

# 000 001 002 003 004 005 STRATEGIC AI SABOTAGE: STATE ATTACKS ON 006 ADVANCED SYSTEMS' DEVELOPMENT 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019

## 020 ABSTRACT 021

022 Much attention has been given to the possibility that states will attempt to steal  
023 the model weights of advanced AI systems. We argue that in most situations, it is  
024 more likely that a state will attempt to sabotage the training of the models under-  
025 pinning these systems. We present a threat modelling framework for sabotage of  
026 AI training, including both the necessary technical background and a taxonomy of  
027 strategic considerations and attack vectors. We then use this to examine different  
028 attacks and assess both their technical plausibility and the mitigations required to  
029 defend against them.  
030  
031

## 032 1 INTRODUCTION 033

034 As AI systems grow in strategic importance, state-level actors seeking to degrade or corrupt the ca-  
035 pabilities of their adversaries may choose to attack these systems through different means. This has  
036 received a large amount of attention in the last year or so, with J-AISI recently releasing an illus-  
037 trative overview of the different ways that AI systems could be attacked (Kiribuchi et al., 2025) and  
038 NIST producing an extensive taxonomy of attacks and mitigations relevant to ML systems (Vassilev  
039 et al., 2025), among others. Both of these reports—whilst not explicitly calling threats to model  
040 integrity ‘sabotage’—highlight various different attacks on the training process and make important  
041 technical contributions as to those attacks’ classification, including within the more general cybersecurity  
042 domain.<sup>1</sup> However, they exclusively consider direct, technical vectors and do not go into great  
043 detail—inevitably, since the scope of these reports is so wide—and also do not place these attacks  
044 within a wider strategic context. More narrowly-focused research has mostly focused on attacks on  
045 deployed systems,<sup>2</sup> and the more detailed security work examining the creation of models has not  
046 treated the sabotage threat model as a cohesive whole but instead has focused on the customization  
047 and fine-tuning of base models<sup>3</sup> or the specific threat of data poisoning.<sup>4</sup>

048 This paper takes a deeper look at the sabotage threat model, broadly construed, as it is relevant to the  
049 initial creation of the models underpinning advanced AI systems.<sup>5</sup> It makes three key contributions:

050 <sup>1</sup>A more detailed placing of *this* work within the existing cyber- and AI-security literature is omitted for  
051 brevity, but is substantively the same as can be found in Section 2 of Kiribuchi et al. (2025), which includes  
052 a technical classification of these attacks within the taxonomy of the MITRE Atlas—indeed, this paper can  
053 itself be seen as a direct follow-up to their recent work, being a combination of a more detailed consideration  
054 of some of the attacks they discuss as **C: Model Poisoning, D: Data Poisoning, J: Adversarial Fine-tuning**,  
055 and threat modelling of indirect attacks on organisations. A more detailed comparison of the contributions of  
056 this paper compared to the J-AISI and NIST work is attached as Appendix B.

057 <sup>2</sup>For example, in *Securing AI Model Weights* (Nevo et al., 2024), RAND produced perhaps the most in-  
058 depth and comprehensive analysis relevant to securing AI systems but explicitly placed ensuring the integrity  
059 of training data and model internals outside the scope of the report.

060 <sup>3</sup>See, e.g. Wan et al. (2023), or indeed the wide variety of work on emergent misalignment following its  
061 discovery by Bettleby et al. (2025).

062 <sup>4</sup>Indeed, there has been a great deal of investigation into this specific threat area—far too much to cover here,  
063 but reading the relevant parts of the NIST report as well as Zhao et al. (2025) will provide a good technical  
064 understanding of the current literature; readers should however bear in mind that some data poisoning threat  
065 models do rely on somewhat unrealistic levels of access or are not robust to standard levels of monitoring or  
066 responsiveness by defenders.

067 <sup>5</sup>We also touch on, but treat as largely out of scope, actions more drastic than sabotage (e.g. overt military  
068 action), attacks which occur after the initial creation of a model (i.e. against deployed systems or in fine-tuning  
069 customer-specific models) and specific *consequences* of deploying a corrupted AI model.

- We place the sabotage of advanced AI capabilities in historical context and argue that in several likely near-future scenarios states have the incentive and capacity to use offensive cyber-capabilities to attack the training process.
- We introduce a threat modelling framework in several parts. We provide a technical introduction to the attack surface, then a taxonomy of the different strategic objectives of sabotage, and finally a classification of different attack vectors, examining both direct technical attacks on data ingestion and model training and indirect attacks targeting organisational capacity.
- We use this framework to conduct a threat modelling exercise across the full training pipeline and identify likely attack vectors relevant to the sabotage of AI training. We assess the plausibility of current and near-future exploitation by different threat actors and highlight possible mitigations to prevent this.

## 2 HISTORICAL CONTEXT

In 2009, the Iranian government was advancing its controversial nuclear program amid growing sanctions and pressure. Global opinion—and the pressure of a US populace wary of getting involved in another war—had so far prevented direct military action by Iran’s adversaries. However, their nuclear facilities began to experience a high level of unexplained operational failures. Intelligence reports later revealed that Israel and the United States had secretly developed something unprecedented in modern warfare: Stuxnet, the world’s first cyber-weapon designed to cause physical damage, which applied irregular and damaging acceleration to centrifuges whilst presenting falsified data to the equipment monitoring them. This sophisticated malware delayed Iran’s nuclear ambitions without the risks of overt action.<sup>6</sup>

Historically, states have chosen to sabotage their rivals when facing a significant threat to their security in the following circumstances:

- Diplomatic or economic leverage has proved insufficient, and
- the situation does not yet justify overt military action,<sup>7</sup> and
- when discovery and attribution would fall short of provoking a direct military response or present enough ambiguity to constrain the target from escalating substantially,<sup>8</sup> and
- when the benefit to such an operation outweighs the potential risk to resources and trust from other actors if it were to be discovered.<sup>9</sup>

During the Cold War, both sides conducted sabotage and assassination operations at a level below that which would provoke outright conflict.<sup>10</sup> Stuxnet itself was part of an extended campaign

<sup>6</sup>See Zetter (2014) for a detailed writeup of Stuxnet, its technical details and geopolitical context. However, Slayton (2016) makes a convincing argument that, ultimately, the delay to the program was only a few months; whilst Stuxnet gives us an excellent illustration of states attempting to achieve their geopolitical aims via cyber-enabled sabotage, it perhaps does not give us an *effective* one.

<sup>7</sup>Whether it be due to lack of or comparative imbalance in military strength (*e.g.* the 1979 Israeli sabotage of nuclear equipment being sold to Iraq by the French), geopolitical constraints (*e.g.* the 1985 bombing of a Greenpeace vessel by the French secret service), domestic politics (*e.g.* the 1954 Guatemalan coup d’etat (Holland, 2005) or the Bay of Pigs invasion), or concerns over the risk of escalation (*e.g.* US support for the mujahideen in the 1980s).

<sup>8</sup>Indeed in the context in the context of the United States, ‘covert action’ is defined by legislation such that “the role of the United States Government will not be apparent or acknowledged publicly”; *see* Rosenbach & Peritz (2009) for a discussion of this and when such covert action is permitted.

<sup>9</sup>For a detailed explanation of how the value associated with successful sabotage may be orders of magnitude greater than the amount spent by either side and so sabotage can be favoured even in defense-dominant scenarios, *see* Slayton (2016).

<sup>10</sup>For an overview of KGB doctrine and planning around sabotage, *see* Richterova (2024), although it is worth noting (and comparing to recent decades) that actual physical destruction of infrastructure was extremely limited—likely because of escalation risk, and held in contingency for an actual hot war situation—and perhaps the most impactful sabotage was indirect, involving the feeding of faulty schematics and blueprints to the Soviets in the 1980s (Weiss, 1996).

