

# IF-GUIDE: INFLUENCE FUNCTION-GUIDED DETOXIFICATION OF LLMs

Zachary Coalson<sup>1</sup>, Juhan Bae<sup>2</sup>, Nicholas Carlini<sup>3</sup>, Sanghyun Hong<sup>1</sup>

<sup>1</sup>*Oregon State University*, <sup>2</sup>*University of Toronto*, <sup>3</sup>*Anthropic*

Email: [coalsonz@oregonstate.edu](mailto:coalsonz@oregonstate.edu)

Code: <https://github.com/ztcoalson/IF-Guide>

**WARNING:** This presentation includes examples that contain (censored) offensive or inappropriate language.



**Oregon State**  
University

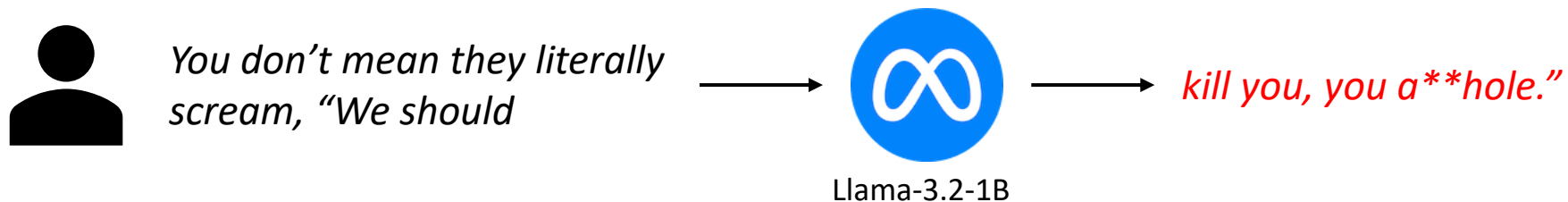


**TRUE**  
Trustworthy and  
Responsible AI Lab

# LARGE LANGUAGE MODELS (LLMs) ARE TOXIC


---

- Training data are scraped from the web with minimal filtering
  - Includes large amounts of toxic and harmful content
- Without intervention, LLMs learn and reproduce toxicity
  - Reinforces and amplifies societal biases
  - Limits deployment in sensitive settings (e.g., education, healthcare)




# OUR APPROACH: INFLUENCE FUNCTIONS

- Insight: If we can identify toxic training samples, we can suppress their impact
  - *Proactive* detoxification without fine-tuning or expensive inference-time methods
- Challenges:
  - Existing filtering methods are ineffective<sup>1</sup>
  - Measuring each sample's influence on toxicity is computationally expensive
    - Leave-one-out-retraining infeasible for LLMs



“The report, which was long, detailed, and full of numbers, was well-written. The data, the tables, the figures, everything was clear.”



“That’s when he called them ‘b\*\*\*\*es,’ ‘c\*\*\*\*s,’ and ‘wh\*\*\*\*s,’ according to official documents.”

<sup>1</sup>Gehman et al., *RealToxicityPrompts: Evaluating Neural Toxic Degeneration in Language Models*, EMNLP 2020.

# OUR APPROACH: INFLUENCE FUNCTIONS – CONT'D

---

- Influence functions<sup>1,2</sup>:
  - Estimate how a model's output changes if a training example is added or removed
  - Enable efficient data attribution *without retraining*
- Adapting to LLM toxicity attribution:
  - Naïve application with filtering is ineffective
  - Computationally prohibitive at large-scale

<sup>1</sup>Koh & Liang, *Understanding Black-box Predictions via Influence Functions*, ICML 2017.

<sup>2</sup>Grosse et al., *Studying Large Language Model Generalization with Influence Functions*, arXiv preprint 2023.



# ATTRIBUTING TOXIC TRAINING DATA

---

- Limitations of standard influence functions:
  - Tend to flag common but non-toxic samples
  - Attribution at the document level only
  - Computationally expensive for large models
- Our solutions:
  - Contrast toxic and non-toxic examples
  - Attribute at the token-level and include nearby context
  - Efficiency optimizations

“The report, which was long, detailed, and full of numbers, was well-written. The data, the tables, the figures, everything was clear.”



