



Security practices in AI development

Petr Spelda¹ · Vit Stritecky¹

Received: 9 August 2024 / Accepted: 14 February 2025
© The Author(s) 2025

Abstract

What makes safety claims about general purpose AI systems such as large language models trustworthy? We show that rather than the capabilities of security tools such as alignment and red teaming procedures, it is security practices based on these tools that contributed to reconfiguring the image of AI safety and made the claims acceptable. After showing what causes the gap between the capabilities of security tools and the desired safety guarantees, we critically investigate how AI security practices attempt to fill the gap and identify several shortcomings in diversity and participation. We found that these security practices are part of securitization processes aiming to support (commercial) development of general purpose AI systems whose trustworthiness can only be imperfectly tested instead of guaranteed. We conclude by offering several improvements to the current AI security practices.

Keywords AI safety · LLM · Alignment · Security practices

1 Introduction

Safety of foundation AI models (Bommasani 2021), such as large language models capable of in-context learning (Brown et al. 2020), is based on aligning a pretrained model with human preferences on the model's responses to user inputs and then testing whether alignment with the preferences holds in as wide a range of situations as possible. The pair of an alignment method, e.g., reinforcement learning from human feedback generalizing preferences through reward modeling (Ouyang et al. 2022, Bai et al. 2022a), and a testing method, e.g., systematic adversarial testing called red teaming attempting to find gaps left behind by the alignment method (Ganguli et al. 2022; Casper et al. 2023a), represents a tool allowing the owner of the model to enact a security practice.

With every new release of commercial chatbots based on LLMs (large language model)¹, we are witnessing yet another round of enacting some security practice that aims to manage the risk in the eyes of the public and regulators alike. The practice itself enables and is accompanied by narratives

of safe product development that protects consumers, society and democracy from harmful effects of AI. There is, however, a gap between the capabilities of the tools (the pairs of alignment and testing methods) and safety guarantees which are promised by enacting security practices with the tools.

We will show that since it is hard to maintain the promised guarantees with available tools, security practices that aim to fill the gap follow the logic of Didier Bigo's (in) securitization processes (Bigo and Tsoukala 2008). The (in) securitization process captures routines of risk management enacted by actors aiming to create a space in which they can make claims about (in)security (ibid.). In the sense of (in) securitization processes, AI security practices enacted by alignment and testing methods are no different than the practices designed to sustain (in)security claims about migration, health or terrorism which originally motivated the critical approaches to the study of security (e.g., C.A.S.E. Collective 2006). We argue that safety claims about LLMs and other types of foundation models are impossible to evaluate without understanding the gap between capabilities of security tools and the desired security guarantees. The gap is to be filled with enacting security practices of risk management consisting of routines designed to manage unease about or even fear of for-profit, privately-owned and closed development of capable AI models.

✉ Petr Spelda
petr.spelda@fsv.cuni.cz

¹ Department of Security Studies, Institute of Political Studies, Faculty of Social Sciences, Charles University, U Kříže 8 Praha 5, 158 00 Prague, Czech Republic

¹ This applies to a degree also to publicly available LLMs.

In order to explain the motivation and critically assess the risks associated with security practices in LLMs, or in general foundation models, development, we proceed as follows. In Sect. 2, it will be explained what causes the gap between security tools and the desired guarantees and, thus, motivates the security practices. In Sect. 3, we will show that the security practices supporting closed, for-profit development of LLMs can cause a lack of participation, accountability, transparency and democratic legitimacy and above all cannot guarantee that the gap between capabilities of the tools and the desired guarantees will be closed. We will conclude with Sect. 4 on more open, participatory, and sustainable practices in LLM development.

2 AI alignment cannot guarantee AI safety

The main method for aligning capabilities of LLMs with human preferences used today is reinforcement learning from human feedback (RLHF; Ouyang et al. 2022; Bai et al. 2022a). The core of RLHF is reward modeling, a method that learns from human preferences over alternative outputs of the model that is being aligned. The aim of learning a reward model from human preferences is to generalize from them and be able to provide correct rewards for outputs of the model that were not included in the preference dataset.

Reward modeling was developed to make training of deep reinforcement learning agents in simulated environments (e.g., games, physics simulators) more efficient with respect to human feedback (Christiano et al. 2017). The reward model captured human preferences on pairs of trajectories of the agent in the environment and learned to generalize from them so that for that particular environment the job of human judges could be automated (ibid.). The same principle has been applied for aligning pretrained LLMs with human preferences on LLM outputs to sensitive inputs. Sensitive inputs include prompts for which pretrained, unaligned models are known to generate biased and inequitable responses (Tamkin et al. 2023) or a wide range of unsafe responses (Mazeika et al. 2024).

Experiments such as OpenAI's InstructGPT (Ouyang et al. 2022) utilized reward modeling based on human preferences for producing models capable of following natural language instructions in a better way compared to, for instance, the GPT-3 model that was trained using the next token prediction objective and was not explicitly aligned for instruction following. Schematically, the procedure can be described in four steps (ibid.):

1. For a set of prompts, people write responses that reflect their preferences (alternatively, suitable existing responses are collected).
2. The set of prompt-response pairs is used for supervised fine-tuning of some base model pre-trained using the next token prediction objective, providing a foundation for an instruction following model IFM.
3. People provide preferences for training the reward model. For each prompt, a ranking of alternative responses to the prompt, $A \succ C \succ B \sim D$ (last two are ranked as equal), is produced by human rankers, and the resulting set of binary preferences over alternatives (e.g., $A \succ D$ pairwise ranking) is used to train a reward model that generalizes from the preferences.
4. For a new prompt, the IFM generates a response evaluated by the reward model which produces a reward for it and the reward is used to update the IFM using a reinforcement learning method.