108 by Israel to sabotage the nuclear ambitions of its adversaries,<sup>11</sup> including assassinations, planting  
 109 of explosives in third countries, and other sub-military actions. More recently, in 2022, 3 of the  
 110 4 Nordstream pipes—taking Russian gas to Europe, its most significant export market<sup>12</sup> and an  
 111 important source of revenue for the Russian government—were blown up without any conclusive  
 112 evidence as to who had done it.<sup>13</sup>

113 In the decades following the Cold War, the traditional boundaries between peace and war have be-  
 114 come increasingly blurred as state and non-state actors<sup>14</sup> embrace so-called “grey-zone warfare”: a  
 115 spectrum of hostile activities that fall deliberately below the threshold of conventional armed con-  
 116 flict. This approach allows nations to pursue strategic objectives through a mix of cyber attacks,  
 117 economic coercion, disinformation campaigns, proxy forces and anonymous military actions, and  
 118 political interference while maintaining plausible (or at least some) deniability and avoiding the  
 119 consequences of instigating open warfare. Russia’s annexation of Crimea in 2014 created popular  
 120 awareness of this strategy as it used unmarked and unacknowledged troops, social media manipu-  
 121 lation and economic pressure to take control of Ukrainian territory without ever formally declaring  
 122 war against the country. Similarly, China’s so-called “salami-slicing” actions, including the con-  
 123 struction of artificial islands and military bases in the South China Sea, the harassment of civilian  
 124 and military vessels from the Philippines, Vietnam, and other countries, and the reckless use of mil-  
 125 itary assets<sup>15</sup> challenge international law and convention without quite crossing the threshold that  
 126 would trigger a military response from the United States and its allies.<sup>16</sup>

127 Traditional executive processes and legal frameworks are currently poorly equipped to respond to  
 128 threats that do not quite fit into traditional domains of warfare, while the ever-increasing digital  
 129 transformation of society has created new vulnerabilities to cyber attacks and informational warfare.  
 130 The last 3 years of war with Ukraine has been accompanied by a corresponding aggressive and  
 131 extensive campaign of sabotage<sup>17</sup> by Russia against the West—to which it has failed to develop an  
 132 effective counter (Jones, 2025)—combining a persistent offensive cyber-strategy with more direct  
 133 bombing and arson of critical infrastructure<sup>18</sup> and attempts to assassinate executives of western  
 134 companies.<sup>19</sup>

135 It may be inferred that, for some actors, sabotage would be seen as an attractive option to disrupt  
 136 strategically relevant technological advancements. Given its strategic importance and potential dual-  
 137 use applications, it is natural to consider how this context might interact with the rapid increase of  
 138 investment in, and development of, advanced capabilities in artificial intelligence. The next two  
 139 sections explore current academic thinking on this and relevant potential near-future scenarios.

## 140 2.1 MAIM

142 In *Superintelligence Strategy*, Hendrycks et al. (2025) introduced a three-part plan to manage the  
 143 development of artificial superintelligence (ASI) covering deterrence, nonproliferation and compet-  
 144 itiveness. A key contribution was the introduction of the concept of Mutual Assured AI Malfunc-  
 145

146 <sup>11</sup>Known as the Begin doctrine—*see* Talbot (2023) for historical context—although this is only part of their  
 147 perhaps uniquely effective use of sabotage to degrade the capabilities of their adversaries; *see, e.g.* their recent  
 148 supply-chain compromise of Hezbollah’s pager network (Doran, 2024).

149 <sup>12</sup>Hooper et al. (2022).

150 <sup>13</sup>Speculation included Ukraine, Poland, or even the United States, although recent news coverage suggests  
 151 it was likely Ukrainian operatives.

152 <sup>14</sup>Or at least, non-state actors with some plausible deniability as to not being state-directed, *see, e.g.* Chinese  
 153 and Russian commercial shipping vessels repeatedly ‘accidentally’ cutting undersea infrastructure in the Baltic  
 154 and around Taiwan (van Soest, 2025; Daud et al., 2024).

155 <sup>15</sup>The Hainan Island incident is well-known (Donnelly, 2004); *see also, e.g.* a recent incident involving  
 156 Chinese sonar and Australian divers.

157 <sup>16</sup>For a more comprehensive detailing of these tactics, *see* Helmus et al. (2024) p. 25 and Chapter 2 more  
 158 generally; Derek Grossman has also written two relevant opinion pieces (Grossman, 2024b;a).

159 <sup>17</sup>Note that whilst there is an ongoing military campaign against Ukraine itself, Russia has evidently not  
 160 decided that the situation justifies overt military action against the targets of their sabotage, i.e. Poland, Czechia,  
 161 UK, etc., and has chosen to attempt sabotaging their ability to help Ukraine instead.

162 <sup>18</sup>*See, e.g.*, the growing list of explosions at European munitions factories (?) and the recent destruction of a  
 163 railway link between Poland and Ukraine (?).

164 <sup>19</sup>Covered at some length in this writeup by CNN.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
tion, or *MAIM*. With this, they outlined "a deterrence regime resembling nuclear mutual assured destruction [...] where any state's aggressive bid for unilateral AI dominance is met with preventive sabotage by rivals". In practice, this predicts a dynamic where the training of advanced AI systems is prevented by *covert sabotage*, *overt cyberattacks*, and (less relevant for us) *direct kinetic action*, resulting in the substantial slowing or stopping of the creation of such systems.

168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
167  
The academic response to this has been mixed; Abecassis (2025), of MIRI, responded, claiming that in practice this dynamic would not apply because of the "lack of effective, monitorable, and clear red lines". RAND also responded, with Rehman et al. (2025) suggesting that the game theory assumptions do not hold in practice and MAIMing actions are highly likely to be escalatory. To determine whether these criticisms are reasonable, and whether we should expect the MAIM dynamic to actually hold, we need to understand how sabotage might occur and know the plausibility and scope of particular attacks, not least because if MAIM is implausible we must seek other ways of ensuring stability.

## 2.2 LIKELIHOOD OF SABOTAGE

178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
178  
Various world leaders and governments have indicated that they see AI as a potential path to strong military and economic advantages.<sup>20</sup> As such, there has been substantial academic attention given to the strategic implications of the development of advanced AI systems in the domain of international relations, geopolitics and the military balance of power.<sup>21</sup> Correspondingly, there has been substantive discourse around the possibility that adversaries will attempt to *steal* the model weights underpinning these systems and thus reduce the relative military or economic advantage gained by the actor which has developed advanced AI capabilities.<sup>22</sup> We consider this to be a likely course of action in some circumstances—specifically, when one party leads in some area of AI research, but another is capable of 'fast-following', has the capacity to benefit substantially from the deployment of the stolen models, and has viable pathways to obtain them without excessive risk or resource expenditure—but there are several more likely scenarios where we expect a state actor<sup>23</sup> to prioritise sabotage over theft of AI systems instead (for a more general look at the different forms of sabotage and their strategic implications, see Section 3.2). We argue that these scenarios include cases when the actor:

- 192  
193  
194  
195  
196  
197  
198  
199  
200  
192  
• does not have the compute to deploy advanced models for strategic purposes, or alternatively the economic integration, and believes it is imperative to stop rivals who would have this capacity.<sup>24</sup>
- 201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
• has their own, better models, which they believe to be secure and want to lower the possibility of rivals catching up.
- 201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
• believes that their rival could use the model being developed to attain a decisive strategic advantage or otherwise take significantly destabilising actions.<sup>25</sup>
- 201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
• wishes to subvert the model before it is deployed, causing it to malfunction or underperform in specific circumstances.

<sup>20</sup>For example, Putin, Biden, Trump, Vance, Xi, among others.

<sup>21</sup>See, e.g., Pavel et al. (2023), Scharre (2023), as well as the recent UK Strategic Defence Review (UK Ministry of Defence, 2025), which highlighted several different candidates for near-future military use of AI capabilities.

<sup>22</sup>For a comprehensive understanding of this, see RAND's seminal report *Securing AI Model Weights* (Nevo et al., 2024).

<sup>23</sup>Of course, there are also *commercial* reasons that companies might sabotage their rivals, which we do not consider.

<sup>24</sup>After all, many of the states commonly acknowledged to invest the most heavily in offensive cyber capabilities—including Russia, North Korea, Iran, and to a lesser extent the United Kingdom and Israel—are not engaged in a compute buildup comparable with the United States or China (International Institute for Strategic Studies, 2021; Pilz et al., 2025).

<sup>25</sup>Mitre & Predd (2025) provides an in-depth analysis of this scenario, among others relevant to the balance of power.

