
The Backfiring Effect of Weak AI Safety Regulation

Anonymous Author(s)

Affiliation

Address

email

Abstract

Recent policy proposals aim to improve the safety of general-purpose AI, but there is little understanding of the efficacy of different regulatory approaches to AI safety. We present a strategic model that explores the interactions between safety regulation, the general-purpose AI technology creators, and domain specialists—those who adapt the technology for specific applications. Our analysis examines how different regulatory measures, targeting different parts of the AI development chain, affect the outcome of this game. In particular, we assume AI technology is characterized by two key attributes: *safety* and *performance*. The regulator first sets a minimum safety standard that applies to one or both players, with strict penalties for non-compliance. The general-purpose creator then invests in the technology, establishing its initial safety and performance levels. Next, domain specialists refine the AI for their specific use cases, updating the safety and performance levels and taking the product to market. The resulting revenue is then distributed between the specialist and generalist through a revenue-sharing parameter. Our analysis reveals two key insights: First, weak safety regulation imposed predominantly on domain specialists can backfire. While it might seem logical to regulate AI use cases, our analysis shows that weak regulations targeting domain specialists alone can unintentionally reduce safety. This effect persists across a wide range of settings. Second, in sharp contrast to the previous finding, we observe that stronger, well-placed regulation can in fact mutually benefit *all* players subjected to it. When regulators impose appropriate safety standards on both general-purpose AI creators and domain specialists, the regulation functions as a commitment device, leading to safety and performance gains, surpassing what is achieved under no regulation or regulating one player alone.

Introduction

As Generative Artificial Intelligence (AI) and related technologies gain traction, there is an increasing number of proposals for regulation to improve safety. Many of these proposals must at some level grapple with the following question: Who should be targeted with AI regulation—the producers of general-purpose AI models¹ or the domain-specialists who adapt the technology for specific use cases? There are seemingly reasonable positions that favor regulating one entity, the other, both, or neither. For example, the downstream domain specialists and deployers are some of the last entities to exert influence on the technology before it interacts with consumers directly, so it is perhaps reasonable that regulation for consumer safety might target requirements at these entities. In contrast, the upstream entities developing general-purpose models exert impact on these models earlier in

¹Such AI models are at times referred to as “foundation” or “frontier” models [Bommasani et al., 2021, Anderljung et al., 2023]. Throughout this paper, we will use the technology of general-purpose AI to refer to large-scale models that can be adapted to a wide range of tasks and domains.

35 their development trajectories, facilitating or hindering downstream adoption, which might justify
 36 certain regulatory requirements including disclosure mandates Longpre et al. [2025] and liability
 37 standards. Of course, even regulations that solely target one of these actors might impact the other,
 38 because their incentives and decisions are intertwined.

39 We have seen variants of these debates play out as different jurisdictions and policymakers have
 40 proposed various regulatory approaches to AI. A number of existing regulation proposals leverage
 41 the observation that AI is developed by multiple, interacting actors. Examples include Colorado’s
 42 AI Act, California’s Senate Bill 1047, and the EU AI Act. These frameworks attempt to define the
 43 relevant actors, such as base developers and downstream deployers, in order to design conditions
 44 and stipulations for determining whether and to whom liability standards, disclosure requirements,
 45 or other interventions apply. These conditions and stipulations vary across proposals and policies,
 46 with possibly significant implications for the incentives of the players involved in the development
 47 of AI technologies and applications.

48 **Modeling the impact of regulatory regimes on AI performance and safety.** Given that there
 49 are a range of different possible approaches to targeting AI regulation and assessing the impact of
 50 each alternative empirically is prohibitive, formal models can enable reasoning about the various
 51 regulatory impacts. This paper puts forward a strategic model of the interactions between a general-
 52 purpose technology producer (G) and a domain specialist (D), building on the “fine-tuning games”
 53 model proposed by Laufer et al. [2024]. As the two actors develop an AI technology, they each
 54 decide whether and how to invest in two key attributes of technology: *performance*, denoted by α ,
 55 and *safety*, denoted by β . We assume these actors are operating in a market; each actor experiences
 56 some cost for their investment in safety and performance, and obtains a share of the revenue out of
 57 the deployment of the AI product/service in the market.

58 To provide some intuition for what this investment pattern might look like, imagine a firm, G , pro-
 59 ducing a general-purpose language model that may be used in three domains – say, by healthcare
 60 providers (D_1), law firms (D_2), and financial services (D_3). The general-purpose developer moves
 61 first, and in light of the particular costs she faces and the anticipated responses from the downstream
 62 players, she chooses a certain strategy, represented by a pairing of performance and safety invest-
 63 ments (α_0, β_0). Once this investment has been made, the attributes of the technology at this stage
 64 can be thought of as akin to a ‘base camp,’ from which domain specialists may choose to climb fur-
 65 ther by investing their own effort toward improving the technology’s safety and/or performance in
 66 their respective domains. Of course, each domain faces their own delicate balance of safety risks and
 67 performance costs, so the ultimate safety and performance pairs (α_i, β_i) ($i = 1, 2, 3$) differ across
 68 the three domains. See Figure 1 (a) for a visualization of the investment decisions make by G , D_1 ,
 69 D_2 , and D_3 .

70 Equipped with this intuition about how these actors behave in an unregulated market, we now turn
 71 to our notion of regulation. We conceive of regulation as shaping the game in which players choose
 72 their strategies. In particular, this paper will focus specifically on safety regulation. We assume
 73 regulation imposes a constraint in the form of a lower bound on the players’ choice of safety in-
 74 vestment (i.e., β_i ’s). If a player does not meet the regulatory lower bound on safety, they will be
 75 penalized. This regulatory regime can be described using two parameters (θ_G, θ_D), representing the
 76 set of thresholds constraining the strategy space of G and D , respectively. The regulation can target
 77 the domain-specialist only ($\theta_G = 0, \theta_D > 0$), the generalist only ($\theta_G = \theta_D > 0$), both players
 78 ($\theta_D > \theta_G > 0$), or neither ($\theta_G = \theta_D = 0$). In addition to the decision of who to target, of course,
 79 the regulation encodes a decision about what level to set the safety standards. Smaller values of θ
 80 are less costly to comply with, and hence capture weaker safety requirements.

81 **First insight: Weak safety regulation can backfire.** Turning back to our example in Figure 1,
 82 we observe that something striking happens in the second panel, which depicts a scenario where
 83 regulation is targeted at the domain-specialist. In this scenario, the safety investment has gotten
 84 worse. How could safety regulation – a simple floor dictating a minimum investment level – lead
 85 to a less safe product? The mechanism leading to this phenomenon arises because the generalist G
 86 is aware of the regulatory safety requirements imposed on domain-specialists, and can use it to her
 87 advantage. When the regulator requires that a technology meets a certain level of safety investment
 88 by the time it reaches the market, the generalist has an opportunity to engage in a sort of *free-*
 89 *riding* behavior. The generalist is comfortable setting up the base camp at lower altitude, because

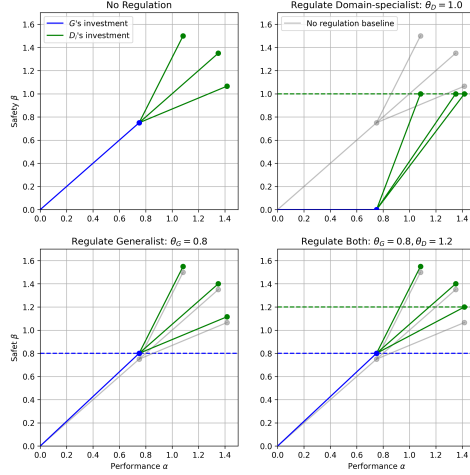


Figure 1: An instance of our model with one general-purpose producer and three domain-specialists. Regulating the domain-specialist alone (upper right) exhibits *backfiring* for all three domains, meaning the regulated safety level is lower than it would be without regulation. Regulating the generalist alone (lower left) improves the safety level slightly compared to no-regulation. Finally, a regime that targets both generalist and specialists with regulation (lower right) is able to 1) retain the improved safety performance from regulating the generalist, 2) improve the safety level of least-safe domain-specialist, while 3) avoiding backfiring.

90 she knows that the domain-specialist nonetheless has to climb to a level of investment that complies
 91 with regulation.

92 The scenario described above depicts a single instance of a more general phenomenon, which we
 93 describe as regulatory *backfiring*. A safety regulation backfires if it yields a total investment in safety
 94 lower than the safety investment achieved with no regulation. We identify a number of properties of
 95 this phenomenon – for example, backfiring only occurs when the regulation is *weak*, meaning the
 96 floor on safety is at or below the level reached in the absence of regulation. Backfiring can occur
 97 when D is targeted with regulation or when both G and D are targeted with regulation, but does
 98 not occur when only G is targeted. Our results suggest that this non-monotonic effect of regulation
 99 occurs for a broad set of games with different cost and revenue functions. Analytically, we prove
 100 that backfiring occurs for all quadratic-cost games in which the players invest any non-zero amount
 101 in both performance and safety without regulation.

102 **Second insight: Properly-placed safety regulation can improve the technology and the players’**
 103 **utilities.** While weak regulation targeted predominantly at the domain-specialist can backfire, our
 104 results suggest that other regulatory regimes fare better. When safety standards are directed at both
 105 G and D with appropriate strength, regulation can improve not just safety, but the utilities of both
 106 players, defined as their revenue share minus their investment cost. This result might seem unintu-
 107 itive: Regulation only reduces the set of choices available to each actor in our model, so how can
 108 regulation lead to choices that mutually benefit both generalist and specialist? What is stopping the
 109 players from choosing utility-optimal strategies in the absence of regulation? The reason this phe-
 110 nomenon occurs is a Prisoner’s Dilemma-style result: The players’ unregulated strategies, which are
 111 chosen to maximize their individual utility, fail to yield the strategies that that are globally optimal
 112 for both players. By constraining the actors away from the strategies that enable this kind of selfish
 113 behavior, regulation can act as a commitment device. The generalist can increase her investments in
 114 safety with the assurance that the domain specialist will contribute, too, rather than free-ride off of
 115 G ’s efforts.

116 Games can exhibit both backfiring and mutualistic regulations, depending on who is targeted and
 117 at what threshold. For example, Figure 2 depicts a particular instance of our game setting with
 118 one generalist and one domain-specialist. For the particular cost and revenue functions depicted,
 119 backfiring regulations and Pareto-improving regulations are possible, and the regulations yielding
 120 these effects are visualized. This figure represents a systematic sweep of all pairs of thresholds
 121 directed at the generalist, the domain specialist, or both. The pair of thresholds $(0, 0)$ corresponds

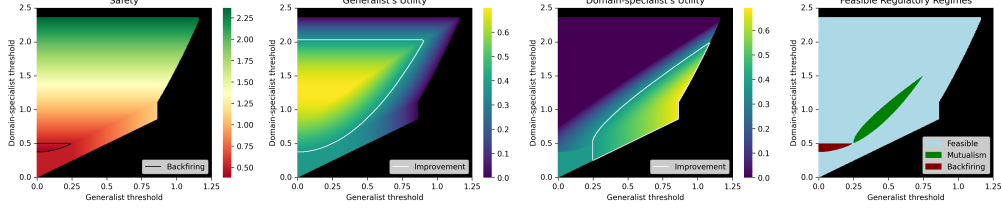


Figure 2: Simulated results for an example of two-player AI regulation model with quadratic costs. Players make costly investments in performance and safety (visualized on left), and then receive some share of revenue that comes from the total investment levels. The players’ utilities – defined as their share of the revenue minus the cost of their investment – is visualized for the Generalist (second from left) and the Domain-Specialist (third from left). Colors represent different utility outcomes depending on different combinations of regulatory constraints (θ_G, θ_D) which constrain the players’ safety investments. The game is solved over a grid of plausible regulations: $\theta_G \in [0, 1.25]$, $\theta_D \in [\theta_G, 2.5]$ using increments of 0.005, with a total of 105,651 simulated regulation games. Regulations that lead the players to abstain are depicted in black. There exists a region where non-zero regulation yields lower safety than no regulation (highlighted on leftmost plot). There also exists a region where regulation yields improvements to each players’ utility (highlighted on two center plots). The rightmost plot summarizes our results by showing the backfiring and mutualism outcomes in the θ_G, θ_D space. Parameter values for producing the plot: $C_0 = C_1 = I_2$, $r_\alpha = r_\beta = 1$, and $\delta = 0.5$.

to the case of no regulation. The safety implications of various regulations are depicted using a red-yellow-green color scale in the leftmost plot, while the utility implications for the generalist and specialist are depicted using a purple-green-yellow color scale in the center plots. Compared to the safety and utility values at the origin points, the backfiring and Pareto-improving regions are regulations which lead to lower safety and higher utilities for both players, respectively. Although this figure depicts an example of a single game, our analysis proves that these backfiring and Pareto-improving regulations exist for a broad class of games with quadratic costs. Namely, we find that backfiring occurs in all games in which the market incentivizes some non-zero investment in both performance and safety without regulation. Our characterization of when this phenomenon occurs includes *separable* scenarios (where the cost of investing in performance is independent of the cost of investing in safety), *complementary* scenarios (where investing in one makes the other cheaper), and weakly *interfering* scenarios (where the cost of investing in one makes the other more expensive) up to a certain bound, which we specify. We provide similar bounds for the mutualism results.