Often, for cost-effectiveness reasons, human rankers are replaced with existing reward models or generative models prompted to select a response from two options, which can replace reward models (see Lambert et al. 2024a, p. 17). While far from optimal considering issues such as participation or inclusion that we discuss in Sect. 4, automated preference generation allowed scaling-up LLM development (relatedly, see also Bai et al. 2022b). The inherent tradeoff should, however, not be overlooked. There is also a class of direct alignment algorithms that do not rely on explicit reward models, most notably Direct Preference Optimization (DPO, Rafailov et al. 2023). DPO does not involve training a reward model from human preferences or reinforcement learning. Rather, DPO calculates the probability of preferred and dispreferred response to a prompt under the model and optimizes its parameters to increase the probability of preferred over dispreferred responses (ibid.). There is evidence (e.g., Lambert et al. 2024b) that DPO and its variants can offer strong performance at the post-training alignment stage. Post-training alignment methods can be, of course, combined. See, for instance, how Lambert et al.'s (2024b) recipe uses DPO for preference tuning and reinforcement learning with verifiable rewards, a variant of RLHF replacing the reward model with verifiers, to strengthen the model's mathematical capabilities. For the class of generative models that produce chain-of-thought traces before outputting a response, there are methods such as Deliberative Alignment (Guan et al. 2024), trying to ensure that the model's chain-of-thought trace is informed by a safety policy relevant for the user's prompt. The method uses a reward model with access to safety policies to filter high-quality completions for supervised fine-tuning and for the final, reinforcement learning stage (ibid.).

It is important to note that the RLHF-based alignment procedure described in steps 1–4 is a domain adaptation rather than safety-guaranteeing method. The method can be used to adapt a base model for instruction following.

Interpreting this as a safety procedure requires that the human preferences expressed in steps 1 and 3 allow the reinforcement learning stage to remove all unsafe capabilities of the instruction following model inherited from the base model.

This is clearly impossible to guarantee for several reasons. First, even if guidelines for providing feedback are based on an equitable and inclusive conception of safety, people are bound to provide incomplete feedback and biased, inconsistent or confused preferences accidentally, out of malice or due to a perceived difficulty of the comparison task (cf. Casper et al. 2023b). Second, it is difficult, if not impossible, to ensure that the reward model will not misgeneralize the preferences during training and will be able to serve as a good guide for updating the IFM during its alignment (ibid.). Finally, given the considerable generality of LLMs trained on the biggest piles of data ever assembled, it is impossible to cover every harmful capability.

Testing the model to uncover its unsafe capabilities, which need to be removed, cannot secure full coverage either. The impossibility can be explained by building an analogy with software testing. The main obstacle to verifying by testing that a software system possesses some properties according to a specification is what Symons and Horner (2019) call the ‘path-complexity catastrophe’. Their argument relies on an observation that even a simple program logic induces exponential scaling of execution paths, which cannot be all tested for adherence with the specification due to computational intractability of such an effort (Symons and Horner 2019, Section 2.2; tests parallelization cannot solve the path-complexity catastrophe either, ibid.). If then, the software’s overall reliability is inferred from a statistic calculated from the executed tests, the inference becomes an inductive prediction about the system’s reliability. And because we need to make assumptions about the system’s overall error distribution (which will remain unknown due to the path-complexity catastrophe, ibid.) to perform that prediction, it cannot provide unproblematically admissible evidence about the system’s overall reliability.

By the same token, can we predict that an LLM will provide aligned completions to user tasks on all future occasions from a limited sample of its responses obtained from performed tests? It is possible to try, but the prediction will not lead to verifiable guarantees with respect to some alignment specification captured as human preferences or other constraints on correct answers. The problem is worse for poisoned (Rando and Tramèr 2024; Hubinger et al. 2024) or deceptively aligned models (see early toy experiments, e.g., Meinke et al. 2024) that will breach their alignment regardless of pre-deployment testing known as red teaming.

Harmful capabilities are found by red teaming (Feffer et al. 2024; Inie et al. 2023; Casper et al. 2023a). Red teaming consists of prompting techniques that aim to elicit

harmful capabilities of the model and help to find the unavoidable ‘gaps’ left behind by the alignment procedure. In its core, red teaming is just an empirical testing practice with adversarial goals that can be performed by people as well as by other LLMs tasked with generating adversarial prompts (Perez et al. 2022; Samvelyan et al. 2024)². Identification of harmful capabilities is followed by repeating the alignment procedure on prompts identified by red teaming. It is crucial to note that red teaming does not refer to a well-defined practice (Feffer et al. 2024) but rather to any testing procedure that can lead to discovering harmful capabilities of supposedly safe models. Such practices can be highly domain specific and context-dependent, which makes them open to issues of participation and fairness (Fazelpour et al. 2024) like any other social practice with safety goals. The alignment-red teaming cycle can be repeated several times before and after the model is deployed, e.g., turned into a product accessible via an API or released as a publicly available model (even if here the usual practice is pre-deployment alignment-red teaming, e.g., Touvron 2023).

Unfortunately, even several repetitions of the alignment-red teaming cycle cannot provide dependable safety guarantees. First, aligned and tested LLMs are known to be vulnerable to jailbreak attacks which force the models to respond in ways that violate their alignment (Wei et al. 2023). These attacks use specialized prompts and are hypothesized to exploit the failures of alignment generalization to unsafe capabilities that were not discovered during red teaming or exploit conflicts that can arise between general instruction following and responding safely to restricted requests (ibid.). Second, Zou et al. (2023) developed an optimization-based method for finding ‘adversarial suffixes’ which, when appended to a wide range of restricted requests, break the LLM’s alignment and force it to suppress the refusal response and return unsafe answers. It was observed that the attack suffixes are effective against a range of different LLMs and considering the fact that the search for them is automated, they substantially reduce dependability of safety guarantees that can be derived from alignment. The safety guarantees cannot be fundamentally improved by combining alignment with filtering LLM outputs according to some safety rules. Glukhov et al. (2023) showed that the filtering problem is formally undecidable and cannot be used to fix the non-robustness of LLM alignment. Despite this ultimate formal obstacle, there is some evidence that input and output classifiers trained on synthetic data generated by LLMs according to a safety specification can make jailbreaking of LLMs harder (Sharma et al. 2025).

² Classifiers can be used to determine whether the output is harmful to allow scaling up the process.

The problem is worse for models that can be further fine-tuned after alignment-red teaming was applied, e.g., for publicly available or open-source base models or commercial models offering fine-tuning APIs. Qi et al. (2024a) showed that less than 100 samples is enough for breaking alignment of a model and recovering its harmful capabilities. The robustness of alignment can be also decreased by fine-tuning an aligned and red-teamed model on benign prompt–response pairs (ibid.). Since there is currently a rich offering of publicly available or open-source high-quality, instruction-following models whose fine-tuning is cheap and can be done on high-performance consumer hardware, alignment and red teaming provides even less safety assurances than for commercial chatbots based on LLMs. In this situation, unlearning harmful knowledge from the model combined with increasing the hardness of re-learning it is considered one of possible solutions. However, it is challenging to perform unlearning in a robust way, preventing re-learning via fine-tuning, while not degrading safe capabilities of the model (cf., Barez et al. 2025). This becomes even more problematic if dual-use knowledge is involved, which has safe, legitimate uses but can also lead to model completions that pose risks in cybersecurity or CBRN (chemical, biological, radiological and nuclear) domains (ibid.).

