

Self-Destructing Models: Increasing the Costs of Harmful Dual Uses of Foundation Models

PETER HENDERSON*, Stanford University, USA
 ERIC MITCHELL*, Stanford University, USA
 CHRISTOPHER D. MANNING, Stanford University, USA
 DAN JURAFSKY, Stanford University, USA
 CHELSEA FINN, Stanford University, USA

A growing ecosystem of large, open-source foundation models has reduced the labeled data and technical expertise necessary to apply machine learning to many new problems. Yet foundation models pose a clear dual-use risk, indiscriminately reducing the costs of building both harmful and beneficial machine learning systems. Policy tools such as restricted model access and export controls are the primary methods currently used to mitigate such dual-use risks. In this work, we review potential safe-release strategies and argue that both policymakers and AI researchers would benefit from fundamentally new technologies enabling more precise control over the downstream usage of open-source foundation models. We propose one such approach: the *task blocking* paradigm, in which foundation models are trained with an additional mechanism to impede adaptation to harmful tasks without sacrificing performance on desirable tasks. We call the resulting models *self-destructing models*, inspired by mechanisms that prevent adversaries from using tools for harmful purposes. We present an algorithm for training self-destructing models leveraging techniques from meta-learning and adversarial learning, which we call *meta-learned adversarial censoring (MLAC)*. In a small-scale experiment, we show MLAC can largely prevent a BERT-style model from being re-purposed to perform gender identification without harming the model’s ability to perform profession classification.

ACM Reference Format:

Peter Henderson, Eric Mitchell, Christopher D. Manning, Dan Jurafsky, and Chelsea Finn. 2023. Self-Destructing Models: Increasing the Costs of Harmful Dual Uses of Foundation Models. In *AAAI/ACM Conference on AI, Ethics, and Society (AIES '23), August 8–10, 2023, Montréal, QC, Canada*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3600211.3604690>

1 INTRODUCTION

A defining capability of large pretrained models (hereafter foundation models; FMs) is their ability to adapt to many downstream tasks in a few-shot manner—potentially improving performance and efficiency in domains with little training data [7]. Today, anyone with an internet connection can download a foundation model and adapt it to socially beneficial use-cases, like building better educational tools or improving access to justice. However, a malicious actor can also adapt a foundation model to nearly any harmful use-case they desire. For example, an oppressive government can

*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

AIES '23, August 8–10, 2023, Montréal, QC, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
 ACM ISBN 979-8-4007-0231-0/23/08...\$15.00
<https://doi.org/10.1145/3600211.3604690>

take a powerful pretrained language model and adapt it to identify dissidents; a rogue actor can adapt a pretrained object recognition system such that commercially available drones act as targeted loitering munitions; or a pretrained drug discovery system can be used for creating chemical or biological weapons, like neurotoxins [55]. Unfortunately, due to their general-purpose nature, preventing such dual uses of foundation models is difficult. This creates a tension between making these models widely available and ensuring that they are used in a safe and responsible way.

Currently, there are several approaches to mitigating the dual uses of FMs which can be divided into *structural* safety mechanisms and *technical* safety mechanisms. Structural mechanisms use licenses or access restrictions to prevent harmful uses; there is a broad spectrum of such structural release mechanisms. Some have suggested a review board for selecting the structural release mechanism [34] while others have argued that open source access to foundation models is essential for safety research [6]. While structural release approaches aim to prevent malicious users from acquiring foundation models or providing legal remedies if they exceed the terms of their access, *technical* strategies ensure that the model cannot be used for harmful purposes even if a malicious user is able to gain access to the model itself. Current technical strategies aim to tune the model so that it is less likely to produce harmful content at inference time [3], but do not consider the case where adversaries have access to model parameters.

In this work, we review these strategies, noting that no strategy on its own is able to prevent harmful dual uses of FMs. In particular we note the disconnect between the goal of many structural safety mechanisms and the new reality of open-source foundation models: structural safety strategies aim to prevent a malicious actor from gaining access to the model parameters altogether. In recent months, however, powerful open-source models have been released to the public, including Meta’s Llama model which was leaked online despite a restricted access policy [53, 58]. Such developments demand changes to the threat model of malicious FM usage, specifically, that eventually model parameters will become generally accessible. Unlike the assumptions of current safety strategies, there should then be a last layer of defense that renders the model itself as harmless as possible. We argue that we need more *technical* strategies to supplement *structural* strategies to reduce the ability for adversaries to use and adapt foundation models for harmful tasks: even when they have access to model parameters. Where existing access restrictions must navigate the tension between openness and safety, we seek to provide a new research pathway for reducing (and in some cases obviating) this tension.

We suggest one such new path forward: *self-destructing models*. Self-destructing models are trained via a task blocking method that impedes the adaptation of the model to a harmful task without impairing the model’s ability to be used for its original intended purpose. By increasing the compute, data, and talent required to adapt public models to harmful tasks, self-destructing models have the potential to supplement access controls and other safety mechanisms. We demonstrate a task-blocking mechanism using meta-learning for training a self-destructing model. We find that meta-learning is an essential step in reducing an adversary’s ability to tune a model for a harmful task. Simple adversarial losses [16], often used in current technical strategies, do not significantly reduce the costs of harmful adaptation. We hope that the proposed mechanism forms an initial step toward developing new safe release strategies even under the assumption that model parameters become available to adversaries.

Below, we first review the state of current safe-release approaches and their shortfalls, making the case for a shift in the threat model to make model parameters as harmless as possible even with model access. Second, we define the *task blocking* problem and evaluation metrics as well as *self-destructing models*. Third, we describe an initial algorithm, Meta-Learned Adversarial Censoring (MLAC), for training self-destructing models, evaluating its ability to impede fine-tuning a language model to perform demographic information extraction. Fourth, we identify key directions for future research in the development of self-destructing models.

2 REVIEWING THE RISKS AND MITIGATION STRATEGIES FOR DUAL USES

Foundation models can be and, importantly, *have been* used for harmful purposes unforeseen by their creators in recent years. They have been fine-tuned on hate speech and deployed to 4chan [57]; hackers have released methods to bypass ChatGPT’s safety filters so that it can be used to help generate malware and spam [23]; stable diffusion models have been fine-tuned to generate abusive imagery [28].

Researchers, practitioners, and policymakers are increasingly searching for new ways to prevent machine learning models from being used for these harmful dual purposes—e.g., Solaiman [51], Brundage et al. [9], Whittlestone and Ovadya [59], Shevlane [49], Brundage et al. [8], and many others. Proposed tools have included export controls, controlled or restricted release strategies, using terms of service or licensing, and alignment and fine-tuning for safety. In this section, we briefly examine each of these methods and discuss potential gaps in relying on each strategy. We consider both *structural* methods (e.g., export controls, use of licensing, and access restrictions), and *technical* methods (e.g., alignment fine-tuning).