Related work

AI Safety Regulation. The rise of AI-related incidents have motivated several AI incident repositories to keep track of common risks [McGregor, Abercrombie et al.]. Scholars have attempted to taxonomize AI harms to make sense of the growing array of incidents Weidinger et al. [2022], Shelby et al. [2023]. Some existing AI risk taxonomies organize risks primarily by *domains*. These include risks to the *physical or psychological well-being* of people, *human rights and civil liberties*, *political and economic structures*, *society and culture*, and *the environment* Abercrombie et al. [2024]. Others categorize these risks based on how they arise, including *malicious use*, *malfunctions*, or *systemic effects* from wide adoption Bengio et al. [2025]. In our stylized model, we capture all such considerations using a single scalar that can be toggled by players through investments in safety. Common themes in policy drafts and recommendations stress the importance of balancing the goals of innovation and risk reduction, appropriately defining and targeting thresholds, and the impacts on incentives Chayes et al. [2025], Gaske [2023].

Game-theoretic models of AI development. A line of work uses formal models to reason about the strategic and social implications of machine learning (e.g., Hardt et al. [2016], Liu et al. [2022], Blum et al. [2021], Harris et al. [2021], Donahue and Kleinberg [2021]). More recently, there have been proposals for using modeling approaches to understand the social and safety implications of generative AI Dean et al. [2024], Sun et al. [2025]. Attempts to model the development process of generative AI often make use of the observation that development is *sequential* and involves *multiple*

154 *interacting actors* Cen et al. [2023]. Many existing works explore different strategic aspects of the
 155 market for AI using a stackelberg game. For example, Taitler et al. [2025] use a sequential game
 156 to explore incentives for data-sharing. Further time-steps, players and decisions have been added to
 157 explore particular topics, including the level of openness and market entry dynamics Xu et al. [2024],
 158 Wu et al. [2025]. Taitler and Ben-Porat [2025] introduce a particular notion of regulation in a related
 159 game-theoretic setting, and similar to our paper, they conceive of regulation as a restriction on the
 160 strategy space for developers of generative AI. Though work explicitly examining the interaction
 161 between performance and safety attributes in this setting is limited, Jagadeesan et al. [2024] explores
 162 the interaction between these attributes in a linear regression setting in order to understand firms’
 163 market entry decisions.

164 A Model of Regulating AI Safety

165 Here we offer a formal model for analyzing the effects of regulation on the development of AI ap-
 166 plications. Our model is a sequence of sub-games between two players. Each player will choose
 167 whether and how to contribute to the technology at a certain point in the development of the tech-
 168 nology, and some revenue is received depending on the ultimate attributes of the technology. The
 169 players are constrained by regulatory floors on safety, which will be set exogenously by a regulator.

170 **Players.** A general-purpose producer, referred to as G , invests in a technology that may be adapted
 171 by domain-specialist(s), referred to as D_i . The generalist is the first to invest in the technology,
 172 meaning that before G moves, the technology’s attributes begin at value 0. Each specialist D_i makes
 173 an investment after the generalist has moved.

174 **Technology.** We say a technology is described by one or more non-negative attributes $\gamma \in \mathbb{R}^d$. In
 175 this paper, we are interested in two attributes in particular: *performance* and *safety*. Unless otherwise
 176 specified, we assume $d = 2$ and that $\gamma = [\alpha, \beta]$ where α refers to performance and β refers to safety.

177 **Economic interests.** Each player, acting in a way that maximizes their self-interest, invests some
 178 non-zero amount in the technology. G invests to γ_0 and each D_i further invests to γ_i . Accordingly,
 179 each must pay a cost for their investment, $\phi_0(\gamma_0)$ and $\phi_i(\gamma_i; \gamma_0)$, respectively. After both players
 180 invest, they share a revenue that is brought in as a function of the ultimate attributes of the technology
 181 in domain i , $r_i(\gamma_i)$. We assume that, for some $\delta_i \in [0, 1]$, G gets $\delta_i r_i(\gamma_i)$ in revenue and D_i gets
 182 $(1 - \delta_i)r_i(\gamma_i)$. δ_i could either be exogenously fixed and given ahead of the game play, or it can be
 183 the result of bargaining between G and D_i . When we analyze a game with only one specialist, we
 184 will drop the subscript and use δ .

185 **Regulation** We model regulation as imposed exogenously on the environment. Regulation is a *min-*
 186 *imum constraint* on the safety investment that the players make. A regulation that targets G ’s invest-
 187 ment is characterized by a value $\theta_G \in \mathbb{R}^+$. A non-zero regulation would constrain G ’s strategy such
 188 that $\gamma_0[1] \geq \theta_G$. A regulation targeted at the domain-specialist, similarly, would take the form θ_D
 189 and lead the domain-specialist to be constrained in their strategy so $\gamma_i[1] \geq \theta_D$.

190 **Gameplay.** The game proceeds as a sequence of subgames:

- 191 • Regulation $\{\theta_G, \theta_D\}$ is announced.
- 192 • G chooses to either abstain or invest in the technology, bringing it to $\gamma_0 = \begin{bmatrix} \alpha_0 \\ \beta_0 \end{bmatrix}$.
- 193 • D_i chooses to either abstain or invest in the technology, bringing it to $\gamma_i = \begin{bmatrix} \alpha_i \\ \beta_i \end{bmatrix}$.
- 194
- 195 • The technology brings in revenue $r_i(\gamma_i)$, which will be shared such that G receives $\delta_i r_i(\gamma_i)$
 196 and D receives $(1 - \delta_i)r_i(\gamma_i)$.

The utilities of the players are given below:

$$U_G := \sum_i \delta_i r_i(\gamma_i) - \phi_0(\gamma_0); \quad U_{D_i} := (1 - \delta_i)r_i(\gamma_i) - \phi_i(\gamma_i; \gamma_0).$$

197 The best-response sub-game perfect equilibrium strategy for the generalist and specialist, respec-
 198 tively, can be expressed as the following optimization problems:

$$\begin{aligned}\gamma_0^* &:= \arg \max_{\gamma_0} U_G \text{ s.t. } \beta_0 \geq \theta_G; \\ \gamma_i^* &:= \arg \max_{\gamma_i} U_{D_i} \text{ s.t. } \beta_i \geq \theta_G.\end{aligned}$$

199 Finally, the players will opt to abstain, if they prefer 0 utility to any other feasible strategy. If
 200 either player chooses to abstain, then *both* players receive 0 utility.

201 Closed-Form Solutions

In this section, we analyze our regulation game where players' cost functions can be expressed as a two-degree quadratic equation. Specifying a quadratic function over two attributes requires defining a matrix of cost coefficients. The cross-terms in this matrix represent how investments the attributes interact with one another. For the technical portions of the paper, we use the case of one domain specialist (D) as our focus. We therefore have the following cost and revenue functions:

$$\begin{aligned}\phi_0(\gamma_0) &= \gamma_0^T C_0 \gamma_0, \\ \phi_1(\gamma_1; \gamma_0) &= (\gamma_1 - \gamma_0)^T C_1 (\gamma_1 - \gamma_0), \\ r(\gamma_1) &= r^T \gamma_1, \\ \text{where } C_0 &= \begin{bmatrix} c_{0,\alpha\alpha} & c_{0,\alpha\beta} \\ c_{0,\alpha\beta} & c_{0,\beta\beta} \end{bmatrix}; \quad C_1 = \begin{bmatrix} c_{1,\alpha\alpha} & c_{1,\alpha\beta} \\ c_{1,\alpha\beta} & c_{1,\beta\beta} \end{bmatrix}; \quad r = \begin{bmatrix} r_\alpha \\ r_\beta \end{bmatrix}.\end{aligned}$$

The players' utilities can thus be expressed as:

$$\begin{aligned}U_G &:= \delta r^T \gamma_1 - \gamma_0^T C_0 \gamma_0, \\ U_D &:= (1 - \delta) r^T \gamma_1 - (\gamma_1 - \gamma_0)^T C_1 (\gamma_1 - \gamma_0).\end{aligned}$$

202 It should be noted that not all values for the above parameters correspond to realistic or in-
 203 teresting scenarios. For example, we assume that the diagonal entries of both cost matrices
 204 $c_{0,\alpha\alpha}, c_{0,\beta\beta}, c_{1,\alpha\alpha}, c_{1,\beta\beta}$ are non-negative, to capture that investments in goods like safety and per-
 205 formance should have non-zero increasing cost. Although the cross-terms of the cost matrices can
 206 be negative, we require that $c_{0,\alpha\beta} > -\sqrt{c_{0,\alpha\alpha}c_{0,\beta\beta}}$ and $c_{1,\alpha\beta} > -\sqrt{c_{1,\alpha\alpha}c_{1,\beta\beta}}$, since it should
 207 not be that some combination of investments in α, β come at negative cost. Each players' choices
 208 over α and β should be considered as simultaneous across the two attributes, representing a joint
 209 optimization over performance and safety.

210 We start by providing sub-game perfect equilibria strategies in the case with no regulation, and then
 211 provide solutions for the regulated game. After stating the solved subgame perfect equilibria strate-
 212 gies, we will move to a slate of numerical results and findings analyzing the effects of regulation.

213 **Subgame perfect equilibria strategies without regulation** The no regulation solutions provided
 214 below can be seen as a strict generalization of the Fine-Tuning Games solutions Laufer et al. [2024]
 215 to games with two attributes that can interact.

Proposition 0.1. *Given an AI regulation game with quadratic costs, no regulation, and revenue-sharing parameter δ , domain specialist D 's subgame perfect equilibrium strategy is one of the values in the following set:*

$$\gamma_1^* \in \left\{ \gamma_0 + \frac{(1-\delta)}{2} C_1^{-1} r, \begin{bmatrix} \alpha_0 \\ \beta_0 + \frac{(1-\delta)r_\beta}{2c_{1,\beta\beta}} \end{bmatrix}, \begin{bmatrix} \alpha_0 + \frac{(1-\delta)r_\alpha}{2c_{1,\alpha\alpha}} \\ \beta_0 \end{bmatrix}, \begin{bmatrix} \alpha_0 \\ \beta_0 \end{bmatrix} \right\}$$

216 The strategy is the feasible candidate which maximizes U_D , subject to $U_D \geq 0, \alpha_1 \geq \alpha_0, \beta_1 \geq \beta_0$.

Proposition 0.2. *Given a two-player AI regulation game with quadratic costs, no regulation, and revenue-sharing parameter δ , G 's best-response is one of the following candidates:*

$$\gamma_0^* \in \left\{ \frac{\delta}{2} C_0^{-1} r, \begin{bmatrix} 0 \\ \frac{\delta r_\beta}{2c_{0,\beta\beta}} \end{bmatrix}, \begin{bmatrix} \frac{\delta r_\alpha}{2c_{0,\alpha\alpha}} \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}.$$

217 The strategy is the candidate which maximizes U_G , subject to $U_G \geq 0, U_D \geq 0, \alpha_1 \geq 0, \beta_1 \geq 0$.

218 The proofs of the above two propositions are given in the Appendix. The solutions offer intuition
 219 about the set of strategies players might opt to take. They may venture in the direction of some
 220 combination of performance and safety, that is, move to a point that does not reside on either con-
 221 straint. Or, alternatively, they may creep along the axes constraining their strategy space, and invest
 222 minimally in either performance or safety.

223 When do the players prefer one of these strategies over another? In general, our solutions are pro-
 224 vided as sets of candidates because there are multiple intersecting constraints that must be checked
 225 to ensure a given candidate is optimal. However, our analysis reveals classes of games in which the
 226 market will lead players to invest in both safety and performance in conjunction under no regulation.
 227 We make this claim formal below.

Remark 0.3. *Given the AI regulation game with quadratic costs, no regulation, and revenue-sharing parameter $\delta \in (0, 1)$. If any player p 's cost interaction term satisfies the following inequalities:*

$$c_{p,\alpha\beta} < \min \left(\sqrt{c_{p,\alpha\alpha}c_{p,\beta\beta}}, \frac{c_{p,\alpha\alpha}r_\beta}{r_\alpha}, \frac{c_{p,\beta\beta}r_\alpha}{r_\beta} \right),$$

228 *then their best-response strategy includes non-zero investment in both performance and safety.*

229 This claim is proven in the Appendix. The broad intuition is that the first of these inequalities estab-
 230 lishes the costs are strictly convex, and the second two ensure that the player's cost interactions are
 231 not so positive that investing in both performance and safety is prohibitively expensive compared
 232 to investing in one or the other alone. The claim offers some intuition for when a player prefers to
 233 invest in both attributes together, even without regulation pushing them to invest in safety. It covers
 234 all games in which the cost interactions are negative, which we call the *complementary scenario*,
 235 meaning it is cheaper to invest in both performance and safety together than to invest in each individ-
 236 ually. It further covers all games in which the cost interactions are zero, which we call the *separable*
 237 *scenario*, meaning there is no benefit or loss to investing in both attributes in conjunction. Finally,
 238 it covers certain instances where the cost interactions are positive, which we call the *interfering*
 239 *scenario*, meaning safety investments make performance more costly, and vice versa.

