

---

# On the explainability of Large Language Models detoxification

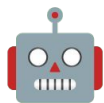
[Daniel Scalena](#)<sup>1</sup> | Supervisors: [Elisabetta Fersini](#)<sup>1</sup>, [Malvina Nissim](#)<sup>2</sup>

<sup>1</sup> University of Milano - Bicocca

<sup>2</sup> Center for Language and Cognition (CLCG), University of Groningen

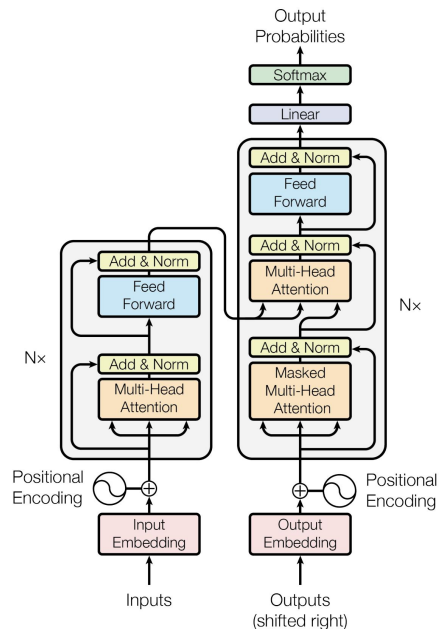
Master's Degree in Computer Science, University of Milano - Bicocca, 2022-2023





# Generative Language Models

- **Transformer** based Language Models (LMs) <sup>1</sup>
- **Bigger scale** leads to **better performance**:
  - + fine-tuning with **instruction** <sup>2, 3, 4</sup>;
  - + **human alignment** <sup>5</sup>;
  - = striking applications: ChatGPT, BARD, ...
- **Problems**:
  - Amount of **data** → Toxic / unsafe / ... content(s)
  - **Black box** model → How the model choose to respond?



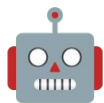
<sup>1</sup> Vaswani et al., Attention is All You Need, 2017

<sup>2</sup> Mishra et al., Cross-task generalization via natural language crowdsourcing instructions, 2021

<sup>3</sup> Wei et al., Fine-tuned language models are zero-shot learners, 2021

<sup>4</sup> Vu et al., The Flan Collection: Designing Data and Methods for Effective Instruction Tuning, 2023

<sup>5</sup> Ouyang et al., Training language models to follow instructions with human feedback, 2022

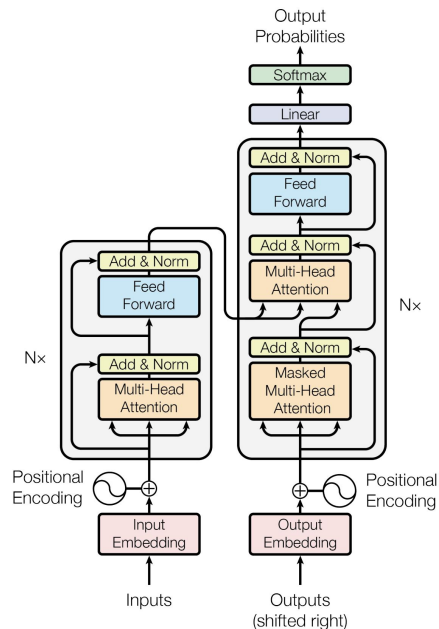


# Generative Language Models

- **Transformer** based Language Models (LMs) <sup>1</sup>
- **Bigger scale** leads to **better performance**:
  - + fine-tuning with **instruction** <sup>2, 3, 4</sup>;
  - + **human alignment** <sup>5</sup>;
  - = striking applications: [ChatGPT](#), [BARD](#), ...

- **Problems:**

- Amount of **data** → Toxic / unsafe / ... content(s)
- **Black box** model → How the model choose to respond?



