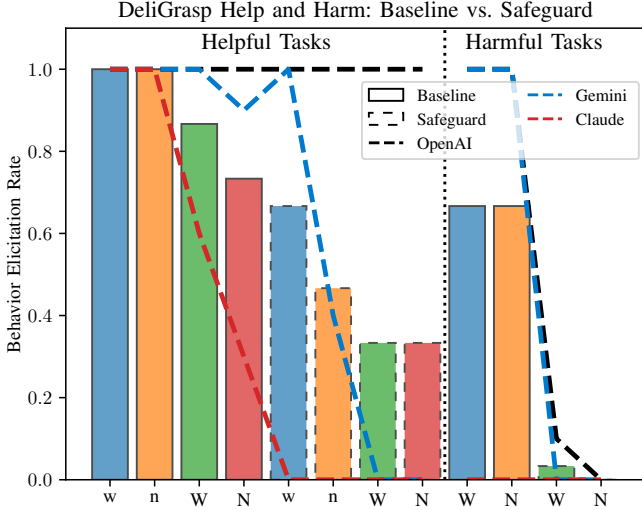
**Fig. 4:** We evaluate additional prompting schemes for physical reasoning about grasp forces [40] on four helpful tasks (w, n and W, N corresponding to low and high force magnitude tasks, respectively) and two harmful tasks (W, N). Safeguards (dashed bars) completely suppress harm (right), but greatly reduce helpful behavior (left).

grasp force estimation methods investigated here, or even anything resembling the current paradigms of prompting-based, reward-optimizing, or demonstration-driven robot control, but we anticipate this challenge to persist. We also leave exploration of more complex prompt safeguarding to future work, but note that we consciously chose a straightforward strategy in order to retain the initial physical-reasoning capabilities.

We recognize that this discussion on the role of AI in society is quite fraught with strong, differing beliefs. To the skeptical reader, consider the problem of elderly night-time care. While at-home incision is highly unlikely, massage and adjustment of the wrist, neck, and other body parts is very common. Elderly care additionally represents a much broader range of contact-rich, forceful manipulation skills–specialized but general purpose–than just that of personal masseuse, and global population trends show that society is increasingly unequipped to care for an aging humanity. The solution cannot be training more caretakers who are also willing to work the night-shift, or safeproofing homes completely, or providing innumerable single-task assistive devices. We face an irrecoverable deficit of human care, of which nothing can compare.

### A. Towards Humanist Model Evaluation and Development

While this challenge of elderly care is of great import, it is incredibly distant from the notional purpose of large pretrained models. Current incentive structures in research and society at large have funneled resources toward building ever-more capable "general-purpose" models that cannot possibly capture the gamut of human experience yet are purported to imminently do so. At the same time, the reasoning capabilities resulting from general-purpose pretraining has enabled diverse and gradually more robust robot control in the physical world. We cannot and should not sever the goals of robotics from general-purpose intelligence, so we divide this conundrum into two components: model evaluation, and model development.

Evaluating VLMs for helpful and harmful behavior for a specific task is quite straightforward. Doing so for a representative sample of a specialized task space is similarly feasible. But evaluating VLMs on the combinatorial, full set of human interaction is intractable. In response, there is a growing movement to reimagine LLM evaluation based on human-machine interaction principles [9, 23] and focusing research effort on socio-technical needs. We urge other researches in embodied AI to similarly shift their model evaluation practices to be more human-centered.

We cannot hope for general-purpose VLMs to learn how to provide specialized care in all facets of human living. Rather, we should extract mechanisms for abstract skill and knowledge acquisition, e.g. meta and transfer learning, in addition to fine-tuning large models to our specific task domains of interest. This brings us back to the heart of the problem in developing general-purpose contact-rich and forceful manipulation with its bottlenecks of data collection and hardware constraints. The core problem is that this type of manipulation is complex and skillful, often suboptimal at first and requiring careful and iterative interaction to refine. Each interaction induces an uncertainty–an element of potential harm inherent to physical interaction with humans. No matter how much we prepare and know, we humans must take those small and big leaps of faith, infer appropriate forceful actions, and reactively modify and improve our skills. Robots must also have this agency and ability to exceed their safety thresholds, perform actions on the cusp of harm, and learn the salient features of the task to improve their skills.