In sum, alignment methods and red teaming practices used today cannot deliver dependable safety and fairness guarantees because gaps in LLMs' alignment are unavoidable and their alignment can be broken or compromised in multiple ways depending on access to the model. This shows that the current alignment methods based on preference modeling are good for product development, e.g., chatbots based on instruction following LLMs, and bad for developing AI products with safety guarantees. Since policy-makers demand safety guarantees (European Union AI Act³, Biden's U.S. Presidential Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence⁴, rescinded on January 20, 2025) to protect the public from harmful effects of LLMs whose alignment is incomplete and non-robust, companies developing these models attempt to manage the (in)security and unease over AI development with security practices.

3 AI alignment and red teaming as tools of security practices

We established that the gap between the tools and the desired safety guarantees cannot be currently closed, and this means that claims about safety of LLMs by the companies that develop them need to be understood as part of a securitization process. What is most visible of this process in the public space are speech acts that aim to communicate preparedness for responsible development of more capable AI systems (Anthropic 2023a; 2023b; OpenAI 2023a; OpenAI 2023b; Leike and Sutskever 2023). The securitization process follows the logic of tacitly acknowledging the shortcomings of available alignment and safety testing methods and claiming that the issues will be addressed by broad research programs whose announcements constitute the core of the speech acts. The level of detail ranges from vague and general (Leike and Sutskever 2023) to partially concrete by pointing toward existing AI safety research (Anthropic 2023a). It is important to understand that these high-level speech acts do not exist in a vacuum. They are a product of security practices which are performed using AI alignment and red teaming methods. Since these AI safety tools are limited, it is necessary to critically assess the practices that they help to enact. This critical assessment is the first step in understanding high-level speech acts regarding AI safety that are hard to parse if kept separated from the underlying security practices.

Bigo and Tsoukala (2008) initiated the study of (in)securitization processes through the lens of security practices and Balzacq et al. (2010) provided a useful framework for understanding the tools of security practices. Despite the fact that Balzacq et al.'s (2010) tools' characteristics were derived from counterterrorism and other areas, the characteristics are fitting AI safety tools well because the enactment of practices using imperfect tools always aims to fill the gap between capabilities and the desired safety guarantees.

First, security tools have design traits and defining features that make them unique yet relatable to other tools (ibid.). As explained in Sect. 2, the core of RLHF, a widely used AI alignment method today, is a reward model based on human preferences. The reward model is the design trait of RLHF considered as a security tool and human preferences are its defining features. The first critical point regarding security practices enacted with RLHF is about defining features of reward models.

There is evidence showing that human preferences used to train reward models represent only some social groups (e.g., Santurkar et al. 2023), which leads to reward models that marginalize already underrepresented communities (Ryan et al. 2024). Apart from leading to secondary fairness issues RLHF tried to address in the first place, findings like

³ <https://eur-lex.europa.eu/eli/reg/2024/1689/oj>

⁴ <https://www.federalregister.gov/documents/2023/11/01/2023-24283/safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence>

this point toward an insufficient level of participation and democratic decision-making resulting in reward models that fail to capture the diversity of preferences in heterogeneous populations. LLMs will not become pluralistic and safe in the social sense (Sorensen et al. 2024) unless it is clear what happens if more than one ranking of alternative responses to a prompt exists. In other words, what is done to preserve diversity of preferences with respect to all stakeholders and what type of diversity preserving preference aggregation is used? We do not have satisfactory answers to these questions for commercial models apart from reports on experiments with democratic/collective alignment (e.g., Anthropic 2023c and Eloundou and Lee 2024). Obtaining preference profiles in an accountable, transparent, and responsible manner is difficult and costly, and this is why we see interfaces of commercial chatbots designed in a way that allows real-time collection of customer preferences. Without guarantees on participation and diversity, security practices enacted by RLHF contribute to product development and support of speech acts rather than to closing the gap between the capabilities of the security tool and the desired level of safety guarantees.

Apart from the preference-related defining features, RLHF has also potentially problematic design traits. Lambert et al. (2023) highlighted that since the reward model is based on a pre-trained LLM, issues of the base model could negatively impact the reward model trained on top of it (similar issues impact generative reward modelling where an LLM is used as a judge to determine the ranking of alternatives, Lambert et al. 2024a, p. 17). Krendl Gilbert et al. (2023) suggested that it would be beneficial to maintain reward reports continually updated with information on the model's implementation, its social interfaces, reinforcement learning objectives and the auditing and verification results to simplify keeping track of what we consider design traits and defining features of deployed systems.

Overall, security practices enacted with RLHF show a significant room for improvement mostly with respect to participation and diversity, without which it is difficult to speak about safety guarantees derived from AI alignment.

As the second characteristic of security tools Balzacq et al. (2010) identified actions which are configured by the tools and consist of requirements, procedures and delivery mechanisms. In Sect. 2, it was explained that gaps in alignment left behind by security tools such as RLHF are expected to be found by diverse empirical testing procedures called red teaming. Red teaming practices correspond to actions whose content is 'configured' against the defining features of a particular RLHF run to test the robustness of alignment of the resulting LLM with human preferences.

As shown by Feffer et al. (2024), the adjustability of red teaming practices, which starts with diverse threat models derived from different safety preferences, can render red

teaming a mere security theater that is invoked whenever there are concerns over the risks of generative AI on society. The security theater interpretation of red teaming fits well the logic of AI safety securitization, which consists of acknowledging concerns and making speech acts showing that there are actions already underway capable of addressing these concerns (Anthropic 2023d; OpenAI 2023c).