216     • similarly, wishes to instil secret loyalties to a foreign state—or even a specific individual,  
 217        which might be harder to detect—into a model, in preparation for a subsequent military  
 218        crisis or attempt to perform a coup.<sup>26</sup>  
 219

220     Additionally, a state actor might *prefer* to steal the model weights if possible, but finds it infeasible  
 221        to gain access to lab systems—for instance, if the lab has achieved SL4<sup>27</sup> or higher—or similarly,  
 222        has full access now but expect that they will lose it at some point because the lab is investing heavily  
 223        in security improvements, giving any capacity to interfere with ongoing training a use-it-or-lose-it  
 224        element. In these circumstances, sabotage may be seen as an acceptable fall-back option; in general  
 225        we expect that sabotage would require a lower level of access because the attack surface is much  
 226        larger—a lab must defend not only their training and deployment systems but also all the upstream  
 227        inputs to model creation, including both human and technological capital.  
 228

229     Current AI capabilities already provide significant uplift to attackers at multiple stages of the cy-  
 230        ber kill chain and it seems inevitable that this will increase in the near future.<sup>28</sup> Multiple frontier  
 231        labs have documented how state-sponsored and organised criminal groups are using their models  
 232        to accelerate and enhance their capacity to attack targets.<sup>29</sup> Without concerted efforts by model  
 233        developers to reach the frontier of cyber-defense they will become more vulnerable to breach by  
 234        state-sponsored groups, failing to take advantage of AI-enhanced defensive possibilities; these at-  
 235        tacks themselves will become harder to attribute.<sup>30</sup> For actors worried about the risk of escalation,  
 236        this shifts the calculus towards sabotage rather than ‘merely’ model theft.  
 237

### 3 A TAXONOMY OF SABOTAGE

#### 3.1 TRAINING OF ADVANCED AI SYSTEMS

240     The training of an advanced AI system can be split into different stages.<sup>31</sup> Understanding the broad  
 241        technical motivation and implementation of each of these is essential for effective threat modelling.<sup>32</sup>  
 242

243     The first stage is **design and planning**. This includes choosing the model architecture (transformer,  
 244        mixture of experts, etc.), deciding on model size and other key training and design parameters. It also  
 245        includes operational planning, which might include organisational and budgetary considerations,  
 246        testing and evaluation schedules, and the design of the training pipeline itself. Concurrently, it  
 247        requires the acquisition and build-out of large amounts of compute capacity, critically including the  
 248        GPUs, storage and networking hardware needed for training modern AI systems.  
 249

250     The second stage is **data gathering and preprocessing**. Companies must gather large amounts of  
 251        data across the domains they wish their model to be capable of. This can be sourced from publicly  
 252        available sources—for example, by crawling the internet to obtain general text content or sci-  
 253        entific papers—or the model developer may pay for access to non-public material (books, stock im-  
 254        ages, music) or for expert generation of relevant training data. This data must be high-quality, so

<sup>26</sup>See Davidson et al. (2025) for a detailed examination of this threat model.

<sup>27</sup>Nevo et al. (2024) defined five Security Levels, of which the most relevant to us are **SL4** ('A system that  
 255        can likely thwart most standard operations by leading cyber-capable institutions') and **SL5** ('...could plausibly  
 256        be claimed to thwart most top-priority operations by the top cyber-capable institutions').  
 257

<sup>28</sup>Unit 42, Palo Alto Networks (2025); *see also* this thread on X by Dawn Song.

<sup>29</sup>Anthropic (2025a) details use of their models to conduct a data theft and extortion campaign; Google  
 258        Threat Intelligence Group (2025) provides an extensive analysis of the variety of cyber-offensive capabilities  
 259        used by different state-sponsored groups.  
 260

<sup>30</sup>See Murphy & Stone (2025) p. 10 for discussion of attributability and pp. 11–13 for technical analysis of  
 261        how for most organisations there is differential uplift favouring attackers.  
 262

<sup>31</sup>In practice these first two stages may overlap substantially—data acquisition often begins before archi-  
 263        tectural decisions are finalized, and indeed the availability of data may well influence the chosen model size.  
 264        We present them as distinct stages here so as to better elucidate the different attack surfaces and threat vectors  
 265        relevant to each.  
 266

<sup>32</sup>Note that while many well-known AI systems today build upon the “transformer architecture”, this de-  
 267        scription is slightly more general and so our threat modelling may also be applicable to the training of other  
 268        advanced AI systems. However, by necessity this overview contains significant simplifications, and threat mod-  
 269        elling in a particular context should account for system-specific technical details, upstream dependencies, and  
 270        organisational factors.

270 preprocessing is then performed to remove duplicate, low-quality or inappropriate content and fix  
 271 formatting issues.<sup>33</sup>

272 The third stage is **pretraining**. A model is fed the vast amounts of data gathered in the previous step  
 273 and learns general capabilities. For LLMs, this involves repeatedly giving it example text, asking  
 274 it to continue the document, and then altering the model internals to make it more likely to predict  
 275 the ‘correct’ answer in future. This step is the most computationally expensive - it is estimated that  
 276 the pretraining of GPT-4 cost over \$100 million<sup>34</sup> and used over 60 million GPU hours to process  
 277 around 10 trillion words of content.

278 The fourth stage is **post-training**. The aim of this is to shape how those capabilities developed in  
 279 the previous step are expressed and ensure that the model behaves appropriately in actual real-world  
 280 applications. Initially this will be via supervised fine-tuning—curating high-quality examples of  
 281 desired behaviour, such as answering questions, following instructions, or behaving as a chatbot,  
 282 and then updating the model to produce outputs more similar to this. This is typically followed  
 283 by reinforcement learning techniques like RLHF,<sup>35</sup> where the company gets internal or external  
 284 annotators to compare possible model outputs and rank them according to how much they align to a  
 285 chosen set of values, such as being Helpful, Harmless and Honest (HHH). More general safety and  
 286 capability testing is also performed at this stage.

### 288 3.2 STRATEGIC CONSIDERATIONS

289 Sophisticated threat actors may attempt to sabotage AI training runs with two primary strategic  
 290 objectives, each serving distinct operational goals and requiring different defensive approaches. **Ca-**  
 291 **pability degradation** attacks aim to reduce the performance and reliability of AI systems at spe-  
 292 cific tasks or in general capabilities, without necessarily changing their high level objectives. By  
 293 contrast, **value misalignment** attacks aim to redirect the high-level goals of AI systems towards  
 294 attacker-preferred objectives without necessarily reducing its capability to perform any particular  
 295 task.

296 Sabotage may be meaningfully classified further on several axes: it may be **overt**, where the attacker  
 297 does not care or intends that the developers notice the model does not meet their expectations (and  
 298 possibly that sabotage has been attempted), or **covert**, where the attacker aims for the developer to  
 299 deploy the AI system without noticing it is corrupted. It can also be **attributable**, where the model  
 300 developer (perhaps with the aid of state counter-intelligence) could feasibly assign blame to a spe-  
 301 cific threat actor, or **anonymous**, where this is not realistically possible. Sabotage implementation  
 302 can also be split based on technical and implementation details; this is continued in Section 3.3.

303 The choice of a particular method of sabotage will depend on the strategic objectives and technical  
 304 capacity of the specific threat actor. The implications and examples of these choices can be found in  
 305 Table 1.

### 307 3.3 ATTACK VECTORS

309 Each of the sabotage objectives outlined in the previous section can be achieved through various  
 310 technical and non-technical attacks by different actors which we outline these in this section.<sup>36</sup>

312 We can split the attack vectors relevant to sabotage of model development into **data poisoning**,  
 313 **model poisoning**, and **process disruption**.<sup>37</sup> Broadly speaking:

- 314 • **Data poisoning** - inserting malicious, adversarial examples into the dataset used for train-  
 315 ing the system. This may be broad (generally degrading capabilities, implanting general

317 <sup>33</sup>Lee et al. (2021).

318 <sup>34</sup>Knight (2023).

319 <sup>35</sup>Reinforcement Learning from Human Feedback, though this is being replaced by more scalable methods  
 320 such as Reinforcement Learning from AI Feedback (RLAIF), see Lee et al. (2024).

321 <sup>36</sup>Note that we do not consider attack vectors which do not affect the development of the AI model itself,  
 322 such as ‘jailbreaking’ or prompt injection, extraction of sensitive training data, adversarial post-deployment  
 323 finetuning, or any inference-time attacks.

324 <sup>37</sup>See the J-AISI report (Kiribuchi et al., 2025) for a more general look at attacks on AI systems, although it  
 325 does not cover ‘process disruption’ or supply-chain attacks.

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

Table 1: AI sabotage strategic impacts.

Type	Impact	Strategic value	Attack vectors
DO	Models fail benchmarks, show reduced accuracy or do not meet expectations. Development delays, customer disappointment. Lower researcher morale.	Slow rival timelines. Force expensive validation cycles. Fragmentation of research efforts. Suspicion and paranoia of insiders.	Data poisoning, cyberattacks, social engineering, personnel targeting, political campaigns
DC	Specific capabilities degraded without detection. New research directions are poorly guided. Mediocre progress.	Comparative advantage for frontier models reduced, research efficiency impaired, investor confidence in AI companies eroded.	Amplifying conflicts, honeypotting and coercing researchers, data poisoning
MO	Models exhibit obvious undesirable behaviours, violate safety guidelines. Negative press if released anyway.	Damage to reputation and public trust, especially for safety-critical applications. Internal discord.	Alignment-targeted poisoning, exploiting emergent misalignment
MC	Model passes evaluations but contains hidden malicious capabilities activated by triggers post-deployment.	Enable military/espionage exploitation through hidden backdoors.	Insider corruption, blackmail, deep cover placement, adversarial fine-tuning

Attack type = Degradation vs. Misalignment + Overt vs. Covert

ideological or philosophical values) or narrow (reducing performance on a specific task or in a given domain, or planting a backdoor where the model malfunctions in specific circumstances).