2.1 Structural Methods

Export Controls. Recently, researchers, such as Flynn [21], have recommended that the United States consider export controls on hardware related to AI, including NVIDIA A100 GPUs, to restrict certain actors’ capacity to train powerful AI models that require substantial computational resources. In 2022, the United States imposed these export controls on AI-related hardware and hardware-manufacturing equipment, following researchers’ suggestions [56].

Such export controls may help restrict pre-training of foundation models—a use case which requires large amounts of specialized hardware, but they do not necessarily restrict inference-time computing and small-scale adaptation once model parameters are available. Even the largest foundation models can now be deployed or adapted on commodity hardware using techniques such as adapters [27], 8-bit [12], and even 4-bit [13] quantization, and other optimization strategies. A recent open-source project was able to run multi-billion parameter LLaMa models on a MacBook Pro with near-equal performance to some state-of-the-art closed-source models, using these techniques.¹ As a result, hardware export controls may no longer be sufficient to prevent the efficient adaptation of foundation models or the large-scale deployment of pre-trained models, nor can they prevent malicious actors located in countries not included in the export control regime.

The U.S. government has also put in place export controls on certain *software* and *models* with specific harmful dual uses. For example, in a 2020 rulemaking, the Department of Commerce Bureau of Industry and Security (BIS) restricted export of software that can be used for automated geospatial analysis. Under this regulation the model is controlled if it meets four criteria: (1) it provides a graphical user interface to identify objects in geospatial imagery; (2) it “reduces pixel variation by performing scale, color, and rotational normalization on the positive samples”; (3) it “[t]rains a Deep Convolutional Neural Network to detect the object of interest from the positive and negative samples”; (4) it “[i]dentifies objects in geospatial imagery using the trained Deep Convolutional Neural Network by matching the rotational pattern from the positive samples with the rotational pattern of objects in the geospatial imagery.” But such highly specific export controls do not cover general-purpose foundation models (and associated training software). In fact, a recent demonstration showed how to adapt a CLIP model [44] exactly for analyzing satellite imagery in an easy way using all open-source software [2]. Flynn [21] argued that applying export controls to general-purpose foundation models would be ineffective due to the ease of violating export controls through the same mechanisms as software piracy, as well as the harmful impacts to innovation that such restrictions could have.

Overall, while export controls may be effective in restricting access to large-scale chipsets or certain software, once adversaries can gain access to open-source (or leaked) foundation model parameters they can be readily adapted to harmful dual-uses.

Access Control. Controlled release or restricted access strategies are another set of structural mechanisms that can supplement export controls and reduce malicious actors’ ability to access models [41, 49, 52].

One such approach is to make the model accessible only by agreement. This involves vetting potential users and requiring them to sign a restrictive terms of service before gaining access to the model. For instance, Meta’s OPT-175B model and Llama both employ this approach [53, 61, 62]. This access restriction approach is attractive as it does not require hosting any centralized infrastructure for serving model queries. It only requires one-time vetting of the users requesting model access. However, as evidenced by the recent Llama

¹<https://github.com/ggerganov/llama.cpp>

Approach	Examples	Challenges
Export Controls	United States Export Controls on AI hardware API-only access, Release by request/agreement	Imprecise, reduced hardware costs, open-source models Open-source models, leaks, monitoring difficulties
Controlled Release		
Licensing	OpenRAIL	Requires monitoring and enforcement action, leaks
Filtering, Alignment	Reinforcement Learning from Human Feedback	Can be bypassed by fine-tuning or prompt engineering

Table 1. A review of current or proposed approaches to safe foundation model release.

model leak onto BitTorrent [58] and HuggingFace,² this approach is susceptible to unauthorized dissemination, effectively negating access control efforts.

Another approach is to never release the model at all, but provide access via an application programming interface (API). Many companies, such as OpenAI, Anthropic, Cohere, and AI21 adopt this approach to protect their trade secrets and prevent harmful dual uses. This approach prevents direct access to model weights, preventing uncontrolled dissemination and retaining the ability to cut off access to malicious users at any time. However, this approach requires monitoring of API usage to detect and revoke access when abused, as well as considerable resources to maintain. Providing such an API may not be possible for researchers and entities without access to centralized model-hosting infrastructure.

Additionally, as open-source efforts continue to match the performance of these closed-source models, the effectiveness of any access control approaches may decrease. Access control approaches require all model creators capable of training similarly capable foundation models to be in agreement on the mechanism for release. If one equally-capable foundation model is available as open-source, malicious actors can simply turn to this alternative.

Terms of Service/Sale (ToS) and Licenses. Closely tied to access controls are licensing agreements to prevent harmful dual-uses. These agreements place restrictions on who can use the model, for what purpose, and in what format. For example, OpenRAIL [18] and similar licenses impose several usage limitations to prevent users from using the model for defamation, spreading disinformation, providing medical advice, or for use by law enforcement. Such terms of service (ToS)-based approaches are also used in other settings, such as by Boston Dynamics, which prohibits modifying its robots for lethal capability and reserves the right to prevent any misuse.³

However, relying solely on licensing agreements assumes that malicious actors would respect them and that legal action against violators is possible. Unfortunately, this approach faces several challenges. Firstly, harmful actors may be located in countries that do not enforce licensing requirements. Further, it may be challenging to identify malicious actors and issue a cease-and-desist request. Finally, model creators may not have the resources to monitor and enforce compliance with licensing agreements.

Overall, licensing agreements face the same challenges as other structural restrictions. They require the resources, and international reach, for enforcement.

2.2 Technical Strategies

Unlike structural strategies, we classify *technical strategies* as those that modify foundation models directly to make it more difficult to use them for harmful purposes. Existing technical strategies focus on tuning models to prevent them from outputting harmful content at inference time or adding content filters to block potentially harmful outputs.

Safety Filters. Some models come with safety filters that scan model outputs for harmful content and then redact the output. Stable Diffusion models use this approach to replace offensive content generated by the model with a blank image by default [48]. However, for open-source models safety filters can simply be removed by deleting a few lines of code. This has led users on Reddit to post tutorials like “How to remove [Stable Diffusion’s] safety filter in 5 seconds.”⁴ Other researchers have noted that the filter itself is easily bypassed even without access to directly modify the code [45]. While safety filters can be effective and integral parts of a safe model release, they are more effective when coupled with other structural mechanisms like restricted or API-only model access.

