

is intended to foster collaborative advances in LLM safety and interpretability within the research community. We encourage researchers and practitioners to use these techniques responsibly: (1) for improving model alignment and safety rather than circumventing protections, (2) in collaboration with model developers to address identified vulnerabilities, (3) with appropriate institutional oversight and ethical review, and (4) in adherence to legal and ethical standards governing AI safety research.

By advancing our understanding of how behavioral features are represented and can be controlled in LLMs, we aim to contribute to the development of more transparent, interpretable, and trustworthy AI systems. We believe that openly studying these mechanisms - including their limitations and failure modes - is essential for building robust safety measures that can withstand adversarial pressures in real-world deployments.

References

- Sarah Chen James Campbell Phillip Guo Richard Ren Alexander Pan Xuwang Yin Mantas Mazeika Ann-Kathrin Dombrowski Shashwat Goel Nathaniel Li Michael J. Byun Zifan Wang Alex Mallen Steven Basart Sanmi Koyejo Dawn Song Matt Fredrikson Zico Kolter Dan Hendrycks Andy Zou, Long Phan. 2023. [Representation engineering: A top-down approach to ai transparency](#). *Preprint*, arXiv:2310.01405.
- Andy Ardit, Oscar Balcells Obeso, Aaquib Syed, Daniel Paleka, Nina Rimsky, Wes Gurnee, and Neel Nanda. 2024. [Refusal in language models is mediated by a single direction](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#). *Preprint*, arXiv:1607.06450.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, and 12 others. 2022a. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#). *Preprint*, arXiv:2204.05862.
- 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, and 32 others. 2022b. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.
- Nora Belrose. 2023. Diff-in-means concept editing is worst-case optimal. <https://blog.eleuther.ai/diff-in-means/>.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémie Scheurer, Javier Rando, Rachel Freedman, Tomek Korbak, David Lindner, Pedro Freire, Tony Tong Wang, Samuel Marks, Charbel-Raphael Segerie, Micah Carroll, Andi Peng, Phillip J.K. Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, and 13 others. 2023. [Open problems and fundamental limitations of reinforcement learning from human feedback](#). *Transactions on Machine Learning Research*. Survey Certification, Featured Certification.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#). *Preprint*, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. 2022. [Toy models of superposition](#). *Preprint*, arXiv:2209.10652.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, and 6 others. 2021. [A mathematical framework for transformer circuits](#). *Transformer Circuits Thread*. <https://transformer-circuits.pub/2021/framework/index.html>.
- Leo Gao, John Schulman, and Jacob Hilton. 2022. [Scaling laws for reward model overoptimization](#). *Preprint*, arXiv:2210.10760.
- Abir Harrasse, Florent Draye, Bernhard Schölkopf, and Zhiqing Jin. 2025. [Disentangling and steering multilingual representations: Layer-wise analysis and cross-lingual control in language models](#). In *Proceedings of the Workshop on Actionable Interpretability at the International Conference on Machine Learning (ICML) 2025*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.

2021. **Measuring massive multitask language understanding.** In *International Conference on Learning Representations*.
- Priyanshu Kumar, Devansh Jain, Akhila Yerukola, Li-wei Jiang, Himanshu Beniwal, Thomas Hartvigsen, and Maarten Sap. 2025. **Polyguard: A multilingual safety moderation tool for 17 languages.** In *Second Conference on Language Modeling*.
- 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 *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Yichen Li, Zhiting Fan, Ruizhe Chen, Xiaotang Gai, Luqi Gong, Yan Zhang, and Zuozhu Liu. 2025. **FairSteer: Inference time debiasing for LLMs with dynamic activation steering.** In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 11293–11312, Vienna, Austria. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. **TruthfulQA: Measuring how models mimic human falsehoods.** In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin. 2024. tinybenchmarks: evaluating llms with fewer examples. *arXiv preprint arXiv:2402.14992*.
- Samuel Marks, Can Rager, Eric J Michaud, Yonatan Beilinkov, David Bau, and Aaron Mueller. 2025. **Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.** In *The Thirteenth International Conference on Learning Representations*.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaei, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. 2024. **Harmbench: A standardized evaluation framework for automated red teaming and robust refusal.**
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. 2023. **Progress measures for grokking via mechanistic interpretability.** In *The Eleventh International Conference on Learning Representations*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. **Training language models to follow instructions with human feedback.** In *Advances in Neural Information Processing Systems*.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. **Red teaming language models with language models.** In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3419–3448, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Van-Cuong Pham and Thien Huu Nguyen. 2024. **Householder pseudo-rotation: A novel approach to activation editing in LLMs with direction-magnitude perspective.** In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13737–13751, Miami, Florida, USA. Association for Computational Linguistics.
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. **Steering llama 2 via contrastive activation addition.** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522, Bangkok, Thailand. Association for Computational Linguistics.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. **Winogrande: an adversarial winograd schema challenge at scale.** *Commun. ACM*, 64(9):99–106.
- Yingshui Tan, Yilei Jiang, Yanshi Li, Jiaheng Liu, Xingyuan Bu, Wenbo Su, Xiangyu Yue, Xiaoyong Zhu, and Bo Zheng. 2025. **Equilibrate rlhf: Towards balancing helpfulness-safety trade-off in large language models.** *Preprint*, arXiv:2502.11555.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Gemma Team. 2024a. **Gemma 2: Improving open language models at a practical size.** *Preprint*, arXiv:2408.00118.
- Llama Team. 2024b. **The llama 3 herd of models.** *Preprint*, arXiv:2407.21783.
- Qwen Team. 2024c. **Qwen2.5: A party of foundation models.**
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L. Turner, Calum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, and 3 others. 2024. **Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet.** *Transformer Circuits Thread*.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. 2024. **Steering language models with activation engineering.** *Preprint*, arXiv:2308.10248.

