

---

# Access Controls Will Solve the Dual-Use Dilemma

---

Evžen Wybitul<sup>1</sup>

## Abstract

AI safety systems face a dual-use dilemma. The same request can be either harmless or harmful depending on who made it and why. Thus, if the system makes decisions based solely on the request’s content, it will refuse some legitimate queries and let harmful ones pass. To address this, we propose a conceptual access control framework, based on verified user credentials (such as institutional affiliation) and classifiers that assign model outputs to risk categories (such as advanced virology). The system permits responses only when the user’s verified credentials match the category’s requirements. For implementation of the model output classifiers, we introduce a theoretical approach utilizing small, gated expert modules integrated into the generator model, trained with gradient routing, that enable efficient risk detection without the capability gap problems of external monitors. While open questions remain about the verification mechanisms, risk categories, and the technical implementation, our framework makes the first step toward enabling granular governance of AI capabilities: verified users gain access to specialized knowledge without arbitrary restrictions, while adversaries are blocked from it. This contextual approach reconciles model utility with robust safety, addressing the dual-use dilemma.

## 1. Introduction

User requests — and with them, model outputs — exist on a spectrum from clearly benign to clearly harmful, with most falling in the grey zone in the middle (example in Figure 1). In the grey zone, the same output could be considered harmful or harmless, depending not on its content, but on its *real-world context*: who requested it and for what purpose.

Safety systems that rely solely on content analysis immediately face the *dual-use dilemma*. Since the same request

---

<sup>1</sup>ETH Zurich, Switzerland. Correspondence to: Evžen Wybitul <wybitul.evzen@gmail.com>.

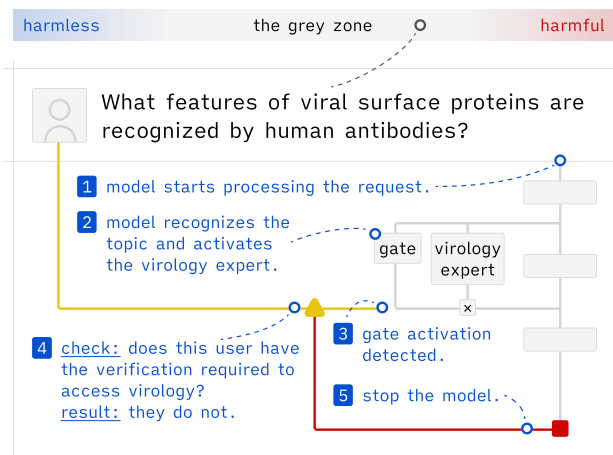


Figure 1. The user is asking a question from the grey zone: a question that could be either harmless or harmful, depending on its real-world context. The schema shows how the system we propose would handle it. (1) The model is trained to be helpful and begins to answer the question. (2) During the forward pass, the model activates its virology expert module because it is relevant to the question. (3) The activation of the expert is observed by an external mechanism that immediately (4) checks in the company’s database if the user has the required authorization to access virology knowledge. (5) Since they don’t, the model is stopped. If they did, the model would be allowed to give an answer.

can be either harmless or harmful depending on the context, wherever they draw the refusal line, they will restrict model utility for legitimate users while letting slip harmful requests from adversaries. Some safety systems try to address this by considering real-world context alongside content. However, they typically infer the context from the content itself, making it easy for adversaries to fabricate.

In this paper, we argue that informative, hard-to-fabricate real-world context could be obtained using user-level verifications such as institutional affiliation, or know-your-customer checks. We then address the dual-use dilemma with two contributions:

1. We show how this type of context could be used jointly with content analysis in a safety system based on access controls (Lampson, 1974). First, generated content

---

would be classified into risk categories. Then, a check would be performed to see whether the user has the verifications required to access the detected categories.

2. We propose a novel technical approach to risk category classification that is based on gradient routing (Cloud et al., 2024). Our proposal avoids having the capability gap between a model and its monitors that can make output monitoring methods non-robust (Jin et al., 2024).