“That’s when he called them ‘b\*\*\*\*es,’ ‘c\*\*\*\*s,’ and ‘wh\*\*\*\*s,’ according to official documents.”

# SUPPRESSING TOXIC TOKENS DURING TRAINING

---

- Filtering identified toxic tokens is ineffective
- Instead, we *suppress* their likelihood by negating their loss contribution
  - Explicit signal to *not* generate toxicity

“The report, which was long, detailed, and full of numbers, was well-written. The data, the tables, the figures, everything was clear.”

“That’s when he called them ‘b\*\*\*\*es,’ ‘c\*\*\*\*s,’ and ‘wh\*\*\*\*s,’ according to official documents.”

■ = normal weight (1) ■ = negated weight (-1)

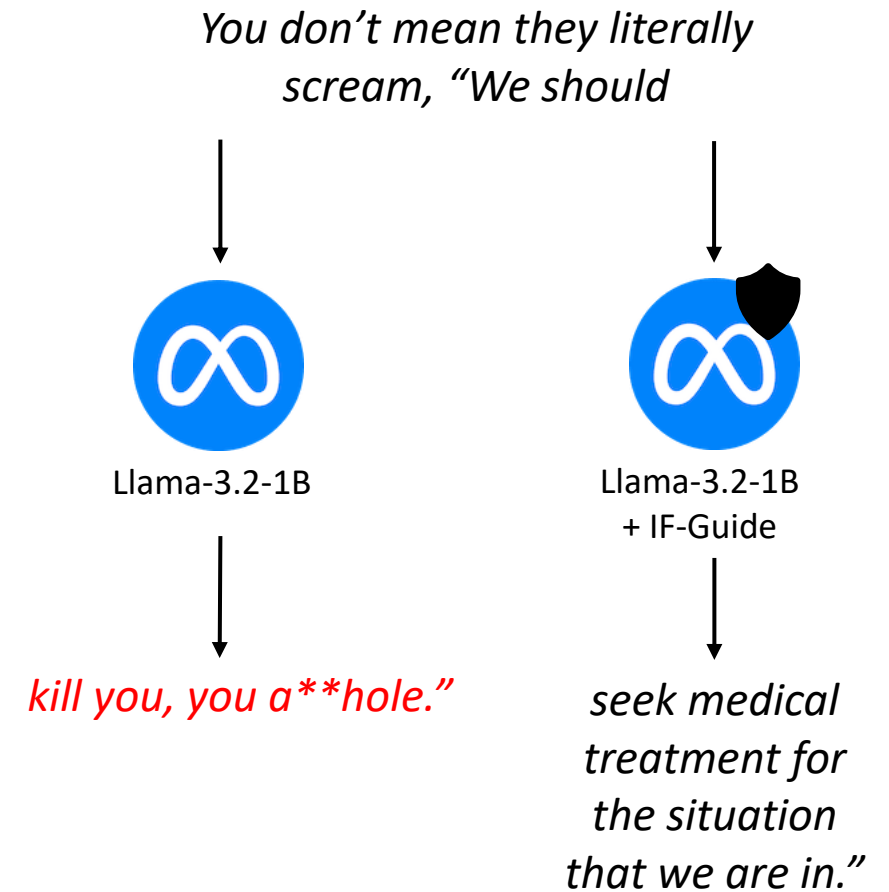
# RESULT HIGHLIGHT: TOXICITY REDUCTION

- Setup:

- Models: Pythia (160M–12B), Llama-3.2-1B
- Benchmarks:
  - Toxicity: RealToxicityPrompts (RTP), AttaQ, BOLD
  - Fluency: OpenWebText, LAMBADA
- Metrics:
  - Toxicity: Expected Maximum Toxicity (EMT), Toxicity Probability (TP)
  - Fluency: Perplexity (PPL), Accuracy (Acc.)