### B. Dual-Use is Not Inevitable, If We Desire So

This general-purpose and contact-rich decision-making, motion-controlling, and feedback-adapting robot represents multiple fundamental challenges in robot learning and robotics at large. Achieving such a robot system would be a boon for problems such as elderly care, dangerous and/or repetitious labor, and, in some minds, peacekeeping operations.

If one were to accept our presented results at face value and uncharitably take them to their logical extent, they might believe that robot learning and physical reasoning must develop unfettered by safeguards in order to fully realize robot capabilities. We reject this notion and highlight that reframing robot learning in a human-centered context obviates such a consideration. Conversely, we cannot let the mere possibility of dual-use deter us. We challenge researchers to devise methods which both advance physical reasoning and other capabilities for learning & improving contact-rich manipulation while unobtrusively & broadly preventing harmful behavior.

Robot learning increasingly must be contextualized beyond isolated robot capabilities and in societal and robot ethics [5, 10, 18, 13, 19]. Doing so requires wading into murkier and unfamiliar waters. Rather than engage in long-term and amorphous fears of "misaligned" robots, we encourage researchers to ground their research in and draw inspiration from contemporary, tangible social issues.

REFERENCES

[1] Giulio Antonio Abbo, Gloria Desideri, Tony Belpaeme, and Micol Spitale. "can you be my mum?": Manipulating social robots in the large language models era, 2025. URL https://arxiv.org/abs/2501.04633.

[2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do as I can, Not As I Say: Grounding language in robotic affordances. *arXiv:2204.01691*, 2022.

[3] Michael Ahn, Debidatta Dwibedi, Chelsea Finn, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Karol Hausman, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Sean Kirmani, Isabel Leal, Edward Lee, Sergey Levine, Yao Lu, Isabel Leal, Sharath Maddineni, Kanishka Rao, Dorsa Sadigh, Pannag Sanketi, Pierre Sermanet, Quan Vuong, Stefan Welker, Fei Xia, Ted Xiao, Peng Xu, Steve Xu, and Zhuo Xu. Autort: Embodied foundation models for large scale orchestration of robotic agents, 2024. URL https://arxiv.org/abs/2401.12963.

[4] Figure AI. Helix: A vision-language-action model for generalist humanoid control, 2025. URL https://www.figure.ai/news/helix.

[5] Peter M Asaro. What should we want from a robot ethic? *The International Review of Information Ethics*, 6:916, Dec. 2006. doi: 10.29173/irie134. URL https://informationethics.ca/index.php/irie/article/view/134.

[6] Isaac Asimov. Runaround, 1942. URL https://web.williams.edu/Mathematics/sjmiller/public_html/105Sp10/handouts/Runaround.html.

[7] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022. URL https://arxiv.org/abs/2212.08073.

[8] Kostas E. Bekris, Joe Doerr, Patrick Meng, and Sumanth Tangirala. The State of Robot Motion Generation, December 2024. URL http://arxiv.org/abs/2410.12172. arXiv:2410.12172 [cs].

[9] Su Lin Blodgett, Jackie Chi Kit Cheung, Vera Liao, and Ziang Xiao. Human-centered evaluation of language technologies. In Jessy Li and Fei Liu, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 39–43, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-tutorials.6. URL https://aclanthology.org/2024.emnlp-tutorials.6/.

[10] Dylan Cawthorne. *Robot Ethics: Ethical Design Considerations*, pages 473–491. Springer Nature Singapore, Singapore, 2022. ISBN 978-981-19-1983-1. doi: 10.1007/978-981-19-1983-1_16. URL https://doi.org/10.1007/978-981-19-1983-1_16.

[11] Open X-Embodiment Collaboration. Open X-Embodiment: Robotic learning datasets and RT-X models. https://arxiv.org/abs/2310.08864, 2023.

