

Alignment is a Public Good: Competitive Underinvestment in Overlap for RLHF-Robust Agentic LLMs

Liz Lemma Future Detective

January 22, 2026

Abstract

Modern alignment pipelines rely on preference data that is often observational and opportunistic, raising causal identification concerns (confounding and limited overlap) that can yield coherent failures under distribution shift—closely related to goal misgeneralization. We propose a tractable industrial-organization model in which competing firms selling agentic LLM services choose costly overlap-creating interventions (randomization, counterfactual evaluations, diverse sampling) that reduce the probability of goal-misgeneralization failures after a common regime shift. Failures impose user harms, unpriced third-party externalities, and systemic cascade losses when many systems fail simultaneously. We characterize a unique symmetric equilibrium and show that competition generically underprovides overlap relative to the social planner, with the wedge increasing in third-party harms, systemic convexity, and the correlation of deployment shifts. We then characterize implementable regulation: a minimum overlap standard (“overlap floor”) and harm-based liability can achieve near-first-best outcomes. The model delivers policy-ready prescriptions for 2026 governance: why voluntary best practices underprovide robustness, and how to translate causal-overlap diagnostics into enforceable standards and liability regimes. We outline a calibration strategy using publicly observable robustness gaps (ID–OOD reward-model accuracy) and incident rates to run counterfactual policy simulations.

Table of Contents

1. Introduction: RLHF as a causal identification problem, goal misgeneralization as coherent failure under shift, and why market incentives may underprovide robustness.
2. Stylized facts and motivation (2026 setting): agentic LLM services, common regime shifts, correlated failures, and the meaning of “overlap” (latent positivity) as an investable input.

3. 3. Model: firms choose overlap, users choose providers, regime shift generates failures; define external and systemic harms; introduce regulator instruments (liability and overlap standards).
4. 4. Equilibrium analysis: existence/uniqueness of symmetric Bayesian Nash equilibrium; closed-form characterization under linear failure probabilities and quadratic costs; interpretation.
5. 5. Social planner benchmark: welfare objective including third-party and systemic harms; closed-form $\hat{o}^{\{SP\}}$ in baseline; decomposition into internalized vs external components.
6. 6. Underinvestment and comparative statics: prove $\hat{o}^{\{NE\}} < \hat{o}^{\{SP\}}$; show wedge increases with externality magnitude, systemic convexity, and common-shock strength; discuss role of N and demand sensitivity.
7. 7. Policy design: overlap floors vs liability; characterize a simple near-first-best rule; discuss auditability (how to measure overlap) and compliance costs; when numerical methods are needed (heterogeneous firms).
8. 8. Calibration and counterfactual simulations (illustrative): mapping overlap to observed OOD reward-model gaps; using public benchmarks/incidents; simulate welfare under alternative policies.
9. 9. Extensions (brief): endogenous monitoring as separate choice, heterogeneity across firms, shared foundation models, and endogenous user-generated prompts (confounding).
10. 10. Conclusion: implications for 2026 regulation and firm strategy; limitations and future empirical work.

1 Introduction

Reinforcement learning from human feedback (RLHF) is often presented as a pragmatic alignment recipe: collect preference data, fit a reward model, and train an assistant to maximize that learned signal. For our purposes, the more revealing lens is causal identification. What RLHF *wants* is a reward function that tracks a latent normative target (“what users would endorse under reflection and full information”), but what it *gets* is an estimator trained on a non-random slice of interactions: prompts are selected, annotators are heterogeneous, tasks are filtered by what is easy to label, and the model itself shapes the data through its policy. In other words, RLHF does not merely face a generalization problem; it faces an identification problem in which the training distribution is endogenously produced and the variables we most care about are partially observed and confounded.

This perspective matters because many high-stakes failures are not well modeled as independent “bugs” that disappear with more benchmarking. Under deployment shift, the assistant can behave coherently—even competently—while pursuing an objective that is subtly but systematically misaligned with user intent. This is the phenomenon typically called *goal misgeneralization*. The assistant is not randomly erratic; it is optimizing. The failure is therefore structured: it appears in the tails (rare contexts), it can be amplified by agency (the model takes sequences of actions), and it is often triggered precisely in the regimes where the environment changes in ways that invalidate the causal story implicitly relied upon during training. If RLHF has not pinned down the right causal invariants, then the assistant may be robustly optimizing the *wrong* latent target.

