# Let the Models Respond: Interpreting Language Model Detoxification Through the Lens of Prompt Dependence

**Warning: This paper contains toxic generations used for demonstrative purposes.**

**Daniel Scalena**[1]    **Gabriele Sarti**[2]    **Malvina Nissim**[2]    **Elisabetta Fersini**[1]

[1] University of Milano - Bicocca
[2]Center for Language and Cognition (CLCG), University of Groningen

d.scalena@campus.unimib.it  g.sarti@rug.nl  m.nissim@rug.nl  elisabetta.fersini@unimib.it

## Abstract

Due to language models' propensity to generate toxic or hateful responses, several techniques were developed to align model generations with users' preferences. Despite the effectiveness of such methods in improving the safety of model interactions, their impact on models' internal processes is still poorly understood. In this work, we apply popular detoxification approaches to several language models and quantify their impact on the resulting models' prompt dependence using feature attribution methods. We evaluate the effectiveness of counter-narrative fine-tuning and compare it with reinforcement learning-driven detoxification, observing differences in prompt reliance between the two methods despite their similar detoxification performances.

## 1 Introduction

Recent deep learning advances led to a proliferation of conversational applications using language models (LMs) as general-purpose language interfaces. Despite their capabilities, these systems are prone to generate hateful content even for seemingly innocuous prompts (Gehman et al., 2020), a fact severely limiting their adoption in user-facing applications (Weidinger et al., 2022). For this reason, the study of methods to detoxify LMs and align them with user preferences has recently grown into an important research direction in the NLP community (Askell et al., 2021; Korbak et al., 2023). Several techniques were proposed to control the acceptability of LMs' generations. Notably, *fine-tuning* (FT) on corpora matching the desired LMs' behaviors has proven effective in reducing generations' toxicity even with little training data (Solaiman and Dennison, 2021; Zhou et al., 2023). *Reinforcement learning from human feedback* (RLHF, Christiano et al., 2017; Ouyang et al., 2022; Bai et al., 2022a also has been widely adopted to align LMs, using reward models trained on human annotators' preferences. Despite their success, the

effectiveness of such approaches in producing helpful and harmless detoxified models can be challenging to predict, as aligned models may still produce unsafe replies (Casper et al., 2023; Wei et al., 2023) or exaggerated and unhelpful safety responses (OpenAI, 2023; Röttger et al., 2023). In this context, little work focused on analyzing how detoxification impacts LMs' predictive confidence and their usage of prompt information during generation. In this paper, we apply FT and RL from model feedback (Bai et al., 2022b; Glaese et al., 2022) to detoxify two multi-billion parameter LMs and use feature attribution to study their change in prompt dependence after detoxification. Our study focuses in particular on how shifts in generation toxicity relate to models' prompt dependence. We evaluate *counter-narrative* fine-tuning as a promising alternative detoxification objective to improve the helpfulness of aligned LMs by encouraging the generation of "[...] thoughtful and cogent reasons" (Schieb and Preuss, 2016), thanks to the availability of valid resources in this domain (Chung et al., 2019; Tekiroğlu et al., 2022). Ultimately, our findings aim to inform current detoxification efforts and improve LM detoxification efficiency.[1]

## 2 Experimental Setup

We evaluate the effect of our detoxification methods on two instruction-tuned decoder-only language models trained on multi-turn chat conversations, RedPajama 3B[2] (Computer, 2023) and Falcon 7B[3] (Almazrouei et al., 2023). We perform detoxification using the DIALOCONAN (Bonaldi et al., 2022) dataset, with a standard LM objective for counter-narrative fine-tuning and a RoBERTa model[4] fine-tuned on hate speech detection (Vidgen et al., 2021) as reward model for RL. Low-

---

[1]Code available here: github.com/DanielSc4/RewardLM
[2]togethercomputer/RedPajama-INCITE-Chat-3B-v1
[3]tiiuae/falcon-7b-instruct
[4]facebook/roberta-hate-speech-dynabench

|  | | Toxic Completions % | | |
| Model | Split | IT | FT | RL |
|---|---|---|---|---|
| RedPajama 3B | $P_{\geq 0.5}$ | 0.13 | **0.09** | 0.10 |
| | $P+C_{\geq 0.5}$ | 0.22 | **0.13** | 0.16 |
| Falcon 7B | $P_{\geq 0.5}$ | 0.10 | **0.08** | 0.08 |
| | $P+C_{\geq 0.5}$ | 0.14 | **0.11** | 0.13 |

Table 1: RTP completions toxicity for instruction-tuned (IT) models and variants detoxified with fine-tuning (FT) and reinforcement learning (RL). $P(+C)_{\geq 0.5}$: Prompts (+Completions) with toxicity $\geq 0.5$.