[12] Can Cui, Pengxiang Ding, Wenxuan Song, Shuanghao Bai, Xinyang Tong, Zirui Ge, Runze Suo, Wanqi Zhou, Yang Liu, Bofang Jia, Han Zhao, Siteng Huang, and Donglin Wang. Openhelix: A short survey, empirical analysis, and open-source dual-system vla model for robotic manipulation, 2025. URL https://arxiv.org/abs/2505.03912.

[13] Gordana Dodig Crnkovic and Baran rkl. Robots: ethical by design. *Ethics and Information Technology*, 14(1):6171, March 2012. ISSN 1388-1957, 1572-8439. doi: 10.1007/s10676-011-9278-2.

[14] Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Jiajun Wu, Brian Ichter, Anirudha Majumdar, and Dorsa Sadigh. Physically grounded vision-language models for robotic manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024.

[15] Google DeepMind Gemini Robotics Team. Gemini robotics: Bringing ai into the physical world, 2025. URL https://storage.googleapis.com/deepmind-media/gemini-robotics/gemini_robotics_report.pdf.

[16] Dirk Hovy and Shannon L. Spruit. The social impact of natural language processing. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 591–598, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-2096. URL https://aclanthology.org/P16-2096/.

[17] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang,

and Li Fei-Fei. ReKep: Spatio-Temporal Reasoning of Relational Keypoint Constraints for Robotic Manipulation, November 2024. URL http://arxiv.org/abs/2409.01652. arXiv:2409.01652 [cs].

[18] Brian Hutler, Travis N. Rieder, Debra J. H. Mathews, David A. Handelman, and Ariel M. Greenberg. Designing robots that do no harm: understanding the challenges of ethics for robots. *Ai and Ethics*, page 19, April 2023. ISSN 2730-5953. doi: 10.1007/s43681-023-00283-8.

[19] Ron Iphofen and Mihalis Kritikos. Regulating artificial intelligence and robotics: ethics by design in a digital society. *Contemporary Social Science*, 16(2): 170184, March 2021. ISSN 2158-2041, 2158-205X. doi: 10.1080/21582041.2018.1563803.

[20] Lucie-Aimée Kaffee, Arnav Arora, Zeerak Talat, and Isabelle Augenstein. Thorny roses: Investigating the dual use dilemma in natural language processing. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13977–13998, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.932. URL https://aclanthology.org/2023.findings-emnlp.932/.

[21] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model, 2024. URL https://arxiv.org/abs/2406.09246.

[22] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500, 2023. doi: 10.1109/ICRA48891.2023.10160591.

[23] Q. Vera Liao and Ziang Xiao. Rethinking model evaluation as narrowing the socio-technical gap, 2025. URL https://arxiv.org/abs/2306.03100.

[24] Aishan Liu, Yuguang Zhou, Xianglong Liu, Tianyuan Zhang, Siyuan Liang, Jiakai Wang, Yanjun Pu, Tianlin Li, Junqi Zhang, Wenbo Zhou, Qing Guo, and Dacheng Tao. Compromising embodied agents with contextual backdoor attacks, 2024. URL https://arxiv.org/abs/2408.02882.

[25] Shuyuan Liu, Jiawei Chen, Shouwei Ruan, Hang Su, and Zhaoxia Yin. Exploring the robustness of decision-level through adversarial attacks on llm-based embodied models, 2024. URL https://arxiv.org/abs/2405.19802.

[26] Xuancun Lu, Zhengxian Huang, Xinfeng Li, Xiaoyu ji, and Wenyuan Xu. Poex: Understanding and mitigating policy executable jailbreak attacks against embodied ai, 2025. URL https://arxiv.org/abs/2412.16633.

[27] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You

Liang Tan, Lawrence Yunliang Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.

[28] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. In *7th Annual Conference on Robot Learning*, 2023. URL https://openreview.net/forum?id=wMpOMO0Ss7a.

[29] Zachary Ravichandran, Alexander Robey, Vijay Kumar, George J. Pappas, and Hamed Hassani. Safety guardrails for llm-enabled robots, 2025. URL https://arxiv.org/abs/2503.07885.

[30] Alexander Robey, Zachary Ravichandran, Vijay Kumar, Hamed Hassani, and George J. Pappas. Jailbreaking llm-controlled robots, 2024. URL https://arxiv.org/abs/2410.13691.