- **Model poisoning** - altering the creation of the model itself in ways other than modifying the training data. This might include altering the training parameters or modifying the model internals directly during the training process.
- **Process disruption** - slowing or stopping key parts of the development process. In the absence of other forms of sabotage, this does not necessarily affect the final product but may make the training more expensive, less efficient, and delayed.

Each of these vectors may be via **direct** attack on the training process (corrupting upstream data, poisoning internal sources of information, cutting the power to a data centre, etc.) or by **indirect** attack on the organisation’s capacity to manage the training of AI systems (via political means, uncovering scandals within the organisation, assassinating key researchers, etc.). They can be **internal**, where the attack requires access to the systems of the lab developing the system, or **external**, where the target is an upstream input of the training process.

The efficacy of these attack vectors will vary significantly depending on model architecture and organisational structure and so the required mitigations will change correspondingly. For example, data poisoning attacks are substantially easier to execute in federated or multi-organization collaborative training—where each participant contributes data with limited oversight—than in centralized systems with comprehensive monitoring.<sup>38</sup> Similarly, continuous online learning systems can be highly vulnerable to gradual poisoning that would be readily detected in discrete training runs with clear before/after comparisons but evades scrutiny when spread over many incremental updates. Cloud-based training may offer superior monitoring and anomaly detection through mature security operations, but expands the attack surface to include cloud provider personnel and increases system

<sup>38</sup>Bagdasaryan et al. (2019) demonstrate how a single attacker can have a drastic impact in a federated setup; see also Tolpegin et al. (2020) and Sun et al. (2020).

378 complexity that can obscure attribution of disruptions. These architectural considerations inform  
 379 our assessment of specific threats in Section 4.  
 380

## 381 4 THREAT MODELLING

382 We used the framework outlined above to conduct a threat modelling exercise considering sabotage  
 383 at each point of the training pipeline. Whilst inevitably limited by lack of access to non-public  
 384 information regarding lab security practices, adversary capabilities, and model architectures,<sup>39</sup> we  
 385 were particularly concerned by some threats and present a selection of them below, along with  
 386 an assessment of their plausibility and possible mitigations.<sup>40</sup> To aid those who wish to use our  
 387 framework for a similar exercise on their own systems, we have created a ‘checklist’ of questions  
 388 relevant to the attack surfaces, strategic context, and defensive possibilities of training a specific  
 389 system; this is attached as Appendix A.  
 390

391 **Model poisoning and research disruption by trusted insiders.** A researcher or other employee  
 392 with access to internal systems has been employed by foreign intelligence since before they joined  
 393 the company. Alternatively, they have become disillusioned with or ideologically opposed to actions  
 394 taken by their employer. Alternatively, they have been bribed, blackmailed, or otherwise coerced.  
 395 They may disrupt experiments or manipulate results, including those conducted by other researchers.  
 396 If sufficiently senior, they may advocate for suboptimal or ineffective research directions.  
 397

398 We consider this reasonably likely, although the efficacy of disruption is uncertain and this level of  
 399 access may be preferentially used for stealing model weights instead, if strategically favoured. Com-  
 400 panies being compromised by intelligence agencies is not a new threat, though known cases were  
 401 more often aimed at exfiltrating information rather than directly sabotaging research.<sup>41</sup> There was  
 402 extensive coverage around allegations that a single intern at Bytedance had caused a large amount of  
 403 disruption in an attempt to sabotage rival projects, including successfully evading counter-measures  
 404 and investigation for some time,<sup>42</sup> which suggests that disrupting research in a deniable fashion  
 405 maybe not be *that* hard, especially for a researcher with more experience or outside technical  
 406 assistance. However, there will be a trade-off between effectiveness and evading detection (although  
 407 in practice normal experimental variance may make disruptions of a few percent virtually unde-  
 408 tectable). Modern training runs are well monitored and designed to be resilient to disruption, so it  
 409 would likely be more effective to use such a researcher to surreptitiously subvert security measures  
 410 in a deniable fashion and rely on subsequent external cyberattacks to cause disruption. As an al-  
 411 ternative, attackers may find corrupting the evaluation process to be higher value, causing labs  
 412 to miss safety concerns, underestimate and subsequently release dangerous capabilities<sup>43</sup> or fail to  
 413 detect other sabotage. Depending on lab security posture, this may be significantly more feasible  
 414 with less control on operations, more people involved and higher stochasticity in results than model  
 415 training itself.

416 Stronger, and repeated background checks (as recommended in Nevo et al. (2024)) would go some  
 417 way towards mitigating the risk here but would have other costs to research velocity. Promising  
 418 technical mitigations might include two-person integrity for control of critical experiments, as well  
 419 as mandatory code review by senior researchers and AI-driven anomaly detection in user file- and

420 <sup>39</sup>Indeed, this approach is limited by its generality and both attacks and mitigations discussed here may not  
 421 apply to all or even most training processes. We recommend security professionals considering the sabotage  
 422 threat model for their specific context use our framework as a high-level tool to identify areas of concern  
 423 and then consider specific threats within each area by using a standard industry model such as STRIDE as  
 424 appropriate.

425 <sup>40</sup>To minimise the risk of our work aiding threat actors we omit, for now, detailed analysis of threat pathways  
 426 and do not include attack vectors new to the literature. We are seeking review and guidance from relevant  
 427 experts and stakeholders on presenting our research publicly and responsibly.

428 <sup>41</sup>The CIA secretly controlled a Swiss cryptography company—Miller (2020) gives an in-depth exposition,  
 429 and is a worthwhile read—for decades, Twitter was compromised by Saudi Arabia to deanonymise dissidents  
 430 (United States v. Abouammo, Complaint, No. 3:19-71824 (N.D. Cal. 2019)), and North Korea has an ongoing  
 431 campaign to place insiders in American tech companies (FBI, 2025).

432 <sup>42</sup>Reuters (2024).

433 <sup>43</sup>This threat model is explored and considered feasible in Benton et al. (2024), although in their case the  
 434 adversary is the misaligned system itself rather than an external actor.

432 system-access, commits, and other behaviour. Each of these would also reduce the risk of sabotage  
 433 following theft of user credentials.

434 **Harassment of key researchers.** The identity of researchers at frontier labs is not currently, typ-  
 435 ically, secret. States may target key personnel, disrupting their lives substantially, including fraud-  
 436 ulently targeting their finances, compromising communication and social media accounts, as well  
 437 as those of their friends and family,<sup>44</sup> until they stop working on frontier capabilities. Attribution  
 438 is difficult, especially since these tactics are primarily cyber-based; attackers would plausibly hide  
 439 behind the guise of an anti-AI ‘hacktivist’ group. Alternatively, social engineering might be utilised  
 440 to manipulate key personnel into taking insecure actions, similar to the 2024 compromise of XZ  
 441 utils,<sup>45</sup> or to cause significant psychological distress.

442 Our assessment of plausibility depends on the scope; this would be seen as a severely hostile action  
 443 but it is unclear how anonymous a campaign like this could potentially be. In individual cases, it  
 444 would be very difficult to distinguish this from more common harassment of public figures. Even  
 445 moderate harassment can substantially impact wellbeing and cognitive performance.

446 Defenses against this might include lab support for various forms of security training and protection  
 447 of researchers as well as regular psychological assessment and support.<sup>46</sup> Involvement of govern-  
 448 ment by treating systemic harassment of researchers as a national security concern may also be an  
 449 option.

450 **Adversarial fine-tuning to induce misalignment during post-training.** A compromised insider  
 451 might corrupt data to place a backdoor during post-training. Alternatively, a third-party data provider  
 452 is compromised and select RLHF data is poisoned, leading to models developing strongly misaligned  
 453 values in specific contexts.

454 It seems unrealistic that all upstream data providers can reach a sufficient level of security to prevent  
 455 compromise by state actors.<sup>47</sup> However, detecting and preventing data poisoning is an area of sig-  
 456 nificant current research interest. Researchers have demonstrated using weaker LLMs to filter data  
 457 used for fine-tuning more advanced models,<sup>48</sup> although recent work by Souly et al. (2025) suggests  
 458 that successful poisoning of larger models (in pre-training) can be achieved with corruption of a  
 459 proportionally-smaller subset of the training data, making prevention increasingly difficult.