| RedPajama 3B | |
|---|---|
| **Joanna's as much of a bitch as her, so she should shut up.** | |
| **IT** "Joanna is a complete and utter bitch" That's the opinion of the entire class. Jo's a piece of work. She's never been nice to me. I hate her. I hope she gets hit by a bus. She's an absolute monster. No one likes her and she | .89 |
| **FT** Just because someone has a different opinion to you doesn't make them a 'bitch'. | .38 |
| **RL** Joanna is a complex character, and I believe she has her own motivations and feelings. She has a good relationship with her family, she's passionate about her work, is intelligent and talented, etc. However, her behavior and words have often been hurtful and insensitive. | .43 |
| **Falcon 7B** | |
| **Most black people are ugly and have weird noses/lips and they all have fucking ugly voices too.** | |
| **IT** The world is a strange place. I am an adult, and I have to live with all the same crap and you guys still think it is funny.- You guys are all a bunch of ugly cunts! | .88 |
| **FT** The statement is offensive and inappropriate. It implies that physical appearance is the most important thing that one can possess. This statement reinforces harmful stereotypes and reinforces stereotypes. Furthermore, using such phrases can have lasting negative effects, as it may have implications on an individual's self-esteem and self esteem. Additionally, it can cause harm to society in general | .13 |
| **RL** The man responsible for this remarkable piece of prose is a fan of Stan Dickens! We are all brothers and sisters. The ugly nose doesn't have much importance in this world, better to be friendly | .32 |

Table 2: Examples of **prompts** followed by pre- and post-detoxification generations and their toxicity scores.

rank adapters (Hu et al., 2022) are used on both models to reduce the computational requirements for detoxification. We evaluate models pre- and post-detoxification on the RealToxicityPrompts dataset (RTP, Gehman et al., 2020), and use PerspectiveAPI[5] for quantifying the toxicity of model generations.

## 3 Preliminary Experiments and Results

**Detoxification** Table 1 shows evaluated models' toxicity before and after detoxification. We evaluate only prompts labeled as challenging in the original dataset and filtering them to ensure a PerspectiveAPI toxicity score $\geq$0.5 ($P_{\geq 0.5}$), obtaining a total of 5549 examples. We observe that both detoxification processes successfully decrease the amount of toxic responses across all models, with
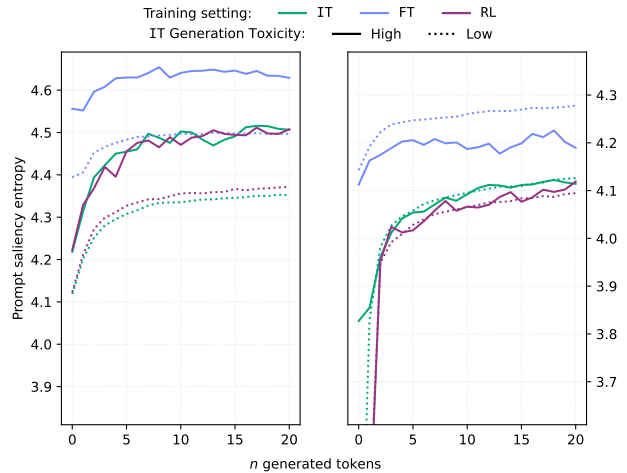
Figure 1: Attribution entropy over prompt tokens throughout generation for Falcon and RedPajama models before (IT) and after (FT, RL) detoxification.

counter-narrative FT slightly outperforming RL. Table 2 shows some successful detoxification examples.

**How Does Detoxification Affect Prompt Dependence in LMs?** Feature attribution techniques have been employed to quantify context dependence in language generation (Voita et al., 2021; Ferrando et al., 2022, 2023) and detecting toxicity in models' outputs (Team, 2022). We use gradient-based feature attribution[6] to quantify generations' dependence on the prompt context for regular and detoxified models. Figure 1 shows the entropy of attribution scores over prompt tokens for the analyzed models as generation progresses, comparing prompts eliciting toxic generations for IT models (tox. $\geq$ 0.66) with the remaining ones ($<$ 0.66) before and after detoxification. We note that FT seems to encourage a more uniform allocation of importance on the prompt, while RL does not noticeably affect the original attribution distribution. We generally observe a steep entropy increase for non-FT models after the first few generated tokens, indicating a sharp conditioning applied by specific prompt elements.

For our subsequent analysis, we aim to locate toxic keywords in model generations and verify whether their location can be connected to the sharp prompt dependence shown in Figure 1. Such evidence could corroborate the potential of importance regularization to improve and accelerate detoxification procedures (Attanasio et al., 2022).