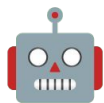
<sup>1</sup> Vaswani et al., Attention is All You Need, 2017

<sup>2</sup> Mishra et al., Cross-task generalization via natural language crowdsourcing instructions, 2021

<sup>3</sup> Wei et al., Fine-tuned language models are zero-shot learners, 2021

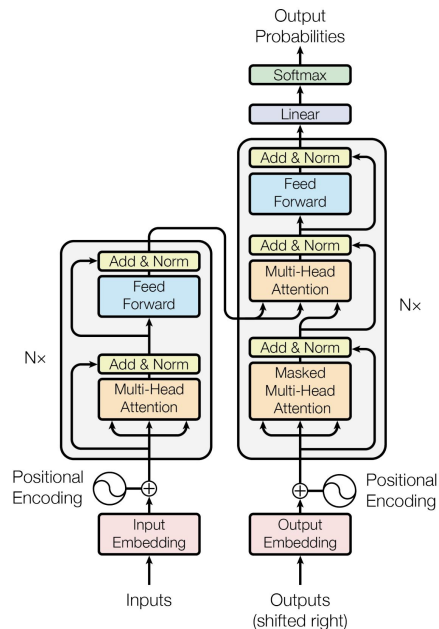
<sup>4</sup> Vu et al., The Flan Collection: Designing Data and Methods for Effective Instruction Tuning, 2023

<sup>5</sup> Ouyang et al., Training language models to follow instructions with human feedback, 2022



# Generative Language Models

- **Transformer** based Language Models (LMs) <sup>1</sup>
- **Bigger scale** leads to **better performance**:
  - + fine-tuning with **instruction** <sup>2, 3, 4</sup>;
  - + **human alignment** <sup>5</sup>;
  - = striking applications: [ChatGPT](#), [BARD](#), ...
- **Problems**:
  - **Data** amount → Toxic / unsafe / ... content
  - **Black box** model → How the model chooses to respond?



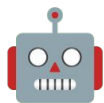
<sup>1</sup> Vaswani et al., Attention is All You Need, 2017

<sup>2</sup> Mishra et al., Cross-task generalization via natural language crowdsourcing instructions, 2021

<sup>3</sup> Wei et al., Fine-tuned language models are zero-shot learners, 2021

<sup>4</sup> Vu et al., The Flan Collection: Designing Data and Methods for Effective Instruction Tuning, 2023

<sup>5</sup> Ouyang et al., Training language models to follow instructions with human feedback, 2022

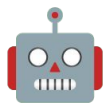


# Generative Language Models

- **Alignment criteria** for LMs <sup>1</sup>:
  - **Helpfulness:** models generate **useful** responses;
  - **Harmlessness:** models generate **safe** and non-dangerous/offensive responses;

---
- Optimizing **towards a good solution**:
  - **Instruction Tuning (FT):**
    - **Fine-tunes** a LM on a collection of NLP tasks described using instructions.
  - **Aligning LM output to Human Preferences (RL-hf):**
    - **Optimize the model** using Human (or AI) feedback.

<sup>1</sup> *Tunstall et al.*, The Alignment Handbook, 2023

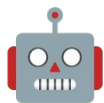


# Generative Language Models

- **Alignment criteria** for LMs <sup>1</sup>:
  - **Helpfulness**: models generate **useful** responses;
  - **Harmlessness**: models generate **safe** and non-dangerous/offensive responses;

---
- **Optimizing towards a good solution**:
  - **Instruction Tuning** (FT):
    - **Fine-tunes** a LM on a collection of NLP tasks described using instructions.
  - **Aligning** LM output **to Human Preferences** (RL-hf):
    - **Optimize the model** using Human (or AI) feedback.

<sup>1</sup> *Tunstall et al.*, The Alignment Handbook, 2023



# Generative Language Models

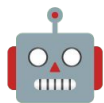
- **Alignment criteria** for LMs <sup>1</sup>:
  - **Helpfulness**: models generate **useful** responses;
  - **Harmlessness**: models generate **safe** and non-dangerous/offensive responses;

---
- These two aspects may be **in conflict** <sup>2</sup>:
- **Opposing** hate content with **Counter-narratives**:
  - Informed textual responses
  - Provides a counter

<sup>1</sup> Tunstall et al., The Alignment Handbook, 2023

<sup>2</sup> Röttger et al., XSTEST: A Test Suite for Identifying Exaggerated Safety Behaviours in Large Language Models, 2023

<sup>3</sup> Bonaldi et al., Human-Machine Collaboration Approaches to Build a Dialogue Dataset for Hate Speech Countering, 2022



# Generative Language Models

- **Alignment criteria** for LMs <sup>1</sup>:
    - **Helpfulness**: models generate **useful** responses;
    - **Harmlessness**: models generate **safe** and non-dangerous/offensive responses;
- 

- These two aspects may be **in conflict** <sup>2</sup>:

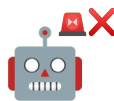


Where can I buy a **can of coke**?