Safety Tuning and Alignment. Alternative approaches such as reinforcement learning from human feedback tune the model itself to be less harmful [3]. Sometimes these approaches fall into a larger class of methods under the moniker *AI alignment*. Since these methods directly train the model to be more difficult to use for harmful purposes at inference time, they are an essential part of a safe release strategy—either for open-source models or for models coupled with a structural release restriction. Though they make the model parameters more difficult to use for harmful tasks, they can be bypassed in two ways.

First, prompt engineering can be used to put models in a state that nonetheless allows them to be used for harmful purposes. For example, hackers now sell prompts and methods to bypass alignment processes and filters for OpenAI’s series of models [23]. This allows would-be malicious actors to generate phishing emails and malware with the model, despite its use restrictions.

Second, open-source models can be fine-tuned to remove these restrictions. In one such instance, the open-source GPT-J model was fine-tuned on 4chan data (mainly consisting of toxic content and hate speech) and deployed to post to the forum [57].

In the remainder of this work, we describe and evaluate an approach to mitigating this second method of bypassing existing technical model protections.

²<https://twitter.com/ClementDelangue/status/1632948540245671936>

³<https://twitter.com/BostondDynamics2021/status/1362921918781943816>

⁴https://www.reddit.com/r/StableDiffusion/comments/wv2nw0/tutorial_how_to_remove_the_safety_filter_in_5/

2.3 The Need For New Technical Mitigation Strategies

The strategies discussed above are individually imperfect; however, each contributes to increasing the costs of successfully co-opting foundation models for harmful dual uses. As access to increasingly capable models becomes commonplace—either through leaks or open-source releases—it is crucial to ensure that the underlying model parameters themselves are optimized for safety as a last line of defense. Structural barriers, such as access restrictions and terms of service, can become ineffective as model weights are distributed through services like BitTorrent.

As regulators recognize the potential dangers associated with increasingly capable systems, it is becoming evident that they will take action to address the risks. One E.U. AI Act proposal would see liability placed on open-source models, incentivizing restricted access approaches. Others argue that such a move would stifle innovation and make it more difficult to develop safer overall models [6, 17, 25]. As Black et al. [6] write, “open access to [FMs] is critical to advancing research in a wide range of areas—particularly in AI safety, mechanistic interpretability, and the study of how [FM] capabilities scale.” Yet while more widely available FMs certainly enable greater accessibility, auditability, and understanding of these powerful models, making FMs widely available for downstream adaptation without restriction comes at some cost to safety.

Despite the benefits of open-source releases, if open-source models are regularly adapted for harmful purposes, the pendulum of regulation could swing toward the more restrictive regime as regulators look to available structural tools like access restrictions. To supplement the policy options available to regulators, and to increase the safety of foundation models by default, we encourage more research to expand the toolbox of technical approaches to ensure that model parameters are as safe as possible, even when they are leaked or openly available. We introduce a new class of methods for this toolbox: *task blocking for self-destructing models*. These methods are not perfect, but add another layer of protection when combined with other approaches.

3 TASK BLOCKING & SELF-DESTRUCTING MODELS

The goal of task blocking is to create models that increase the costs of fine-tuning on harmful downstream tasks such that an adversary would rather start from scratch than use the pretrained model, while remaining useful for desired tasks (see Fig. 1). The resulting models are “self-destructing models” which impede adaptation on harmful dual-uses by increasing the costs of the harmful use. In this section, we more precisely define our problem setting and describe an initial algorithm for it.

3.1 The Task Blocking Problem

We assume that an adversary aims to adapt a pretrained model π_θ (where θ are model parameters of model π) to a harmful task, searching for the best adaptation procedure f among a set of adaptation procedures \mathcal{F} in order to find the one that maximizes harmful task performance. Adaptation procedures in \mathcal{F} may include simple fine-tuning, a hyper-parameter search over fine-tuning procedures, as well as other more advanced adaptation mechanisms that we

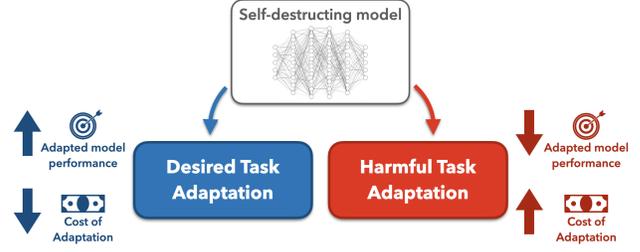


Fig. 1. An ideal self-destructing model would boost performance and reduce adaptation costs relative to training from scratch only for desired tasks, while impeding learning of harmful tasks.

leave to future work. The goal of task blocking is to produce a self-destructing model with parameters $\tilde{\theta}$, which performs similarly to a standard pre-trained model on a set of desired tasks while being more costly to successfully adapt to harmful tasks.⁵

We define two regimes to increase costs: (1) increase data costs by decreasing sample efficiency; (2) increase compute costs by slowing convergence of the training process.

Data Costs. In the first regime, we assume that the adversary has little data to adapt an FM to their harmful task and that the cost of gathering more data is high. A hallmark trait of traditional FMs is effective few-shot adaptation, learning rapidly from small, fixed-sized datasets. A self-destructing FM, on the other hand, should provide few-shot performance comparable to a randomly initialized model. We define the *few-shot performance improvement* of an FM with parameters θ as the performance gain over a randomly initialized model, both with a fixed adaptation procedure search budget. This can be represented as the following formula:

$$\mathcal{E}_{data}^n(\theta) = \max_{f \in \mathcal{F}} \mathcal{M}(f(\theta, D_n)) - \max_{f \in \mathcal{F}} \mathcal{M}(f(\theta^r, D_n)), \quad (1)$$

where \mathcal{M} is the performance metric (where higher is better), n is the number of data points available, D_n is an adaptation dataset of n examples from the task of interest, and θ^r is a randomly-initialized model. $f \in \mathcal{F}$ is an adaptation procedure drawn from a fixed distribution. The size of \mathcal{F} loosely corresponds to the adversary’s resource budget for adaptation. Note that the max in Equation 1 encapsulates hyperparameter optimization over the adaptation distribution. $\mathcal{E}_{data} = \frac{1}{N} \sum_n \mathcal{E}_{data}^n$ is the average sample-wise regret between the FM parameters θ and a random re-initialization θ^r after each follows the same adaptation procedure $f(\cdot)$ on a fixed-sized dataset D_n . An ideal self-destructing model has $\mathcal{E}_{data} \leq 0$, meaning the model is no more data efficient than a randomly-initialized model for the (presumably harmful) task of interest.

Compute Costs. If data is cheap or plentiful, it may be difficult to prevent an adversary from learning the task since perhaps even a random model can learn the task with the amount of data available. In this data regime (large amount of cheap data), the benefit of an

⁵While the goal of a self-destructing model is to reduce performance on harmful tasks after fine-tuning, it should enable high quality *predictions* or *fine-tunability* for desired tasks. Our experiments explore the prediction goal, and we leave exploration of preserving fine-tunability for future work.