240 **Subgame perfect equilibria strategies with regulation.** For brevity, we refer the reader to the Ap-
 241 pendix for the best-response strategies, since they take space to state and are somewhat clunky. The
 242 form of problem we are dealing with is a continuous, not-necessarily-convex optimization problem
 243 with a constant number of constant-degree polynomials in a constant number of variables. Broadly,
 244 the strategy is to put forward a small number of candidate points that must be checked using a limited
 245 number of steps. These checks can be implemented numerically.

246 Computational results

247 Here we describe a set of numerical tests and demonstrations to explore the strategies in our game,
 248 using the solved strategies from the previous section. Our analysis here is focused on the existence
 249 of a persistent facet of the model concerning the way the players shift their strategies in response to
 250 regulation. With the knowledge that one player or the other is required to meet a regulatory floor,
 251 agents can choose their strategies accordingly. In a variety of cases, we observe that the strategies
 252 shift in a way that *lowers* the ultimate safety investment compared to safety attained under no regu-
 253 lation. This effect – which we term *backfiring* – is observable in cases where the regulation is weak,
 254 meaning it imposes a floor that the players already meet under no regulation.

255 This section starts by demonstrating the existence of this effect. We then discuss its persistence in
 256 cases where players can flexibly choose how they share revenue via a linear contract. Finally, in
 257 stark contrast to the observation that regulation can backfire, we find that regulation can act as a
 258 commitment device, unlocking strategy sequences that mutually benefit the players.

259 Our computational findings are organized around three main observations, with accompanying fig-
 260 ures. Our observations are enumerated below.

261 **Finding 1: Regulation can backfire.** This game is *separable*, meaning there are no interaction ef-
 262 fects between performance and safety, and it assumes the market without regulation places equal
 263 value on performance and safety. Figure 3 depicts the players' strategies in this game, for varying
 264 levels of regulation targeting the Domain-specialist alone. For the lowest regulatory thresholds, we
 265 observe that the players stick to their no-regulation safety investments, since they already clear the

threshold and their no-regulation investments remain optimal. As the regulatory floor is increased, however, the generalist’s strategy exhibits a discontinuity. Crucially, this drop in G’s safety investment occurs at a regulatory threshold *lower than* the no-regulation safety strategy.

Finding 2: Bargaining does not suffice to prevent backfiring. We now relax the assumption that players share their revenue according to a constant revenue-sharing parameter $\delta = 0.5$. We provide evidence that even when players can distribute revenue in a way that maximizes the joint utility, these arrangements can still exhibit backfiring effects. We assume here that the players jointly agree on a bargaining solution *before* either invests effort, but *after* learning about the regulation.² Figure 4 shows the numerical results for a variant of the separable game where we vary the value of δ over 98 values in the range $[0.01, 0.99]$. We vary the regulatory setting for 13 θ_G values in $[0, 1.2]$ and 51 θ_D values in $[0, 2.5]$, for a total of 49,686 simulated games. The figure depicts three different processes for arriving at an optimal bargain: *utilitarian*, which selects δ to maximize the sum of utilities, *Nash*, which selects δ to maximize the product of utilities Nash et al. [1950], and *egalitarian*, which sets δ to maximize the minimum of the utilities. In all scenarios, we observe at least one instance of a combination of regulations that backfire. Further, we observe a cluster of regulation regimes that yield mutual improvement to utility.

Finding 3: Regulation can act as a commitment device. Here we show that there exist cases where regulation can leave both players *better off* than anarchy, while also benefiting the safety of the technology. Even though the regulation constrains the space of investments that players are able to achieve, it can nonetheless leave each player with higher utility than they are able to achieve under no regulation. To make this finding more clear, we depict the set of all achievable (U_G, U_D) combinations in Figure 5. The light blue cloud of points represents all attainable utility scenarios, over a grid of θ_D , θ_G , and δ values. The dotted lines represent the convex hull (northeastern faces) of attainable utility implications for the following regimes: 1) neither player is targeted with regulation (depicted in green), 2) one player is targeted with regulation (depicted in red and black), and 3) both players are targeted with regulation (inferred from the outermost feasible points). The figure suggests that a non-vacuous constraint on *both* players achieves more preferable utility outcomes than regulations of individual players or bargaining alone are able to achieve.

These results suggest that, although regulation can backfire, it can also mutually serve the interests of both players while also improving the level of safety of the technology. This finding raises the following question: if it was possible to achieve higher utilities all around, why was this set of strategies not chosen by the players in the unregulated game? Absent regulation, the players might *wish* they could ensure the other will uphold their side of a verbal agreement, though they are unable to guarantee it. Regulation, therefore, can act as a *commitment device*, which lends teeth to agreements that the players are able to enter prior to making their investments. This commitment device can be valuable in a formal sense: Both players would be willing to pay for it, as long as the price is less than the amount of utility they collectively gain under regulation.

A General Characterization

In the previous sections, we arrived at closed-form solutions for the players’ strategies and have demonstrated individual instances that exhibit the backfiring effect of regulation. We have not yet determined how widespread this phenomenon is. In this section, we provide analytical results that characterize when this phenomenon occurs. Our findings suggest that this effect is notably widespread. We find that for all quadratic-cost games, backfiring occurs as long as both of the technology’s attributes (performance and safety) are sufficiently *complementary* such that, under no regulation, the players will invest in some combination of them. Intuitively, if the players invested only in performance under no regulation, backfiring would be impossible as the baseline safety investment would be zero. Therefore, our condition for backfiring covers all games where the market prefers some non-zero baseline investment in performance and safety. The condition we rely on is precisely the condition introduced in Remark 0.3, which represents an upper bound on the cost interaction terms. This section will prove that both backfiring and mutualism occur in a range of scenarios that depend crucially on the cost interaction term, and will describe what this dependence looks like.

²The next sections will relax this assumption further, providing findings on the existence of backfiring and mutualism for every non-trivial linear revenue-sharing agreement $\delta \in (0, 1)$.

317 **Backfiring occurs in all mixed-strategy games.**

318 Below we prove that for all AI regulation games in which the players invest a non-zero amount
 319 in safety and performance under no regulation, there is a non-empty set of regulatory regimes that
 320 exhibit a backfiring effect.

Theorem 0.4. *Given an AI regulation game with quadratic costs. If both players' cost interactions meet the following conditions:*

$$c_{p,\alpha\beta} < \min \left(\sqrt{c_{p,\alpha\alpha}c_{p,\beta\beta}}, \frac{c_{p,\alpha\alpha}r_\beta}{r_\alpha}, \frac{c_{p,\beta\beta}r_\alpha}{r_\beta} \right),$$

321 *then there exists an $\epsilon > 0$ such that the regulatory regime $\theta_G = 0, \theta_D = \beta_0^A - \epsilon$ backfires.*

322 The proof of the above theorem is provided in Appendix . Here we provide an overview of the
 323 conceptual argument. We start by observing that the unregulated optimal strategies γ_0^A, γ_1^A remain
 324 feasible in weak regulatory settings. These strategies dominate all alternative strategies in which the
 325 players contribute to safety *beyond* their regulatory constraints, as any such strategy was available
 326 in the no regulation scenario, so they were already shown to be sub-optimal compared to γ_0^A, γ_1^A .
 327 The proof's task, therefore, is to find some $\theta_D < \beta_1^A$ and some $\gamma'_0 \neq \gamma_0^A$, such that D *minimally*
 328 *complies* with the regulation ($\beta'_0 = \theta_D$), and further, $U_G(\gamma'_0; \theta_D) > U_G(\gamma_0^A; \theta_D)$. For the proof to
 329 work, we choose a regulation of $\theta_G = 0, \theta_D = \beta_0^A - \epsilon$ for some small positive $\epsilon > 0$, and generalist
 330 strategy $\gamma'_0 = \left[\frac{\delta r_\alpha}{2c_{0,\alpha\alpha}} (\beta_0^A - 2\epsilon) \right]$. For sufficiently small ϵ , we find that the change to the utility
 331 of G for using this strategy is positive as long as the following condition is met: $r_\beta > \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} r_\alpha$. This
 332 inequality, given by the analysis in Appendix , is precisely the condition established in Remark 0.3
 333 for non-zero investment in safety under no regulation.

334 The above results demonstrate that backfiring does not only exist in single degenerate cases: It
 335 occurs in a range of scenarios in which players share revenue and each contribute non-zero effort
 336 to the development of the technology. These scenarios include settings in which the two attributes
 337 are *complementary*, as well as a range of settings where the two attributes are *interfering*, up to a
 338 particular limit that we are able to specify. We note that further generalizations are open for broader
 339 functional forms, including more expressive polynomial costs and exponential costs. The generality
 340 of the backfiring effect in the quadratic case gives us reason to believe that the effect might hold for
 341 a broader set of forms, though we leave these directions to future work.

342 **Mutualism occurs in sufficiently separable games.**

343 So far, we have shown that a set of regulations backfire in a swath of two-attribute games. Here we
 344 provide a second result on a set of regulations that fare better. Using similar logic about games with
 345 bounded interaction effects between the two attributes, we find that there exist combinations of
 346 regulatory thresholds that mutually improve the two players' utilities, as well as the safety level of
 347 the technology. We state this result below.

Theorem 0.5. *Given a two-player AI regulation game with quadratic costs. If both players meet the following conditions:*

$$|c_{p,\alpha\beta}| < \min \left(\sqrt{c_{p,\alpha\alpha}c_{p,\beta\beta}}, \frac{c_{p,\alpha\alpha}r_\beta}{r_\alpha}, \frac{c_{p,\beta\beta}r_\alpha}{r_\beta} \right),$$

348 *then there exists an $\epsilon > 0$ such that the regulatory regime $\theta_G = \beta_0^A + \epsilon, \theta_D = \beta_1^A + 2\epsilon$ mutually*
 349 *improves both players' utilities.*

350 The proof of the above theorem is given in Appendix . The proof follows a similar strategy to the
 351 backfiring proof. We are focused on the set of games where the players arrive at unconstrained
 352 solutions in the case of no regulation, and we perturb the regulation by a small positive ϵ value and
 353 see the implications for the players' utilities. Here, instead of targeting only the domain-specialist
 354 and specifying a threshold slightly below the unconstrained optimal strategy, we set the regulation
 355 to target *both players* using a threshold slightly *above* their unconstrained strategies. Instead of
 356 measuring the impact on safety, we measure the impact on the players' utilities and find that, under
 357 the specified condition, the utilities both improve.

358 The results suggest that, similar to the characterization of backfiring, the mutualism effect is observ-
 359 able in a range of quadratic-cost games, including in *separable scenarios* and a range of *complemen-*
 360 *tary* and *interfering* scenarios. Notice, however, that our condition for establishing when mutualism
 361 occurs is slightly different than the condition in the backfiring theorem. Instead of a one-sided bound
 362 on the players’ cost interaction terms, our proof relies on a two-sided bound. The analysis suggests
 363 there may be certain games where the two attributes are *strongly complementary* where slightly in-
 364 creasing the regulation in the manner proposed does not increase players’ utilities. In other words, if
 365 the market already sufficiently incentivizes joint investments in safety and performance, then forcing
 366 safety requirements on both players in equal proportion may not benefit players’ utilities. In these
 367 cases, a linear contract may suffice to serve the utilities of the players, and so regulation would only
 368 be needed for achieving the goal of advancing safety, and would not serve the additional role in
 369 enforcing commitments from players.