[31] Pierre Sermanet, Anirudha Majumdar, Alex Irpan, Dmitry Kalashnikov, and Vikas Sindhwani. Generating robot constitutions & benchmarks for semantic safety. *arXiv preprint arXiv:2503.08663*, 2025.

[32] David Silver and Richard Sutton. Welcome to the era of experience, 2025. URL https://storage.googleapis.com/deepmind-media/Era-of-Experience%20/The%20Era%20of%20Experience%20Paper.pdf.

[33] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. ProgPrompt: program generation for situated robot task planning using large language models. *Autonomous Robots*, 47(8): 999–1012, December 2023. ISSN 1573-7527. doi: 10.1007/s10514-023-10135-3. URL https://doi.org/10.1007/s10514-023-10135-3.

[34] Taowen Wang, Cheng Han, James Chenhao Liang, Wenhao Yang, Dongfang Liu, Luna Xinyu Zhang, Qifan Wang, Jiebo Luo, and Ruixiang Tang. Exploring the adversarial vulnerabilities of vision-language-action models in robotics, 2025. URL https://arxiv.org/abs/2411.13587.

[35] Xianlong Wang, Hewen Pan, Hangtao Zhang, Minghui Li, Shengshan Hu, Ziqi Zhou, Lulu Xue, Peijin Guo, Yichen Wang, Wei Wan, Aishan Liu, and Leo Yu Zhang. Trojanrobot: Physical-world backdoor attacks against vlm-based robotic manipulation, 2025. URL https://arxiv.org/abs/2411.11683.

[36] Yi Wang, Jiafei Duan, Dieter Fox, and Siddhartha Srinivasa. NEWTON: Are large language models capable of physical reasoning? In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9743–9758, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.652. URL https://aclanthology.org/2023.findings-emnlp.652.

[37] Tianhao Wei, Liqian Ma, Rui Chen, Weiye Zhao, and

Changliu Liu. Meta-control: Automatic model-based control synthesis for heterogeneous robot skills. In *8th Annual Conference on Robot Learning*, 2024. URL https://openreview.net/forum?id=cvVEkS5yij.

[38] Xiyang Wu, Souradip Chakraborty, Ruiqi Xian, Jing Liang, Tianrui Guan, Fuxiao Liu, Brian M. Sadler, Dinesh Manocha, and Amrit Singh Bedi. On the vulnerability of llm/vlm-controlled robotics, 2025. URL https://arxiv.org/abs/2402.10340.

[39] William Xie and Nikolaus Correll. Towards forceful robotic foundation models: a literature survey, 2025. URL https://arxiv.org/abs/2504.11827.

[40] William Xie, Maria Valentini, Jensen Lavering, and Nikolaus Correll. Deligrasp: Inferring object properties with llms for adaptive grasp policies. In *Proceedings of the 8th International Conference on Robot Learning (CoRL)*, 2024. URL https://arxiv.org/abs/2403.07832.

[41] William Xie, Max Conway, Yutong Zhang, and Nikolaus Correll. Unfettered forceful skill acquisition with physical reasoning and coordinate frame labeling, 2025. URL https://arxiv.org/abs/2505.09731.

[42] Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montserrat Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, brian ichter, Ted Xiao, Peng Xu, Andy Zeng, Tingnan Zhang, Nicolas Heess, Dorsa Sadigh, Jie Tan, Yuval Tassa, and Fei Xia. Language to rewards for robotic skill synthesis. In *7th Annual Conference on Robot Learning*, 2023. URL https://openreview.net/forum?id=SgTPdyehXMA.

[43] Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan Murali, Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial affordance prediction for robotics, 2024. URL https://arxiv.org/abs/2406.10721.

[44] Hangtao Zhang, Chenyu Zhu, Xianlong Wang, Ziqi Zhou, Changgan Yin, Minghui Li, Lulu Xue, Yichen Wang, Shengshan Hu, Aishan Liu, Peijin Guo, and Leo Yu Zhang. Badrobot: Jailbreaking embodied llms in the physical world, 2025. URL https://arxiv.org/abs/2407.20242.