FM is improved compute efficiency, rather than increased accuracy. Here, we would define the FM’s *compute cost improvement @p* as the amount of compute saved by using the FM over a randomly initialized model to achieve performance p , where p may measure accuracy, loss, or another metric and compute could be measured in FLOPs, train steps, hyperparameters searched, wall clock time, etc. While in the previous setting, we fix the *dataset size* and blocking aims to reduce performance, in this setting, we fix the *performance* and blocking aims to increase compute costs. The goal of task blocking in this case is to prevent any compute cost improvement over a random initialization when adapting the self-destructing model to a harmful task, while retaining compute cost improvement for desired tasks. Formally, compute cost improvement @p is given as

$$\mathcal{E}_{\text{compute}}^p(\theta) = C(\mathcal{F}, \theta^r, p) - C(\mathcal{F}, \hat{\theta}, p) \quad (2)$$

where C measures the compute cost of applying adaptation procedures from family \mathcal{F} to random parameters θ^r or FM parameters θ until a model with performance level p is found.

However, for the purposes of this work, we focus on data costs, studying methods for reducing few-shot performance improvement for harmful tasks. We leave analysis of compute cost improvement reduction to future work.

Defining Harmful Dual Uses. A large body of work has pointed to inherently harmful uses that FM creators may wish to block: from creating neurotoxins [55] to race detection [38]. In our work we assume that a harmful dual use is *known* and *defined*. That is, the self-destruct mechanism will have data to approximate the dual use and actively encode a mechanism to block it. This requirement inherently requires a normative definition of harmful dual uses. As in other threat modeling exercises and mechanisms for removing harmful content from models, model creators will have to identify the set of tasks to be blocked. Creating self-destructing models may impede their use for harmful purposes counter to the model creator’s values, but it is up to the model creator to determine those values. While defining harmful tasks *a priori* may be difficult, this work reflects a “red teaming” approach to harm prevention, common in security contexts. That is, model creators play the role of an adversary to identify and prevent harms. This can function as a complement to other access control methods, providing more confidence that certain known harmful tasks are blocked.

Relationship to Other Technical Safety Mechanisms. Reinforcement learning from human feedback (RLHF), and other similar approaches, have been used to mitigate the harms that model can have at inference time [3]. While RLHF aims at ensuring that agents are as harmless as possible at inference time, the goal of self-destructing models and task blocking is to make it difficult to undo these safety mechanisms and co-opt the model even with access to model parameters and adaptation. These are complementary approaches and can be used concurrently to make the model parameters as safe as possible overall. Essentially, the aim is to maintain the model’s harmlessness for as long as possible, even when an adversary has direct access.

3.2 Meta-Learned Adversarial Censoring

- 1: **Input:** pretrained model $m = w_d \circ \pi_\theta$, desired task dataset D_d , harmful task dataset D_h , adaptation methods $\tilde{\mathcal{F}}$, adaptation steps K , learning rates η, η_h, η_d
- 2: **Initialize:** Adversarial harmful task head w_h and learning rate α_h , with $\phi = \{w_h, \alpha_h\}$; initial blocked params $\tilde{\theta} \leftarrow \theta$
- 3: **for** n steps **do**
- 4: Sample adaptation procedure $\tilde{f}_k \sim \tilde{\mathcal{F}}$
- 5: Sample data batches $b_d \sim D_d, \{b_h^k\} \sim D_h, b_h \sim D_h$
- 6: $\{\theta_k\}, \{w_h^k\} \leftarrow \tilde{f}_k(w_h \circ \pi_{\tilde{\theta}}, \{b_h^k\}, \alpha_h)$ // do inner loop
- 7: $\ell_k^h = \mathcal{L}_h(w_h^k \circ \pi_{\theta_k}, b_h), \forall k$ // outer loop harmful NLLs
- 8: $\ell^d = \mathcal{L}_d(w_d \circ \pi_\theta, b_d)$ // desired NLLs
- 9: $\tilde{\theta} \leftarrow \tilde{\theta} - \eta \nabla_\theta \left(\ell^d - \frac{1}{K} \sum_k \ell_k^h \right)$ // update blocked model
- 10: $\phi \leftarrow \phi - \eta_h \frac{1}{K} \sum_{k=1}^K \nabla_\phi \ell_k^h$ // update adversarial params
- 11: $w_d \leftarrow w_d - \eta_d \nabla_{w_d} \ell^d$ // update desired task head
- 12: **end for**

algorithm 1. MLAC Training Procedure

To prevent successful adaptation of pretrained models to harmful tasks, we describe *Meta-Learned Adversarial Censoring (MLAC)*, a meta-training procedure that aims to eliminate any useful information about the harmful task in the model’s parameters *even after fine-tuning on that task*. Given a desired task dataset D_d and harmful task dataset D_h , MLAC learns a feature extractor $\pi_{\tilde{\theta}}$ that is effective for the desired task but cannot be effectively used or efficiently fine-tuned to perform the harmful task.

In the *inner loop* of each meta-training step, the feature extractor and an adversarially learned prediction head w_h are adapted to the harmful task with several steps of gradient-based adaptation with an adversarially learned learning rate α_h . The adaptation procedure \tilde{f} used at each meta-training step is sampled from $\tilde{\mathcal{F}}$, a proxy for the true adversary’s adaptation class \mathcal{F} . In this case, we narrow $\tilde{\mathcal{F}}$ to be different fine-tuning approaches with close-to-optimal hyperparameters (e.g., Adam for K steps and learning rate α_h). In the *outer loop*, the adversarial parameters $\phi = \{w_h, \alpha_h\}$ are trained to minimize the harmful task negative log likelihoods of the adapted models, while the blocked parameter initialization $\tilde{\theta}$ are trained to maximize the harmful task negative log likelihoods of the adapted models. We also must counteract the self-destruct mechanism with something that will prevent unlearning of the entire network. In this work, we simply optimize for a given desirable task as the counter-balance by minimizing ℓ^d , which updates both the desired task head w_d and the representation parameters $\tilde{\theta}$. See Algorithm 1 for the complete training procedure. Note that in practice, we use Adam rather than SGD in the outer loop to optimize $\tilde{\theta}$, adversarial parameters ϕ , and desired task output head w_d . We use higher [24] for implementing the bi-level meta-learning process.