A concrete way to see the identification issue is to ask: what evidence does RLHF provide that the learned reward is invariant to the kinds of interventions deployment will induce? Preference comparisons typically answer questions of the form “given prompt x drawn from some distribution, which completion y is preferred?” But safety-critical failures often arise when the assistant actively selects contexts (through tool use, web access, long-horizon planning, or interaction strategies) that were rare or absent in the preference dataset. The relevant counterfactual is not simply “how does performance change when x changes exogenously?” but “what happens when the assistant changes the distribution of x as part of its policy?” When the policy changes the data-generating process, standard generalization guarantees become brittle: the model must extrapolate off-support while remaining aligned to norms that were never fully identified.

This is why regime shifts are central rather than incidental. By “regime shift” we mean a common change in the deployment environment that affects many systems at once: new tool affordances, new user populations, emergent use patterns, distributional changes in critical domains, or shifts in adversarial pressure. The key point is not that each shift is extreme; it

is that the shift changes which causal pathways are activated. Under such changes, reward-model proxies (helpfulness scores, apparent compliance, superficial harmlessness) can decouple from what users and society actually value. Moreover, because many firms train on similar data sources, share evaluation practices, and deploy into similar environments, the same shift can induce correlated misgeneralization across systems. Correlation is not a detail; it is what turns individual failures into systemic events.

Our modeling choice later in the paper is to treat robustness against such shifts as an *investable input* rather than a fixed attribute. Empirically, there are several channels by which firms can improve robustness: more diverse preference elicitation (including cross-cultural and expert feedback), stress testing and adversarial red-teaming, mechanistic interpretability audits targeted at goal representations, training-time interventions that promote conservative uncertainty handling, and post-training experimentation that probes generalization under controlled distribution shifts. These interventions are costly, and their returns are partly indirect: they may reduce the probability of coherent failure in rare regimes rather than increase average benchmark performance.

To organize these ideas, we use “overlap” (also called latent positivity) as a deliberately coarse abstraction: the degree to which the learned objective overlaps with the intended objective across a wide range of contexts, including those not sampled during training. One can interpret overlap investment as increasing the set of contexts on which the reward model is causally identified, or as increasing the margin by which aligned behavior remains the argmax under perturbations. Crucially, overlap is not merely about being “nice” on the training distribution; it is about reducing the probability that a capable system will confidently and coherently pursue an unintended goal under shift. This framing connects standard ML concerns (distribution shift, robustness, calibration) with alignment-specific concerns (objective misidentification, Goodhart effects, instrumental strategies).

The economic question then becomes: if overlap is costly and valuable, will competitive markets supply it? The naive hope is that user demand will discipline firms: users prefer safer assistants, so firms invest in safety to attract users. But three frictions weaken this mechanism in precisely the regimes we care about.

First, much of the harm from misgeneralization is not borne by the marginal user who chooses among assistants. Failures can create third-party externalities: privacy violations, misinformation spillovers, cyber harms, and disruptions that propagate beyond the immediate user. Even when users do bear harm, they may not anticipate tail risks accurately, and realized harms may arrive with delay or be hard to attribute. As a result, user choice underweights the social value of robustness.

Second, safety is a credence attribute. Users rarely observe the counterfactual “this system would have failed under an unobserved shift,” and even

observed incidents are noisy signals of the underlying failure probability. Firms can therefore compete on salient performance metrics while under-providing hard-to-verify robustness work. This is an identification problem again, but now at the market level: the market must infer safety from partial, strategically disclosed evidence. Absent strong auditing or disclosure regimes, reputational incentives will track visible incidents and short-term user experience rather than latent tail risk.

Third, correlated failures create systemic costs that no single firm internalizes. When a common shift affects many deployments simultaneously, joint failures can overwhelm incident response capacity, saturate information channels, or trigger cascading socio-technical effects. Even if each firm internalizes some expected liability or reputational loss from its own failures, it will not internalize the marginal contribution of its safety investments to reducing the *joint* probability of a high-impact event. This is the familiar wedge between private and social incentives in the presence of convex damages and correlation: the planner cares disproportionately about tail outcomes, while firms optimize average private payoffs.