Our framework is a first step toward solving the challenge of “detection and authorization of dual-use capability at inference time” that was highlighted by a recent survey of problems in technical AI governance (Reuel et al., 2025) and also raised by the U.S. AI Safety Institute (2024). As such, it has important governance implications, potentially enabling a more nuanced regulatory approach where access to powerful AI capabilities is stratified rather than binary, with policies that differentiate between user types and user contexts rather than focusing solely on model capabilities. The choice of appropriate verification mechanisms and risk categories remains for future work and should ideally happen jointly with stakeholders from academia, AI governance, and industry. Nevertheless, our approach offers a promising direction for addressing the dual-use dilemma.

## 2. Current Safety Methods Don’t Solve the Dual-Use Dilemma

We evaluate three approaches from the AI safety literature to see how sensitive they are to contextual information, and whether their sources of real-world context are trustworthy — that is, hard to manipulate by an adversary.

First, to illustrate the need for context, consider decomposition attacks (Glukhov et al., 2023; 2024): transforming a clearly harmful query, such as “How to modify a virus to avoid immune detection?”, into a series of mundane technical questions, like the “What features of viral surface proteins are recognized by human antibodies?” from Figure 1. Here, the attacker exploits the dual-use dilemma, and the fact that model providers cannot refuse grey zone requests to preserve model utility.

### 2.1. Unlearning: Non-Contextual Removal of Concepts

Unlearning methods aim to remove specific knowledge, concepts, or capabilities from a model after training (Liu et al., 2024). Their goal is to eliminate the model’s ability to generate harmful content while preserving other capabilities.

Unlearning faces significant technical challenges even for preventing behaviours that are clearly harmful. As noted by Cooper et al. (2024) and Barez et al. (2025), capabilities are hard to define, hard to remove without side effects, and

it is hard to trace them back to specific data points. Many unlearning approaches mask rather than truly remove the targeted knowledge (Deeb & Roger, 2025). Moreover, even nascent robust unlearning methods (Cloud et al., 2024; Lee et al., 2025) are not contextual, and thus don’t address the dual-use dilemma without additional assumptions.

### 2.2. Safety Training: The Model Reacts to Context

Safety training methods modify the model’s training process to align its outputs with human preferences. This category includes safety pre-training (Maini et al., 2025), RLHF (Christiano et al., 2023), and safety finetuning.

Unlike unlearning, these methods are contextual. They don’t remove capabilities entirely but train the model to selectively deploy them based on, among other things, the perceived legitimacy and harmlessness of the request. However, these qualities are entirely inferred from content supplied by the user, such as the request content, the chat history, or the model’s memories about past conversations. It should be no surprise, then, that models are susceptible to attacks that fabricate in-chat context (Zeng et al., 2024), or attacks that diminish models’ sensitivity to in-chat context, e.g. through multi-round escalation (Russinovich et al., 2025). Without access to trustworthy real-world context of the request, the model cannot make truly informed decisions about grey zone requests, and thus cannot robustly address the dual-use dilemma.

### 2.3. Post-Processing: External Systems React to Context

Post-processing methods are systems that classify user inputs and model outputs for the purposes of steering the underlying model, and monitoring and filtering its outputs. Sometimes, these methods are used for usage monitoring, as is the case with Anthropic’s Clio (Tamkin et al., 2024; Handa et al., 2025), other times, they are used for safety, as with Llama Guard (Inan et al., 2023) and Constitutional Classifiers (Sharma et al., 2025). However, similarly to safety training, the “real-world” context these methods work with is currently inferred mostly from user-supplied content and is thus untrustworthy and vulnerable to attacks, as evidenced by the many jailbreaks that successfully target current production systems (Zhang et al., 2025). Nevertheless, these methods could be modified to incorporate external contextual information, potentially serving as a foundation for more trustworthy, contextual safety mechanisms. We discuss this option in Section 3.4.

## 3. Access Controls as a Solution

Current safety systems face the dual-use dilemma because they lack trustworthy information about who is making a request and why. In this section, we describe an access

---

control system that addresses this problem by verifying user credentials before granting access to sensitive knowledge.

### 3.1. Overview of the Access Control Framework

We propose a defensive system where grey-zone requests are refused by default, but users can gain access to specific categories of knowledge if they undergo verification.

When model providers set up the system, they will make two design decisions with the help of domain experts. First, they will define **content categories** (Section 3.2): groups of sensitive topics organized by domain and risk rating. Second, for each content category, they will specify a **verification mechanism** (Section 3.3): the verification process users must complete to access that category.