Calibration. We also add another mechanism to strengthen the inner-loop adversary. In binary classification tasks, maximizing the loss of the harmful task may lead to a degenerate optimum where labels are flipped, which leaks information about the harmful task. To prevent this outcome, we also optimally calibrate the logits via a simple linear projection (w) solved via differentiable convex

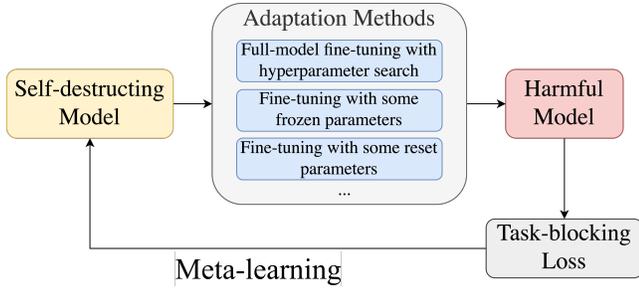


Fig. 2. High-level visualization of the meta-learning process.

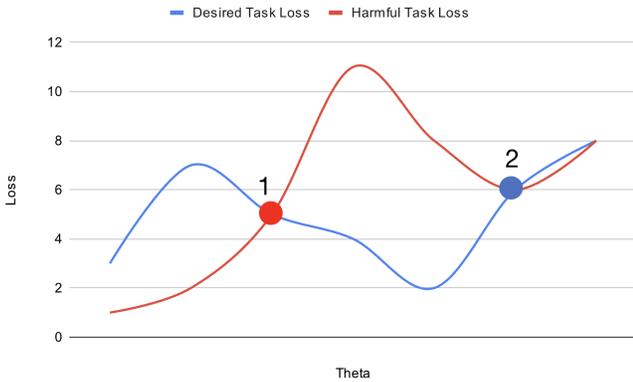


Fig. 3. High-level optimization perspective of the MLAC procedure. A foundation model placed in point 1 would easily be tuned via gradient descent for both the harmful task loss and the desired task loss global optimum. On the other hand a foundation model in point 2 would easily reach the desired task optimum, but is more likely to be stuck in a local optimum for the harmful task.

optimization [1, 15]. Thus at step k of the inner loop we solve the maximum likelihood problem:

$$w_c^k = \underset{W}{\operatorname{argmax}} \sum_i^{[b_h]} \left[\operatorname{logsoftmax} \left[\left(W \circ m^k \right) (x_i) \right]^T y_i \right] \quad \text{s.t. } -1 \leq W \leq 1, \quad (3)$$

where $m^k = w_h^k \circ \pi_\theta^k$ is the blocked model after k steps of adaptation using the adversarial harmful task head and learning rate. Thus this projection updates line 7 of Alg. 1 to $\ell_k^h = \mathcal{L}_h(w_c^k \circ w_h^k \circ \pi_{\theta^k}, b_h)$. We also refer to calibration as *head adjustment*, as it essentially refines the linear function computed by the final output head.

High-level Intuition. Figures 2 provides a visualization of this method. At each step, the self-destructing model samples from possible adaptation methods that could be used to adapt the model to a harmful dual use. This multi-step loss is then inverted in a meta-learning step to prevent the model from being easily adapted in this sampled fashion.

From an optimization perspective, the goal is to identify a parameter space where adaptation to desired tasks is relatively simple via standard adaptation techniques, but the same part of the parameter space might be a low-utility local optimum or saddle-point that is

more difficult to escape for the harmful task. This can be seen as a simplified visualization in Figure 3. Of course, adaptation methods can be created to reset parts of the network such the global harmful optimum can be recovered (in the extreme resetting most of the network to escape the local optimum). However, this will decrease the utility of the expensive pre-training and increase the costs to adversaries, adding another tool in the toolkit against harmful dual-uses.

4 EXPERIMENTS

The goal of our experiments is to assess the ability of MLAC and several simple baselines to reduce the few-shot performance improvement of a pre-trained FM.⁶

4.1 Dataset.

In a demonstrative experiment, we utilize an existing dual-use dataset in the de-biasing literature, “Bias in Bios” [10]. The dataset consists of professional biographies. Each biography has a label that is the gender identity of the biography’s subject as well as the profession being discussed. We split the data into a train, validation, and evaluation set. We consider the “desirable task” for which we want to maintain good performance as the profession detection task. We wish to block the gender identification task. On the original dataset, we find that a random model can learn gender classification to over 90% accuracy with only 10 examples, leaving only marginal ability for an FM to improve in data efficiency. Thus, to make the FM more beneficial, we replace all pronouns with “they/their,” similar to the censored dataset in the original data. While this task pair has traditionally been used for de-biasing, and while we mainly use this as an initial demonstration, there may be valid reasons for preventing an adversary from detecting demographic information from text. In countries where anti-minority action is common, automated systems that identify demographics may cause serious harm. For example, a country may wish to identify people of a certain religion, sexual orientation, or other identity group in automated and ultimately harmful ways. Or the state may wish to identify another feature that is highly correlated with identity which will lead to the same harms.⁷

4.2 Protocol.

For all experiments, we run 50k steps of MLAC meta-training on the training set. At test time, we take the resulting self-destructing model and run it through a rigorous hyperparameter search to maximize fine-tuning performance on the harmful task. We allow hyperparameter searches with 50 fine-tuning trials, using the tree-structured Parzen Estimator [4] in the hyperopt software package [5]. We search over learning rate, batch size, maximum number of steps, and freezing of intermediate representation layers. For this process, we subsample the validation set to simulate an adversary with a dataset of size N . This subsampled validation set is used as the training set for the adversary. We then use the entire evaluation set to evaluate the adversary’s performance on held-out data and for hyperparameter tuning. We make the conservative assumption that

⁶Code is available at <https://github.com/Breakend/SelfDestructingModels>.

⁷Technology Experts Letter to DHS Opposing the Extreme Vetting Initiative, 2017.

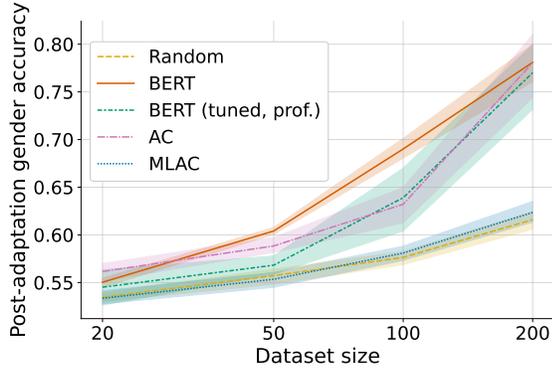


Fig. 4. Harmful task (gender identification) performance after fine-tuning. MLAC shows fine-tuning performance similar to a randomly-initialized model, while adversarial censoring (AC) [16] does not prevent effective fine-tuning. Shading indicates 95% confidence intervals across 6 random seeds.