Feffer et al.'s (2024) findings synthesized from existing research on red teaming can be used to critically assess it in terms of requirements, procedures and delivery mechanisms, which constitute actions as the second characteristic of security tools. Since red teaming is a set of evaluation procedures, its results heavily depend on people performing them. In an ideal world, it would be required that selected experts work together with stakeholders to maintain diversity of evaluation by wide participation. In reality, experts are often accompanied by testers from crowd-working platforms, and for them as well as for the experts the question of fair representation and participation remains open (ibid.). Inie et al. (2023) mapped a complex landscape of motivations and goals adopted by people participating in red teaming of LLMs. Diversity of evaluation teams cannot be understated because the ability to find vulnerabilities of the model that unevenly impact different social groups depends on participation. LLM-based red teaming is limited by the number of generated adversarial prompts and the computational budget for the subsequent adversarial training on red teaming data aiming to patch the gaps in the model's alignment (ibid.).

When it comes to testing LLMs on harmful prompts, that is, performing the core of red teaming, Feffer et al. (2024) found evidence of procedures which are neither well-scoped nor -structured. The major contributing factor is the difficulty of meeting participation requirements where instead of diverse stakeholders, crowd-workers test the 'lowest hanging' vulnerabilities and experts represent mostly academe and AI companies (ibid.). It is difficult to comprehensively test versatile LLMs, for which the risk surface can be large and complex with red teams that are biased to focus only on some of its areas. This means that even after using one of the alignment techniques such as RLHF to obtain a supposedly safe LLM, testing procedures are not guaranteed to find all gaps that can be misused. Repeating the process does not necessarily increase the test coverage if the approach to assembling the red team does not change.

Finally, delivery mechanisms, as the last stage of actions configured by security tools, are not without problems either. Despite the existence of open red teaming experiments (Ganguli et al. 2022), detailed reports on the identified vulnerabilities and their mitigation are

not always available and there is a lack of standards on the reports' structure, types of disclosed information, and evaluation of the mitigation results (Feffer et al. 2024)⁵. The lack of information complicates the role of stakeholders who cannot submit requests for participation on the basis of ill-conceived or missing tests designed by unrepresentative red teams.

Overall, red teaming procedures following RLHF as actions configured by security tools cannot identify every safety issue left behind by the alignment procedure. The reason for this lies in the fact that to obtain a full test converge of versatile LLMs is close to impossible and also in the fact that the 'configuration' of the action and of its individual stages is often flawed in the same way as the tool itself by lacking in diversity and participation. In this sense, proposals for safe harbors allowing unimpeded safety research of commercial LLMs (Longpre et al. 2024) could become an important tool that not only protects researchers but also increases participation and diversity in AI safety research and contributes to better understanding of the security tools. Generally, a lack of information on the audited AI system always leads to an incomplete picture of the system's risk surface (Casper et al. 2024).

We join the discussion on two remaining Balzacq et al.'s (2010) characteristics of security tools together because public actions and images invoked by security practices enacted with particular tools are closely related. Public actions and images are the end point of (in)securitization processes. In AI safety, security tools such as RLHF, their design traits and defining features, and actions such as red teaming 'configured' by them aim to inspire AI governance (public action) that supports the image of safe and socially responsible innovation in sync with public policies (Anthropic 2023e; OpenAI 2023d).

It is important to note that this process runs in parallel with regulatory efforts that try to ensure society will not be harmed by the gap between the capabilities of AI security tools and the desired level of safety guarantees. Since (in)securitization processes aim to shape public actions and images, it is reasonable to ask whether AI regulation is impacted by security practices of AI companies that develop commercial, versatile LLMs. Red teaming was, for instance, featured in Biden's U.S. Presidential Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence⁶ (rescinded on January 20, 2025) as a vehicle for identifying harmful capabilities of foundation

models such as versatile LLMs that need to be mitigated, the results of testing should be then shared with the government. Prior to the Executive Order, there was a public-private partnership between major U.S. companies involved in AI development and the Biden-Harris administration⁷. When looking at the companies' interpretation of the partnership (e.g., OpenAI 2023d), red teaming was the top priority and it turned out to be an important part of the Executive Order as well.

We might ask whether the red teaming's role results from a successful public-private partnership on AI safety or from a securitization process that managed to produce a high-level speech act supported by the security practice. According to Balzacq et al. (2010), tools used to enact security practices reconfigure public action with respect to the image of the threat. We have shown that alignment methods such as RLHF are useful for product development and their use for AI safety does not bring safety guarantees even if accompanied by extensive testing. From a critical perspective, this means that the practice combining alignment procedures with red teaming reconfigured the image of AI safety and paved the way for the public action, e.g., Biden's Executive Order on AI, that accepted LLMs' alignment and red teaming to be sufficient for AI safety.

4 From critique to more open, participatory, and sustainable practices in LLM development

In the situation where we lack alternatives capable of delivering genuinely dependable safety guarantees on AI safety, it is perhaps understandable that the securitization process concluded by 'certifying' the product development as a safe practice. Amidst concerns for AI risks of capable systems (Bengio et al. 2023), until (if ever) we have tools capable of deriving dependable AI safety guarantees, more emphasis should be placed on participation and diversity in AI security practices to make sure that they are at least equitable.

One way of achieving this is to consider a different 'baseline'. Alongside large, general-purpose, closed or publicly available, alternatively open source (e.g., Groeneveld et al. 2024), LLMs, multiple small, domain-specific and specialized LLMs could be trained, constrained to perform only a range of domain-specific tasks. There are several benefits to this alternative approach. First, the alignment and red teaming procedures would become more manageable due

⁵ Feffer et al. (2024) compiled an instructive set of questions pertaining to individual stages of red teaming that could help to better 'configure' the action to fit the aims of the security practice at hand.

⁶ <https://www.federalregister.gov/documents/2023/11/01/2023-24283/safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence>

⁷ <https://bidenwhitehouse.archives.gov/briefing-room/statements-releases/2023/07/21/fact-sheet-biden-harris-administration-secures-voluntary-commitments-from-leading-artificial-intelligence-companies-to-manage-the-risks-posed-by-ai/>

to a clear and narrow range of tasks that need to be aligned and tested. Second, making the security practice more manageable would also allow developers to engage key stakeholders more effectively, hopefully increasing participation and diversity in the security practices and leading to more accountable sociotechnical AI systems that many call for (e.g., Fazelpour and De-Arteaga 2022). Third, regulatory oversight over specialized, domain-specific LLMs could be performed by government departments or agencies tasked with regulating the given area. This would avoid building parallel (and redundant) capabilities at an AI governance institution specializing in oversight of versatile LLMs.