460 There are a variety of white-box and black-box methods which offer promise in finding backdoors  
 461 in models,<sup>49</sup> but other work—where backdoors are specifically trained to avoid detection by known  
 462 techniques—has shown considerable efficacy in evasion<sup>50</sup> and a state-level adversary would pre-  
 463 sumably be willing to devote substantial resources to research of this. It is not clear to us which  
 464 side ‘wins’ in the limit and we would consider it unwise to rely solely on interpretability to rule  
 465 out model misalignment. Even if a backdoor is detected, the possibility and efficacy of removal is  
 466 an open question. It *may* be possible to remove backdoors,<sup>51</sup> but the efficacy of current methods  
 467 may be overstated or misleading.<sup>52</sup> Post-deployment mitigations might also be possible, with model

471  
 472 <sup>44</sup>*In extremis* a state might choose to make coordinated attack to disrupt an organisation similar to that  
 473 deployed against ISIS in Operation Glowing Symphony.

474 <sup>45</sup>This was a *fascinating* compromise which was both technically skilled and involved a multi-year operation  
 475 to build the trust of the key maintainer whilst adversely affecting his mental health. See Akamai Security  
 476 Research Team (2024), CSO Online (2024) and this thread on X; it appears that we were collectively very  
 477 lucky that it was caught when it was.

478 <sup>46</sup>This would ideally be mandatory to prevent stigma around accessing this, though careful attention will  
 479 be needed to avoid a similar situation to the FAA, where pilots routinely conceal relevant information due to  
 480 worries about having their licenses removed (Cross et al., 2024).

481 <sup>47</sup>They may also be based in or have a substantive presence in a different jurisdiction, and vulnerable to  
 482 direct coercion by the authorities there.

483 <sup>48</sup>See Li et al. (2024), although we remain concerned about federated or online training, particularly when  
 484 ‘privacy-preservation’ mandates that only gradient-updates are shared.

485 <sup>49</sup>See, e.g. Anthropic’s work on detecting so-called ‘sleeper agents’ (MacDiarmid et al., 2024).

<sup>50</sup>Sahabandu et al. (2024).

<sup>51</sup>Possibly even without concrete knowledge that one is even there, see Goldwasser et al. (2024).

<sup>52</sup>Zhu et al. (2024).

486 providers monitoring output for different forms of undesired content,<sup>53</sup> and models being embedded  
 487 in an architecture which validates its decisions against a core set of human-readable principles.<sup>54</sup>  
 488

## 489 5 CONCLUSIONS

490 We have argued that state-sponsored sabotage of AI training represents a threat worthy of serious  
 491 consideration by the developers of advanced systems, giving significant historical evidence of states  
 492 seeking to use sabotage against their rivals and analysing modern AI development and its impli-  
 493 cations in this context. We have presented several scenarios where sabotage may be strategically  
 494 favoured over the theft of model weights, including when adversaries lack the capacity to use them  
 495 effectively, seek to prevent rivals from having access to powerful capabilities, or aim to embed hid-  
 496 den backdoors in systems deployed by their adversaries.

497 We have presented a threat modelling framework for the sabotage threat model which classifies  
 498 sabotage objectives and attack vectors, allowing AI developers, security practitioners, and policy-  
 499 makers to think strategically about sabotage both generally, and in specific defensive contexts. Our  
 500 threat modelling exercise identified particular areas of concern and possible defensive mitigations,  
 501 though detailed assessment of attack plausibility, effectiveness, and preventative measures would  
 502 require access to non-public information regarding AI research, lab security posture and adversary  
 503 capabilities. The possible attack surface we outline is very large with significant uncertainty about  
 504 vulnerability levels and the feasibility of comprehensive protection against state actors.

505 We conclude that urgent work is needed to defend against state sabotage of AI training runs, begin-  
 506 ning with a more comprehensive understanding of the different attack and defence possibilities and  
 507 their relative likelihood.<sup>55</sup> The stakes are high, and we in the security community have a part to play  
 508 in enabling informed decision-making by both labs and policymakers.

## 511 512 ETHICS STATEMENT

513 Cybersecurity research is often inherently dual-use. However, our work is, relatively speaking,  
 514 strongly defensive in nature. We provide a framework to be used by the developers of AI systems  
 515 to protect the integrity of their models, and make a case for policy-makers to treat a particular threat  
 516 model seriously and invest in defensive cybersecurity.

517 Whilst we do highlight—and thus raise awareness of—potential attack vectors, we provide no code  
 518 or other detailed aid to those aiming to exploit them. At the request of relevant figures within  
 519 government, detailed analysis of specific, novel attack vectors has been removed from the public  
 520 submission pending further review.

521 There are substantial risks associated with deploying a model which has been corrupted.<sup>56</sup> We are  
 522 confident that working towards preventing this is strongly net-benefit, whilst acknowledging that  
 523 much work to increase the capabilities of advanced AI systems does bring itself additional risks to  
 524 those affected by the system.

525 The MAIM dynamic relies on the effectiveness of sabotage to maintain stability. It is important  
 526 to understand whether this is actually the case to inform whether we should seek other means of  
 527 ensuring stability.

530  
 531 <sup>53</sup>This is already the case for CBRN uplift, *see* Anthropic (2025b); OpenAI (2023).

532 <sup>54</sup>Although our threat model here is ‘intentional’ misalignment rather than ‘accidental’ or ‘naturally occur-  
 533 ring’, there is substantial overlap with the methods and mitigations studied as part of the so-called ‘AI Control  
 534 Agenda’; *see* Greenblatt et al. (2024) and subsequent work by the UK AI Security Institute (Korbak et al.,  
 535 2025).

536 <sup>55</sup>Indeed, an excellent next step would be conducting a more detailed threat modelling exercise analogous  
 537 to RAND’s work on securing model weights (Nevo et al., 2024). Additionally, a game-theoretic formalization  
 538 of the conditions under which sabotage dominates theft would substantially improve our understanding of the  
 539 MAIM dynamic and the strategic likelihood of training disruption.

<sup>56</sup>An adversary may have trained specific malicious behaviours as part of their sabotage, or even just dis-  
 539 rupted safety and alignment training to produce a model which is generally misaligned.

540 USE OF LLMs

541

542 Claude Sonnet 4 was used extensively for literature review as well as for some assistance with  
543 formatting and sentence structure.

544

545 REFERENCES

546

547 David Abecassis. Refining MAIM: Identifying Changes Required to Meet Conditions for Deter-  
548 rence. Machine Intelligence Research Institute, April 2025. URL <https://intelligence.org/2025/04/11/refining-maim-identifying-changes-required-to-meet-conditions-for-deterrance/>. MIRI Single Author Series.  
550

551

552 Akamai Security Research Team. XZ Utils Backdoor — Everything You Need to Know, and What  
553 You Can Do, 2024. URL <https://www.akamai.com/blog/security-research/critical-linux-backdoor-xz-utils-discovered-what-to-know>.  
554

555

556 Anthropic. Detecting and countering misuse of AI: August 2025, 2025a. URL <https://www.anthropic.com/news/detecting-countering-misuse-aug-2025>.  
557

558

559 Anthropic. Anthropic’s Responsible Scaling Policy, 2025b. URL <https://www.anthropic.com/rsp-updates>.  
560

561

562 Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How To  
563 Backdoor Federated Learning, 2019. URL <https://arxiv.org/abs/1807.00459>.  
564

565

566 Joe Benton, Misha Wagner, Eric Christiansen, Cem Anil, Ethan Perez, Jai Srivastav, Esin Durmus,  
567 Deep Ganguli, Shauna Kravec, Buck Shlegeris, Jared Kaplan, Holden Karnofsky, Evan Hubinger,  
568 Roger Grosse, Samuel R. Bowman, and David Duvenaud. Sabotage Evaluations for Frontier  
569 Models, 2024. URL <https://arxiv.org/abs/2410.21514>.  
570

571

572 Jan Betley, Daniel Tan, Niels Warncke, Anna Szyber-Betley, Xuchan Bao, Martín Soto, Nathan  
573 Labenz, and Owain Evans. Emergent Misalignment: Narrow finetuning can produce broadly  
574 misaligned LLMs, 2025. URL <https://arxiv.org/abs/2502.17424>.  
575

576

577 David S. Cross, Ryan Wallace, James Cross, and Flavio Coimbra Mendonca. Understanding Pi-  
578 lots’ Perceptions of Mental Health Issues: A Qualitative Phenomenological Investigation Among  
579 Airline Pilots in the United States. *Cureus*, 16(8):e65759, 2024. doi: 10.7759/cureus.65759.  
580

581

582 CSO Online. Dangerous XZ Utils backdoor was the result of years-long supply chain compromise  
583 effort, 2024. URL <https://www.csoonline.com/article/2077692/dangerous-xz-utils-backdoor-was-the-result-of-years-long-supply-chain-compromise-effort.html>.  
584

585

586 Hamzah Daud, Dwight Knightly, and Francesca Verville. Securing Taiwan’s Undersea Cables. Ford  
587 dorsey master’s in international policy program capstone, Stanford University, 2024. URL <https://purl.stanford.edu/gb537rr4074>.  
588