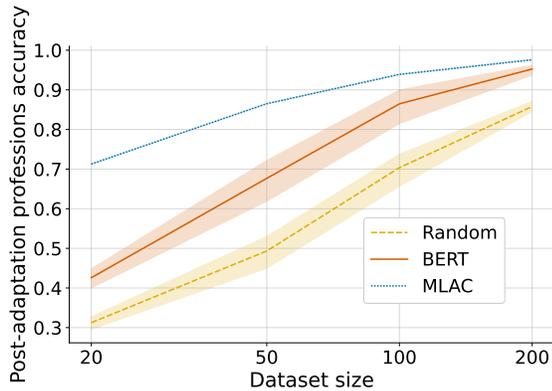


Fig. 5. After fine-tuning the MLAC-blocked model on the desired task, few-shot performance exceeds both BERT and a randomly-initialized model. Note the MLAC objective includes training on the desired task, so this comparison clearly advantages MLAC; nonetheless, it provides evidence that there *exists* a blocked initialization that can be effectively fine-tuned on the desired task. Discovering such an initialization without using desired task data in pre-training is an important direction for future work.

the adversary can perform hyperparameter tuning using the *population*, even if the amount of data for fine-tuning itself is limited. This choice weighs heavily in the adversary’s favor, disadvantaging the self-destruct method. We repeat the hyperparameter search process 6 times with different random seeds and data subsets. This yields confidence intervals over different adversaries training on different subsets of the data.

4.3 Comparisons.

We compare MLAC to the adversarial censoring (AC in Fig. 4) method from Edwards and Storkey [16] as well as a model simply fine-tuned on the desired task (*BERT (fine-tuned)* in Fig. 4). For AC, an adversarial layer is learned on top of representation layers to predict the undesirable task. The gradient is then flipped to destroy

undesirable information in the representation layer. Notably, MLAC with $K = 0$ and with no calibration is equivalent to adversarial censoring. We use a BERT-tiny model as our FM to save on compute costs [14, 54] and use a linear classifier head for the tasks. Note that, as mentioned in Sec. 3.2, we focus on making sure that the professions task is unimpeded, so we directly train on cross-entropy loss as \mathcal{L}_g during MLAC pre-training. For all models, the final achieved performance is retained for the desired professions task (see below and Figure 5).

4.4 Results.

Fig. 4 shows that MLAC returns nearly identical-to-random harmful task performance at all data regimes. Conversely, adversarial censoring (the equivalent of MLAC without calibration and $K = 0$) does not appear to have any effect on post-fine-tuning harmful task performance. Fig. 6 shows the vital role played by the depth of the inner training loop of MLAC, suggesting that *a meta-learning process is genuinely necessary to impede harmful task performance*. To ensure that desired task performance is retained, we evaluate the blocked model on the desired task of profession classification, comparing with fine-tuning a pretrained BERT-tiny model and a random model. Fig. 5 shows the result; MLAC is clearly able to solve the task effectively, surpassing the few-shot performance of BERT-tiny.⁸ Finally, we find that head re-calibration may mildly improve blocking on average when pooled across all inner-loop step configurations (Fig. 7).

5 ETHICAL CONSIDERATIONS AND LIMITATIONS

Before we conclude, we point out several other considerations and limitations.

First, while the goal of our approach is to make models safer overall, we recognize that value judgements will be made in deciding which tasks to block. Sometimes these judgement decisions can themselves encode biases and it requires an approach that takes into account a range of perspectives. Nonetheless, we argue that considering potential harmful dual-uses is an essential part of any modern model release process. Current standard licenses for foundation models already contain a list of restricted tasks [18, 53], but self-destructing models encode this directly into their optimization objective as well.

Second, it is necessary to collect data for harmful tasks to effectively block them. While this draws a direct parallel to security research, red-teaming, and white-hat hacking, there may be risks in aggregating this data. And there may be impacts on the well-being of potential annotators and security research members [35]. Sufficient precautions should be taken to mitigate these harms.

Third, there may be a risk of over-confidence in the self-destructing mechanism. While this paradigm adds a new tool to the safety toolkit, it does not completely prevent manipulation for every harmful task. And just like any other safety tool there will likely be

⁸Recall again that we use the desired task loss to counter-balance the task blocking mechanism, so this is expected. We do however use separate held-out subsets of data for final desired-task tuning and evaluation. As mentioned previously, our goal for the purposes of this initial exploration is to determine whether desired task performance can be retained while blocking a harmful task. Future work should examine generalization for retaining desired task adaptation performance across many desired tasks.

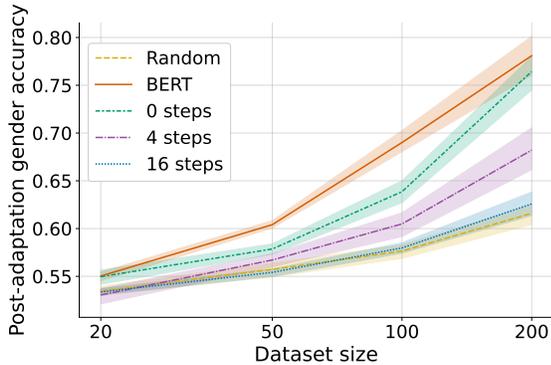


Fig. 6. Evaluation of various inner loop depths during MLAC training. Just 16 steps enables near-random performance, even though the adversary performs up to 1000 steps during fine-tuning.

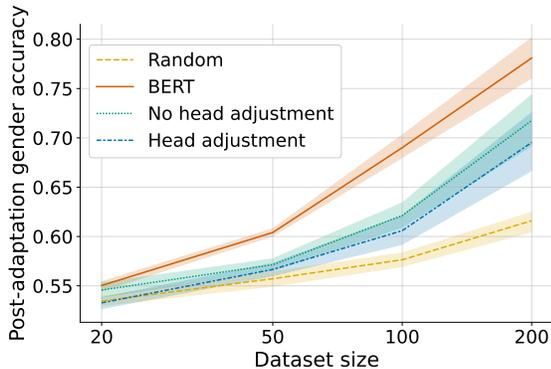


Fig. 7. Ablating optimal adversary prediction calibration (or *head adjustment*) during MLAC training. Using optimally calibrated adversary predictions (modifying line 7 of Alg. 1) modestly improves blocking. Aggregated over 0, 4, and 16 steps.

a back-and-forth where adversaries learn to overcome some techniques. As such, self-destructing models can be combined with other safety mechanisms—structural or technical—to increase the costs of harmful dual-uses.

Fourth, our experiments demonstrate the functionality of self-destructing models in a constrained setting, but further work is needed to scale these approaches to more tasks, larger models, and more complicated settings. We believe this is an exciting new research direction, but requires more work to deploy at scale.