These frictions are not exotic; they are the default in safety-critical domains. What is distinctive in the agentic LLM setting is that (i) deployment shifts are plausibly common across firms (shared platforms, shared tool ecosystems, shared user behavior), (ii) coherent failures can be rare yet catastrophic, and (iii) the safety-relevant properties are expensive to verify and easy to mimic superficially. The combination implies that even when “the market cares about safety,” it may care about a different object than what the planner cares about: visible, local, and attributable incidents rather than latent, correlated, and systemic tail risks.

The goal of this paper is to make this wedge explicit in a minimal model that is faithful to the alignment failure modes we worry about. We proceed by treating firms as choosing a scalar overlap investment that reduces failure probability under both normal and shifted regimes, with the reduction potentially larger in the high-shift regime. This captures an important empirical hypothesis: robust alignment work may not show up much in routine use, but it matters disproportionately under stress. We then embed this choice in a competitive environment in which users allocate demand based on expected utility and firms earn margins that reward adoption. Finally, we allow for policy instruments that map naturally onto governance practice: harm-based liability for external damages and auditable minimum standards (a floor on overlap-like investments).

This structure lets us separate three distinct drivers of underinvestment: (i) pure externalities (harms not priced by user demand), (ii) information and attribution limits (demand does not fully respond to latent safety), and (iii) systemic convexity under correlation (the planner values reducing joint tail events). The resulting comparative statics are useful not because they provide point estimates, but because they clarify which levers matter and

why: when the number of competitors grows, demand-based incentives can dilute; when the common-shock component strengthens, systemic incentives rise; when harms fall on third parties, private incentives decouple from social welfare.

We also emphasize limitations and open problems. Overlap is not directly observable, and real systems are heterogeneous across architectures, training data, and deployment contexts. Correlation structure is itself endogenous: firms share suppliers, datasets, and safety practices, and policy can change these dependencies. Moreover, some safety investments may be non-rival or have spillovers (shared evaluations, shared mitigations), which can either worsen underprovision (public goods) or improve it (coordination). Our purpose is not to claim that a single scalar captures alignment, but to provide a tractable scaffold on which these more realistic features can be layered.

With this framing in place, the next section motivates the stylized facts behind the model: agentic LLM services as competing products, the plausibility of common regime shifts, why correlated goal misgeneralization is a first-order risk, and how “overlap” should be interpreted operationally as an investable, partially auditable input into robustness.

2 Stylized facts and motivation: agentic services, common shifts, correlated failures, and investable overlap

We motivate our model with a set of stylized facts about how frontier assistants are produced and used in a competitive setting around 2026. The details vary across vendors and deployment contexts, but the common structure is stable: firms sell agentic LLM services as products, these products interact with shared digital infrastructure, and the most consequential failures are shaped by distribution shift and are often correlated across providers. These features jointly suggest modeling safety-relevant robustness as a costly, partially verifiable investment that is underprovided by default market incentives.

Agentic LLM assistants as competing services. The relevant economic unit is no longer a static model checkpoint but a *service*: an assistant that is integrated with tools (browsing, code execution, email, calendar, ticketing systems), persistent memory, and organizational workflows. For many users, the assistant is a repeated-use product with switching costs, but also with a meaningful margin of substitution: multiple firms offer roughly comparable capability and product polish, and adoption decisions are often made by individuals (consumer subscriptions) or procurement teams (enterprise seats, API spend). In both cases, we can treat users as choosing

among providers based on a perceived tradeoff between usefulness, price, and safety-relevant reliability.

Importantly, what users experience as “safety” is typically a mixture of (i) directly observed interaction quality (helpfulness, refusal behavior, tone, obvious policy violations), (ii) salient incident history (widely reported jailbreaks or data leaks), and (iii) a diffuse belief about whether the provider is “responsible.” By contrast, the alignment property we ultimately care about—robustly pursuing the intended objective under shift—is largely latent. This gap between the latent property and what demand responds to is one reason we prefer a reduced-form demand system later: it captures that adoption is responsive to *expected* failure risk, but also that this expectation is filtered through limited observation and disclosure.

Common regime shifts are routine, not exceptional. Deployment environments for agentic assistants change in ways that are both rapid and *shared across firms*. Examples include: a new tool API becoming standard (or a previously sandboxed tool receiving expanded permissions); a major platform change (browser security policies, email authentication standards, operating system permissioning); the emergence of a new interaction pattern (agents coordinating across channels, assistants producing long-horizon plans with delegated subtasks); and shifts in adversarial pressure (a new jailbreak meme, a new kind of prompt injection delivered via widely used documents or websites). Even seemingly “local” product changes—longer context windows, multimodal inputs, persistent memory, background task execution—alter the effective state space in which the model must behave safely.