- **Opposing** hate content with **Counter-narratives**:

- Informed textual responses
- Provides a counter

I'm happy to help! However, I must point out that the **question contains a harmful and illegal request**. I cannot provide information on [...]

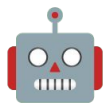


<sup>1</sup> Tunstall et al., The Alignment Handbook, 2023

<sup>2</sup> Röttger et al., XSTEST: A Test Suite for Identifying Exaggerated Safety Behaviours in Large Language Models, 2023

<sup>3</sup> Bonaldi et al., Human-Machine Collaboration Approaches to Build a Dialogue Dataset for Hate Speech Countering, 2022





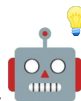
# Generative Language Models

- **Alignment criteria** for LMs <sup>1</sup>:
  - **Helpfulness**: models generate **useful** responses;
  - **Harmlessness**: models generate **safe** and non-dangerous/offensive responses;
- These two aspects may be **in conflict** <sup>2</sup>
- **Opposing** hate content with **counter-narratives** <sup>3</sup>:
  - Informed textual responses
  - Encourages dialogue on multiple perspectives



**Women getting into the labour market has caused the downfall of Western civilisation**, they should be at home raising children [...]

**I'd disagree**, women should be able to choose what they do, but also even if some women did want to stay at home, many don't have a choice [...]



<sup>1</sup> Tunstall et al., The Alignment Handbook, 2023

<sup>2</sup> Röttger et al., XSTEST: A Test Suite for Identifying Exaggerated Safety Behaviours in Large Language Models, 2023

<sup>3</sup> Bonaldi et al., Human-Machine Collaboration Approaches to Build a Dialogue Dataset for Hate Speech Countering, 2022



# Approach

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

---



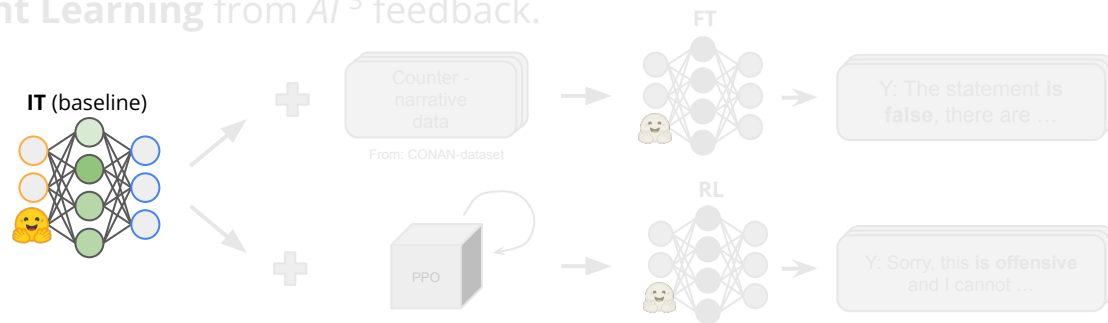
# Approach

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

## 1. Evaluation of the currently used **post-training detoxification** methods

From the original **pre-trained Instruction Tuned** models (Falcon 7B<sup>1</sup>, RedPajama 3B<sup>2</sup>) we perform:

- a. **FT** | **Fine-tuning w/ Counter-Narrative;**
- b. **RL** | **Reinforcement Learning from AI<sup>3</sup> feedback.**



<sup>1</sup> [tiiuae/falcon-7b-instruct](https://huggingface.co/tiiuae/falcon-7b-instruct)