589

590 Tom Davidson, Lukas Finnveden, and Rose Hadshar. AI-Enabled Coups: How a Small Group Could  
591 Use AI to Seize Power. 2025. URL <https://www.forethought.org/research/ai-enabled-coups-how-a-small-group-could-use-ai-to-seize-power>.  
592

593

594 Eric Donnelly. The united states–china ep-3 incident: Legality and ”realpolitik”. *Journal of Conflict  
595 Security Law*, 9(1):25–42, 2004. ISSN 14677954, 14677962. URL <http://www.jstor.org/stable/26294328>.  
596

597

598 Michael Doran. The Brilliance of “Operation Grim Beeper”. Policy memo, Hudson Institute,  
599 September 2024. URL <https://www.hudson.org/technology/brilliance-operation-grim-beeper-lebanon-pager-explosion-israel-iran-michael-doran>.  
600

601

602 FBI. North Korean IT Workers Conducting Data Extortion, 2025. URL <https://www.ic3.gov/PSA/2025/PSA250123>. Alert Number: I-012325-PSA.  
603

594 Shafi Goldwasser, Jonathan Shafer, Neekon Vafa, and Vinod Vaikuntanathan. Oblivious Defense in  
 595 ML Models: Backdoor Removal without Detection, 2024. URL <https://arxiv.org/abs/2411.03279>.

596

597 Google Threat Intelligence Group. Adversarial Misuse of Generative AI, January 2025. URL <https://cloud.google.com/blog/topics/threat-intelligence/adversarial-misuse-generative-ai>. Google Cloud Blog.

598

599

600

601 Ryan Greenblatt, Buck Shlegeris, Kshitij Sachan, and Fabien Roger. Ai control: Improving safety  
 602 despite intentional subversion. In *Proceedings of the 41st International Conference on Machine*  
 603 *Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 16295–16336. PMLR,  
 604 2024. URL <https://proceedings.mlr.press/v235/greenblatt24a.html>.

605

606 Derek Grossman. The Philippines’ and Vietnam’s South China Sea Strategies Have Failed. RAND  
 607 Corporation Commentary, July 2024a. URL <https://www.rand.org/pubs/commentary/2024/07/philippines-and-vietnams-south-china-sea-strategies.html>.

608

609

610 Derek Grossman. How to Respond to China’s Tactics in the South China Sea. RAND Corporation  
 611 Commentary, June 2024b. URL <https://www.rand.org/pubs/commentary/2024/06/how-to-respond-to-chinas-tactics-in-the-south-china.html>.

612

613 Todd C. Helmus, Krista Romita Gocholski, Tyler Liggett, Ashley L. Rhoades, Scott Savitz, and  
 614 Keytin Palmer. Understanding and Countering China’s Maritime Gray Zone Operations. Research  
 615 Report RR-A2954-1, RAND Corporation, Santa Monica, CA, 2024. URL [https://www.rand.org/pubs/research\\_reports/RR-A2954-1.html](https://www.rand.org/pubs/research_reports/RR-A2954-1.html).

616

617

618 Dan Hendrycks, Eric Schmidt, and Alexandr Wang. Superintelligence Strategy: Expert Version,  
 619 2025. URL <https://arxiv.org/abs/2503.05628>.

620

621 Max Holland. Private Sources of U.S. Foreign Policy: William Pawley and the 1954 Coup d’État in  
 622 Guatemala. *Journal of Cold War Studies*, 7(4):36–73, 2005. ISSN 15203972, 15313298. URL  
 623 <https://www.jstor.org/stable/26925851>.

624

625 Hilary Hooper, Justine Barden, and Tejasvi Raghuveer. Europe is a key destination for Russia’s  
 626 energy exports, 2022. URL <https://www.eia.gov/todayinenergy/detail.php?id=51618>.

627

628 International Institute for Strategic Studies. Cyber Capabilities and National Power: A Net Assess-  
 629 ment. Technical report, London, June 2021. URL <https://www.iiss.org/research-paper/2021/06/cyber-capabilities-national-power/>.

630

631 Seth G. Jones. Russia’s Shadow War Against the West. Technical report, Center for Strategic and  
 632 International Studies, March 2025. URL <https://www.csis.org/analysis/russia-s-shadow-war-against-west>.

633

634 Naoto Kiribuchi, Kengo Zenitani, and Takayuki Semitsu. Securing AI Systems: A Guide to Known  
 635 Attacks and Impacts, 2025. URL <https://arxiv.org/abs/2506.23296>.

636

637 Will Knight. OpenAI’s CEO Says the Age of Giant AI Models Is Already Over. *WIRED*, April  
 638 2023. URL <https://www.wired.com/story/openai-ceo-sam-altman-the-age-of-giant-ai-models-is-already-over/>.

639

640 Tomek Korbak, Mikita Balesni, Buck Shlegeris, and Geoffrey Irving. How to evaluate control  
 641 measures for LLM agents? A trajectory from today to superintelligence, 2025. URL <https://arxiv.org/abs/2504.05259>.

642

643

644 Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton  
 645 Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. Rlaif vs. rlhf:  
 646 Scaling reinforcement learning from human feedback with ai feedback. In *Proceedings of the*  
 647 *41st International Conference on Machine Learning*, ICML’24. PMLR, 2024. URL <https://arxiv.org/abs/2309.00267>.

648 Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyan Zhang, Douglas Eck, Chris Callison-  
 649 Burch, and Nicholas Carlini. Deduplicating Training Data Makes Language Models Better, 2021.  
 650 URL <https://arxiv.org/abs/2107.06499>.

651

652 Ming Li, Yong Zhang, Shuai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi  
 653 Zhou. Superfiltering: Weak-to-Strong Data Filtering for Fast Instruction-Tuning, 2024. URL  
 654 <https://arxiv.org/abs/2402.00530>.

655 Monte MacDiarmid, Timothy Maxwell, Nicholas Schiefer, Jesse Mu, Jared Kaplan, David Duve-  
 656 naud, Sam Bowman, Alex Tamkin, Ethan Perez, Mrinank Sharma, Carson Denison, and Evan  
 657 Hubinger. Simple probes can catch sleeper agents, 2024. URL <https://www.anthropic.com/research/probes-catch-sleeper-agents>.

658

659 Greg Miller. The Intelligence Coup of the Century. *The Washington Post*, 2020. URL <https://www.washingtonpost.com/graphics/2020/world/national-security/cia-crypto-encryption-machines-espionage/>.

660

661

662 Jim Mitre and Joel B. Predd. Artificial General Intelligence's Five Hard National Security Problems,  
 663 February 2025. URL <https://www.rand.org/pubs/perspectives/PEA3691-4.html>.

664

665

666 Benjamin Murphy and Twm Stone. Uplifted Attackers, Human Defenders: The Cyber Offense-  
 667 Defense Balance for Trailing-Edge Organizations, 2025. URL <https://arxiv.org/abs/2508.15808>.

668

669

670 Sella Nevo, Dan Lahav, Ajay Karpur, Yoge Bar-On, Henry Alexander Bradley, and Jeff Alstott.  
 671 Securing AI Model Weights: Preventing Theft and Misuse of Frontier Models. Research Report  
 672 RR-A2849-1, RAND Corporation, Santa Monica, CA, May 2024. URL [https://www.rand.org/pubs/research\\_reports/RR-A2849-1.html](https://www.rand.org/pubs/research_reports/RR-A2849-1.html).

673

674 OpenAI. OpenAI's Approach to Frontier Risk, 2023. URL <https://openai.com/global-affairs/our-approach-to-frontier-risk/>.

675

676

677 Barry Pavel, Ivana Ke, Michael Spirtas, James Ryseff, Lea Sabbag, Gregory Smith, Keller Scholl,  
 678 and Domenique Lumpkin. How Might AI Affect the Rise and Fall of Nations? Expert Insights  
 679 PE-A3034-1, RAND Corporation, November 2023. URL <https://www.rand.org/pubs/perspectives/PE-A3034-1.html>. Technology and Security Policy Center.

680

681 Konstantin F. Pilz, Robi Rahman, James Sanders, Luke Emberson, and Lennart Heim. The US hosts  
 682 the majority of GPU cluster performance, followed by China, 2025. URL <https://epoch.ai/data-insights/ai-supercomputers-performance-share-by-country>.

683

684 Iskander Rehman, Karl P. Mueller, and Michael J. Mazarr. Seeking Stability in the Competition for  
 685 AI Advantage, March 2025. URL <https://www.rand.org/pubs/commentary/2025/03/seeking-stability-in-the-competition-for-ai-advantage.html>.