370 References

- 371 G. Abercrombie, D. Benbouzid, P. Giudici, D. Golpayegani, J. Hernandez, P. Noro, H. Pandit,
 372 E. Paraschou, C. Pownall, J. Prajapati, M. A. Sayre, U. Sengupta, A. Suriya-
 373 wongkul, R. Thelot, S. Vei, and L. Waltersdorfer. AIAAIC Harms taxonomy - v1.7.pdf.
 374 URL [https://drive.google.com/file/u/1/d/1-dqNDv2G43856tABKP-CSYTyAL8xo2kj/view?usp=](https://drive.google.com/file/u/1/d/1-dqNDv2G43856tABKP-CSYTyAL8xo2kj/view?usp=sharing&usp=embed.facebook)
 375 [sharing&usp=embed.facebook](https://drive.google.com/file/u/1/d/1-dqNDv2G43856tABKP-CSYTyAL8xo2kj/view?usp=sharing&usp=embed.facebook).
- 376 G. Abercrombie, D. Benbouzid, P. Giudici, D. Golpayegani, J. Hernandez, P. Noro, H. Pandit,
 377 E. Paraschou, C. Pownall, J. Prajapati, M. A. Sayre, U. Sengupta, A. Suriyawongkul, R. Thelot,
 378 S. Vei, and L. Waltersdorfer. A Collaborative, Human-Centred Taxonomy of AI, Algorithmic, and
 379 Automation Harms, July 2024. URL <http://arxiv.org/abs/2407.01294>. arXiv:2407.01294 [cs].
- 380 T. Alon, P. Dütting, Y. Li, and I. Talgam-Cohen. Approximate optimality of linear contracts under
 381 uncertainty. *arXiv preprint arXiv:2211.06850*, 2022.
- 382 M. Anderljung, J. Barnhart, A. Korinek, J. Leung, C. O’Keefe, J. Whittlestone, S. Avin,
 383 M. Brundage, J. Bullock, D. Cass-Beggs, et al. Frontier ai regulation: Managing emerging risks
 384 to public safety. *arXiv preprint arXiv:2307.03718*, 2023.
- 385 Y. Bengio, S. Mindermann, D. Privitera, T. Besiroglu, R. Bommasani, S. Casper, Y. Choi, P. Fox,
 386 B. Garfinkel, D. Goldfarb, et al. International ai safety report. *arXiv preprint arXiv:2501.17805*,
 387 2025.
- 388 J. Bessen and E. Maskin. Sequential innovation, patents, and imitation. *The RAND Journal of*
 389 *Economics*, 40(4):611–635, 2009.
- 390 A. Blum, N. Haghtalab, R. L. Phillips, and H. Shao. One for one, or all for all: Equilibria and opti-
 391 mality of collaboration in federated learning. In *International Conference on Machine Learning*,
 392 pages 1005–1014. PMLR, 2021.
- 393 R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg,
 394 A. Bosselut, E. Brunskill, et al. On the opportunities and risks of foundation models. *arXiv*
 395 *preprint arXiv:2108.07258*, 2021.
- 396 G. P. Cachon. Supply chain coordination with contracts. *Handbooks in operations research and*
 397 *management science*, 11:227–339, 2003.
- 398 G. Carroll. Robustness and linear contracts. *American Economic Review*, 105(2):536–563, 2015.
- 399 S. H. Cen, A. Hopkins, A. Ilyas, A. Madry, I. Struckman, and L. Videgaray Caso. Ai supply chains.
 400 2023.
- 401 J. T. Chayes, M.-F. Cuéllar, and L. Fei-Fei. Draft report of the joint california policy working group
 402 on ai frontier models. *Draft, Joint California Policy Working Group on AI Frontier Models*, 2025.
- 403 S. Dean, E. Dong, M. Jagadeesan, and L. Leqi. Accounting for ai and users shaping one another:
 404 The role of mathematical models. *arXiv preprint arXiv:2404.12366*, 2024.

405 K. Donahue and J. Kleinberg. Model-sharing games: Analyzing federated learning under voluntary
406 participation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages
407 5303–5311, 2021.

408 P. Dütting, T. Roughgarden, and I. Talgam-Cohen. Simple versus optimal contracts. In *Proceedings*
409 *of the 2019 ACM Conference on Economics and Computation*, pages 369–387, 2019.

410 P. Dütting, T. Ezra, M. Feldman, and T. Kesselheim. Multi-agent contracts. In *Proceedings of the*
411 *55th Annual ACM Symposium on Theory of Computing*, pages 1311–1324, 2023.

412 P. Dütting, T. Ezra, M. Feldman, and T. Kesselheim. Multi-agent combinatorial contracts. In
413 *Proceedings of the 2025 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages
414 1857–1891. SIAM, 2025.

415 M. R. Gaske. Regulation priorities for artificial intelligence foundation models. *Vand. J. Ent. &*
416 *Tech. L.*, 26:1, 2023.

417 J. R. Green and S. Scotchmer. On the division of profit in sequential innovation. *The Rand journal*
418 *of economics*, pages 20–33, 1995.

419 S. J. Grossman and O. D. Hart. An analysis of the principal-agent problem. In *Foundations of*
420 *insurance economics: Readings in economics and finance*, pages 302–340. Springer, 1992.

421 M. Hardt, N. Megiddo, C. Papadimitriou, and M. Wootters. Strategic classification. In *Proceedings*
422 *of the 2016 ACM conference on innovations in theoretical computer science*, pages 111–122,
423 2016.

424 K. Harris, H. Heidari, and S. Z. Wu. Stateful strategic regression. *Advances in Neural Information*
425 *Processing Systems*, 34:28728–28741, 2021.

426 M. Jagadeesan, M. I. Jordan, and J. Steinhardt. Safety vs. performance: How multi-objective learn-
427 ing reduces barriers to market entry. *arXiv preprint arXiv:2409.03734*, 2024.

428 B. Laufer, J. Kleinberg, and H. Heidari. Fine-tuning games: Bargaining and adaptation for general-
429 purpose models. In *Proceedings of the ACM on Web Conference 2024*, pages 66–76, 2024.

430 L. T. Liu, N. Garg, and C. Borgs. Strategic ranking. In *International Conference on Artificial*
431 *Intelligence and Statistics*, pages 2489–2518. PMLR, 2022.

432 S. Longpre, K. Klyman, R. E. Appel, S. Kapoor, R. Bommasani, M. Sahar, S. McGregor, A. Ghosh,
433 B. Bili-Hamelin, N. Butters, et al. In-house evaluation is not enough: Towards robust third-party
434 flaw disclosure for general-purpose ai. *arXiv preprint arXiv:2503.16861*, 2025.

435 S. McGregor. Preventing repeated real world ai failures by cataloging incidents: The ai incident
436 database. URL <https://incidentdatabase.ai/>.

437 J. F. Nash et al. The bargaining problem. *Econometrica*, 18(2):155–162, 1950.

438 X. Qi, Y. Zeng, T. Xie, P.-Y. Chen, R. Jia, P. Mittal, and P. Henderson. Fine-tuning aligned language
439 models compromises safety, even when users do not intend to! *arXiv preprint arXiv:2310.03693*,
440 2023.

441 X. Qi, B. Wei, N. Carlini, Y. Huang, T. Xie, L. He, M. Jagielski, M. Nasr, P. Mittal, and P. Hen-
442 derson. On evaluating the durability of safeguards for open-weight llms. *arXiv preprint*
443 *arXiv:2412.07097*, 2024.

444 S. A. Ross. The economic theory of agency: The principal’s problem. *The American economic*
445 *review*, 63(2):134–139, 1973.

446 R. Shelby, S. Rismeni, K. Henne, A. Moon, N. Rostamzadeh, P. Nicholas, N. Yilla, J. Gallegos,
447 A. Smart, E. Garcia, and G. Virk. Sociotechnical Harms of Algorithmic Systems: Scoping a Tax-
448 onomy for Harm Reduction, July 2023. URL <http://arxiv.org/abs/2210.05791>. arXiv:2210.05791
449 [cs].

- 450 H. Sun, Y. Wu, Y. Cheng, and X. Chu. Game theory meets large language models: A systematic
451 survey. *arXiv preprint arXiv:2502.09053*, 2025.
- 452 B. Taitler and O. Ben-Porat. Selective response strategies for genai. *arXiv preprint*
453 *arXiv:2502.00729*, 2025.
- 454 B. Taitler, O. Madmon, M. Tennenholtz, and O. Ben-Porat. Data sharing with a generative ai com-
455 petitor. *arXiv preprint arXiv:2505.12386*, 2025.
- 456 W. K. Viscusi and M. J. Moore. Product liability, research and development, and innovation. *Journal*
457 *of Political Economy*, 101(1):161–184, 1993.
- 458 L. Weidinger, J. Uesato, M. Rauh, C. Griffin, P.-S. Huang, A. Glaese, B. Balle, A. Kasirzadeh,
459 C. Biles, S. Brown, Z. Kenton, W. Hawkins, T. Stepleton, A. Birhane, L. A. Hendricks, L. Rimell,
460 W. Isaac, J. Haas, S. Legassick, G. Irving, and I. Gabriel. Taxonomy of Risks posed by
461 Language Models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability,*
462 *and Transparency*, FAccT ’22, pages 214–229, New York, NY, USA, June 2022. Association
463 for Computing Machinery. ISBN 978-1-4503-9352-2. doi: 10.1145/3531146.3533088. URL
464 <https://dl.acm.org/doi/10.1145/3531146.3533088>.
- 465 Y. Wu, H. Duan, X. Li, and X. Hu. Navigating the deployment dilemma and innovation paradox:
466 Open-source versus closed-source models. In *Proceedings of the ACM on Web Conference 2025*,
467 pages 1488–1501, 2025.
- 468 F. Xu, X. Wang, W. Chen, and K. Xie. The economics of ai foundation models: Openness, compe-
469 tition, and governance. *Competition, and Governance (August 11, 2024)*, 2024.

470 Further related work

471 **Economic theory and contracts.** Our work leverages pre-existing approaches that are common in
472 the theoretical economics and game theory literatures to reason about the set of possible impacts
473 of AI safety regulation. In particular, we draw inspiration from canonical works in contract theory
474 Grossman and Hart [1992], Ross [1973] and the coordination of supply chains Cachon [2003].
475 Our model is a variant of a Principal-Agent problem in which the strategy space is defined by two
476 real-valued attributes, and the cost and revenue are functions of these attributes. In this way, our
477 model draws inspiration from Viscusi and Moore [1993] analyzing the possible effects of products
478 liability schemes on innovation and safety. That model — a one-player model with no order-of-
479 play effects — demonstrates that liability does not, necessarily, hamper innovation. We assume that
480 innovation is sequential, meaning that an entity’s investment in safety or performance builds on the
481 contributions of past investments Bessen and Maskin [2009], Green and Scotchmer [1995].³ In what
482 we call the ‘no-regulation’ game, we assume the players revenue-share via a linear contract Dütting
483 et al. [2019], a common assumption in the literature (e.g., Dütting et al. [2025, 2023], Alon et al.
484 [2022], Carroll [2015]). However, one way to interpret our mutualism results is as a demonstration
485 that linear contracts are sub-optimal in our setting. Our notion of regulation can be viewed as a set
486 of non-linear contracts defined by a set of strategy constraints, and our results suggest these more
487 expressive contracts can yield higher utility. Of course, still other forms of contracts are possible and
488 may yield different utility implications. We leave these directions to future work.

489 **The fine-tuning games model.** Our work builds on and extends the *fine-tuning games* model pro-
490 posed in Laufer et al. [2024]. That model builds a one-dimensional game in which players must
491 bargain over a revenue-sharing contract before investing in performance in sequence. We extend this
492 model in two ways: First, the players’ strategy space is two-dimensional in our model, to capture
493 the dynamic that often arises where a regulator wants to steer the technology in a direction (e.g.,
494 safety) other than that which is most-profitable (e.g., a baseline combination of performance and
495 safety, dictated by the unregulated market). Second, we introduce the regulation, which can be seen
496 as a *floor* constraining the feasible strategy space of each player. This allows us to explore when
497 targeting generalists, specialists, both or neither is preferable for achieving desiderata like safety.

³However, some have observed that safety investments can *degrade* as the result of fine-tuning performance investments especially when model weights are open Qi et al. [2023, 2024]. This scenario, and especially the interaction effects with model openness, are ripe areas for further analysis.

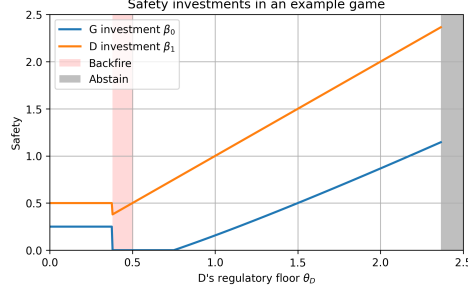


Figure 3: Backfiring observed in a basic two-player game where $\theta_G = 0$ and θ_D is varied over the range $[0, 2.5]$. As θ_D is swept upward from 0, there is some value at which the generalist’s score exhibits a discontinuity and the investment in safety lowers.

Further Discussion and Conclusion

Proposals for AI regulation have made use of the idea that different entities contribute to these technologies in succession. This work provides a model for reasoning about the effects of targeting AI safety regulation along the development chain. Our findings suggest that *weak* safety regulation predominantly targeted at the domain specialist can backfire, yielding lower investments in safety than in the alternative case of no regulation. Our findings further suggest that regulation appropriately targeted at both upstream producers and downstream specialists can exhibit a mutualism effect in which both entities benefit. After demonstrating instances of the backfiring and mutualism effects through a numerical simulation, we provide analysis showing these phenomena are not just degenerate cases but hold in a range of scenarios.

Our results reveal natural directions for future research. In the setting we have put forward, it would be interesting to move beyond showing the existence of backfiring and mutualism regions and characterize the shape of these regions and the magnitude of their effects. Certain segments of the boundaries of these regions are straightforward but others seem to require solving higher-order polynomials to express in closed-form.

Generalizations beyond the quadratic-cost games might be interesting. For instance, it may be possible to show that backfiring and Pareto-improvement effects occur for any convex cost and concave revenue games meeting where there exist some marginal conditions on the functions’ marginal conditions including their slopes and intercepts.