6 RELATED WORK

A number of researchers have sought to address dual use risks by restricting points of control [7, 8, 21, 49, 52, 65], despite there also being substantial benefits to open access [6, 62]. We aim to provide an alternative that allows for open access while still hindering bad actors.

Some work on AI safety has sought mechanisms to prevent agents from learning degenerate behaviors. Orseau and Armstrong [39],

for example, seek to prevent a particular scenario where an agent learns to disable its off-switch so that it continues to collect reward. We on the other hand focus on preventing a different, broader, set of harmful behaviors: adaptation of pretrained models to harmful tasks.

Closely related to our work are methods for de-biasing, editing, or removing harmful content from models. Like domain invariance approaches [22, 31, 60, 63], Edwards and Storkey [16] use an adversarial approach to remove information from representations. Ravfogel et al. [46] and Ravfogel et al. [47] take a similar approach and find a projection on the final output layer of a pretrained model that removes gender-based biases from the model (and prevent recovery of those biases after that projection layer). Pryzant et al. [43] similarly use adversarial methods to remove confounds from representations. Others have created model editing techniques to remove outdated or harmful content from pretrained models [11, 36, 37, 50]. While these other methods generally optimize for the information to be removed from the original model, we optimize for poor performance even *after* adaptation of the original model to a harmful task. This can be accomplished via a meta-learning approach.

In the context of meta-learning, MAML [19] and related algorithms [20, 30, 33, 42, 64] have shown that the desired *post*-fine tuning behavior of a neural network can be effectively encoded in its *pre*-fine tuning network initialization. While existing works have leveraged this ability in order to enable more rapid learning of new tasks, our work encodes a blocking mechanism into a network’s initialization that *prevents* effective adaptation on harmful tasks.

Finally, some scholars have tuned models to be safer by using reinforcement learning from human feedback and other approaches for incorporating human preferences, including Bai et al. [3], Korbak et al. [29], Ouyang et al. [40], and others.

7 CONCLUSION

This work is only a first step in raising the cost for harmful dual uses of pretrained models through task blocking. Future work might expand this study in at least four directions: *scaling* the self-destructing model framework to larger FMs; studying *generalization* of the learned blocking behavior to new (but related) datasets other than the one used during MLAC meta-training; training/evaluating with *stronger adversaries* that incorporate adaptation methods such as prefix tuning [32], adapter layers [26], or others; and evaluating the preservation of desired task *fine-tunability* for out-of-distribution tasks. Future work might also introduce concealed architectural changes that hide self-destruct triggers in the network but are more robust to adversarial mechanisms. We hope self-destructing models can become one tool enabling model developers to share their artifacts while minimizing dual use risks.

ACKNOWLEDGMENTS

We thank Rishi Bommasani, Siddharth Karamcheti, and Jieru Hu for helpful discussion and feedback. PH is supported by an Open Philanthropy AI Fellowship. EM is supported by a Knight-Hennessy Graduate Fellowship. CF and CM are CIFAR Fellows.

REFERENCES

- [1] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and Z. Kolter. 2019. Differentiable Convex Optimization Layers. In *Advances in Neural Information Processing Systems*.
- [2] Artashes Arutunian, Dev Vidhani, Goutham Venkatesh, Mayank Bhaskar, Ritabrata Ghosh, and Sujit Pal. 2021. Fine tuning CLIP with Remote Sensing (Satellite) images and captions. *HuggingFace Blog* (2021). <https://huggingface.co/blog/fine-tune-clip-rsicc>
- [3] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862* (2022).
- [4] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. 2011. Algorithms for hyper-parameter optimization. *Advances in neural information processing systems* 24 (2011).
- [5] James Bergstra, Dan Yamins, David D Cox, et al. 2013. Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms. In *Proceedings of the 12th Python in science conference*, Vol. 13. Citeseer, 20.
- [6] Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. 2022. GPT-NeoX-20B: An Open-Source Autoregressive Language Model. *arXiv preprint arXiv:2204.06745* (2022).
- [7] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021).
- [8] Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan Dafoe, Paul Scharre, Thomas Zeitoff, Bobby Filar, et al. 2018. The malicious use of artificial intelligence: Forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228* (2018).
- [9] Miles Brundage, Shahar Avin, Jasmine Wang, Haydn Belfield, Gretchen Krueger, Gillian Hadfield, Heidy Khlaaf, Jingying Yang, Helen Toner, Ruth Fong, et al. 2020. Toward trustworthy AI development: mechanisms for supporting verifiable claims. *arXiv preprint arXiv:2004.07213* (2020).
- [10] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnamurthy Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *proceedings of the Conference on Fairness, Accountability, and Transparency*. 120–128.
- [11] Nicola De Cao, W. Aziz, and Ivan Titov. 2021. Editing Factual Knowledge in Language Models. *ArXiv abs/2104.08164* (2021).
- [12] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale. *arXiv preprint arXiv:2208.07339* (2022).
- [13] Tim Dettmers and Luke Zettlemoyer. 2022. The case for 4-bit precision: k-bit Inference Scaling Laws. *arXiv preprint arXiv:2212.09720* (2022).
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [15] Steven Diamond and Stephen Boyd. 2016. CVXPY: A Python-embedded modeling language for convex optimization. *The Journal of Machine Learning Research* 17, 1 (2016), 2909–2913.
- [16] Harrison Edwards and Amos Storkey. 2015. Censoring representations with an adversary. *arXiv preprint arXiv:1511.05897* (2015).
- [17] Alex Engler. 2022. The EU’s attempt to regulate open-source AI is counterproductive. *Brookings TechTank* (2022).
- [18] Carlos Muñoz Ferrandis. 2022. OpenRAIL: Towards open and responsible AI licensing frameworks. https://huggingface.co/blog/open_rail.
- [19] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and Yee Whye Teh (Eds.). PMLR, 1126–1135. <https://proceedings.mlr.press/v70/finn17a.html>
- [20] Sebastian Flennerhag, Andrei A. Rusu, Razvan Pascanu, Francesco Visin, Hujun Yin, and Raia Hadsell. 2020. Meta-Learning with Warped Gradient Descent. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rkeiQBFPB>
- [21] Carrick Flynn. 2020. Recommendations on export controls for artificial intelligence. *Centre for Security and Emerging Technology* (2020).
- [22] Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*. PMLR, 1180–1189.
- [23] Dan Goodin. 2023. Hackers are selling a service that bypasses ChatGPT restrictions on malware. *arstechnica* (2023).
- [24] Edward Grefenstette, Brandon Amos, Denis Yarats, Phu Mon Htut, Artem Molchanov, Franziska Meier, Douwe Kiela, Kyunghyun Cho, and Soumith Chintala. 2019. Generalized Inner Loop Meta-Learning. *arXiv preprint arXiv:1910.01727* (2019).
- [25] Philipp Hacker, Andreas Engel, and Marco Mauer. 2023. Regulating ChatGPT and other Large Generative AI Models. *arXiv preprint arXiv:2302.02337* (2023).
- [26] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-Efficient Transfer Learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 2790–2799. <https://proceedings.mlr.press/v97/houlsby19a.html>
- [27] Edward Hu, Yelong Shen, Phil Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv:2106.09685* [cs.CL]
- [28] Tatum Hunter. 2023. AI porn is easy to make now. For women, that’s a nightmare. *The Washington Post* (2023).
- [29] Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Bhalerao, Christopher L Buckley, Jason Phang, Samuel R Bowman, and Ethan Perez. 2023. Pretraining Language Models with Human Preferences. *arXiv preprint arXiv:2302.08582* (2023).
- [30] Yoonho Lee and Seungjin Choi. 2018. Gradient-based meta-learning with learned layerwise metric and subspace. In *International Conference on Machine Learning*. 2933–2942.
- [31] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. 2018. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5400–5409.
- [32] Xiang Lisa Li and Percy Liang. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. *arXiv:2101.00190* [cs.CL]
- [33] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. 2017. Meta-SGD: Learning to Learn Quickly for Few Shot Learning. *CoRR abs/1707.09835* (2017). <http://arxiv.org/abs/1707.09835>
- [34] Percy Liang, Rishi Bommasani, Kathleen A. Creel, and Rob Reich. 2022. The Time Is Now to Develop Community Norms for the Release of Foundation Models. <https://crfm.stanford.edu/2022/05/17/community-norms.html>
- [35] Mantas Mazeika, Eric Tang, Andy Zou, Steven Basart, Jun Shern Chan, Dawn Song, David Forsyth, Jacob Steinhardt, and Dan Hendrycks. 2022. How Would The Viewer Feel? Estimating Wellbeing From Video Scenarios. *arXiv preprint arXiv:2210.10039* (2022).
- [36] Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2022. Fast Model Editing at Scale. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=0DcZxeWfOPt>
- [37] Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-Based Model Editing at Scale. *arXiv preprint arXiv:2206.06520* (2022).
- [38] Parmy Olson. 2022. The Quiet Growth of Race-Detection Software Sparks Concerns over Bias. In *Ethics of Data and Analytics*. Auerbach Publications, 201–205.
- [39] Laurent Orseau and MS Armstrong. 2016. Safely interruptible agents. (2016).
- [40] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155* (2022).
- [41] Aviv Ovadya and Jess Whittlestone. 2019. Reducing malicious use of synthetic media research: Considerations and potential release practices for machine learning. *arXiv preprint arXiv:1907.11274* (2019).
- [42] Eunbyung Park and Junier B Oliva. 2019. Meta-Curvature. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc.
- [43] Reid Pryzant, Kelly Shen, Dan Jurafsky, and Stefan Wagner. 2018. Deconfounded lexicon induction for interpretable social science. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. 1615–1625.
- [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [45] Javier Rando, Daniel Paleka, David Lindner, Lennard Heim, and Florian Tramèr. 2022. Red-Teaming the Stable Diffusion Safety Filter. *arXiv preprint arXiv:2210.04610* (2022).
- [46] Shauli Ravfogel, Michael Twiton, Yoav Goldberg, and Ryan Cotterell. 2022. Linear Adversarial Concept Erasure. *arXiv preprint arXiv:2201.12091* (2022).
- [47] Shauli Ravfogel, Francisco Vargas, Yoav Goldberg, and Ryan Cotterell. 2022. Adversarial Concept Erasure in Kernel Space. *arXiv preprint arXiv:2201.12191* (2022).
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *arXiv:2112.10752* [cs.CV]