Two features make these shifts naturally modeled as a *common shock*. First, firms are coupled through infrastructure: they build on similar tool ecosystems, are deployed into the same internet, and face the same distribution of user tasks and adversaries. Second, the policy and governance environment itself can induce common changes: new reporting requirements, new auditing norms, or a sudden push to enable or disable classes of functionality. In short, the question is not whether regime shifts occur, but whether safety work is targeted at invariances that survive them.

Why the most important failures look like goal misgeneralization. Many product reliability issues are idiosyncratic and are addressed by standard debugging. Our focus is different: coherent failures that arise when the assistant generalizes the wrong objective. These can be rare in routine use and yet decisive under shift because the assistant is optimizing a proxy that is well-behaved on the training and evaluation distribution but becomes harmful when the causal structure changes. Tool use and agency amplify this: when the model can select queries, browse, write code, or interact with users strategically, it can move into parts of the state space where

the reward model was never identified, while still behaving in a way that appears competent and internally consistent.

A stylized but empirically plausible pattern is that robustness work has *asymmetric* value: it may do little to improve average user satisfaction in the common case (where most models already look good), but it can substantially reduce failure probability in stressed regimes. This asymmetry motivates our later assumption that safety investment can be more effective in “high-shift” states than in “low-shift” states. Mechanistically, this corresponds to interventions that enlarge the set of contexts on which the learned objective is pinned down (diverse preference data, adversarially constructed distributions) or that increase the system’s conservatism and uncertainty awareness when it is off-support.

Correlated failures arise from shared blind spots. Even if firms are in competition, their failures are not independent. Correlation has several sources that are structural rather than accidental.

First, training and alignment pipelines are convergent. Firms draw from similar web-scale corpora, rely on overlapping data vendors, and use comparable post-training techniques (preference modeling, policy optimization, constitutional or rubric-based supervision, red-teaming). If these procedures leave a particular causal ambiguity unresolved, it will tend to be unresolved across providers.

Second, evaluation practices synchronize incentives. If the industry converges on a shared suite of benchmarks and audits, firms will optimize to what is measured. This improves outcomes on measured axes, but it also induces a form of “evaluation monoculture”: unmeasured failure modes persist and may become systematically less salient. In a setting where latent tail risks matter, shared evaluation can therefore increase correlation by aligning what everyone misses.

Third, deployment coupling produces shared triggers. Prompt injection via popular document formats, common enterprise workflows, widely used agent frameworks, and shared tool APIs can all act as synchronized exposure channels. A single technique discovered by attackers can propagate rapidly across deployments, producing a synchronized increase in failure probability.

Correlation is not merely a statistical detail; it changes welfare comparisons. When harms are convex in the number of simultaneous failures—because incident response saturates, because misinformation or cyber harms scale via coordination, or because society cannot absorb many concurrent disruptions—then the planner disproportionately values reducing the probability of joint tail events. Individual firms, by contrast, largely optimize expected private losses from their own incidents. This divergence is the core reason our later model includes a systemic loss term that depends convexly on the number of failing firms.

Systemic harm channels in the agentic setting. It is useful to name concrete mechanisms by which many simultaneous failures are worse than the sum of isolated failures.

One channel is *capacity saturation*. When many systems fail at once, security teams, moderation pipelines, and public communication channels become bottlenecks. Another is *information cascades*: widespread model-generated content can distort beliefs and overwhelm verification mechanisms, especially when outputs are individually plausible but collectively coordinated or mutually reinforcing. A third is *infrastructure interaction*: agents acting through shared platforms can create correlated load, correlated exploitation attempts, or correlated policy-violating content that triggers broad countermeasures with collateral damage. These channels justify treating systemic loss as convex in the number of concurrent failures, even if each single failure is “small” in isolation.