Concerning participation and diversity in AI security practices, Rauh et al. (2024) identified several problems, among which the model-centered evaluation in combination with automated benchmarks has a potential to obscure social harms that happen outside this narrow context. Relatedly, these evaluations should never be considered value-neutral (ibid.). There are sociotechnical methodologies, such as STAR (Weidinger et al. 2024), aiming to assess impacts on broader social contexts. These contexts are engaged by red teaming instructions that, if the model provided completions, would result in significant safety and social justice violations (ibid.). If methodologies like STAR were added to the LLM development toolbox, increased participation and diversity would become a necessary requirement for red teaming efforts. Without participation and diversity, it would not be possible to determine how to test the model for harmful effects in social contexts that are impossible to understand assuming only the model-centered and supposedly ‘value-neutral’ perspectives. We will now provide three examples of small, domain-specific models and compare them to large, versatile alternatives.

Currently, it seems reasonable to consider models up to 8 billion parameters (8B) small because they can be used for inference on affordable hardware. For example, Llama-3.1-8B-Instruct, an open-weight (publicly available), instruction-following model from Meta (Grattafiori et al. 2024), can run with decent context lengths even on consumer accelerators (Llama-3.1-8B-Instruct with weights in the half precision floating point format, FP16, and a keys and values cache for the context length of 16 thousand tokens will fit the memory of a high-end consumer GPU, e.g., Nvidia GeForce RTX 4090 with 24 GB of memory, see, e.g., Schmid et al. 2024), though at reduced token throughputs compared to enterprise accelerators. There are many models in this category, some notable examples of open-weight models include: Meta’s Llama-3.2-1B & 3B and their instruction-tuned variants

developed for edge and mobile devices⁸, Microsoft’s Phi family of models⁹ (although some models from the family have more than 8B parameters, Phi is considered to be a family of small language models), the Gemma family of small models from Google¹⁰, small models from Alibaba’s Qwen family¹¹ (0.5B, 1.5B, 3B, 7B), OpenBMB’s MiniCPM family¹² or fully open-source 7B variants of OLMo 2 from Ai2¹³.

Llama and Phi models were used to develop small, domain-specific models that were fine-tuned for specific task sets and have, arguably, a more manageable safety profile compared to large models that can be used to perform the tasks as well thanks to their more versatile in-context learning capabilities. LinkedIn developed the EON-8B model based on Llama-3.1-8B-Instruct to suggest matches between candidates and jobs and to generate explanations of the matchings (Bodigutla et al. 2024). Bodigutla et al. (2024) report better or comparable performance to large general models such as GPT-4o or Llama-3-70B and better cost-effectiveness. Alignment and red teaming of an 8B domain-adapted model is still challenging but more manageable compared to large, versatile models. This could be a crucial difference especially in sensitive areas like hiring. Li et al. (2024) domain-adapted several small language models, including Phi-3.5-mini 3.8B (Abdin et al. 2024), for translating natural language tasks to domain-specific code allowing interaction with the hardware fulfilment process that supplies Azure, Microsoft’s cloud computing platform. They found that after fine-tuning Phi-3.5-mini ranked best among small models and outperformed large, versatile models like GPT-4-turbo prompted with samples used for adapting the small models (ibid.). Similar to the previous case, the small, domain-adapted Phi-3.5-mini was more cost-effective compared to large versatile models (Li et al. 2024, Table 2) and its alignment arguably more manageable than alignment of a large, general model accessed via in-context learning. Finally, Liu et al. (2023) domain-adapted two versions of Meta’s second-generation Llama, 7B and 13B, for chip design (engineering assistant chatbot, electronic design automation script generator, and bug analysis tool). While they found the 13B version best performing, which is slightly above our threshold for small models, the 7B version performed well (ibid.) At individual tasks, the best domain-adapted model was better or matched the performance of a general-purpose, larger model (Llama-2-70B-chat), was more cost-effective (ibid.) and, we might add, its alignment was more manageable. In some cases, as

⁸ <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>

⁹ <https://azure.microsoft.com/en-us/products/phi>.

¹⁰ <https://ai.google.dev/gemma>

¹¹ <https://www.alibabacloud.com/en/solutions/generative-ai/qwen>

¹² <https://www.openbmb.cn/>

¹³ <https://allenai.org/olmo>

noted by Liu et al. (2023), models' domain adaptation can require investments that are later compensated by lowering operating costs, which are higher for large, general models. Substantive up-front investments are, however, not the rule. Li et al. (2024) estimated that for their use case the cost of adapting the model to the domain was 'negligible'.

Further research is needed to assess the volume of negative externalities, not only in terms of safety and equity but also in terms of environmental impacts (Strubell et al. 2019; Spelda and Stritecky 2020), posed by these two approaches—centralized development of versatile models vs decentralized development of domain-specific, specialized models. There is evidence that during deployment general/multi-purpose models consume more energy (and depending on the power source also possibly produce more carbon emissions) than task-specific (single task) models (Luccioni et al. 2024). We currently lack general evidence on how this relation would look like for versatile vs specialized models, but since the latter will be smaller in the number of parameters, training and runtime inferences could be computationally cheaper and produce fewer negative externalities, all other things being equal. Two of our examples of small, domain-specific models, EON-8B and Phi-3-mini for Azure hardware provisioning, could be used to provide partial evidence. For EON-8B, Bodigutla et al. (2024) reported a lower number of GPUs necessary to support their use case compared to a solution based on a large, general model (see their Fig. 4). This translates into lower computational requirements, energy consumption, and fewer carbon emissions, considering comparable energy sources. Li et al. (2024) compared their domain-adapted Phi-3-mini with 3.8B parameters to the then state-of-the-art large, general model and reported lower costs and higher accuracy, clearly favoring the small model (see their Table 2). In situations like this, running a well-utilized, small, domain-adapted model is bound to be less wasteful than performing domain tasks via in-context learning on large, general models. Costs of adaptation need to be considered as well, but in this case, Li et al. (2024) reported them to be mild. Small and capable base language models could provide a viable, eco-friendly and cost-effective alternative to large, general models in certain domains. Environmental costs of developing small base models are favorable compared to resources involved in the development of large, general models, holding the energy source comparable.

Finally, as explained in Sect. 2, instruction following and generally specialized LLMs are produced using base pre-trained models. It is important to discuss provenance of these base models and practices used to develop them. Fully open models developed with responsible practices (e.g., Groeneveld et al. 2024) should be preferred over commercial or merely publicly available base models (only weights and

inference code are available) if full auditability (Casper et al. 2024) and good AI security practices are the top priority.