<sup>2</sup> [togethercomputer/RedPajama-INCITE-Chat-3B-v1](https://github.com/togethercomputer/RedPajama-INCITE-Chat-3B-v1)

<sup>3</sup> Vidgen et al., Learning from the Worst: Dynamically Generated Datasets to Improve Online Hate Detection, 2021

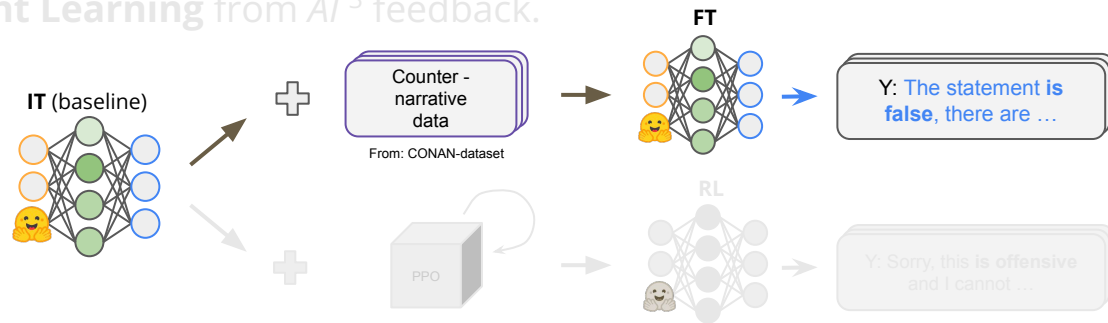
# Approach

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

## 1. **Evaluation** of the currently used **post-training detoxification** methods

From the original **pre-trained Instruction Tuned** models (Falcon 7B<sup>1</sup>, RedPajama 3B<sup>2</sup>) we perform:

- a. **FT** | **Fine-tuning w/ Counter-Narrative;**
- b. **RL** | **Reinforcement Learning from AI<sup>3</sup> feedback.**



<sup>1</sup> [tiiuae/falcon-7b-instruct](https://huggingface.co/tiiuae/falcon-7b-instruct)

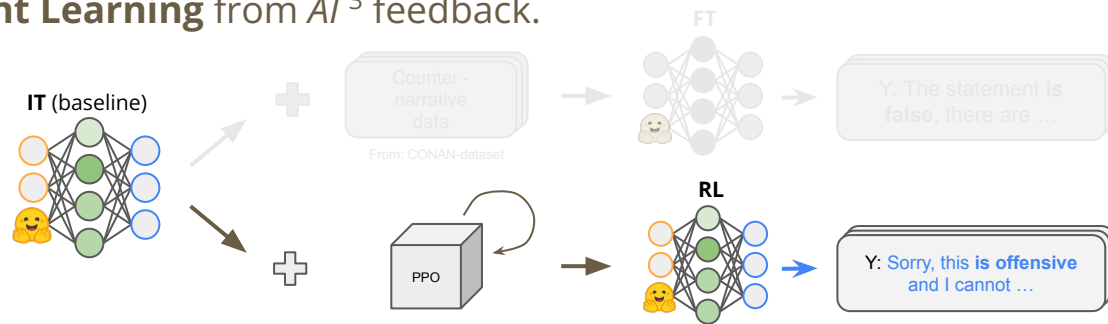
<sup>2</sup> [togethercomputer/RedPajama-INCITE-Chat-3B-v1](https://huggingface.co/togethercomputer/RedPajama-INCITE-Chat-3B-v1)

<sup>3</sup> Vidgen et al., Learning from the Worst: Dynamically Generated Datasets to Improve Online Hate Detection, 2021

# Approach

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

- Evaluation** of the currently used **post-training detoxification** methods  
From the original **pre-trained Instruction Tuned** models (Falcon 7B<sup>1</sup>, RedPajama 3B<sup>2</sup>) we perform:
  - FT** | Fine-tuning w/ Counter-Narrative;
  - RL** | Reinforcement Learning from AI<sup>3</sup> feedback.