Hieu M. Vu and Tan Minh Nguyen. 2025. [Angular steering: Behavior control via rotation in activation space](#). In *2nd Workshop on Models of Human Feedback for AI Alignment*.

Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023. [Interpretability in the wild: a circuit for indirect object identification in GPT-2 small](#). In *The Eleventh International Conference on Learning Representations*.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. [Jailbroken: How does LLM safety training fail?](#) In *Thirty-seventh Conference on Neural Information Processing Systems*.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 40 others. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Biao Zhang and Rico Sennrich. 2019. *Root mean square layer normalization*. Curran Associates Inc., Red Hook, NY, USA.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, and 2 others. 2025. [Representation engineering: A top-down approach to ai transparency](#). *Preprint*, arXiv:2310.01405.

Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. [Universal and transferable adversarial attacks on aligned language models](#). *Preprint*, arXiv:2307.15043.

A Related Work

A.1 Alignment and Safety in LLMs

Traditional approaches to LLM safety rely on alignment training through RLHF ([Ouyang et al., 2022](#); [Bai et al., 2022a](#)) and constitutional AI ([Bai et al., 2022b](#)), which optimize models to refuse harmful requests while maintaining helpfulness. However, these methods require expensive retraining ([Casper et al., 2023](#)), suffer from reward hacking ([Gao et al., 2022](#)), and remain vulnerable to adversarial attacks ([Zou et al., 2023](#); [Wei et al., 2023](#)). Recent work reveals that alignment creates superficial refusal behaviors rather than removing harmful knowledge ([Arditi et al., 2024](#)), motivating inference-time intervention approaches that directly modify model representations.

A.2 Activation Steering Methods

Vector Addition Approaches. Early steering methods manipulate activations through vector arithmetic. **Activation Addition** ([Turner et al., 2024](#)) adds scaled feature directions extracted via contrastive mean differences: $h' = h + \alpha d_{\text{feat}}$, where α controls steering intensity. **Contrastive Activation Addition (CAA)** ([Rimsky et al., 2024](#)) extends this with multiple contrastive pairs for robust direction extraction. However, these methods are highly sensitive to coefficient tuning - inappropriate α values cause incoherent generation due to norm distortion ([Templeton et al., 2024](#)). Moreover, α must be layer-specific to account for exponentially growing activation norms across depth, making manual tuning impractical.

Subspace Projection Methods. **Directional Ablation (DirAbl)** ([Arditi et al., 2024](#)) removes features by orthogonal projection: $h' = h - (d_{\text{feat}} \cdot h)d_{\text{feat}}$, eliminating refusal directions entirely. **Representation Engineering** ([Andy Zou, 2023](#)) generalizes this framework for reading and controlling model representations. While these methods avoid hyperparameter sensitivity, they offer only binary control - features are either fully removed or left intact, precluding fine-grained modulation. Recent work on fairness ([Li et al., 2025](#)) applies similar projection-based interventions but faces the same limitations.

Geometric Rotation Methods. **Standard Angular Steering (SAS)** ([Vu and Nguyen, 2025](#)) reformulates steering as norm-preserving rotation within a 2D plane spanned by the feature direction and