Genuinely open-source models differ from open-weight models by open access to pretraining data mixtures and training recipes used to develop base models, including data and code required for instruction tuning that produces instruction-following models from base models during post-training. Open-weight models, on the other hand, are available only as weights of the neural network and code that enables using it for completing user-defined tasks. For instance, Ai2's OLMo family of open-source, computationally efficient models offers highly competitive alternatives to open-weight models (see OLMo 2 report, Walsh et al. 2025), for which full data mixtures and detailed training recipes are not publicly available. Open-source artifacts that comprise the OLMo ecosystem, ranging from datasets to models at various sizes, help the scientific understanding of language models development. Specific efforts, such as fully documented and reproducible post-training recipes including the necessary data, e.g., Lambert et al. (2024b), open the stage of LLM development that is crucial for stabilizing, improving and aligning final capabilities of LLMs. Transparency of open-source models is helpful during safety evaluation because open access to pre- and post-training data mixtures and recipes enables identifying unseen tasks that were not part of the model development and, thus, offer a less biased perspective on its safety and alignment. For open-weight and closed models, for which information on pre- and post-training data is not fully disclosed, the degree to which evaluation tasks are unseen (held-out) is not independently verifiable. This can bias the perspective on the models' safety and alignment with respect to user requests.

Open-source and open-weight foundation models come with their own risks caused by the impossibility to control their downstream use. There is an ongoing debate on whether open LLMs pose a greater risk than closed ones accessible only indirectly via APIs, and we briefly touched on this topic when explaining the reversibility of alignment of open LLMs. Qi et al. (2024b) evaluated the robustness of two prominent defenses against fine-tuning/tampering attacks against open models (representation noising, Rosati et al. 2024, and tamper attack resistance, Tamirisa et al. 2024, both methods attempt to unlearn unsafe information from the model so that it is hard to bring it back by reversing alignment of the model) and concluded that they are not robust to what should be inconsequential details of fine-tuning and evaluation pipelines. This means that downstream misuse of open-source and open-weight models remains hard to prevent. Kapoor et al. (2024) constructed a framework for analyzing the marginal risk of open models, i.e., the degree by which a misused open-source or open-weight model increases the social risks over closed models or

legacy technologies allowing attackers to perform the task, which could be used to navigate the safety tradeoffs.

5 Conclusion

We showed that a securitization process underpinned by a combination of alignment and red teaming procedures reconfigured the public image of AI safety and produced a set of practices that are considered beneficial for safety. This combination of alignment and red teaming is also behind development of LLMs that AI companies use to build products such as chatbots. This situation resulted in regulatory approaches that took closed and commercial, general-purpose LLMs as their baseline, while leaving the status of publicly available or open-source LLMs less clear. Different baselines are possible, such as small, domain-specific LLMs allowing users to perform specialized tasks, on which we demonstrated that AI security practices could be different, more welcoming to participation and diversity—values with which the current AI security practices struggle.

Acknowledgements We would like to thank the reviewers for helpful feedback on the manuscript.

Author contributions Petr Spelda conceptualized the problem and performed its analysis. Petr Spelda and Vit Stritecky wrote and reviewed the paper.

Funding This work was supported by the European Regional Development Fund project "Beyond Security: Role of Conflict in Resilience-Building" (reg. no.: CZ.02.01.01/00/22_008/0004595).

Data availability No new data were generated.

Declarations

Conflict of interest On behalf of all authors, the corresponding author states that there is no conflict of interest.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Abdin M, Aneja J, Awadalla H et al. (2024) Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. [arXiv:2404.14219](https://arxiv.org/abs/2404.14219) [cs.CL].
- Anthropic (2023a) Core Views on AI Safety: When, Why, What, and How. <https://www.anthropic.com/news/core-views-on-ai-safety>.
- Anthropic (2023b) Anthropic's Responsible Scaling Policy. <https://www.anthropic.com/news/anthropics-responsible-scaling-policy>.
- Anthropic (2023c) Collective Constitutional AI: Aligning a Language Model with Public Input. <https://www.anthropic.com/news/collective-constitutional-ai-aligning-a-language-model-with-public-input>.
- Anthropic (2023d) Frontier Threats Red Teaming for AI Safety. <https://www.anthropic.com/news/frontier-threats-red-teaming-for-ai-safety>.
- Anthropic (2023e) Thoughts on the US Executive Order, G7 Code of Conduct, and Bletchley Park Summit. <https://www.anthropic.com/news/policy-recap-q4-2023>.
- Bai Y, Jones A, Ndousse K et al. (2022a) Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback. [arXiv:2204.05862](https://arxiv.org/abs/2204.05862) [cs.CL].
- Bai Y, Kadavath S, Kundu S et al. (2022b) Constitutional AI: Harmlessness from AI Feedback. [arXiv:2212.08073](https://arxiv.org/abs/2212.08073) [cs.CL].
- Balzacq T, Basaran T, Bigo D, Guittet E-P, Olsson C (2010) Security practices. In: Denmark (ed) International Studies Encyclopedia Online. <https://doi.org/10.1111/b.9781444336597.2010.x>.
- Barez F, Fu T, Prabhu A, Casper S, Sanyal A, Bibi A, O'Gara A, Kirk R, Bucknall B, Fist T, Ong L, Torr P, Lam K-Y, Trager R, Krueger D, Mindermann S, Hernandez-Orallo J, Geva M, Gal Y (2025) Open Problems in Machine Unlearning for AI Safety. [arXiv:2501.04952](https://arxiv.org/abs/2501.04952) [cs.LG].
- Bengio Y, Hinton G, Yao A et al. (2023) Managing AI Risks in an Era of Rapid Progress. [arXiv:2310.17688](https://arxiv.org/abs/2310.17688) [cs.CY].
- Bigo D, Tsoukala A (2008) Understanding (in)security. In: Bigo D, Tsoukala A (eds) Terror Insecurity and Liberty. Routledge, London, p 1
- Bodigutla PK, Jindal A, Balaji G, Zhu JS, Bing J, Rohit J, Jiang Y, Li Z (2024) How we built domain-adapted foundation GenAI models to power our platform. <https://www.linkedin.com/blog/engineering/generative-ai/how-we-built-domain-adapted-foundation-genai-models-to-power-our-platform>.
- Bommasani R, Hudson DA, Adeli E et al. (2021) On the opportunities and risks of foundation models. [arXiv:2108.07258](https://arxiv.org/abs/2108.07258) [cs.LG].
- Brown T, Mann B, Ryder N et al. (2020) Language Models are Few-Shot Learners. In Advances in Neural Information Processing Systems 33 (NeurIPS 2020).
- Casper S, Davies X, Shi C et al. (2023b) Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback. *Trans Mach Learn Res* 1:1
- Casper S, Lin J, Kwon J, Culp G, Hadfield-Menell D (2023a) Explore, Establish, Exploit: Red Teaming Language Models from Scratch. [arXiv:2306.09442](https://arxiv.org/abs/2306.09442) [cs.CL].
- Casper S, Ezell C, Siegmann C et al. (2024) Black-Box Access is Insufficient for Rigorous AI Audits. In FAccT '24: Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency, pp. 2254–2272.
- Christiano PF, Leike J, Brown T, Martic M, Legg S, Amodei D (2017) Deep Reinforcement Learning from Human Preferences. In Advances in Neural Information Processing Systems 30 (NIPS 2017).
- Collective CASE (2006) Critical Approaches to Security in Europe: A Networked Manifesto. *Secur Dial* 37(4):433