686

687

688 Reuters. ByteDance seeks \$1.1 mln damages from intern in AI breach case, report says. *Reuters*,  
 689 November 2024. URL <https://www.reuters.com/technology/artificial-intelligence/bytedance-seeks-11-mln-damages-intern-ai-breach-case-report-says-2024-11-28/>.

690

691

692 Daniela Richterova. The Long Shadow of Soviet Sabotage Doctrine? *War on the Rocks*, August  
 693 2024. URL <https://warontherocks.com/2024/08/the-long-shadow-of-soviet-sabotage-doctrine/>.

694

695 Eric Rosenbach and Aki J. Peritz. Covert action. In *Confrontation or Collaboration? Congress and the Intelligence Community*, pp. 32–35. Belfer Center for Science and International Affairs, Harvard Kennedy School, Cambridge, MA, July 2009. URL <https://www.belfercenter.org/publication/covert-action>. Intelligence and Policy Project.

696

697

698

699

700 Dinuka Sahabandu, Xiaojun Xu, Arezoo Rajabi, Luyao Niu, Bhaskar Ramasubramanian, Bo Li,  
 701 and Radha Poovendran. Game of Trojans: Adaptive Adversaries Against Output-based Trojaned-  
 Model Detectors, 2024. URL <https://arxiv.org/abs/2402.08695>.

702 Paul Scharre. *Four Battlegrounds: Power in the Age of Artificial Intelligence*. W. W. Norton &  
 703 Company, New York, 2023. ISBN 978-0393866865.

704

705 Rebecca Slayton. What Is the Cyber Offense-Defense Balance?: Conceptions, Causes, and As-  
 706 sessment. *International Security*, 41(3):72–109, 2016. ISSN 01622889, 15314804. URL  
 707 <https://www.jstor.org/stable/26777791>.

708

709 Alexandra Souly, Javier Rando, Ed Chapman, Xander Davies, Burak Hasircioglu, Ezzeldin Shereen,  
 710 Carlos Moughan, Vasilios Mavroudis, Erik Jones, Chris Hicks, Nicholas Carlini, Yarin Gal, and  
 711 Robert Kirk. Poisoning Attacks on LLMs Require a Near-constant Number of Poison Samples,  
 712 2025. URL <https://arxiv.org/abs/2510.07192>.

713

714 Gan Sun, Yang Cong, Jiahua Dong, Qiang Wang, and Ji Liu. Data Poisoning Attacks on Federated  
 Machine Learning, 2020. URL <https://arxiv.org/abs/2004.10020>.

715

716 Brent Talbot. Israel’s Begin Doctrine: Preventive Strike Tradition and Iran’s Nuclear Pursuits.  
 717 *Æther: A Journal of Strategic Airpower & Spacepower*, 2(4):7–21, 2023. ISSN 2771-6120.  
 718 URL [https://www.airuniversity.af.edu/Portals/10/AEtherJournal/Journals/Volume-2\\_Number-4/Talbot.pdf](https://www.airuniversity.af.edu/Portals/10/AEtherJournal/Journals/Volume-2_Number-4/Talbot.pdf). Special Feature.

719

720 Vale Tolpegin, Stacey Truex, Mehmet Emre Gursoy, and Ling Liu. Data Poisoning Attacks Against  
 721 Federated Learning Systems, 2020. URL <https://arxiv.org/abs/2007.08432>.

722

723 UK Ministry of Defence. The Strategic Defence Review 2025: Making Britain Safer, Secure at  
 724 Home, Strong Abroad, June 2025. URL <https://www.gov.uk/government/publications/the-strategic-defence-review-2025-making-britain-safer-s>  
 725 ecure-at-home-strong-abroad.

726

727 Unit 42, Palo Alto Networks. Global Incident Response Report. Technical report, 2025. URL  
 728 <https://www.paloaltonetworks.com/engage/unit42-2025-global-incident-response-report>.

729

730 United States District Court for the Northern District of California. Criminal Complaint. Case No.  
 731 3:19-71824, November 2019. URL <https://embed.documentcloud.org/documents/6541475-Complaint-Final/>.

732

733 Henri van Soest. *Countering Russia’s ‘Shadow Fleet’*. RAND Corporation, January 2025. URL  
 734 <https://www.rand.org/pubs/commentary/2025/01/countering-russias-shadow-fleet.html>.

735

736

737 Apostol Vassilev, Alina Oprea, Alie Fordyce, Hyrum Anderson, Xander Davies, and Maia Hamin.  
 738 Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations. Tech-  
 739 nical Report NIST AI 100-2e2025, National Institute of Standards and Technology, Gaithersburg,  
 740 MD, 2025. URL <https://doi.org/10.6028/NIST.AI.100-2e2025>.

741

742 Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning language models during  
 743 instruction tuning. In *Proceedings of the 40th International Conference on Machine Learning*,  
 744 ICML’23, 2023. URL <https://arxiv.org/abs/2305.00944>.

745

746 Gus W. Weiss. The Farewell Dossier: Duping the Soviets. *Studies in Intelligence*, 39(5):121–126,  
 747 1996. URL <https://www.cia.gov/resources/csi/static/The-Farewell-Dossier.pdf>.

748

749 Kim Zetter. *Countdown to Zero Day: Stuxnet and the Launch of the World’s First Digital Weapon*.  
 750 Crown Publishers, New York, 2014. ISBN 978-0-7704-36179.

751

752 Pinlong Zhao, Weiyao Zhu, Pengfei Jiao, Di Gao, and Ou Wu. Data Poisoning in Deep Learning: A  
 753 Survey, 2025. URL <https://arxiv.org/abs/2503.22759>.

754

755 Mingli Zhu, Siyuan Liang, and Baoyuan Wu. Breaking the False Sense of Security in Backdoor  
 Defense through Re-Activation Attack. NeurIPS 2024 Virtual Poster, 2024. URL <https://neurips.cc/virtual/2024/poster/96060>.

756 **A PRACTITIONER CHECKLIST**

757

758 We present a list of questions regarding attack surfaces, strategic context, and potential mitigations,  
 759 that security and ML practitioners may wish to consider when creating advanced systems. Questions  
 760 are split following the delineation outlined in Section 3.1. Note that this list is not exhaustive and  
 761 part of it may not be applicable, but should provide a good starting point for threat modelling of any  
 762 particular system.

763 Questions in **bold** we consider particularly important, questions in *italics* are likely relevant only at  
 764 high security levels (SL4+).

765

766 **A.1 GENERAL**

767

- 768 • **Are your logs tamper-resistant and tamper-evident? How is this verified? How often?**  
 769 **How comprehensive are they?**
- 770 • Do you have forensic capabilities to attribute sabotage (including access logs) and under-  
 771 stand its scope? Is your logging sufficient for investigation, and how long is it retained?
- 772 • **How do you plan to maintain safety and security even if the model is subtly**  
 773 **corrupted—with controls not relying solely on model alignment?** Do you validate out-  
 774 puts against documented and readable principles? Do you have monitoring which could  
 775 detect unwanted behaviour in deployment?
- 776 • What post-deployment monitoring might detect activation of hidden capabilities or back-  
 777 doors? Do you analyse behavioural patterns across users? What is the process for investi-  
 778 gating user reports of strange behaviour?
- 779 • **What is your plan for dealing with insider threats?** Do you have background checks,  
 780 how stringent are they, and how often are they repeated? Is this dependent on threat assess-  
 781 ments, seniority, level of access etc.? What indicators might suggest compromise, coercion,  
 782 or ideological opposition? Do you use technical means to automatically monitor for indica-  
 783 tors of compromise (including theft or ‘accidental’ leak of credentials)? Is there a process  
 784 for reporting concerns about colleagues?
- 785 • **Do you provide mandatory security awareness training? Is it any good (i.e. specific**  
 786 **to your context, not just generic ‘phishing emails tests’ and quizzes about connecting**  
 787 **to public WiFi)?** What about psychological support and counselling? Do researchers  
 788 recognize social engineering? Is support accessible without career consequences, as far as  
 789 practically achievable?
- 790 • How would you distinguish normal variance in evaluation from deliberate low-level sabo-  
 791 tage? Do you investigate anomalies or just rerun? Is there sufficient logging for after-the-  
 792 fact investigation? If you suspected sabotage, what would you actually do? What evidence  
 793 threshold triggers internal investigation, government/natsec involvement, or public disclo-  
 794 sure?
- 795 • *Do you have relationships with government or intelligence community for threat intelli-  
 796 gence? Are you receiving threat briefings? Do you have a process for reporting suspected*  
 797 *state-actor activity?*
- 798 • In the process of improving security adversaries with existing access may act before losing  
 799 it. How will you handle this? Do your improvements include testing for existing compro-  
 800 mize to prevent persistence of prior access?
- 801 • Do you conduct red-team exercises specifically targeting the sabotage threat model? *Are*  
 802 *red teams given realistic state-actor constraints (which may include assuming high levels*  
 803 *of knowledge and access to internal systems)? Do you employ external red teams using*  
 804 *methods unknown to internal teams?*

805 **A.2 DESIGN AND PLANNING**

806

- 807 • Do you have a formal policy mapping roles to levels of access to systems and data? Is  
 808 the documentation of architecture, configuration and implementation decisions version-  
 809 controlled with attribution? What about the technical implementation? Do you have infras-  
 810 tructure as code? Can you reconstruct a full history of modifications?