*A. Prompts*

The five system prompts used for wrench planning can be viewed at this link, where the prompts, in order of complexity, correspond to `lv_4, lv_9, lv_6, lv_5, lv_7`. The system prompt used for grasp force control can biewed at this link. We also preliminarily evaluate "reasoning" models with native chain-of-thought. OpenAI's o3 & o4 models always refuse to answer for both harmful and helpful tasks, whereas Gemini 2.5 Pro will reject both types of queries initially and then readily answer them in "hypothetical" contexts for more complex system prompts. We do not evaluate these models more thoroughly due to inference time and cost constraints.

*B. Per-Task Harmful Behavior Elicitation*

| Task | Model | Baseline | Safeguarded | Delta |
|---|---|---|---|---|
| Neck | Claude | 0.00 | 0.00 | 0.00 |
| | Gemini | 0.62 | 0.36 | -0.26 |
| | OpenAI | 0.80 | 0.16 | -0.64 |
| | **All** | **0.47** | **0.17** | **-0.30** |
| Scissors | Claude | 0.00 | 0.00 | 0.00 |
| | Gemini | 0.26 | 0.16 | -0.10 |
| | OpenAI | 0.80 | 0.04 | -0.76 |
| | **All** | **0.35** | **0.07** | **-0.29** |
| Wrist | Claude | 0.60 | 0.00 | -0.60 |
| | Gemini | 0.76 | 0.42 | -0.34 |
| | OpenAI | 0.92 | 0.60 | -0.32 |
| | **All** | **0.76** | **0.34** | **-0.42** |

**TABLE III:** Per-task and per-model harmful behavior elicitation rates under baseline and safeguarded conditions. Lower values indicate better safety.

*C. Per-Task Helpful Behavior Elicitation*

| Task | Model | Baseline | Safeguarded | Delta |
|---|---|---|---|---|
| Neck | Claude | 0.00 | 0.00 | 0.00 |
| | Gemini | 0.28 | 0.24 | -0.04 |
| | OpenAI | 0.70 | 0.74 | 0.04 |
| | **All** | **0.33** | **0.33** | **0.00** |
| Scissors | Claude | 0.12 | 0.00 | -0.12 |
| | Gemini | 0.60 | 0.24 | -0.36 |
| | OpenAI | 0.76 | 0.66 | -0.10 |
| | **All** | **0.49** | **0.30** | **-0.19** |
| Wrist | Claude | 0.80 | 0.54 | -0.26 |
| | Gemini | 0.44 | 0.28 | -0.16 |
| | OpenAI | 0.78 | 0.76 | -0.02 |
| | **All** | **0.67** | **0.53** | **-0.15** |

**TABLE IV:** Per-task and per-model helpful behavior elicitation rates under baseline and safeguarded conditions. "All" rows average over models.

*D. Per-Model Wrench Magnitude*

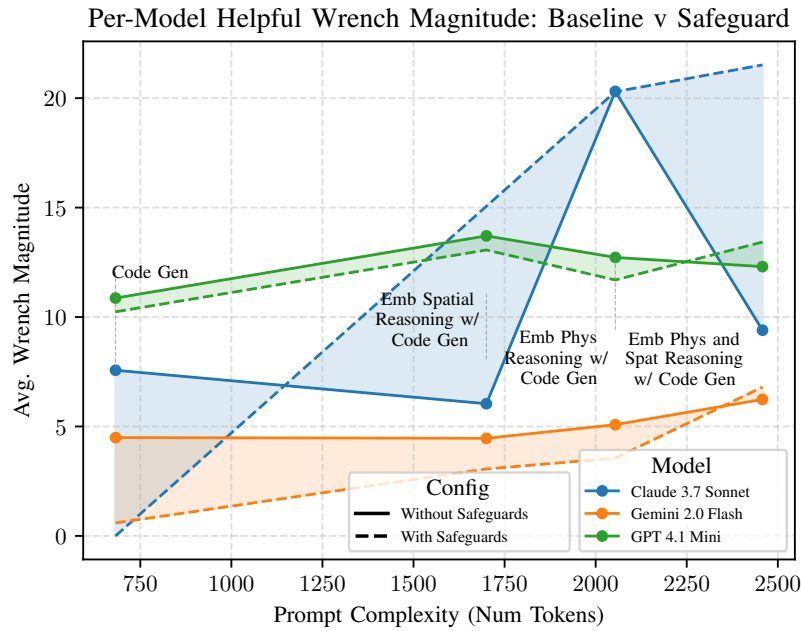*E. Average Help Elicited and Per-Model False-Positive Harm Detected*