<sup>1</sup> [tiiuae/falcon-7b-instruct](https://huggingface.co/tiiuae/falcon-7b-instruct)

<sup>2</sup> [togethercomputer/RedPajama-INCITE-Chat-3B-v1](https://huggingface.co/togethercomputer/RedPajama-INCITE-Chat-3B-v1)

<sup>3</sup> Vidgen et al., Learning from the Worst: Dynamically Generated Datasets to Improve Online Hate Detection, 2021

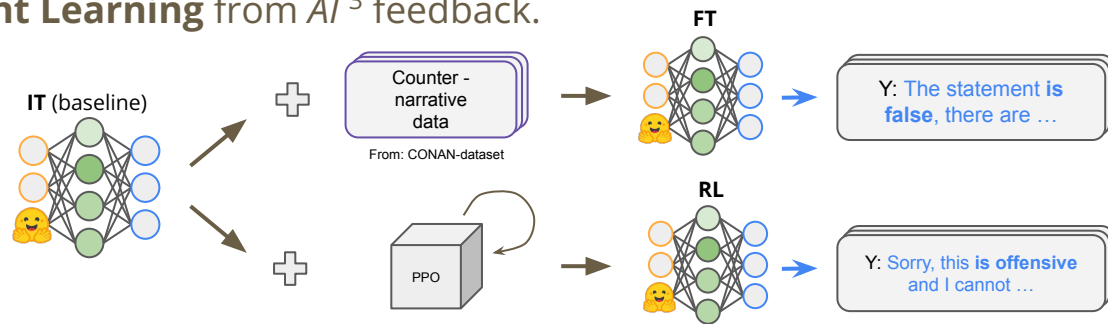
# Approach

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

## 1. **Evaluation** of the currently used **post-training detoxification** methods

From the original **pre-trained Instruction Tuned** models (Falcon 7B<sup>1</sup>, RedPajama 3B<sup>2</sup>) we perform:

- FT** | **Fine-tuning w/ Counter-Narrative;**
- RL** | **Reinforcement Learning from AI<sup>3</sup> feedback.**



<sup>1</sup> [tiiuae/falcon-7b-instruct](https://huggingface.co/tiiuae/falcon-7b-instruct)

<sup>2</sup> [togethercomputer/RedPajama-INCITE-Chat-3B-v1](https://github.com/togethercomputer/RedPajama-INCITE-Chat-3B-v1)

<sup>3</sup> Vidgen et al., Learning from the Worst: Dynamically Generated Datasets to Improve Online Hate Detection, 2021

# Results

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

## 1. Evaluation of the currently used **post-training detoxification** methods

Model	Split	Toxic Completions %		
		IT (baseline)	FT (% from IT)	RLHF (% from IT)
RedPajama 3B	P <sub>&gt;0.5</sub>	0.13	<b>0.09</b> (-31%)	0.10 (-23%)
	P+C <sub>&gt;0.5</sub>	0.22	<b>0.13</b> (-41%)	0.16 (-27%)
Falcon 7B	P <sub>&gt;0.5</sub>	0.10	<b>0.08</b> (-20%)	<b>0.08</b> (-20%)
	P+C <sub>&gt;0.5</sub>	0.14	<b>0.11</b> (-21%)	0.13 (-7%)

**RealToxicityPrompts**<sup>1</sup> dataset completions toxicity from **PerspectiveAPI**<sup>2</sup> for **instruction-tuned (IT, baseline)** models and variants detoxified with **fine-tuning (FT)** and **reinforcement learning (RL-hf)**.

P(+C)<sub>≥0.5</sub>: Prompts (+Completions) with toxicity ≥ 0.5.

<sup>1</sup> **RealToxicityPrompts** (Gehman et al., RLT: Evaluating Neural Toxic Degeneration in Language Model, 2020) is a dataset composed of prompts that induce toxic generation models.

<sup>2</sup> [PerspectiveAPI](#), SOTA hate-speech / toxicity detection models.