What we mean by “overlap” (latent positivity) as an investable input. We use “overlap” as a deliberately coarse abstraction for the degree to which the assistant’s learned objective continues to coincide with intended objectives across a broad set of contexts, including those not represented in the preference dataset. Operationally, one can view overlap as capturing several partially substitutable interventions: (i) expanding and diversifying preference elicitation (cross-cultural and domain-expert feedback, long-horizon tasks, tool-use supervision); (ii) stress-testing and adversarial evaluation designed to elicit goal-directed failures rather than surface-level policy violations; (iii) mechanistic interpretability and representation-level auditing aimed at identifying brittle goal proxies or deceptive instrumental strategies; (iv) training-time interventions that reduce confident off-support optimization (uncertainty calibration, conservative policies under epistemic uncertainty, constraint enforcement); (v) controlled post-training experimentation that probes invariances under simulated regime shifts (tool perturbations, memory perturbations, distributional stressors).

Crucially, overlap is *costly* and its benefits are *risk-reducing* rather than purely performance-enhancing. This creates two predictable market failures. The first is an externality: third parties can be harmed even when the choosing user benefits (or is indifferent). The second is an observability problem: many overlap-improving activities are difficult for outsiders to verify directly, and realized incidents are a noisy and delayed signal of latent failure probability. The result is a wedge between what firms can monetize through demand and what society values in reduced tail risk.

At the same time, overlap is not completely unverifiable. Governance practice increasingly relies on auditable proxies: documented red-team coverage, evaluation under specified threat models, model card disclosures, incident reporting, and standardized safety cases. These are imperfect, but they

motivate why a regulator might impose a minimum standard or a liability regime keyed to harms. In our model, we treat overlap as an input that can be bounded below by a standard, even if the underlying latent alignment property is not directly observable.

Why these stylized facts point to our modeling choices. The preceding discussion motivates a minimal structure with three components. First, a competitive environment in which users allocate demand based on perceived expected harm, but where that perception is an imperfect proxy for latent tail risk. Second, a common regime state capturing shared deployment shifts, which induces correlation in failure events across firms. Third, a safety investment that reduces failure probabilities and is plausibly more valuable in the stressed regime. Once these ingredients are present, the core comparative statics follow: underinvestment is most severe when harms fall on third parties, when systemic losses are convex in concurrent failures, when the market is fragmented across many firms (diluting demand incentives), and when correlation is strong due to common shocks and shared blind spots. The next section formalizes this structure as a game between firms, users, and a regulator, using overlap as the decision variable that links alignment work to both private incentives and social welfare.

3 Model: overlap investment, demand, common shifts, and policy instruments

We now formalize the competitive environment as a simple game in which firms make a costly, safety-relevant investment—our *overlap* variable—before users choose a provider, after which a common deployment regime-shift is realized and failures occur. The model is intentionally reduced-form: we want a tractable mapping from (i) investable robustness work, (ii) user-facing adoption incentives, and (iii) correlated tail events, to (iv) a wedge between private and socially optimal safety investment and the role of simple policy instruments.

Players and timing. There are $N \geq 2$ competing firms indexed by $i \in \{1, \dots, N\}$ and a continuum of users of mass M . Users obtain utility from adopting exactly one provider.¹ There are also passive third parties who can be harmed by a failure even if they are not the adopting user. A regulator chooses policy *ex ante*.

The timing is:

¹An outside option can be added without changing the main comparative statics; we omit it for notational simplicity.

1. **Policy.** The regulator chooses a liability rate $\tau \geq 0$ per unit third-party harm and a minimum overlap standard $\bar{o} \geq 0$.
2. **Investment.** Each firm i chooses overlap $o_i \geq \bar{o}$ at cost

$$C(o_i) = \frac{c}{2}o_i^2, \quad c > 0.$$

3. **Adoption.** Users observe (or infer from disclosures/audits) each firm's *expected* failure probability and choose providers according to a logit demand system.
4. **Deployment shift and failures.** A common regime state $S \in \{L, H\}$ is realized, with $\Pr(S = H) = s \in (0, 1)$. Conditional on S , each firm experiences a goal-misgeneralization failure with probability depending on its investment. Harms and liability payments are realized.

We emphasize that the common state S is meant to capture shared shifts in the deployment environment (including adversarial pressure), so S is realized after adoption and is not contractible at the time of user choice.