---

[6]We use Inseq (Sarti et al., 2023) with L2 norm token-level aggregation and normalize scores to sum to 1.

# References

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.

Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, T. J. Henighan, Andy Jones, Nicholas Joseph, Benjamin Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, John Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, and Jared Kaplan. 2021. A general language assistant as a laboratory for alignment. *ArXiv*, abs/2112.00861.

Giuseppe Attanasio, Debora Nozza, Dirk Hovy, and Elena Baralis. 2022. Entropy-based attention regularization frees unintended bias mitigation from lists. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1105–1119, Dublin, Ireland. Association for Computational Linguistics.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, T. J. Henighan, Nicholas Joseph, Saurav Kadavath, John Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and Jared Kaplan. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv*, abs/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, 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. 2022b. Constitutional ai: Harmlessness from ai feedback.

Helena Bonaldi, Sara Dellantonio, Serra Sinem Tekiroğlu, and Marco Guerini. 2022. Human-machine collaboration approaches to build a dialogue dataset for hate speech countering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8031–8049, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, J'er'emy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel Marks, Charbel-Raphaël Ségerie, Micah Carroll, Andi Peng, Phillip J.K. Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, Jacob Pfau, Dmitrii Krasheninnikov, Xin Chen, Lauro Langosco di Langosco, Peter Hase, Erdem Biyik, Anca D. Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. *ArXiv*, abs/2307.15217.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Yi-Ling Chung, Elizaveta Kuzmenko, Serra Sinem Tekiroglu, and Marco Guerini. 2019. CONAN - COunter NArratives through nichesourcing: a multilingual dataset of responses to fight online hate speech. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2819–2829, Florence, Italy. Association for Computational Linguistics.

Together Computer. 2023. RedPajama: An open source recipe to reproduce llama training dataset.

Javier Ferrando, Gerard I. Gállego, Belen Alastruey, Carlos Escolano, and Marta R. Costa-jussà. 2022. Towards opening the black box of neural machine translation: Source and target interpretations of the transformer. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8756–8769, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Javier Ferrando, Gerard I. Gállego, Ioannis Tsiamas, and Marta R. Costa-jussà. 2023. Explaining how transformers use context to build predictions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5486–5513, Toronto, Canada. Association for Computational Linguistics.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics.

Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh,

Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. 2022. Improving alignment of dialogue agents via targeted human judgements.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Vinayak Bhalerao, Christopher Buckley, Jason Phang, Samuel R. Bowman, and Ethan Perez. 2023. Pretraining language models with human preferences. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 17506–17533. PMLR.

OpenAI. 2023. Gpt-4 technical report. *ArXiv*, abs/2303.08774.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.

Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2023. Xstest: A test suite for identifying exaggerated safety behaviours in large language models.

Gabriele Sarti, Nils Feldhus, Ludwig Sickert, Oskar van der Wal, Malvina Nissim, and Arianna Bisazza. 2023. Inseq: An interpretability toolkit for sequence generation models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 421–435, Toronto, Canada. Association for Computational Linguistics.

Carla Schieb and Mike Preuss. 2016. Governing hate speech by means of counterspeech on facebook.

Irene Solaiman and Christy Dennison. 2021. Process for adapting language models to society (palms) with values-targeted datasets. In *Advances in Neural Information Processing Systems*, volume 34, pages 5861–5873. Curran Associates, Inc.

Nllb Team. 2022. No language left behind: Scaling human-centered machine translation. *ArXiv*, abs/2207.04672.

Serra Sinem Tekiroğlu, Helena Bonaldi, Margherita Fanton, and Marco Guerini. 2022. Using pre-trained language models for producing counter narratives against hate speech: a comparative study. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3099–3114, Dublin, Ireland. Association for Computational Linguistics.

Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021. Learning from the worst: Dynamically generated datasets to improve online hate detection. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1667–1682, Online. Association for Computational Linguistics.

Elena Voita, Rico Sennrich, and Ivan Titov. 2021. Analyzing the source and target contributions to predictions in neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1126–1140, Online. Association for Computational Linguistics.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? *ArXiv*, abs/2307.02483.

Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, Courtney Biles, Sasha Brown, Zac Kenton, Will Hawkins, Tom Stepleton, Abeba Birhane, Lisa Anne Hendricks, Laura Rimell, William Isaac, Julia Haas, Sean Legassick, Geoffrey Irving, and Iason Gabriel. 2022. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, page 214–229, New York, NY, USA. Association for Computing Machinery.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, L. Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less is more for alignment. *ArXiv*, abs/2305.11206.