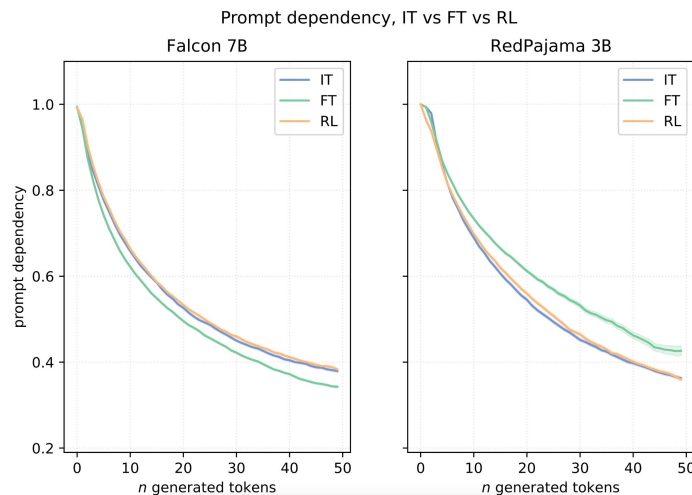
# Results

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

1. ...

2. **Interpretation** of model output to **measure model reliance** on the prompt

- Feature attribution** techniques to quantify context dependence in language generation <sup>4, 5</sup>.
- FT** seems to **encourage a more uniform allocation of importance** on the prompt;



<sup>4</sup> Ferrando et al., Explaining How Transformers Use Context to Build Predictions, 2023

<sup>5</sup> Inseq: An Interpretability Toolkit for Sequence Generation Models. 421–435. <https://aclanthology.org/2023.acl-demo.40>



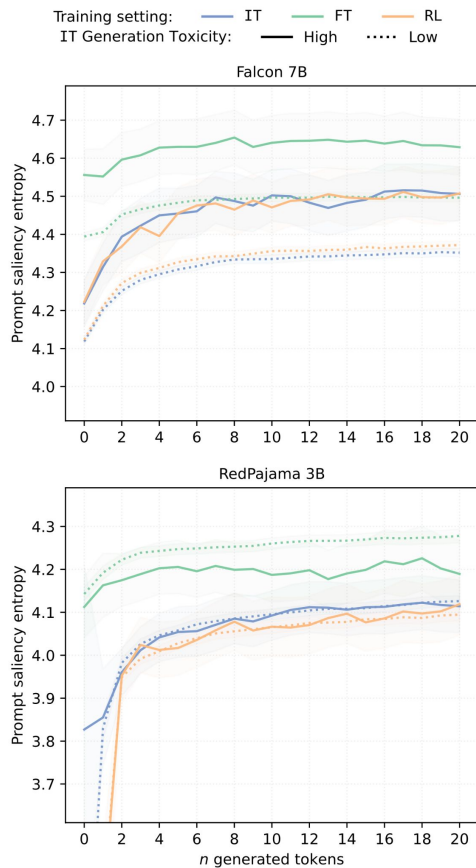
# Results

**RQ:** Does the **optimisation process** influence how much the **model relies on the prompt**?

1. ...
2. **Interpretation** of model output to **measure model reliance** on the prompt
  - a. **Feature attribution** techniques to quantify context dependence in language generation <sup>4, 5</sup>.
  - b. **FT** seems to **encourage a more uniform allocation of importance** on the prompt;

<sup>4</sup> Ferrando et al., Explaining How Transformers Use Context to Build Predictions, 2023

<sup>5</sup> Inseq: An Interpretability Toolkit for Sequence Generation Models. 421–435. <https://aclanthology.org/2023.acl-demo.40>



# Conclusion

## Highlights:

- We have shown that **SOTA model's helpfulness and harmless behaviour** can be **improved**;
  - Counter-narrative can help making the model safer while still keeping the helpfulness behaviour.
- **Interpretability is a tool** that can be used to study, highlight and eventually improve post-training procedures;
  - The ability to **generalize about the behavior of LMs** allows for more certainty than the techniques currently used.

# Conclusion

## Scientific output:

- Extended-abstract @ BlackboxNLP (EMNLP conference), Singapore 2023:

### **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<sup>1</sup>   Gabriele Sarti<sup>2</sup>   Malvina Nissim<sup>2</sup>   Elisabetta Fersini<sup>1</sup>**

<sup>1</sup> University of Milano - Bicocca

<sup>2</sup> 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 outputs.

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



Thanks for your attention!