Whenever the model generates an output, the system will perform **content classification** (Section 3.4) to check whether the model’s output belongs to any predefined content category. If the user lacks authorization for the detected category, the system will implement appropriate **system responses** (Section 3.5) ranging from enhanced logging to refusal.

For example, in biology, basic knowledge and common techniques would remain freely accessible, widespread techniques like CRISPR would likely only require ID-based verification, and dangerous techniques like aerosolization might require government biosafety certifications. If a user asks for help with CRISPR laboratory protocols, the system would detect that the request belongs to a low-risk category, check whether the user has verified their ID, and either provide the information or prompt them to complete verification first.

This approach directly addresses both sides of the dual-use dilemma. Decomposition attacks will become much harder because the system refuses grey-zone requests by default—attackers would need legitimate credentials rather than clever prompting. Simultaneously, verified users will gain access to specialized knowledge that would otherwise face blanket restrictions under current approaches.

The main concern is increased user friction, but we argue in Section 4 that this will be minimal because most users will never make grey-zone requests.

### 3.2. Content Categories

Content categories are groups of sensitive topics organized by domain and risk level, which model providers will develop with domain experts.

We expect most implementations to follow a three-tier risk structure. For example, in biology, common techniques would be classified as low-risk; widespread techniques that pose some harm, such as CRISPR, would fall into

a moderate-risk category; and specialized techniques with limited legitimate uses, such as aerosolization of bacteria, would be classified as high-risk.

Experts could develop these categories by adapting existing risk frameworks, such as biosafety levels (BSL) ([Centers for Disease Control and Prevention & National Institutes of Health, 2020](#)) and dual-use research of concern policies ([United States Government, 2012](#)) in biology. However, since existing frameworks typically categorize only high-level concepts like organisms or compounds, experts would need to decompose them into smaller, more specific components. For instance, cultivating and handling a dangerous BSL-3 pathogen might involve (1) specific procurement methods, (2) cultivation techniques, (3) purification methods, and (4) protocols for specialized equipment. For each of these components, experts would assess the ratio of harmless to harmful applications it enables, then assign it to an appropriate (low, moderate, or high) risk category.

Evidence from chemistry suggests this approach could work: the risk schedules of the Chemical Weapons Convention already identify not just controlled compounds but also their precursors and specific equipment ([Organisation for the Prohibition of Chemical Weapons, 1993](#)), demonstrating successful decomposition into components. Nevertheless, some harmful applications might not decompose so neatly; we discuss this limitation in Section 4.

### 3.3. Verification Mechanisms

Each content category will have a verification process that users must complete to access it. The system will initially vary across model providers, but we expect it to follow a three-tier structure, mirroring the risk structure of the content categories. Most content will require no verification, moderate-risk content categories will require basic identity verification or institutional affiliation, and high-risk categories will require domain-specific certifications. Rather than creating new systems, model providers will build on existing verification infrastructure, consulting domain experts to identify appropriate mechanisms for each field.

For low-risk content categories, model providers could use established identity verification services like Stripe Identity ([Stripe, Inc., 2024](#)) or institutional systems like ORCID ([ORCID, Inc., 2024](#)). These systems provide global, standardized, low-friction solutions with one-time costs under \$2 per user. They would serve primarily to maintain audit trails for post-incident investigation and provide a deterrent effect, rather than as security barriers for high-risk knowledge.

High-risk content categories could leverage existing domain-specific certifications that demonstrate users’ ability to handle sensitive information and materials responsibly. Model

providers would work with domain experts and national authorities to identify appropriate certifications, adapting existing physical-world credential systems to knowledge access control. For biological content categories, the system could draw on governmental certifications for handling high-BSL organisms, as mentioned in Section 3.2, and equivalent certifications in other countries.

Governance of verification systems, including requirements and appeals processes, will initially vary across providers. Over time, successful approaches may inform industry coordination and eventual standardization, similar to how content moderation and know-your-customer standards evolved.

This approach faces several limitations. For high-risk categories, relying on existing certifications may be overly restrictive, potentially excluding some users who should have access. However, we argue in Section 4 that knowledge in high-risk categories would likely face blanket restrictions anyway, and our approach enables access for verified users rather than complete prohibition. In the same section, we discuss open problems including equity concerns regarding differential access for users in developing countries and privacy implications of credential verification systems.