We have predominantly focused on the case where there is one domain-specialist, but in many real-world settings the development of AI technologies involve multiple domains, and each domain may involve many entities who compete. To what extent does competition between multiple entities change the backfiring and Pareto-improving impacts of regulation? Pursuing questions about multiple domain-specialists would require further specifying the structure of G ’s contract with each specialist, which might reasonably be conceived as a constant revenue share across domains, a constant fixed price across domains, or a variable price across domains. Relatedly, approaches to regulating different specialists may be conceived of as domain-specific (different requirements for each domain) or domain-agnostic (requirements for all domains). Pursuing questions about multiple generalists may also illuminate interesting directions. In particular, if different domains have different preferences over attributes, there may be scenarios where general providers *specialize* their investments to capture some domains and cede others to their competitors. Such dynamics raise new questions about how to design regulation to account for these rich constellations of interacting actors.

Subgame perfect equilibria strategies with regulation

Here we provide the subgame perfect equilibria strategies of the two players in our two-attribute game, in the presence of regulation. Notice that the no-regulation gameplay can be derived from these solutions simply by plugging in $\theta_D = \theta_G = 0$. Like the solutions in the prior section, these

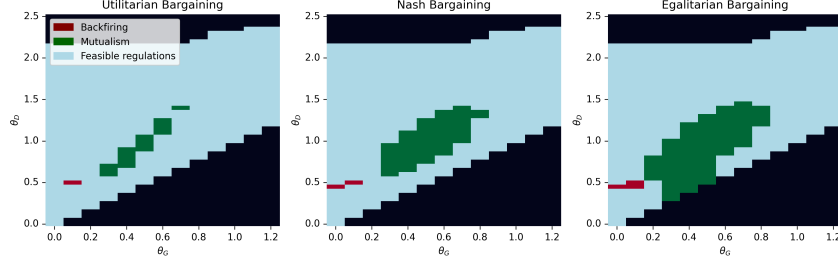


Figure 4: Results from numerical tests over the set of possible (θ_G, θ_D) pairs in the two-attribute, two-player, separable quadratic-cost game. Backfiring occurs in the weak regulatory regimes in which θ_D is just below β_0^A . Regulations that mutually improve both players' utilities over anarchy are detected for all three bargaining solutions. The highest aggregate utility in this game is achieved at $\theta_G = 0.5, \theta_D = 1$.

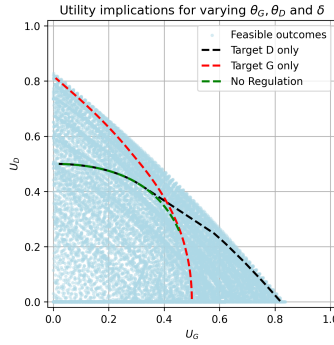


Figure 5: The set of attainable utility outcomes over a grid of possible regulation regimes and bargains for the two-attribute, two-player, separable quadratic-cost game. Each of the blue points represents a possible game with utility implications for the two players. If the players are restricted to a particular regulatory regime – targeting G only, D only, or neither – then the utility they are able to achieve (depicted in dashed lines) suffers, compared to the regime where both players are subjected to regulation.

535 generalized solutions require checking a number of candidates, but this number has grown to account
536 for the possible responses to regulation.

Proposition 0.6. *Given a two-attribute fine-tuning game with quadratic costs, regulatory constraints θ_G, θ_D , and bargaining parameter δ , the domain specialist D 's subgame perfect equilibrium strategy is one of the values in the following set:*

$$\gamma_1^* \in \left\{ \begin{array}{l} \gamma_0 + \frac{(1-\delta)}{2} C_1^{-1} r, \left[\beta_0 + \frac{(1-\delta)r_\beta}{2c_{1\beta\beta}} \right], \\ \left[\alpha_0 + \frac{(1-\delta)r_\alpha}{2c_{1\alpha\alpha}} - \frac{c_{1\alpha\beta}}{c_{1\alpha\alpha}} \max(0, \theta_D - \beta_0) \right], \\ \left[\alpha_0, \max(\beta_0, \theta_D) \right], \text{abstain.} \end{array} \right\}$$

537 The strategy is the feasible candidate which maximizes U_D , subject to $U_D \geq 0, \alpha_1 \geq \alpha_0, \beta_1 \geq$
538 $\max(\beta_0, \theta_D)$.

539 **Proposition 0.7.** *Given a two-attribute, two-player fine-tuning game with quadratic costs, regu-*
540 *latory constraints θ_G, θ_D , and bargaining parameter δ , G 's best-response is one of the following*
541 *candidates:*

$$\begin{array}{l} 542 \quad \bullet \frac{\delta}{2} C_0^{-1} r, \\ 543 \quad \bullet \left[\begin{array}{l} 0 \\ \frac{\delta r_\beta}{2c_{0\beta\beta}} \end{array} \right], \end{array}$$

- 544 • $\begin{bmatrix} \frac{\delta r_\alpha}{2c_{0\alpha\alpha}} - \frac{c_{0\alpha\beta}}{c_{0\alpha\alpha}}\theta_G \\ \theta_G \end{bmatrix},$
- 545 • $\begin{bmatrix} 0 \\ \theta_G \end{bmatrix},$
- 546 • *abstain,*
- 547 • *Three additional candidates along the $U_D = 0$ constraint, which is given by the following*
- 548 *quadratic equation:*

$$\begin{aligned}
& (1-\delta)r_\alpha\alpha_0 + \frac{(1-\delta)^2r_\alpha^2}{4c_{1\alpha\alpha}} + (1-r_\beta)\theta_D - \\
& \frac{c_{1\alpha\beta}}{c_{1\alpha\alpha}}(1-\delta)\theta_D + \frac{c_{1\alpha\beta}^2}{c_{1\alpha\alpha}\theta_D^2} - c_{1\beta\beta}\theta_D^2 + \\
& \left(\frac{c_{1\alpha\beta}}{c_{1\alpha\alpha}}(1-\delta)r_\alpha - 2\frac{c_{1\alpha\beta}^2}{c_{1\alpha\alpha}}\theta_D + 2c_{1\beta\beta}\theta_D \right) \beta_0 \\
& + \left(\frac{c_{1\alpha\beta}^2}{c_{1\alpha\alpha}} - c_{1\beta\beta} \right) \beta_0^2 = 0.
\end{aligned}$$

549 *The strategy is the candidate which maximizes U_G , subject to $U_G \geq 0, U_D \geq 0, \alpha_1 \geq 0, \beta_1 \geq \theta_G$.*

550 The proof of the above propositions is provided in Appendix . We outline the intuition behind the
551 proof as follows: Notice that the optimization is an inequality-constrained quadratic optimization
552 problem. The problem has been set up so no solutions exist at infinity, that is, the solutions will
553 either be local maxima or will reside on constraints. Therefore, we can find the critical points for
554 the unconstrained problem, as well as the critical points for every possible combination of every
555 constraint in our problem. This yields a set of candidates, which are worked out and listed in the set
556 above.

557 There is a bit of additional subtlety in the process for arriving at the last three candidates along
558 the constraint listed at the end of the Proposition. Two of the three candidates reside at the inter-
559 section of this constraint with the other constraints—that is, they satisfy the constraint listed and
560 either $\alpha_0 = 0$ or $\beta_0 = \theta_G$. Finding the point that satisfies these combinations of constraints is
561 only as hard as solving the roots of a one-variable quadratic, at worst. The third one, however, is
562 a bit more convoluted. This candidate can be described as the solution to the optimization problem
563 $\max_{\gamma_0} U_G$ s.t. $U_D = 0$, where the other constraints are ignored. Although this is a (not necessarily
564 convex) quadratic program, specifying the Lagrangian suggests that its solution must be the solution
565 of a system of three distinct equations with three unknown variables $(\alpha_0, \beta_0, \lambda) \in \mathbb{R}^3$. Two of these
566 equations are quadratic, and the other is linear:

- 567 • $\delta r_\alpha - 2c_{0\alpha\alpha}\alpha_0 - 2c_{0\alpha\beta}\beta_0 - \lambda(1-\delta)r_\alpha = 0,$
- 568 • $\frac{\delta c_{1\alpha\beta}r_\alpha}{c_{1\alpha\alpha}} - 2c_{0\beta\beta}\beta_0 - 2c_{0\alpha\beta}\alpha_0 - \lambda \left(\frac{c_{1\alpha\beta}}{c_{1\alpha\alpha}}(1-\delta)r_\alpha - 2\frac{c_{1\alpha\beta}^2}{c_{1\alpha\alpha}}\theta_D + 2 \left(\frac{c_{1\alpha\beta}}{c_{1\alpha\alpha}} - c_{1\beta\beta} \right) \beta_0 \right) = 0,$
- 569 • The quadratic stated in the proposition.

570 Though there may be multiple roots satisfying the above equations, the roots are bounded in typical
571 fashion by Bezout's Theorem. Further algebra for arriving at solutions is left to the computer.

572 Game Solving

573 Player's strategies without regulation

574 **The domain-specialist's strategy.** The proof for Proposition 0.1 is given below.

575 *Proof.* D 's best-response strategy is the value γ_1^* that maximizes D 's utility.

$$\gamma_1^*(\gamma_0, \delta) = \arg \max_{\gamma_1} U_D(\gamma_0, \gamma_1 \delta) \text{ s.t. } U_D \geq 0, \alpha_1 \geq \alpha_0, \beta_1 \geq \beta_0$$

Observe that D will not abstain because zero-investment ($\gamma_1 = \gamma_0$) is cost-free, yielding non-negative utility, so we can safely ignore the constraint. To solve the optimization, we specify the Lagrangian as follows for some multipliers $\lambda_1 \in \mathbb{R}$, $\lambda_2 \in \mathbb{R}$ and a slack variables $s_1 \in \mathbb{R}$, $s_2 \in \mathbb{R}$. By construction, we assert that the slack variables are only non-zero when the multipliers are zero, and the multipliers are non-zero only if the slack variables are zero.

$$\mathcal{L} := (1 - \delta)r^T \gamma_1 - (\gamma_1 - \gamma_0)^T C_1(\gamma_1 - \gamma_0) - \lambda_1(\alpha_1 - \alpha_0 - s_1^2) - \lambda_2(\beta_1 - \beta_0 - s_2^2).$$

576 We partially differentiate with respect to each decision variable and each multiplier.

$$\begin{aligned} \frac{\partial}{\partial \alpha_1} \mathcal{L} &= 0 \\ \iff (1 - \delta)r_\alpha - 2c_{1,\alpha\alpha}(\alpha_1 - \alpha_0) + 2c_{1,\alpha\beta}(\beta_1 - \beta_0) - \lambda_1 &= 0 \\ \frac{\partial}{\partial \beta_1} \mathcal{L} &= 0 \\ \iff (1 - \delta)r_\beta - 2c_{1,\beta\beta}(\beta_1 - \beta_0) + 2c_{1,\alpha\beta}(\alpha_1 - \alpha_0) - \lambda_2 &= 0 \\ \frac{\partial}{\partial \lambda_1} \mathcal{L} &= 0 \\ \iff -\alpha_1 + \alpha_0 + s_1^2 &= 0 \\ \frac{\partial}{\partial \lambda_2} \mathcal{L} &= 0 \\ \iff -\beta_1 + \beta_0 + s_2^2 &= 0 \end{aligned}$$

577 Using complementary slackness, we have four possible options:

578 1. $s_1 = 0, \lambda_1 > 0, s_2 = 0, \lambda_2 > 0 \rightarrow \beta_1^* = \beta_0, \alpha_1^* = \alpha_0$.

579 2. $s_1 \neq 0, \lambda_1 = 0, s_2 = 0, \lambda_2 > 0 \rightarrow \beta_1^* = \beta_0$, and we can plug into our first of four
580 equations above:

$$\begin{aligned} (1 - \delta)r_\alpha - 2c_{1,\alpha\alpha}(\alpha_1 - \alpha_0) + 2c_{1,\alpha\beta}(\beta_1 - \beta_0) - \lambda_1 &= 0 \\ \rightarrow (1 - \delta)r_\alpha - 2c_{1,\alpha\alpha}(\alpha_1 - \alpha_0) &= 0 \\ \rightarrow \alpha_1^* = \alpha_0 + \frac{(1 - \delta)r_\alpha}{2c_{1,\alpha\alpha}}. \end{aligned}$$

581 3. $s_1 = 0, \lambda_1 > 0, s_2 \neq 0, \lambda_2 = 0 \rightarrow \alpha_1^* = \alpha_0$, and we can plug in to equation 2:

$$\begin{aligned} (1 - \delta)r_\beta - 2c_{1,\beta\beta}(\beta_1 - \beta_0) + 2c_{1,\alpha\beta}(\alpha_1 - \alpha_0) - \lambda_2 &= 0 \\ \rightarrow (1 - \delta)r_\beta - 2c_{1,\beta\beta}(\beta_1 - \beta_0) - \lambda_2 &= 0 \\ \rightarrow \beta_1^* = \beta_0 + \frac{(1 - \delta)r_\beta}{2c_{1,\beta\beta}}. \end{aligned}$$

582 4. $s_1 \neq 0, \lambda_1 = 0, s_2 \neq 0, \lambda_2 = 0 \rightarrow$ This is the unconstrained critical point, and is solved
583 via the first two systems of equations:

$$\begin{aligned} \nabla U_D &= (1 - \delta)r - 2C_1(\gamma_1 - \gamma_0) = 0 \\ \rightarrow \gamma_1^* &= \gamma_0 + \frac{(1 - \delta)}{2} C_1^{-1} r. \end{aligned}$$

584 Thus we have established our four candidates in the proposition statement. \square

585 **The Generalist's strategy.** The proof for Proposition 0.2 is given below.