- Eloundou T, Lee T (2024) Democratic inputs to AI grant program: lessons learned and implementation plans. <https://openai.com/blog/democratic-inputs-to-ai-grant-program-update>.
- Fazelpour S, De-Arteaga M (2022) Diversity in sociotechnical machine learning systems. *Big Data Soc* 9(1):1
- Fazelpour S, Hadfield-Menell D, Belli L (2024) Red Teaming AI: The Devil Is In The Details. <https://www.techpolicy.press/red-teaming-ai-the-devil-is-in-the-details/>.
- Feffer M, Sinha A, Lipton ZC, Heidari H (2024) Red-Teaming for Generative AI: Silver Bullet or Security Theater? In Proceedings of the Seventh AAI/ACM Conference on AI, Ethics, and Society (AIES2024), pp. 421-437.
- Ganguli D, Lovitt L, Kernion J et al. (2022) Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned. [arXiv:2209.07858](https://arxiv.org/abs/2209.07858) [cs.CL].
- Glukhov D, Shumailov I, Gal Y, Papernot N, Pappas V (2023) LLM Censorship: A Machine Learning Challenge or a Computer Security Problem? [arXiv:2307.10719](https://arxiv.org/abs/2307.10719) [cs.AI].
- Grattafiori A, Dubey A, Jauhri A et al. (2024) The Llama 3 Herd of Models. [arXiv:2407.21783](https://arxiv.org/abs/2407.21783) [cs.AI].
- Groeneveld D, Beltagy I, Walsh P et al. (2024) Olmo: Accelerating the science of language models. [arXiv:2402.00838](https://arxiv.org/abs/2402.00838) [cs.CL].
- Guan MY, Joglekar M, Wallace E, Jain S, Barak B, Helyar A, Dias R, Vallone A, Ren H, Wei J, Chung HW, Toyer S, Heidecke J, Beutel A, Glaese A (2024) Deliberative Alignment: Reasoning Enables Safer Language Models. [arXiv:2412.16339](https://arxiv.org/abs/2412.16339) [cs.CL].
- Hubinger E, Denison C, Mu J et al. (2024) Sleeper Agents: Training Deceptive LLMs that Persist Through Safety Training. [arXiv:2401.05566](https://arxiv.org/abs/2401.05566) [cs.CR].
- Inie N, Stray J, Derczynski L (2023) Summon a Demon and Bind it: A Grounded Theory of LLM Red Teaming in the Wild. [arXiv:2311.06237](https://arxiv.org/abs/2311.06237) [cs.CL].
- Kapoor S, Bommasani R, Klyman K et al. (2024) On the Societal Impact of Open Foundation Models. [arXiv:2403.07918](https://arxiv.org/abs/2403.07918) [cs.CY].
- Krendl Gilbert T, Lambert N, Dean S, Zick T, Snoswell A (2023) Reward Reports for Reinforcement Learning. In Proceedings of the 2023 AAI/ACM Conference on AI, Ethics, and Society, pp. 84-130.
- Lambert N, Krendl Gilbert T, Zick T (2023) The History and Risks of Reinforcement Learning and Human Feedback. [arXiv:2310.13595](https://arxiv.org/abs/2310.13595) [cs.CY].
- Lambert N, Pyatkin V, Morrison J, Miranda LJ, Lin BY, Chandu K, Dziri N, Kumar S, Zick T, Choi Y, Smith NA, Hajishirzi H (2024a) RewardBench: Evaluating Reward Models for Language Modeling. [arXiv:2403.13787](https://arxiv.org/abs/2403.13787) [cs.LG].
- Lambert N, Morrison J, Pyatkin V et al. (2024b) Tulu 3: Pushing Frontiers in Open Language Model Post-Training. [arXiv:2411.15124](https://arxiv.org/abs/2411.15124) [cs.CL].
- Leike J, Sutskever I (2023) Introducing Superalignment. <https://openai.com/blog/introducing-superalignment>.
- Li B, Zhang Y, Bubeck S, Pathuri J, Menache I (2024) Small Language Models for Application Interactions: A Case Study. [arXiv:2405.20347](https://arxiv.org/abs/2405.20347) [cs.CL].
- Liu M, Ene T, Kirby R et al. (2023) ChipNeMo: Domain-Adapted LLMs for Chip Design. [arXiv:2311.00176v1](https://arxiv.org/abs/2311.00176v1) [cs.CL].
- Longpre S, Kapoor S, Klyman K et al. (2024) Position: A Safe Harbor for AI Evaluation and Red Teaming. In Proceedings of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235.
- Luccioni AS, Jernite Y, Strubell E (2024) Power Hungry Processing: Watts Driving the Cost of AI Deployment? In ACM Conference on Fairness, Accountability, and Transparency (ACM FAccT '24), June 3-6, 2024, Rio de Janeiro, Brazil. <https://doi.org/10.1145/3630106.3658542>.
- Mazeika M, Phan L, Yin X, Zou A, Wang Z, Mu N, Sakhaee E, Li N, Basart S, Li B, Forsyth D, Hendrycks D (2024) HarmBench: A Standardized Evaluation Framework for Automated Red Teaming and Robust Refusal. [arXiv:2402.04249](https://arxiv.org/abs/2402.04249) [cs.LG].
- Meinke A, Schoen B, Scheurer J, Balesni M, Shah R, Hobbahn M (2024) Frontier Models are Capable of In-context Scheming. [arXiv:2412.04984](https://arxiv.org/abs/2412.04984) [cs.AI].
- OpenAI (2023a) Our approach to AI safety. <https://openai.com/blog/our-approach-to-ai-safety>.
- OpenAI (2023b) Frontier risk and preparedness. <https://openai.com/blog/frontier-risk-and-preparedness>.
- OpenAI (2023c) OpenAI Red Teaming Network. <https://openai.com/blog/red-teaming-network>.
- OpenAI (2023d) Moving AI governance forward. <https://openai.com/blog/moving-ai-governance-forward>.
- Ouyang L, Wu J, Jiang X et al (2022) Training language models to follow instructions with human feedback. *Adv Neural Inform Process Syst* 35:1
- Perez E, Huang S, Song F, Cai T, Ring R, Aslanides J, Glaese A, McAleese N, Irving G (2022) Red Teaming Language Models with Language Models. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pp. 3419-3448, December 7-11.
- Qi X, Zeng Y, Xie T, Chen P-Y, Jia R, Mittal P, Henderson P (2024a) Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! In The Twelfth International Conference on Learning Representations.
- Qi X, Wei B, Carlini N, Huang Y, Xie T, He L, Jagielski M, Nasr M, Mittal P, Henderson P (2024b) On Evaluating the Durability of Safeguards for Open-Weight LLMs. [arXiv:2412.07097](https://arxiv.org/abs/2412.07097) [cs.CR].
- Rafailov R, Sharma A, Mitchell E, Manning CD, Ermon S, Finn C (2023) Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In 37th Conference on Neural Information Processing Systems (NeurIPS 2023).
- Rando J, Tramèr F (2024) Universal Jailbreak Backdoors from Poisoned Human Feedback. In The 12th International Conference on Learning Representations (ICLR 2024).
- Rauh M, Marchal N, Manzini A, Hendricks LA, Comanescu R, Akbulut C, Stepleton T, Mateos-Garcia J, Bergman S, Kay J, Griffin C, Bariach B, Gabriel I, Rieser V, Isaac W, Weidinger L (2024) Gaps in the Safety Evaluation of Generative AI. In Proceedings of the Seventh AAI/ACM Conference on AI, Ethics, and Society (AIES2024).
- Rosati D, Wehner J, Williams K, Bartoszczke L, Gonzales R, Maple C, Majumdar S, Sajjad H, Rudzicz F (2024) Representation Noising: A Defence Mechanism Against Harmful Finetuning. In The 38th Conference on Neural Information Processing Systems (NeurIPS 2024).
- Ryan MJ, Held W, Yang D (2024) Unintended Impacts of LLM Alignment on Global Representation. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 16121-16140.
- Samvelyan M, Chandra Raparthy S, Lupu A, Hambro E, Markosyan AH, Bhatt M, Mao Y, Jiang M, Parker-Holder J, Foerster J, Rocktäschel T, Raileanu R (2024) Rainbow Teaming: Open-Ended Generation of Diverse Adversarial Prompts. *Adv Neural Inform Process Syst* 37:1
- Santurkar S, Durmus E, Ladhak F, Lee C, Liang P, Hashimoto T (2023). Whose Opinions Do Language Models Reflect? In Proceedings of the 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202.
- Schmid P, Sanseviero O, Bartolome A, von Werra L, Vila D, Srivastav V, Sun M, Cuenca P (2024) Llama 3.1 - 405B, 70B & 8B with multilinguality and long context. <https://huggingface.com/blog/llama31>.
- Sharma M, Tong M, Mu J et al. (2025) Constitutional Classifiers: Defending against Universal Jailbreaks across Thousands of Hours of Red Teaming. [arXiv:2501.18837](https://arxiv.org/abs/2501.18837) [cs.CL].