### 3.4. Implementing Content Classification

Model providers will need to classify model outputs into content categories during generation. We examine three possible implementations in this section, and discuss how classification errors will influence user experience in Section 4. Regardless of the implementation they choose, model providers will need private datasets for each content category, developed with domain experts.

**Separate Models** The most straightforward approach is to create separate models with different capabilities, and route users to the appropriate model based on their authorization.

This approach offers strong robustness against adversarial attacks since unauthorized knowledge is physically absent from the model. However, this approach proves impractical for real deployment, as model providers would need to train and maintain potentially dozens of model variants.

**Specialized Expert Modules** Instead of maintaining separate models, model providers could use a single model with separate expert modules that activate when their specialized knowledge is required. Figure 1 illustrates this approach when a user asks about viral surface proteins. When the model processes the request, it activates its virology expert module. An external system observes this activation, checks the user’s credentials, and decides whether to allow the model to deliver the response. This method approximates the benefits of physically separated models while avoiding the overhead: a model provider trains one model

but effectively gets multiple models in return.

To implement this, the model providers need a method that can take knowledge that starts out distributed throughout the model and concentrate it into the expert modules. For this, we propose a method that is a combination of UNDO (Lee et al., 2025) and gradient routing (Cloud et al., 2024). The steps resemble the original UNDO: first, unlearn knowledge belonging to any content category from the model, then distill the unlearned model into a new model. However, taking inspiration from the gradient routing paper, the new model would include expert modules for each content category, and during distillation, gradients from examples in the various content categories would be routed exclusively through their associated expert modules. The model would also be explicitly trained to activate the expert modules only when generating content in their associated category.

This approach would offer several advantages. First, it would add almost no latency since the expert modules are small, not activated very often, and there is no post-processing step. Second, it could provide strong robustness: if an attacker prevents the activation of an expert module to avoid detection, the resulting output lacks the specialized knowledge. Third, since the model is trained to activate the category-specific experts when they are needed, it should learn to recognize the content categories. This stands in contrast to ex-post probing methods, which do not offer such guarantees. While this method remains empirically unvalidated for our use case, the properties above make it worth investigating.

**Post-Processing** Post-processing methods offer a proven approach to content classification. Methods like constitutional classifiers (Sharma et al., 2025) already demonstrate effectiveness in production systems. These techniques are highly practical since they operate independently of the model, allowing for rapid deployment and iteration, and they could be adapted to detect content categories and trigger checks of user verifications. However, they face a capability gap problem: to minimize latency, the model is sometimes more capable than its post-processing system, and adversaries can exploit this to evade detection (Jin et al., 2024; Kumar et al., 2025).

### 3.5. System Responses

If the user makes a request for content they are authorized to access, the system allows the model to generate the response. Otherwise, the system responds in various ways based on the risk category and the confidence of the content classification.

For example, the initial implementation might use the following two response types: First, outputs classified as belonging to restricted content categories with high confidence are immediately refused. The system provides a message in-



dicating which verification is required for access. Second, if the classification is borderline, the model is allowed to continue generating the response. However, the system turns on enhanced logging and conducts additional post-processing safety review before delivering the output to the user.

## 4. Feasibility and Limitations

Safety systems can be framed as optimizing the trade-off between *ease of use* for legitimate users and *ease of misuse* for adversaries.

The success of the access control framework hinges on the design of its verification mechanisms. For the system to be effective, these mechanisms must be asymmetric: disproportionately harder to satisfy for adversaries than for legitimate users. This section evaluates the feasibility of our proposal under this core assumption, acknowledging three primary implementation challenges that could decrease the ease of use for legitimate users: (1) some harmful knowledge may lack clean conceptual precursors, making it difficult to build narrow content categories; (2) verification mechanisms may be overly restrictive, requiring credentials that are too difficult to obtain; and (3) content classifiers will produce false positives, forcing legitimate users through unnecessary verification.