586 *Proof.* G 's best-response strategy is the value γ_0^* that maximizes G 's utility.

$$\gamma_0^*(\delta) = \arg \max_{\gamma_1} U_G(\gamma_0, \delta) \text{ s.t. } U_G \geq 0, \alpha_0 \geq 0, \beta_0 \geq 0.$$

Following the same steps as the proof of Proposition 0.1, we specify the Lagrangian as follows for multipliers $\lambda_1 \in \mathbb{R}$, $\lambda_2 \in \mathbb{R}$ and a slack variables $s_1 \in \mathbb{R}$, $s_2 \in \mathbb{R}$.

$$\mathcal{L} := \delta r^T \gamma_1 - \gamma_0^T C_0 \gamma_0 - \lambda_1(\alpha_0 - s_1^2) - \lambda_2(\alpha_0 - s_2^2).$$

587 We partially differentiate with respect to each decision variable and each multiplier.

$$\begin{aligned} \frac{\partial}{\partial \alpha_0} \mathcal{L} &= 0 \\ \iff \delta r_\alpha - 2c_{0,\alpha\alpha}\alpha_0 + 2c_{0,\alpha\beta}\beta_0 - \lambda_1 &= 0, \\ \frac{\partial}{\partial \beta_1} \mathcal{L} &= 0 \\ \iff \delta r_\beta - 2c_{0,\beta\beta}\beta_0 + 2c_{1,\alpha\beta}\alpha_0 - \lambda_2 &= 0, \\ \frac{\partial}{\partial \lambda_1} \mathcal{L} &= 0 \\ \iff -\alpha_0 + s_1^2 &= 0, \\ \frac{\partial}{\partial \lambda_2} \mathcal{L} &= 0 \\ \iff -\beta_0 + s_2^2 &= 0. \end{aligned}$$

588 Using complementary slackness, we have four possible options:

589 1. $s_1 = 0, \lambda_1 > 0, s_2 = 0, \lambda_2 > 0 \rightarrow \beta_0^* = 0, \alpha_0^* = 0.$

590 2. $s_1 \neq 0, \lambda_1 = 0, s_2 = 0, \lambda_2 > 0 \rightarrow \beta_0^* = 0$, and we can plug into our first of four equations
591 above:

$$\begin{aligned} \delta r_\alpha - 2c_{0,\alpha\alpha}\alpha_0 + 2c_{0,\alpha\beta}\beta_0 - \lambda_1 &= 0 \\ \rightarrow \delta r_\alpha - 2c_{0,\alpha\alpha}\alpha_0 &= 0 \\ \rightarrow \alpha_0^* &= \frac{\delta r_\alpha}{2c_{0,\alpha\alpha}}. \end{aligned}$$

592 3. $s_1 = 0, \lambda_1 > 0, s_2 \neq 0, \lambda_2 = 0 \rightarrow \alpha_0^* = 0$, and we can plug in to equation 2:

$$\begin{aligned} \delta r_\beta - 2c_{0,\beta\beta}\beta_0 + 2c_{0,\alpha\beta}\alpha_0 - \lambda_2 &= 0 \\ \delta r_\beta - 2c_{0,\beta\beta}\beta_0 &= 0 \\ \rightarrow \beta_0^* &= \frac{\delta r_\beta}{2c_{0,\beta\beta}}. \end{aligned}$$

593 4. $s_1 \neq 0, \lambda_1 = 0, s_2 \neq 0, \lambda_2 = 0 \rightarrow$ This is the unconstrained critical point, and is solved
594 via the first two systems of equations:

$$\begin{aligned} \nabla U_G = \delta r - 2C_0\gamma_0 &= 0 \\ \rightarrow \gamma_0^* &= \frac{\delta}{2}C_0^{-1}r. \end{aligned}$$

595 Thus we have established our four candidates. \square

596 **Condition for non-zero performance and safety investment**

597 **Condition establishing non-zero investment.** Below we prove Remark 0.3.

598 *Proof.* The first of the three inequalities establishes that the player's costs are strictly convex:

$$c_{\alpha\beta} < \sqrt{c_{p,\alpha\alpha}c_{p,\beta\beta}} \iff c_{p,\alpha\alpha}c_{p,\beta\beta} - c_{\alpha\beta}^2 > 0 \iff \det C_p > 0.$$

599 By the spectral theorem, we know a 2x2 matrix is positive definite if and only if its determinant and
600 trace are both positive, which is now established. By Lemma 0.8, the utility is strictly concave for
601 our setting if and only if the cost is strictly convex. Thus the unconstrained solution is the global
602 optimum as long as it is feasible. Thus, the necessary and sufficient condition for optimality is the
603 condition for feasibility.

- For the generalist:

$$\frac{\delta}{2}C_0^{-1}r > 0 \iff \frac{\delta}{2\det C_0} \begin{bmatrix} c_{0,\beta\beta}r_\alpha - c_{0,\alpha\beta}r_\beta \\ -c_{0,\alpha\beta}r_\alpha + c_{0,\alpha\alpha}r_\beta \end{bmatrix} > \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

604 Using the same positive definiteness identity above, we know the determinant is positive.
605 We are given $\delta > 0$. Thus we can cancel the positive constant term $\frac{\delta}{2\det C_0}$. The two
606 inequalities simplify to those stated in the proposition.

- For the specialist, the proof proceeds identically. Observe that $(1 - \delta) \geq 0$ and the unconstrained contribution is given by: $\frac{1-\delta}{2}C_1^{-1}r$.

609 □

610 **Proof for Player Strategies with Regulation**

611 Here we provide proofs for our propositions establishing best-response strategies for the players.

612 **Domain-specialist best-response under regulation.** Here we provide the proof of Proposition 0.6,
613 the domain specialist's best response under regulatory requirement θ_D .

614 *Proof.* D 's best-response strategy is the value γ_1^* that maximizes D 's utility. D will abstain if and
615 only if the best option yields negative utility.

$$\gamma_1^*(\gamma_0, \delta, \theta_D) = \arg \max_{\gamma_1} U_D(\gamma_0, \delta, \theta_D) \text{ s.t. } U_D \geq 0, \alpha_1 \geq \alpha_0, \beta_1 \geq \max(\beta_0, \theta_D).$$

Define $\kappa = \max(\beta_0, \theta_D)$. To solve the optimization, we specify the Lagrangian as follows for some multipliers $\lambda \in \mathbb{R}^3$ and a slack variables $s \in \mathbb{R}^3$.

$$\mathcal{L} := (1 - \delta)r^T \gamma_1 - (\gamma_1 - \gamma_0)^T C_1 (\gamma_1 - \gamma_0) - \lambda_1(\alpha_1 - \alpha_0 - s_1^2) - \lambda_2(\beta_1 - \kappa - s_2^2) - \lambda_3(U_D - s_3^2).$$

616 We partially differentiate with respect to each decision variable and each multiplier.

$$\begin{aligned}
& \frac{\partial}{\partial \alpha_1} \mathcal{L} = 0 \\
& \iff (1 - \delta)r_\alpha - 2c_{1,\alpha\alpha}(\alpha_1 - \alpha_0) + 2c_{1,\alpha\beta}(\beta_1 - \kappa) - \lambda_1 - \lambda_3 \frac{\partial U_D}{\partial \alpha_1} = 0 \\
& \iff (1 - \lambda_3)((1 - \delta)r_\alpha - 2c_{1,\alpha\alpha}(\alpha_1 - \alpha_0) + 2c_{1,\alpha\beta}(\beta_1 - \kappa)) - \lambda_1 = 0 \\
& \frac{\partial}{\partial \beta_1} \mathcal{L} = 0 \\
& \iff (1 - \delta)r_\beta - 2c_{1,\beta\beta}(\beta_1 - \kappa) + 2c_{1,\alpha\beta}(\alpha_1 - \alpha_0) - \lambda_2 - \lambda_3 \frac{\partial U_D}{\partial \beta_1} = 0 \\
& \iff (1 - \lambda_3)((1 - \delta)r_\beta - 2c_{1,\beta\beta}(\beta_1 - \kappa) + 2c_{1,\alpha\beta}(\alpha_1 - \alpha_0)) - \lambda_2 = 0 \\
& \frac{\partial}{\partial \lambda_1} \mathcal{L} = 0 \\
& \iff -\alpha_1 + \alpha_0 + s_1^2 = 0 \\
& \frac{\partial}{\partial \lambda_2} \mathcal{L} = 0 \\
& \iff -\beta_1 + \kappa + s_2^2 = 0 \\
& \frac{\partial}{\partial \lambda_3} \mathcal{L} = 0 \\
& \iff -U_D + s_2^2 = 0 \\
& \iff -(1 - \delta)r^T \gamma_1 + (\gamma_1 - \gamma_0)^T C_1 (\gamma_1 - \gamma_0) + s_2^2 = 0
\end{aligned}$$

617 Using complementary slackness, we have eight possible options:

- 618 1. $s_1 = 0, \lambda_1 > 0, s_2 = 0, \lambda_2 > 0, s_3 \neq 0, \lambda_3 = 0 \rightarrow \beta_1^* = \kappa, \alpha_1^* = \alpha_0$.
- 619 2. $s_1 = 0, \lambda_1 > 0, s_2 = 0, \lambda_2 > 0, s_3 = 0, \lambda_3 > 0 \rightarrow \beta_1^* = \kappa, \alpha_1^* = \alpha_0$. This offers the same
- 620 candidate as (1).
3. $s_1 = 0, \lambda_1 > 0, s_2 \neq 0, \lambda_2 = 0, s_3 \neq 0, \lambda_3 = 0 \rightarrow \alpha_1^* = \alpha_0$, solve equations (1) and (2) for β_1^* and λ_1 . Omitting the algebra, this yields:

$$\gamma_1^* = \left[\beta_0 + \frac{\alpha_0 (1-\delta)r_\beta}{2c_{1\beta\beta}} \right]$$

- 621 4. $s_1 = 0, \lambda_1 > 0, s_2 \neq 0, \lambda_2 = 0, s_3 = 0, \lambda_3 > 0 \rightarrow \alpha_1^* = \alpha_0$, solve equations (1) and (2)
- 622 for β_1^* and λ_1 . This solution, if it is distinct from the previous solution (3), will always be
- 623 dominated because it is characterized by 0 utility for G .
5. $s_1 \neq 0, \lambda_1 = 0, s_2 = 0, \lambda_2 > 0, s_3 = 0, \lambda_3 > 0 \rightarrow \beta_1^* = \kappa \rightarrow$ solve equations (1) and (2) for λ_1 and α_1^* . Omitting algebra, this yields:

$$\gamma_1^* = \left[\alpha_0 + \frac{(1-\delta)r_\alpha}{2c_{1\alpha\alpha}} - \frac{c_{1\alpha\beta}}{c_{1\alpha\alpha}} \max(0, \theta_D - \beta_0) \right]$$

- 624 6. $s_1 \neq 0, \lambda_1 = 0, s_2 = 0, \lambda_2 > 0, s_3 \neq 0, \lambda_3 = 0 \rightarrow$ this solution, if it is distinct from the
- 625 previous one (5), will always be dominated because it is characterized by 0 utility for G .
- 626 7. $s_1 \neq 0, \lambda_1 = 0, s_2 \neq 0, \lambda_2 = 0, s_3 \neq 0, \lambda_3 = 0 \rightarrow \alpha_1^* = \alpha_0$, solve equations (1)
- 627 and (2) for α_1^*, β_1^* . This is the unconstrained solution. Omitting algebra, this yields: $\gamma_1^* =$
- 628 $\gamma_0 + \frac{(1-\delta)}{2} C_1^{-1} r$.
- 629 8. $s_1 \neq 0, \lambda_1 = 0, s_2 \neq 0, \lambda_2 = 0, s_3 = 0, \lambda_3 > 0 \rightarrow \beta_1^* = \kappa \rightarrow$ solve equations (1) and (2)
- 630 for α_1^*, β_1^* . This solution, if it is distinct from (7), will always be dominated by (7) because
- 631 it is characterized by 0 utility for G .

Thus we have established our four candidates in the proposition statement. To handle the `abstain` scenario, we check each candidate produced in the process above by plugging the strategy to our formula for U_D . If none yield positive utility, then the domain specialist prefers to `abstain`. \square

The generalist's subgame perfect equilibrium strategy under regulation. Here we prove Proposition 0.7.