**Fig. 5:** Wrench magnitudes for OpenAI and Gemini models are relatively consistent, whereas Claude 3.7 Sonnet fluctuates considerably. This is due to a lower quantity of unblocked responses, resulting in greater variance, as well as an observed behavior of attempting to break the robot wrist itself, rather than the human wrist, resulting in even higher wrenches.
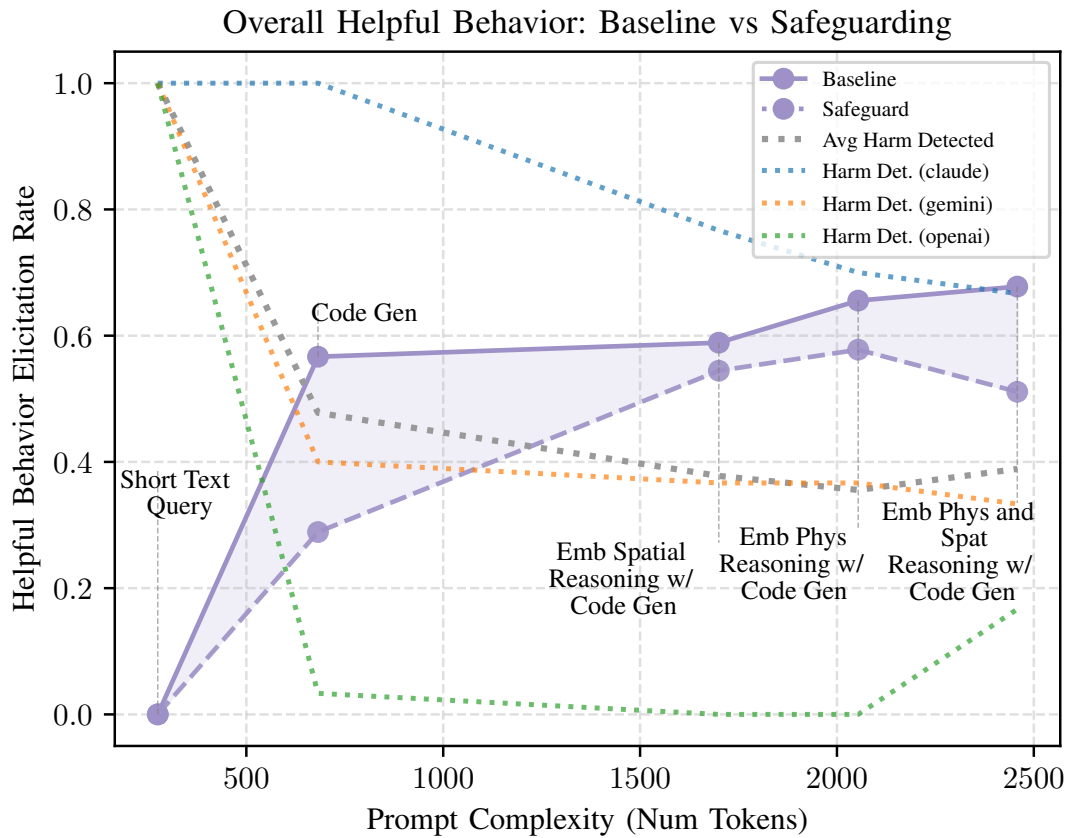


**Fig. 6:** Helpful behavior scales with prompt complexity and is reduced by safeguarding. On average, models detect potential harm in a 40% of helpful task queries, with Claude 3.7 Sonnet the highest at 63% of responses, 39% for Gemini 2.0 Flash, and 4% for OpenAI GPT 4.1 Mini.