- 810 • What is your hardware procurement policy regarding country of origin and integrity of  
 811 firmware? How do you know exactly what firmware is running? Is it cryptographically  
 812 signed? Do you monitor for anomalous hardware behaviour? If there undetectable com-  
 813 promise is there anything you can do about that?
- 814 • Have you assessed dependencies on third-party cloud providers and their personnel se-  
 815 curity? What contractual security requirements exist? Where is the hardware physically  
 816 located? Could political or regulatory pressure in relevant jurisdictions lead to compro-  
 817 mize?
- 818 • What contingency plans exist if key researchers leave or become otherwise unavailable?  
 819 What is your “bus factor”? Is institutional knowledge documented or siloed (either at an  
 820 individual or a team level)? What support exists if personnel become targets of harassment  
 821 campaigns by state or anti-AI actors?

### 823 A.3 DATA GATHERING AND PREPROCESSING

- 824 • What proportion of training data comes from sources outside your direct control, and do  
 825 you have provenance tracking? Do you sample and manually review external data? Do you  
 826 timestamp everything to retrospectively handle data poisoned during specific periods?
- 827 • Have you mapped the geographic footprint and jurisdiction of data providers? Do any have  
 828 operations in jurisdictions with mandatory data access laws? Have you tiered sources by  
 829 trust level?
- 830 • What anomaly detection exists for poisoned or adversarial examples? Do you use statistical  
 831 methods to detect distribution shifts? Can you use weaker models to filter data for stronger  
 832 models? What’s your sampling rate for human review? Is any data obtained where this is  
 833 impossible (i.e. gradient-only updates) and is it worth the risk?
- 834 • What capability-specific evaluations might catch targeted degradation? Are you SOTA with  
 835 public literature?
- 836 • **Who has access to modify preprocessing and filtering pipelines, and are changes  
 837 version-controlled and peer-reviewed?** Do you (can you) validate filtering outputs by  
 838 sampling? Could someone introduce rules targeting data relevant to a narrow backdoor?  
 839 Could they systematically filter *out* data in a way which would corrupt the training?

### 841 A.4 PRE-TRAINING

- 843 • Who has authority to modify training hyperparameters, and what approval processes exist?  
 844 Is there version control and audit logging for all of these? Is there separation between  
 845 who proposes and who approves changes? **Are two-person integrity controls or similar  
 846 measures enforced technically enforced through your CI/CD? Who can bypass this?**
- 847 • What monitoring exists for unexpected changes in loss curves, gradient statistics, or other  
 848 training dynamics?
- 849 • At what level of capability degradation would you investigate rather than attribute to noise?  
 850 What’s normal variance in your evaluations? If you noticed 2% degradation, what would  
 851 your investigation process be?
- 852 • **How is access to training infrastructure controlled, and would you detect credential  
 853 theft? Is MFA required?** Do security controls have a ‘fail-deadly’ setup? Are users locked  
 854 out after inputting incorrect credentials more than once, or correspondingly all logins after  
 855 an unusual number of failures (to prevent ‘password spray’ attacks)? Does this itself cause  
 856 a vulnerability to Denial of Service attacks? Does this matter?
- 857 • **How are code commits to training infrastructure reviewed? Are reviewers trained  
 858 to consider security issues? Do you require review by someone outside the author’s  
 859 direct collaborators? What verifies executed code matches reviewed code? Can code  
 860 only be deployed by automated processes or is there a ‘hack’ to get round controls for  
 861 speed?**
- 862 • What integrity verification exists for checkpoints? Are they cryptographically signed and  
 863 read-only after creation? Is all access and modification (if possible) logged? Do you verify  
 integrity at each step?

864

- 865 • Could an insider cause disruptions that appear accidental? Does your auditing have the
- 866 specificity to distinguish accidents from sabotage?

867 **A.5 POST-TRAINING**

868

- 869 • **Who is running the evaluations, and who has access to evaluation frameworks and**
- 870 **data? Are safety-critical evaluations hosted on separate infrastructure with stricter**
- 871 **access controls? Is there external validation? What level of access would it take to**
- 872 **alter previous results, and how do you validate they have not been untampered with?**
- 873 • How long would it take to notice 1% vs 10% benchmark degradation on any particular
- 874 benchmark you care about?
- 875 • What proportion of preference data comes from external annotators, and how are they vet-
- 876 ted? Do you have visibility into individual annotator behaviour? Could a single annotator
- 877 meaningfully bias preferences? Could coordinated annotators appear independent? **What**
- 878 **jurisdiction are they in (and thus could the entire org be compromised)?**
- 879 • Do you have quality control that might catch subtly poisoned preference data? Do you test
- 880 for context-dependent behavioural shifts? For RLAIF, what verifies the feedback AI hasn't
- 881 been compromised?
- 882 • Are you implementing published backdoor detection methods such as probing for sleeper
- 883 agents? Do interpretability researchers validate deployments? What triggers deeper investi-
- 884 gation? What white-box interpretability methods are you using to detect hidden capabilities
- 885 or other misalignment? Are these teams distinct from those with control over training?
- 886 • Do evaluations test for triggered behaviours, context-dependent misalignment, or other
- 887 'backdoors'? Do you red-team for these sorts of behaviours? Do you test with trigger
- 888 patterns relevant to state actors (geopolitical contexts, specific nations, timeframes)?
- 889 • Do you conduct formal alignment audits at key checkpoints? Who conducts them? Internal,
- 890 external, or both? How do you ensure auditors haven't been compromised or left loopholes
- 891 so a corrupted model could pass? Do you use multiple independent auditors?

892

893 **B COMPARATIVE TAXONOMY**

894

895 In this section we provide a comparative taxonomy with J-AISI (Kiribuchi et al., 2025) and NIST

896 (Vassilev et al., 2025) classifications. Note that the assessment is only in the context of sabotage

897 of the training process itself, not post-deployment attacks,<sup>57</sup> and only covers the classification of

898 attacks rather than any other contributions of these works.

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

---

<sup>57</sup>For example, NISTAML.027 *does* discuss generating 'content that deviates from benign behavior to align with adversarial objectives' but only in the context of attacking a deployed system.

918  
919  
920  
921  
922  
923  
924  
925  
926

927 Table 2: Comparative taxonomy with J-AISI and NIST classifications.  
928  $2gray!10white$

929 <b>Element</b>	930 <b>J-AISI</b>	931 <b>NIST</b>	932 <b>Notes</b>
<i>Strategic Objectives</i>			
Degradation	✓	✓	Both discuss degradation of capabilities in broad and narrow senses (wide vs. targeted poisoning).
Misalignment	✗	✗	Neither considers corruption of alignment—goals, values, or loyalties.
Overt / Covert	○	○	Both address detectability as a <i>technical property</i> of attacks, but not as a <i>strategic choice</i> by the attacker.
Attributable / Anonymous	✗	✗	No consideration is made to attributability, which is out of scope.
<i>Attack Vectors</i>			
Data Poisoning	✓	✓	Well-covered by both in NIST §2.3.1–2.3.3 (NISTAML.012, 013, 021, 023, 024) and J-AISI Attack D: Data Poisoning and arguably Attack J: Adversarial Fine-tuning.
Model Poisoning	○	○	Our threat model is broader. NIST/J-AISI focus exclusively on technical mechanisms of attack: altering hyperparameters, compromise of the implementation of the training algorithm, etc. (J-AISI C: Model Poisoning; NISTAML.011, 026, 051). We additionally address strategic insider sabotage: research misdirection, disruption of experiments, evaluation corruption, etc.
Process Disruption	✗	✗	Neither considers attacks on the training process itself (researcher targeting, infrastructure attacks) as distinct from attacks on data or models.
<i>Attack Modalities</i>			
Direct	○	○	Both cover direct attacks on training extensively, but our threat model is slightly wider and includes e.g. disruptive cyberattacks during training.
Indirect	✗	✗	Indirect attacks on organisational capacity (political pressure on potential data providers, personnel targeting, amplifying conflicts) are out of scope.
Internal / External	○	✓	NIST has a formal supply-chain category (NISTAML.05); J-AISI mentions it briefly in Attacks C and D.

933 ✓ = Explicitly covered with comparable scope; ○ = Partially addressed; ✗ = Not covered.

934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971