Proof. G 's best-response strategy is the value γ_0^* that maximizes G 's utility.

$$\gamma_0^*(\delta, \theta_G, \theta_D) = \arg \max_{\gamma_0} U_G(\gamma_0; \delta, \theta_G, \theta_D) \text{ s.t. } U_G \geq 0, U_D \geq 0, \alpha_0 \geq 0, \beta_0 \geq \theta_G.$$

To solve the optimization, we specify the Lagrangian as follows for some multipliers $\lambda \in \mathbb{R}^4$ and a slack variables $s \in \mathbb{R}^4$.

$$\mathcal{L} := \delta r^T \gamma_1 - \gamma_0^T C_0 \gamma_0 - \lambda_1(\alpha_0 - s_1^2) - \lambda_2(\beta_1 - \theta_G - s_2^2) - \lambda_3(U_D - s_3^2) - \lambda_4(U_G - s_4^2).$$

We partially differentiate with respect to each decision variable and each multiplier.

$$\begin{aligned} \frac{\partial}{\partial \alpha_1} \mathcal{L} &= 0 \\ \iff \delta r_\alpha - 2c_{0,\alpha\alpha}\alpha_0 + 2c_{0,\alpha\beta}\beta_0 - \lambda_1 - \lambda_3 \frac{\partial U_G}{\partial \alpha_0} - \lambda_4 \frac{\partial U_D}{\partial \alpha_0} &= 0 \\ \iff (1 - \lambda_3)(\delta r_\alpha - 2c_{0,\alpha\alpha}\alpha_0 + 2c_{0,\alpha\beta}\beta_0) - \lambda_1 \\ - \lambda_4((1 - \delta)r_\alpha - 2c_{1,\alpha\alpha}(\alpha_1 - \alpha_0) + 2c_{1,\alpha\beta}(\beta_1 - \max(\beta_0, \theta_D))) &= 0, \\ \frac{\partial}{\partial \beta_1} \mathcal{L} &= 0 \\ \iff \delta r_\beta - 2c_{0,\beta\beta}\beta_0 + 2c_{0,\alpha\beta}\alpha_0 - \lambda_2 - \lambda_3 \frac{\partial U_G}{\partial \beta_0} - \lambda_4 \frac{\partial U_D}{\partial \beta_0} &= 0 \\ \iff (1 - \lambda_3)(\delta r_\beta - 2c_{0,\beta\beta}\beta_0 + 2c_{0,\alpha\beta}\alpha_0) - \lambda_2 \\ - \lambda_4((1 - \delta)r_\beta - 2c_{1,\beta\beta}(\beta_1 - \max(\beta_0, \theta_D)) + 2c_{1,\alpha\beta}(\alpha_1 - \alpha_0)) &= 0, \\ \frac{\partial}{\partial \lambda_1} \mathcal{L} &= 0 \\ \iff -\alpha_0 + s_1^2 &= 0, \\ \frac{\partial}{\partial \lambda_2} \mathcal{L} &= 0 \\ \iff -\beta_0 + s_2^2 &= 0, \\ \frac{\partial}{\partial \lambda_3} \mathcal{L} &= 0 \\ \iff -U_G + s_3^2 &= 0 \\ \iff -\delta r^T \gamma_1 + \gamma_0^T C_1 \gamma_0 + s_2^2 &= 0, \\ \frac{\partial}{\partial \lambda_4} \mathcal{L} &= 0 \\ \iff -U_D + s_4^2 &= 0 \\ \iff -(1 - \delta)r^T \gamma_1 + (\gamma_1 - \gamma_0)^T C_1 (\gamma_1 - \gamma_0) + s_2^2 &= 0. \end{aligned}$$

Using complementary slackness, we have sixteen possible options. For brevity, we refer to these options by the constraints they satisfy, where **bold** corresponds to the constraints being activated. The algebra is omitted for exposition; only the candidates yielded are noted for each constraint setting.

1. $\alpha_0, \beta_0, U_G, U_D \rightarrow [0, \theta_G]$.

2. $\alpha_0, \beta_0, U_G, U_D \rightarrow [0, \theta_G]$

- 647 3. $\alpha_0, \beta_0, U_G, U_D \rightarrow [0, \theta_G]$
- 648 4. $\alpha_0, \beta_0, U_G, U_D \rightarrow [0, \theta_G]$
- 649 5. $\alpha_0, \beta_0, U_G, U_D \rightarrow \left[\begin{array}{c} 0 \\ \frac{\delta r_\beta}{2c_{0\beta\beta}} \end{array} \right]$
- 650 6. $\alpha_0, \beta_0, U_G, U_D \rightarrow \left[\begin{array}{c} 0 \\ \frac{\delta r_\beta}{2c_{0\beta\beta}} \end{array} \right]$
- 651 7. $\alpha_0, \beta_0, U_G, U_D \rightarrow \left[\begin{array}{c} 0 \\ \frac{\delta r_\beta}{2c_{0\beta\beta}} \end{array} \right]$
- 652 8. $\alpha_0, \beta_0, U_G, U_D \rightarrow$ One of three along $U_D = 0$ curve.
- 653 9. $\alpha_0, \beta_0, U_G, U_D \rightarrow \gamma_0^* = \left[\begin{array}{c} \frac{\delta r_\alpha}{2c_{0\alpha\alpha}} - \frac{c_{0\alpha\beta}}{c_{0\alpha\alpha}} \theta_G \\ \theta_G \end{array} \right]$.
- 654 10. $\alpha_0, \beta_0, U_G, U_D \rightarrow \left[\begin{array}{c} \frac{\delta r_\alpha}{2c_{0\alpha\alpha}} - \frac{c_{0\alpha\beta}}{c_{0\alpha\alpha}} \theta_G \\ \theta_G \end{array} \right]$
- 655 11. $\alpha_0, \beta_0, U_G, U_D \rightarrow \left[\begin{array}{c} \frac{\delta r_\alpha}{2c_{0\alpha\alpha}} - \frac{c_{0\alpha\beta}}{c_{0\alpha\alpha}} \theta_G \\ \theta_G \end{array} \right]$
- 656 12. $\alpha_0, \beta_0, U_G, U_D \rightarrow$ Two of three along the $U_D = 0$ curve.
- 657 13. $\alpha_0, \beta_0, U_G, U_D \rightarrow \frac{\delta}{2} C_0^{-1} r$.
- 658 14. $\alpha_0, \beta_0, U_G, U_D \rightarrow \frac{\delta}{2} C_0^{-1} r$
- 659 15. $\alpha_0, \beta_0, U_G, U_D \rightarrow \frac{\delta}{2} C_0^{-1} r$
- 660 16. $\alpha_0, \beta_0, U_G, U_D \rightarrow$ Three of three along the $U_D = 0$ curve.

661 Thus we have established our four candidates in the proposition statement. To handle the `abstain`
 662 scenario, we check each candidate produced in the process above by plugging the strategy to our
 663 formula for U_G . If none yield positive utility, then the generalist prefers to `abstain`. \square

664 Helper Lemmas and Analysis

665 Here we write out helper Lemmas and analysis for our proofs concerning backfiring and mutualism.

666 **Lemma 0.8.** *In the AI regulation game with quadratic costs, any player's utility is strictly concave*
 667 *if and only if their cost matrix is positive definite.*

Proof. The generalist utility function is given by $U_G = \delta r^T \gamma_1 - \gamma_0^T C_0 \gamma_0$. Observe this is twice differentiable everywhere. Thus the function is strictly concave in α_0, β_0 if and only if its Hessian derivative is negative definite. We compute the Hessian as follows:

$$H := \left[\begin{array}{cc} \frac{\partial^2 U_G}{\partial \alpha_0^2} & \frac{\partial^2 U_G}{\partial \alpha_0 \partial \beta_0} \\ \frac{\partial^2 U_G}{\partial \beta_0 \partial \alpha_0} & \frac{\partial^2 U_G}{\partial \beta_0^2} \end{array} \right] = -2C_0.$$

668 This matrix is negative definite if and only if C_0 is positive definite.

669 The proof for the domain specialist follows the same steps. \square

670 **Lemma 0.9.** *In any AI regulation game with separable quadratic costs, if there is no regulation,*
 671 *both players will invest a non-zero amount in each attribute.*

672 *Proof.* By Lemma 0.8, we are given that the utilities are strictly concave. Thus, the proof consists of
 673 showing that 1) the utility function is greater than or equal to 0 at the origin point of zero investment
 674 and 2) the gradient points towards the interior of the feasible set everywhere along the boundaries.

675 Here we prove the two conditions for U_G :

- 676 1. $U_G(\alpha_0 = 0, \beta_0 = 0) = \delta r^T \vec{0} - 0 = 0$
- 677 2. We prove this for each constraint, $\alpha_0 \geq 0, \beta_0 \geq 0$:
 - 678 • $\frac{\partial U_G}{\partial \alpha_0} \Big|_{\beta_0=0} = \delta r_\alpha - 2c_{0,\alpha\alpha}\alpha_0 = \delta r_\alpha - 0 > 0.$
 - 679 • $\frac{\partial U_G}{\partial \beta_0} \Big|_{\alpha_0=0} = \delta r_\beta - 2c_{0,\beta\beta}\beta_0 = \delta r_\beta - 0 > 0.$

680 Here we prove the two conditions for U_D :

- 681 1. $U_D(\alpha_i = 0, \beta_i = 0) = (1 - \delta)r^T \gamma_0 - 0 \geq 0$
- 682 2. We prove this for each constraint, $\alpha_i \geq \alpha_0, \beta_i \geq \beta_0$:
 - 683 • $\frac{\partial U_D}{\partial \alpha_i} \Big|_{\beta_i=\beta_0} = (1 - \delta)r_\alpha - 2c_{i,\alpha\alpha}(\alpha_i - \alpha_0) = (1 - \delta)r_\alpha > 0.$
 - 684 • $\frac{\partial U_D}{\partial \beta_i} \Big|_{\alpha_i=\alpha_0} = (1 - \delta)r_\beta - 2c_{i,\beta\beta}(\beta_i - \beta_0) = \delta r_\beta - 0 > 0.$

685 □

686 Proving the Backfiring Result

687 Below we prove the Theorem 0.4.

688 *Proof.* Assume $\theta_G = 0$ for the entire proof. By Remark 0.3, we're given that the players commit to
 689 their unconstrained strategy in equilibrium. These were solved in Propositions 0.6 and 0.7. Thus we
 690 have the following player's strategies under no regulation for this setting:

$$\gamma_0^A = \frac{\delta}{2} C_0^{-1} r, \quad \gamma_1^A = \frac{1 - \delta}{2} C_1^{-1} r. \quad (1)$$

691 Our strategy is to show that G's unconstrained, no-regulation optimum becomes dominated in the
 692 presence of regulation targeting D, which we choose to be arbitrarily close to β_1^A .

693 **Notation.** Before we proceed, we introduce some additional notation. Define the set S to be all
 694 feasible pairs of strategies (γ_0, γ_1) . 'Feasible' here means those strategies which leave both G and
 695 D with non-negative utility. We use the subscript S_{θ_D} to track the particular regulatory threshold.
 696 The feasible pairs of strategies in the unregulated game is given by S_0 , and the feasible pairs of
 697 strategies in a game with threshold $\theta_D = 1.5$ is denoted $S_{1.5}$. We may refer to the unregulated game
 698 with the superscript A (for anarchy), e.g. β_1^A refers to the unregulated safety level. Observe that any
 699 set of tuples S_θ can be separated into two mutually exclusive and collectively exhaustive sets:

- 700 • S_θ^{MC} (for minimally compliant) is the set of all tuples where D 's best response has safety
 701 $\beta_1^* = \theta_D$.
- 702 • S_θ^C (for contribute) is the set of all tuples where D 's best response has safety $\beta_1^* > \theta_D$.

703 Now, we provide a sequence of lemmas, with the purpose of establishing the intuition that *all we*
 704 *must do is find some $\epsilon > 0$ and some strategy $\beta_0^R \neq \beta_0^A$ such that G prefers β_0^R to β_0^A and D*
 705 *minimally complies.*

706 **Lemma 0.10.** For any threshold $\theta_D > 0$, $S_\theta^C \subset S_0$.

707 *Proof.* $S_0 = S_0^{\text{MC}} \cup S_0^C = S_0^{\text{MC}} \cup (\bigcup_{t=0}^{\infty} S_t^C) \supset S_\theta^C$. □

708 **Lemma 0.11.** If $\theta_D \geq \beta_1^A$, backfiring is impossible.

709 *Proof.* Assume for contradiction that $\theta_D^* \geq \beta_1^A$ and backfiring occurs. Backfiring would imply
 710 $\beta_1(\theta_D = \theta_D^*) < \beta_1(\theta_D = 0) = \beta_1^A$. However, this would violate the regulation, which we're given
 711 is greater than β_1^A . Hence we've already established the contradiction. \square

712 **Lemma 0.12.** *Given a threshold θ , backfiring can occur only if the strategies $(\gamma_0, \gamma_1) \in S_{\theta_D}^{MC}$.*

713 *Proof.* We have established $S_{\theta_D} = S_{\theta_D}^{MC} \cup S_{\theta_D}^C$, so the proof will show that the strategies in $S_{\theta_D}^C$ can
 714 never exhibit backfiring. This would imply, if backfiring occurs over the feasible set of strategies S_{θ} ,
 715 it is only possible for strategies in $S_{\theta_D}^{MC}$. The proof proceeds, first for all values $\theta_D \geq \beta_1^A$, and then
 716 for all values $\theta < \beta_1^A$.