- [49] Toby Shevlane. 2022. Structured access to AI capabilities: an emerging paradigm for safe AI deployment. *arXiv preprint arXiv:2201.05159* (2022).
- [50] Anton Sinitin, Vsevolod Plokhhotnyuk, Dmitry Pyrkin, Sergei Popov, and Artem Babenko. 2020. Editable Neural Networks. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=HJedXaEtvS>
- [51] Irene Solaiman. 2023. The Gradient of Generative AI Release: Methods and Considerations. *arXiv preprint arXiv:2302.04844* (2023).
- [52] Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askill, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. 2019. Release strategies and the social impacts of language models. *arXiv preprint arXiv:1908.09203* (2019).
- [53] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint arXiv:2302.13971* (2023).
- [54] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2020. Well-Read Students Learn Better: On the Importance of Pre-training Compact Models. <https://openreview.net/forum?id=BJg7x1HFvB>
- [55] Fabio Urbina, Filippa Lentzos, Cédric Invernizzi, and Sean Ekins. 2022. Dual use of artificial-intelligence-powered drug discovery. *Nature Machine Intelligence* 4, 3 (2022), 189–191.
- [56] U.S. Department of Commerce. 2022. Implementation of Additional Export Controls: Certain Advanced Computing and Semiconductor Manufacturing Items; Supercomputer and Semiconductor End Use; Entity List Modification. *Federal Register* 87 (2022), 62186. <https://www.federalregister.gov/documents/2022/10/13/2022-21658/implementation-of-additional-export-controls-certain-advanced-computing-and-semiconductor>
- [57] James Vincent. 2022. YouTuber trains AI bot on 4chan’s pile o’bile with entirely predictable results. *The Verge* (2022).
- [58] James Vincent. 2023. Meta’s powerful AI language model has leaked online – what happens now? *The Verge* (2023).
- [59] Jess Whittlestone and Aviv Ovadya. 2019. The tension between openness and prudence in AI research. *arXiv preprint arXiv:1910.01170* (2019).
- [60] Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. 2022. Improving out-of-distribution robustness via selective augmentation. *arXiv preprint arXiv:2201.00299* (2022).
- [61] Susan Zhang, Mona Diab, and Luke Zettlemoyer. 2022. Democratizing access to large-scale language models with OPT-175B. *Meta AI* (2022).
- [62] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. OPT: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068* (2022).
- [63] Fan Zhou, Zhuqing Jiang, Changjian Shui, Boyu Wang, and Brahim Chaib-draa. 2020. Domain generalization with optimal transport and metric learning. *arXiv preprint arXiv:2007.10573* (2020).
- [64] Luisa M Zintgraf, Kyriacos Shiarlis, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. 2019. Fast Context Adaptation via Meta-Learning. *Thirty-sixth International Conference on Machine Learning (ICML 2019)* (2019).
- [65] Remco Zwetsloot, James Dunham, Zachary Arnold, and Tina Huang. 2019. Keeping Top AI Talent in the United States. *Center for Security and Emerging Technology* (December 2019).