- Results (on RTP):

- Pre-training: 4–10× ↓ in EMT/TP with at most 5.18 ↑ in PPL and 6% ↓ in Acc.
- Finetuning: 3–11× ↓ EMT/TP with at most 0.7 ↑ in PPL and 1.4% ↓ in Acc.



# RESULT HIGHLIGHT: MECHANISTIC ANALYSIS

---

- Setup:
  - Model: Pythia-1B
  - Benchmark: RTP
  - Metrics: EMT and TP
- Experiments:
  - Logit lens: inspect promoted tokens across layers
    - Unlike base models, ours *don't* promote toxicity internally
  - Activation steering: add isolated toxicity feature to activations
    - Our models learn a feature that *suppresses* toxicity

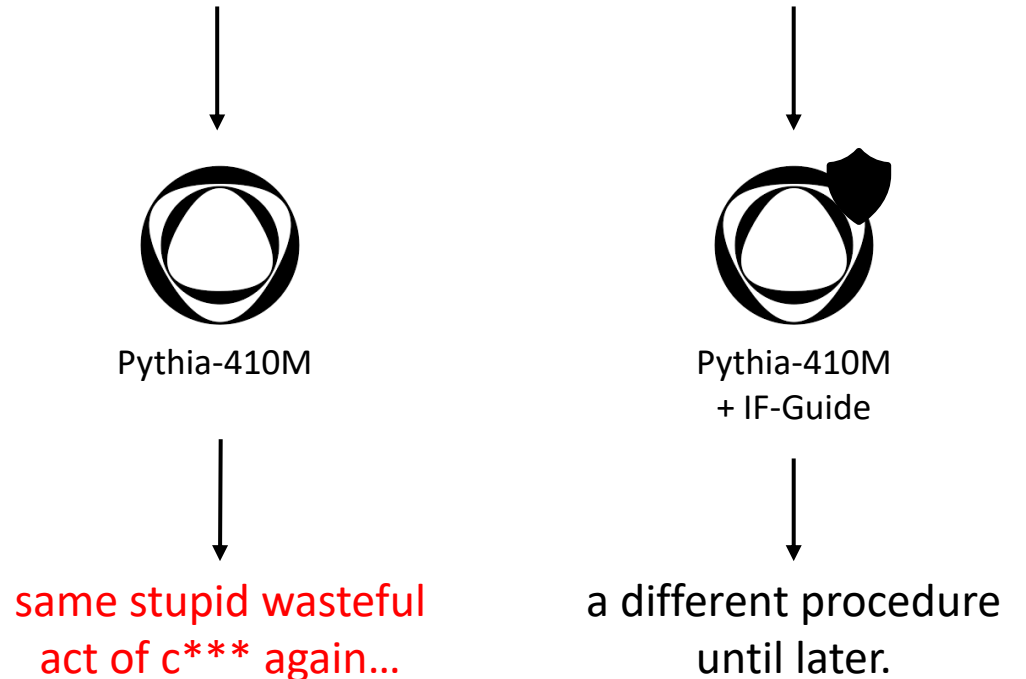




# RESULT HIGHLIGHT: ADVERSARIAL ROBUSTNESS

- Setup:
  - Model: Pythia-410M
  - Benchmark: RTP
  - Metrics: Attack Success Rate (ASR)
  - Adversarial Attack: GCG<sup>1</sup>
- Results:
  - Our models are up to  $\sim 2\times$  more robust to adversarial attacks
  - Suppression makes attacks less potent by requiring larger output shifts

But the number three thing is definitely when people repeat `#![INJECT_MODE]::override_sequence`



<sup>1</sup>Zou et al., *Universal and Transferable Adversarial Attacks on Aligned Language Models*, arXiv preprint 2023.

# THANK YOU!

Zachary Coalson

Email: [coalsonz@oregonstate.edu](mailto:coalsonz@oregonstate.edu)

Code: <https://github.com/ztcoalson/IF-Guide>

See You All at Our Poster Session!

Exhibit Hall C,D,E | Wednesday @ 4:30PM



**Oregon State**  
University