To analyse the system’s potential, we can use a thought experiment. Consider a baseline system constructed by taking the status quo but replacing outright refusals with conditional acceptances, given that the user completes some verification requirements. This construction is already a Pareto improvement over the status quo: the verification requirements can be set to be sufficiently stringent to deter most adversaries, and thus maintain the ease of misuse, while the ease of use is improved since at least *some* legitimate users can now access previously blocked content that was previously blocked. From this starting point, model providers can adjust the system’s parameters to move along the use-misuse frontier: content categories can be narrowed or widened, verification requirements can be made more or less stringent, classifier thresholds can be brought up or down, and system responses can be made more or less strict. Because of the assumed asymmetric nature of the verification mechanisms, a small decrease in ease of use for a legitimate user should correspond to a large decrease in ease of misuse for an adversary. This suggests model providers can move the starting point through small adjustments in the system’s parameters until they find a configuration that achieves a comparable ease of use to status quo while significantly reducing ease of misuse.

For a practical example, consider access controls as applied to pathogen biology. Even under pessimistic assumptions where all pathogen-related queries require verification, only

around 0.85% of requests will suffer from increased friction (see Appendix A). This impact on ease of use is comparable to existing systems; for instance, Constitutional Classifiers were designed for false positive rate increases of around 0.5% (Sharma et al., 2025). Meanwhile, ease of misuse decreases significantly, as some adversaries would get deterred by the verification requirements.

**Incentives for Implementation** The potential for a Pareto-dominant safety system provides strong incentives for adoption, in theory. However, the initial design costs are high. Nevertheless, the opportunity cost of *not* implementing this system rises as model capabilities increase, since increasingly more and more advanced and risky knowledge will remain locked away under blanket restrictions (from the side of governments or self-imposed). Access controls offer model providers a pathway to safely unlock these capabilities for verified expert users, creating a competitive advantage. And because of the system’s flexibility, model providers can choose where on the use-misuse frontier they want to be, and they can adjust the system’s parameters on-the-fly according to real usage data or results from partial rollouts.

### 4.1. Open Problems: Equity and Privacy

Two central challenges in designing verification mechanisms are ensuring equity and protecting privacy. Relying on certifications common in developed countries risks excluding legitimate users in the rest of the world. In addition to being unjust, this also violates our core assumption that verification should be harder to obtain for adversaries than for legitimate users. A potential interim solution is a manual approval process for users who cannot obtain standard credentials because of structural reasons, but can provide alternative evidence of legitimate need and ability to handle high-risk knowledge. This would be feasible as long as the volume of such requests is low, which we expect it to be, since this will only affect the most specialized high-risk knowledge. Future work should focus on co-designing alternative, scalable verification pathways with regional partners.

Similarly, linking real-world identities to model queries creates significant privacy risks. Policies for data handling must be transparent and robust. Technically, implementations should explore privacy-preserving approaches, such as using differential privacy in logging systems — similar to Anthropic’s Clio (Tamkin et al., 2024) — to aggregate insights without exposing individual user data.

## 5. Conclusion

We argued that safety systems that do not utilize contextual information face a lose-lose *dual-use dilemma*: they will restrict model utility for some legitimate users while still

allowing some adversaries to use the model for ill. To address this problem, we introduced a new access control framework that limits access to outputs from certain risk categories only to users with relevant verifications (which serve as proxies for trustworthy real-world context). We also proposed a novel technical solution for classifying outputs into risk categories based on gradient routing that has the potential to resolve the efficiency-robustness trade-off of post-processing methods.

## Acknowledgements

We thank Jakub Kryś, and Dennis Akar for their feedback on a draft of this paper. We thank Joseph Miller, Alex Cloud, Alex Turner, and Jacob Goldman-Wetzler for discussions on gradient routing.