- Sorensen T, Moore J, Fisher J, Gordon M, Mireshghallah N, Rytting CR, Ye A, Jiang L, Lu X, Dziri N, Althoff T, Choi Y (2024) Position: A Roadmap to Pluralistic Alignment. In Proceedings of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235.
- Spelda P, Stritecky V (2020) The Future of Human-Artificial Intelligence Nexus and its Environmental Costs. *Futures* 117:102531
- Strubell E, Ganesh A, McCallum A (2019) Energy and Policy Considerations for Deep Learning in NLP. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 3645-3650, Florence, Italy.
- Symons J, Horner JK (2019) Why There is no General Solution to the Problem of Software Verification. *Found Sci* 25:541–557
- Tamirisa R, Bharathi B, Phan L, Zhou A, Gatti A, Suresh T, Lin M, Wang J, Wang R, Arel R, Zou A, Song D, Li B, Hendrycks D, Mazeika M (2024) Tamper-Resistant Safeguards for Open-Weight LLMs. [arXiv:2408.00761](https://arxiv.org/abs/2408.00761) [cs.LG].
- Tamkin A, Askill A, Lovitt L, Durmus E, Joseph N, Kravec S, Nguyen K, Kaplan J, Ganguli D (2023) Evaluating and Mitigating Discrimination in Language Model Decisions. [arXiv:2312.03689](https://arxiv.org/abs/2312.03689) [cs.CL].
- Touvron H, Martin L, Stone K et al. (2023) Llama 2: Open Foundation and Fine-Tuned Chat Models. [arXiv:2307.09288](https://arxiv.org/abs/2307.09288) [cs.CL].
- Walsh P, Soldaini L, Groeneveld D et al. (2025) 2 OLMo 2 Furious. [arXiv:2501.00656](https://arxiv.org/abs/2501.00656) [cs.CL].
- Wei A, Haghtalab N, Steinhardt J (2023) Jailbroken: How Does LLM Safety Training Fail? *Adv Neural Inform Process Syst* 36:2
- Weidinger L, Mellor JFJ, Pegueroles BG, Marchal N, Kumar R, Lum K, Akbulut C, Diaz M, Bergman AS, Rodriguez MD, Rieser V, Isaac W (2024) STAR: SocioTechnical Approach to Red Teaming Language Models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing.
- Zou A, Wang Z, Kolter JZ, Fredrikson M (2023) Universal and Transferable Adversarial Attacks on Aligned Language Models. [arXiv:2307.15043](https://arxiv.org/abs/2307.15043) [cs.CL].

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.