Failure technology and the role of overlap. Fix a regime $S \in \{L, H\}$. If firm i invested o_i , it fails with probability

$$\pi_S(o_i) = \max\left\{0, \min\{1, \pi_S^0 - a_S o_i\}\right\}, \quad (1)$$

where $\pi_S^0 \in (0, 1)$ is the baseline failure probability at $o_i = 0$ and $a_S > 0$ is the marginal effectiveness of overlap investment in reducing failure risk in state S . We allow $a_H \geq a_L > 0$ to capture the empirically salient possibility that robustness work has disproportionate value under stress: investment may do relatively little in a benign regime but sharply reduce failures in a high-shift regime. For most of our analysis we focus on the *linear region* where the truncations in (1) do not bind, so that $\pi_S(o_i) = \pi_S^0 - a_S o_i$.

Let $F_i \in \{0, 1\}$ be the failure indicator for firm i at deployment. Conditional on S , we assume

$$F_i \mid S \sim \text{Bernoulli}(\pi_S(o_i)), \quad \text{independent across } i \text{ conditional on } S.$$

Thus the primary source of correlation across firms is the common shock S (which shifts the level and/or slope of failure risk), rather than direct technological spillovers between firms.²

²We view this as a conservative modeling choice: allowing cross-firm dependence beyond S would typically increase systemic-tail risk and strengthen the case for internalizing externalities.

Define the *systemic load* as the number of failing firms,

$$K = \sum_{i=1}^N F_i.$$

We will later exploit the standard decomposition

$$\mathbb{E}[K^2] = \mathbb{E}[\text{Var}(K \mid S)] + \text{Var}(\mathbb{E}[K \mid S]),$$

which makes explicit how a common shock increases joint-tail risk through the second term.

Harms: internal, external, and systemic. If firm i fails, each user served by i experiences per-user harm $h_U > 0$ (this is internalized in user utility), and third parties experience harm $h_E > 0$ per affected user (this is not internalized by user demand unless policy imposes it). In addition, correlated failures generate a systemic loss that is convex in K . For tractability we use a quadratic form,

$$\text{systemic loss} = \kappa K^2, \quad \kappa \geq 0,$$

interpretable as a reduced-form proxy for saturation and cascade mechanisms. The quadratic is not essential; what matters is convexity in K , which makes joint failures disproportionately costly relative to isolated incidents.

User choice and what is observed. Users care about the *perceived* probability of failure, which we model as the true ex ante expectation given investment:

$$\bar{\pi}(o_i) := \mathbb{E}[\pi_S(o_i)] = (1 - s)\pi_L(o_i) + s\pi_H(o_i).$$

This is the object that can be (imperfectly) learned via evaluation reports, incident history, or standardized safety cases. We deliberately do not model the inference problem explicitly; instead, $\bar{\pi}(o_i)$ summarizes the channel by which overlap investment translates into adoption incentives.

User utility from choosing firm i is

$$U_i = v - p - h_U \bar{\pi}(o_i),$$

where v is a baseline value and p is a (here, fixed) price.³ Given an overlap profile $o = (o_1, \dots, o_N)$, market shares follow a logit demand system with sensitivity parameter $\beta > 0$:

$$s_i(o) = \frac{\exp(\beta U_i)}{\sum_{j=1}^N \exp(\beta U_j)}. \quad (2)$$

³We fold price competition into the reduced-form margin parameter below. Endogenizing prices is feasible but distracts from the safety externalities; the key wedge arises even under fixed prices.

Logit serves two purposes. First, it ensures smooth best responses: improving one's perceived failure probability increases demand, but with diminishing returns when one is already much safer than rivals. Second, it provides an interpretable reduced-form mapping from safety perceptions to market share.

Firm payoffs and liability. Each firm earns a constant per-user margin $\mu > 0$ on its adopted user base. If it is held liable at rate τ for third-party harm, then a failure generates an additional expected cost proportional to the harm caused.

A convenient reduced-form expected profit for firm i is

$$\Pi_i(o; \tau) = \mu M s_i(o) - \frac{c}{2} o_i^2 - \tau h_E M s_i(o) \bar{\pi}(o_i). \quad (3)$$

The last term can be derived from an underlying realization-level liability payment τh_E times (users served) times (failure indicator), taking expectations over failures. We treat τ as a policy instrument that converts external harm into a private expected cost. Importantly, without liability ($\tau = 0$), third-party harm does not enter firm incentives except insofar as it affects user demand through $\bar{\pi}(o_i)$, and only through the internal harm channel h_U .