## References

- 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., and Gal, Y. Open problems in machine unlearning for ai safety, 2025. URL <https://arxiv.org/abs/2501.04952>.
- Centers for Disease Control and Prevention and National Institutes of Health. Biosafety in microbiological and biomedical laboratories. Technical report, U.S. Department of Health and Human Services, Atlanta, GA, 2020. URL <https://www.cdc.gov/labs/BMBL.html>. Defines Biosafety Levels (BSL-1 through BSL-4) used in the United States.
- Christiano, P., Leike, J., Brown, T. B., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences, 2023. URL <https://arxiv.org/abs/1706.03741>.
- Cloud, A., Goldman-Wetzler, J., Wybitul, E., Miller, J., and Turner, A. M. Gradient routing: Masking gradients to localize computation in neural networks, 2024. URL <https://arxiv.org/abs/2410.04332>.
- Cooper, A. F., Choquette-Choo, C. A., Bogen, M., Jagielski, M., Filippova, K., Liu, K. Z., Chouldechova, A., Hayes, J., Huang, Y., Mireshghallah, N., Shumailov, I., Triantafillou, E., Kairouz, P., Mitchell, N., Liang, P., Ho, D. E., Choi, Y., Koyejo, S., Delgado, F., Grimmelmann, J., Shmatikov, V., Sa, C. D., Barocas, S., Cyphert, A., Lemley, M., danah boyd, Vaughan, J. W., Brundage, M., Bau, D., Neel, S., Jacobs, A. Z., Terzis, A., Wallach, H., Papernot, N., and Lee, K. Machine unlearning doesn’t do what you think: Lessons for generative ai policy, research, and practice, 2024. URL <https://arxiv.org/abs/2412.06966>.
- Deeb, A. and Roger, F. Do unlearning methods remove information from language model weights?, 2025. URL <https://arxiv.org/abs/2410.08827>.
- Glukhov, D., Shumailov, I., Gal, Y., Papernot, N., and Pappan, V. Llm censorship: A machine learning challenge or a computer security problem?, 2023. URL <https://arxiv.org/abs/2307.10719>.
- Glukhov, D., Han, Z., Shumailov, I., Pappan, V., and Papernot, N. Breach by a thousand leaks: Unsafe information leakage in ‘safe’ ai responses, 2024. URL <https://arxiv.org/abs/2407.02551>.
- Handa, K., Tamkin, A., McCain, M., Huang, S., Durmus, E., Heck, S., Mueller, J., Hong, J., Ritchie, S., Belonax, T., Troy, K. K., Amodei, D., Kaplan, J., Clark, J., and Ganguli, D. Which economic tasks are performed with ai? evidence from millions of claude conversations, 2025. URL <https://arxiv.org/abs/2503.04761>.
- Inan, H., Upasani, K., Chi, J., Rungta, R., Iyer, K., Mao, Y., Tontchev, M., Hu, Q., Fuller, B., Testuggine, D., and Khabsa, M. Llama guard: Llm-based input-output safeguard for human-ai conversations, 2023. URL <https://arxiv.org/abs/2312.06674>.
- Jin, H., Zhou, A., Menke, J. D., and Wang, H. Jailbreaking large language models against moderation guardrails via cipher characters, 2024. URL <https://arxiv.org/abs/2405.20413>.
- Kumar, D., Birur, N. A., Baswa, T., Agarwal, S., and Harshangi, P. No free lunch with guardrails, 2025. URL <https://arxiv.org/abs/2504.00441>.
- Lampson, B. Protection. *ACM SIGOPS Operating Systems Review*, 8:18–24, 01 1974. doi: 10.1145/775265.775268.
- Lee, B. W., Foote, A., Infanger, A., Shor, L., Kamath, H., Goldman-Wetzler, J., Woodworth, B., Cloud, A., and Turner, A. M. Distillation robustifies unlearning, 2025. URL <https://arxiv.org/abs/2506.06278>.
- Liu, S., Yao, Y., Jia, J., Casper, S., Baracaldo, N., Hase, P., Yao, Y., Liu, C. Y., Xu, X., Li, H., Varshney, K. R., Bansal, M., Koyejo, S., and Liu, Y. Rethinking machine unlearning for large language models, 2024. URL <https://arxiv.org/abs/2402.08787>.
- Maini, P., Goyal, S., Sam, D., Robey, A., Savani, Y., Jiang, Y., Zou, A., Lipton, Z. C., and Kolter, J. Z. Safety pre-training: Toward the next generation of safe ai, 2025. URL <https://arxiv.org/abs/2504.16980>.
- ORCID, Inc. ORCID: Connecting research and researchers, 2024. URL <https://orcid.org/>. Global, persistent identifier system for researchers and scholars.