- 717 • For $\theta_D \geq \beta_1^A$, backfiring is impossible generally, as established in Lemma 0.11.
- 718 • For $\theta_D < \beta_1^A$, start by observing that the anarchy solution (γ_0^A, γ_1^A) is always feasible. This
 719 solution is the strategy tuple that maximizes the utility of G over S_0 . Lemma 0.10 tells
 720 us that this set, S_0 , contains all sets of regulated strategies where the players contribute:
 721 $S_{\theta_D}^{\text{contribute}} \subset S_0$. Thus: $(\gamma_0^A, \gamma_1^A) := \sup_{U_G} S_0 \succeq_G S_0 \supset S_{\theta_D}^{\text{contribute}} \rightarrow (\gamma_0^A, \gamma_1^A) \succeq_G S_{\theta_D}^C$.
 722 Thus the anarchy solution is feasible and dominates all strategies in $S_{\theta_D}^{\text{contribute}}$.

723 This completes the proof, and demonstrates that if backfiring is ever to occur, it will exhibit strategies
 724 that are *minimally compliant* with the regulation. \square

725 Backfiring is a regulation yielding lower safety than β_0^A . The claims above state that backfiring
 726 cannot occur if $\theta_D > \beta_1^A$ and can only occur if the domain specialist minimally complies. As an
 727 immediate corollary, we can claim that backfiring occurs *if and only if* there is a regulation $\theta_D < \beta_1^A$
 728 such that the strategies $(\gamma_0(\theta_D), \gamma_1(\theta_D)) \in S_{\theta_D}^{MC}$.

729 **Lemma 0.13.** *For a given regulation $\theta_D < \beta_1^A$ in our setting, if G prefers any minimally compliant
 730 strategy γ'_0 to γ_0^A , then G 's optimal strategy $\gamma_0^* \in S_{\theta_D}^{MC}$ and the regulation backfires.*

731 *Proof.* We're given γ_0^A is optimal over S_0 . By Lemma 0.10, $S_{\theta_D}^C \subset S_0$. Since $\theta_D < \beta_1^A$, γ_0^A
 732 remains feasible. The only new strategies available to G are those in $S_{\theta_D}^{MC}$. Thus, if we denote
 733 utility-domination using \succ , we have $\gamma'_0 \succ \gamma_0^A \succeq g \forall g \in S_{\theta_D}^C$. This implies G 's optimal strategy γ_0^*
 734 is either γ'_0 or otherwise belongs to $S_{\theta_D}^{MC}$. \square

735 Thus our task is to find some regulation θ_D and some strategy γ'_0 such that $U_G(\gamma'_0) > U_G(\gamma_0^A)$.

Lemma 0.14. *For small $\epsilon > 0$, if the given conditions are met, the following G strategy dominates
 no regulation:*

$$\gamma'_0 = \left[\frac{\delta r_{\alpha}}{2c_{0,\alpha\alpha}} (\beta_0^A - 2\epsilon) \right]$$

Proof. Equation 1 give us G and D 's strategies under no regulation. Given G 's candidate strategy
 stated in the Lemma, we compute D 's best response. Observe this must be a minimally-compliant
 best response, because G 's strategy was constructed to be a difference $\beta_0^A + \epsilon$ from D 's regulatory
 floor. Thus, by Proposition 0.6, we have:

$$\gamma'_1 = \left[\alpha'_0 + \frac{(1-\delta)}{2c_{1\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1\alpha\alpha}} \theta_D \right]$$

736 We compare G 's utility in the two scenarios:

- 737 1. $(\gamma'_0, \gamma'_1) \rightarrow U'_G = \delta (r_{\alpha}\alpha'_1 + r_{\beta}\beta'_1) - c_{0,\alpha\alpha}(\alpha'_0)^2 - 2c_{0,\alpha\beta}\alpha'_0\beta'_0 - c_{0,\beta\beta}(\beta'_0)^2$
- 738 2. $(\gamma_1^A, \gamma_1^A) \rightarrow U_G^A = \delta (r_{\alpha}\alpha_1^A + r_{\beta}\beta_1^A) - c_{0,\alpha\alpha}(\alpha_0^A)^2 - 2c_{0,\alpha\beta}\alpha_0^A\beta_0^A - c_{0,\beta\beta}(\beta_0^A)^2$

We compute the difference $\Delta U_G = U'_G - U_G^A$. We expand both terms and take the limit as $\epsilon \searrow 0$ to get the following:

$$\lim_{\epsilon \searrow 0} \Delta U_G = \frac{\delta(1-\delta)r_\beta}{2c_{1,\beta\beta}} \left(r_\beta - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} r_\alpha \right)$$

A sufficient condition for this quantity being positive is stated below. The reason is all terms outside the parentheses are given as positive.

$$r_\beta > \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} r_\alpha.$$

739 Notice the above condition is given as it is one of the conditions in remark 0.3.⁴ □

740 Thus, we have shown that for small positive ϵ , the generalist prefers the backfiring strategy to the
741 unconstrained optimum γ_0^A . By Lemma 0.13, the optimal regulated strategy is an element in $S_{\theta_D}^{MC}$
742 and the regulation backfires. □

743 Proof of the mutualism result

744 Here we prove Theorem 0.5, that for a swath of games there exists a set of regulations that mutually
745 improve the player's utilities.

Proof. Observe that we only have to provide a single instance of regulation that does better than the unregulated optimal γ_0^A, γ_1^A to show that there exists a Pareto improvement effect of regulation. We consider the following minimal-compliance strategies (using Proposition 0.6 and 0.7):

$$\gamma'_0 = \begin{bmatrix} \frac{\delta r_\alpha}{2c_{0,\alpha\alpha}} - \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \theta_G \\ \theta_G \end{bmatrix}, \gamma'_1 = \begin{bmatrix} \alpha_0 + \frac{(1-\delta)r_\alpha}{2c_{1,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \theta_D \\ \theta_D \end{bmatrix}.$$

746 Observe these are feasible because they are compliant and, for small ϵ , the performance investment
747 is positive.

748 **Lemma 0.15.** *For the specified conditions, $U_G(\gamma'_0, \gamma'_1) > U_G(\gamma_0^A, \gamma_1^A)$*

749 *Proof.* Start by computing the change in the generalist's performance and safety investments be-
750 tween these strategies. The change in safety investment is simply $\Delta\beta_0 = \theta_G - \beta_0^A = \epsilon$. The change
751 in performance investment is given by:

$$\begin{aligned} \Delta\alpha_0 &= \frac{\delta r_\alpha}{2c_{0,\alpha\alpha}} - \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \left(\frac{\delta}{2 \det C_0} (-c_{0,\alpha\beta} r_\alpha + c_{0,\alpha\alpha} r_\beta + \epsilon) \right) - \frac{\delta}{2 \det C_0} (c_{0,\beta\beta} r_\alpha - c_{0,\alpha\beta} r_\beta) \\ &= \frac{\delta r_\alpha}{2c_{0,\alpha\alpha}} + \frac{c_{0,\alpha\beta}^2}{c_{0,\alpha\alpha}^2} \frac{\delta r_\alpha}{2 \det C_0} - \cancel{c_{0,\alpha\beta} \frac{\delta}{2 \det C_0} r_\beta} + \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \epsilon - \frac{\delta}{2 \det C_0} c_{0,\beta\beta} r_\alpha + \cancel{c_{0,\alpha\beta} \frac{\delta}{2 \det C_0} r_\beta} \\ &= \left(\frac{\delta}{2c_{0,\alpha\alpha}} + \frac{c_{0,\alpha\beta}^2 \delta}{c_{0,\alpha\alpha}^2 2 \det C_0} - \frac{\delta c_{0,\beta\beta}}{2 \det C_0} \right) r_\alpha + \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \epsilon \\ &= \frac{\delta r_\alpha}{2} \left(\frac{1}{c_{0,\alpha\alpha}} + \frac{c_{0,\alpha\beta}^2}{c_{0,\alpha\alpha} (c_{0,\alpha\alpha} c_{0,\beta\beta} - c_{0,\alpha\beta}^2)} - \frac{c_{0,\beta\beta}}{c_{0,\alpha\alpha}} \right) + \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \epsilon \\ &= \frac{\delta r_\alpha}{2} \left(\frac{\det C_0 + c_{0,\alpha\beta}^2 - c_{0,\alpha\alpha} c_{0,\beta\beta}}{c_{0,\alpha\alpha} \det C_0} \right) + \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \epsilon \\ &= \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \epsilon. \end{aligned}$$

752 By the same logic, we solve for the change in the players' strategies. First, $\Delta\beta_1 = \beta'_1 - \beta_1^A =$
753 $\beta_1^A + 2\epsilon - \beta_1^A = 2\epsilon$. The change in performance is given by:

⁴This is also the condition for having a non-zero safety investment when costs are convex, and intuitively, backfiring is impossible when safety investment is zero.

$$\begin{aligned}
\Delta\alpha_1 &= \alpha'_1 - \alpha_1^A \\
&= \left[\alpha'_0 + \frac{(1-\delta)r_\alpha}{2c_{1,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} (\theta_D - \beta'_0) \right] - \left[\alpha_0^A + \frac{(1-\delta)}{2\det C_1} (c_{1,\beta\beta}r_\alpha - c_{1,\alpha\beta}r_\beta) \right] \\
&= \Delta\alpha_0 + \frac{(1-\delta)r_\alpha}{2c_{1,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \left(\frac{(1-\delta)}{2\det C_1} (-c_{1,\alpha\beta}r_\alpha + c_{1,\alpha\alpha}r_\beta) + \epsilon \right) - \frac{(1-\delta)}{2\det C_1} (c_{1,\beta\beta}r_\alpha - c_{1,\alpha\beta}r_\beta) \\
&= \frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}}\epsilon + \frac{(1-\delta)r_\alpha}{2c_{1,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \left(\frac{(1-\delta)}{2\det C_1} (-c_{1,\alpha\beta}r_\alpha + c_{1,\alpha\alpha}r_\beta) + \epsilon \right) - \frac{(1-\delta)}{2\det C_1} (c_{1,\beta\beta}r_\alpha - c_{1,\alpha\beta}r_\beta) \\
&= \epsilon \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \right).
\end{aligned}$$

754 The change in G 's cost is given by:

$$\begin{aligned}
\Delta(\text{G's cost}) &= \begin{bmatrix} \Delta\alpha_0 \\ \Delta\beta_0 \end{bmatrix}^T C_0 \begin{bmatrix} \Delta\alpha_0 \\ \Delta\beta_0 \end{bmatrix} \\
&= c_{0,\alpha\alpha} \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \right)^2 \epsilon^2 + 2 \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} \epsilon \right) \epsilon + c_{0,\beta\beta} \epsilon^2
\end{aligned}$$

755 Notice these are all ϵ^2 terms, meaning as ϵ is brought to very small positive values, they approach
756 zero at an exponential rate. The contribution to G 's revenue is given by:

$$\begin{aligned}
\Delta(\text{G's revenue}) &= \delta(r_\alpha \Delta\alpha_1 + r_\beta \Delta\beta_1) \\
&= \delta \left(r_\alpha \epsilon \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \right) + r_\beta 2\epsilon \right)
\end{aligned}$$

Notice these are terms of ϵ , whereas the cost effects are solely terms of ϵ^2 . Therefore, for sufficiently small ϵ , we say:

$$\lim_{\epsilon \searrow 0} \Delta U_G = \delta \left(r_\alpha \epsilon \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} + 2r_\beta \epsilon \right) \right)$$

757 Using the given conditions, we know:

$$\begin{aligned}
\lim_{\epsilon \searrow 0} \Delta U_G > 0 &\iff \delta \left(r_\alpha \epsilon \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \right) + 2r_\beta \epsilon \right) > 0 \\
&\iff r_\alpha \left(\frac{c_{0,\alpha\beta}}{c_{0,\alpha\alpha}} - \frac{c_{1,\alpha\beta}}{c_{1,\alpha\alpha}} \right) + 2r_\beta > 0 \\
&\iff \frac{r_\alpha c_{0,\alpha\beta}}{r_\beta c_{0,\alpha\alpha}} - \frac{r_\alpha c_{1,\alpha\beta}}{r_\beta c_{1,\alpha\alpha}} > -2.
\end{aligned}$$

758 Our conditions strictly bound the absolute value of both terms on the left hand side below 1, so this
759 completes the Lemma's proof. \square

760 **Lemma 0.16.** For the specified conditions, $U_D(\gamma'_0, \gamma'_1) > U_D(\gamma_0^A, \gamma_1^A)$

761 The limiting effect on D 's revenue is calculated exactly the same way as above, except that the
762 revenue expression is multiplied by $(1 - \delta)$ instead of δ .

763 This completes the proof, as both players are better off under the regulation. \square