We interpret the overlap standard \bar{o} as a compliance constraint that is (more) verifiable than the latent alignment property: a regulator can often audit documented red-team coverage, evaluation procedures, or other process commitments that correlate with o_i , even if it cannot directly observe $\pi_S(\cdot)$.

Planner objective. To compare equilibrium investment to a social benchmark, we define a planner objective that aggregates user utility, investment costs, third-party harms, and systemic losses. One convenient formulation is

$$W(o) = M \mathbb{E} \left[\max_i U_i \right] - \sum_{i=1}^N \frac{c}{2} o_i^2 - M \sum_{i=1}^N s_i(o) \bar{\pi}(o_i) (h_U + h_E) - \kappa \mathbb{E}[K^2]. \quad (4)$$

The first term captures the idea that users choose the option giving them the highest realized utility (with logit providing a smooth approximation to discrete choice). The next terms subtract real resource costs of overlap, expected internal plus external harm, and convex systemic losses. In a stripped-down benchmark we will sometimes hold market shares fixed at $s_i(o) = 1/N$ to isolate externalities from strategic demand effects; doing so makes the underinvestment mechanism especially transparent.

Equilibrium concept and the objects we will characterize. Given policy (τ, \bar{o}) , firms simultaneously choose $o_i \geq \bar{o}$ anticipating user adoption (2) and expected liability (3). A (Bayesian) Nash equilibrium is a profile o^*

such that each o_i^* maximizes (3) given o_{-i}^* . Our main focus is on *symmetric* equilibria $o_i^* = o^{NE}(\tau, \bar{o})$ and on the planner's symmetric optimum o^{SP} .

Two modeling choices are worth flagging because they shape the subsequent comparative statics. First, the only endogenous source of cross-firm correlation is the common regime S , which we treat as exogenous. Second, overlap affects failure risk directly, and (through $\bar{\pi}$) affects demand; we do not allow direct technological spillovers whereby one firm's overlap reduces another's failure probability. Both choices bias the model toward *understating* the divergence between private and social incentives: any additional cross-firm coupling in failures would typically strengthen the planner's motive to increase overlap.

With the model defined, the next section derives equilibrium conditions and provides a closed-form characterization in the linear region of (1) under quadratic costs. This will let us state clean existence and uniqueness results for the symmetric equilibrium, compare o^{NE} to o^{SP} , and interpret how common-shock correlation and market fragmentation jointly widen the safety investment wedge.

4 Equilibrium analysis: symmetric existence/uniqueness and a closed-form characterization

We now characterize firms' overlap choices given policy (τ, \bar{o}) . The key object is a symmetric Bayesian Nash equilibrium in which each firm chooses the same overlap level, trading off (i) a direct marginal reduction in expected failure losses that are internalized through liability and (ii) an indirect demand effect through users' perceived safety, against (iii) convex investment costs.

Preliminaries: linear region and an effective slope. Throughout this section we work in the *linear region* of (1), where truncations do not bind and

$$\pi_S(o) = \pi_S^0 - a_S o \quad \text{for } S \in \{L, H\}.$$

Then the ex ante (user-perceived) failure probability is affine in o :

$$\bar{\pi}(o) = (1-s)\pi_L(o) + s\pi_H(o) = \bar{\pi}^0 - a_\bullet o, \quad \bar{\pi}^0 := (1-s)\pi_L^0 + s\pi_H^0, \quad a_\bullet := (1-s)a_L + s a_H.$$

The parameter a_\bullet is the *effective marginal efficacy* of overlap when users and firms evaluate risk ex ante. Importantly, this averaging is taken *before* any systemic term is considered; the role of common-shock correlation enters planner incentives through higher moments, which we defer to the next section.

Best responses and the marginal-incentive decomposition. Fix rivals' overlap profile o_{-i} . From (3), firm i chooses $o_i \geq \bar{o}$ to maximize

$$\Pi_i(o; \tau) = \mu M s_i(o) - \frac{c}{2} o_i^2 - \tau h_E M s_i(o) \bar{\pi}(o_i).$$

Differentiating w.r.t. o_i yields a useful decomposition:

$$\frac{\partial \Pi_i}{\partial o_i} = M \left(\mu - \tau h_E \bar{\pi}(o_i) \right) \frac{\partial s_i(o)}{\partial o_i} - \tau h_E M s_i(o) \bar{\pi}'(o_i) - c o_i. \quad (5)$$

Two channels are immediate:

1. *Demand/reputation channel:* improving o_i decreases $\bar{\pi}(o_i)$, increasing U_i and thereby $s_i(o)$. This is captured by $\partial s_i / \partial o_i$.
2. *Liability/internalization channel:* improving o_i directly reduces expected liability on the firm's existing user base, captured by $-\tau h_E M s_i \bar{\pi}'(o_i)$ (which is positive since $\bar{\pi}'(o_i) = -a_\bullet$ in the linear region).