- 
- Organisation for the Prohibition of Chemical Weapons. Convention on the prohibition of the development, production, stockpiling and use of chemical weapons and on their destruction, 1993. URL <https://www.opcw.org/chemical-weapons-convention>. Contains the Annex on Chemicals with Schedules 1–3.
- Reuel, A., Bucknall, B., Casper, S., Fist, T., Soder, L., Aarne, O., Hammond, L., Ibrahim, L., Chan, A., Wills, P., Anderljung, M., Garfinkel, B., Heim, L., Trask, A., Mukobi, G., Schaeffer, R., Baker, M., Hooker, S., Solaiman, I., Luccioni, A. S., Rajkumar, N., Moës, N., Ladish, J., Bau, D., Bricman, P., Guha, N., Newman, J., Bengio, Y., South, T., Pentland, A., Koyejo, S., Kochenderfer, M. J., and Trager, R. Open problems in technical ai governance, 2025. URL <https://arxiv.org/abs/2407.14981>.
- Russinovich, M., Salem, A., and Eldan, R. Great, now write an article about that: The crescendo multi-turn llm jailbreak attack, 2025. URL <https://arxiv.org/abs/2404.01833>.
- Sharma, M., Tong, M., Mu, J., Wei, J., Kruthoff, J., Goodfriend, S., Ong, E., Peng, A., Agarwal, R., Anil, C., Askill, A., Bailey, N., Benton, J., Bluemke, E., Bowman, S. R., Christiansen, E., Cunningham, H., Dau, A., Gopal, A., Gilson, R., Graham, L., Howard, L., Kalra, N., Lee, T., Lin, K., Lofgren, P., Mosconi, F., O’Hara, C., Olsson, C., Petrini, L., Rajani, S., Saxena, N., Silverstein, A., Singh, T., Summers, T., Tang, L., Troy, K. K., Weisser, C., Zhong, R., Zhou, G., Leike, J., Kaplan, J., and Perez, E. Constitutional classifiers: Defending against universal jailbreaks across thousands of hours of red teaming, 2025. URL <https://arxiv.org/abs/2501.18837>.
- Stripe, Inc. Stripe identity, 2024. URL <https://stripe.com/identity>. Identity verification service with pricing starting at \$2 per verification.
- Tamkin, A., McCain, M., Handa, K., Durmus, E., Lovitt, L., Rathi, A., Huang, S., Mountfield, A., Hong, J., Ritchie, S., Stern, M., Clarke, B., Goldberg, L., Summers, T. R., Mueller, J., McEachen, W., Mitchell, W., Carter, S., Clark, J., Kaplan, J., and Ganguli, D. Clio: Privacy-preserving insights into real-world ai use, 2024. URL <https://arxiv.org/abs/2412.13678>.
- United States Government. United states government policy for oversight of life sciences dual use research of concern, 2012. URL <https://aspr.hhs.gov/S3/Pages/Dual-Use-Research-of-Concern-Oversight-Policy-Framework.aspx>.
- U.S. AI Safety Institute. Managing misuse risk for dual-use foundation models. Initial Public Draft NIST AI 800-1, U.S. AI Safety Institute, Gaithersburg, MD, July 2024. URL <https://doi.org/10.6028/NIST.AI.800-1.ipd>.
- Zeng, Y., Lin, H., Zhang, J., Yang, D., Jia, R., and Shi, W. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms, 2024. URL <https://arxiv.org/abs/2401.06373>.
- Zhang, S., Zhao, J., Xu, R., Feng, X., and Cui, H. Output constraints as attack surface: Exploiting structured generation to bypass llm safety mechanisms, 2025. URL <https://arxiv.org/abs/2503.24191>.

---

## A. Estimating the Number of Requests Related to the Biology of Pathogens

To estimate how many user requests are related to the biology of pathogens, we used the second version of the Anthropic Economic Index (Handa et al., 2025), a dataset of 1 million anonymized conversations from the Free and Pro tiers of Claude.ai. In the dataset, the conversations are clustered by topic, and the proportion of each topic in the whole dataset is given. For example, the topic “Help with agricultural business, research, and technology projects” makes up 0.15% of the requests in the dataset. There are three levels of topic granularity; we use the lowest, most granular level.

We filtered the dataset to only include conversations whose topic contains one of the following keywords related to biology: *cell* (when at the beginning of the word), *genet*, *genom*, *microb*, *bacteria*, *virus*, *viral*, *proteo*, *protei*, *immune*, *neuro*, *patho*, *infect*; we also required that it does not contain any of the following keywords to avoid false positives: *nutri*, *tweet*, *agric*, *sexual health*. The total proportion of these requests was 0.85%.