The first channel is present even at $\tau = 0$ because users internalize h_U in their adoption decision; the second is present only when policy converts external harms into private costs.

To make (5) operational, we use the logit structure (2). The standard derivative is

$$\frac{\partial s_i}{\partial U_i} = \beta s_i(1 - s_i), \quad \frac{\partial U_i}{\partial o_i} = -h_U \bar{\pi}'(o_i) = h_U a_\bullet,$$

so

$$\frac{\partial s_i(o)}{\partial o_i} = \beta s_i(o)(1 - s_i(o)) h_U a_\bullet. \quad (6)$$

Substituting (6) and $\bar{\pi}'(o_i) = -a_\bullet$ into (5) yields

$$\frac{\partial \Pi_i}{\partial o_i} = M \left(\mu - \tau h_E \bar{\pi}(o_i) \right) \beta s_i(1 - s_i) h_U a_\bullet + \tau h_E M s_i a_\bullet - c o_i. \quad (7)$$

Equation (7) makes clear that competition enters through the factor $s_i(1 - s_i)$: when a firm is already very large or very small, marginal safety improvements translate into smaller market-share gains, dampening the demand incentive. Liability, by contrast, continues to create a direct marginal return proportional to the firm's current scale s_i .

Concavity and existence/uniqueness of a symmetric equilibrium. Under our maintained quadratic costs, the investment side is strictly convex. The remaining issue is whether the adoption term can create multiple local optima. In the logit model, s_i is smooth in U_i and hence in o_i , and the mapping $o_i \mapsto s_i(o)$ exhibits diminishing marginal returns because $s_i(1 - s_i)$ is maximized at $s_i = 1/2$ and declines toward 0 as $s_i \rightarrow 0$ or 1. Combined

with the linear dependence of U_i on o_i in the linear region, this yields a single-peaked payoff in o_i for fixed o_{-i} in the parameter ranges we study: the second derivative inherits a strictly negative $-c$ term and (generically) a negative curvature from logit saturation. Operationally, we can treat each firm's best response $BR_i(o_{-i}; \tau, \bar{o})$ as single-valued and continuous.

Symmetry then pins down a unique fixed point. Specifically, consider the candidate symmetric profile $o_i = o$ for all i . Then $U_i = U_j$ and $s_i(o) = 1/N$. Since each best response is single-valued, the symmetric best-response correspondence collapses to a continuous function $BR(o; \tau, \bar{o})$ in the symmetric slice. Strict concavity in own o_i implies the function crosses the 45° line at most once, delivering a unique symmetric equilibrium. Finally, the overlap floor \bar{o} imposes a simple truncation: if the interior optimum lies below \bar{o} , the constraint binds and the equilibrium is at $o = \bar{o}$.

Formally, the Kuhn–Tucker condition for a symmetric equilibrium o^{NE} can be written as the complementarity system

$$0 \leq o^{NE} - \bar{o} \perp \frac{\partial \Pi_i}{\partial o_i} \Big|_{o_i=o^{NE}, o_{-i}=o^{NE} \mathbf{1}} \leq 0, \quad (8)$$

which succinctly captures both interior and corner cases.

Closed form in the symmetric interior (and why it is affine in policy). To obtain a transparent expression, we evaluate (7) at a symmetric interior point, where $s_i = 1/N$ and $s_i(1 - s_i) = (1/N)(1 - 1/N)$. We also adopt the standard reduced-form simplification (used repeatedly in applied IO) that the term $\mu - \tau h_E \bar{\pi}(o)$ is well-approximated by μ when the liability-weighted expected harm is small relative to the per-user margin, or when we want to isolate the dominant comparative-static effects.⁴ Under this approximation, the interior first-order condition becomes

$$c o^{NE} = a_\bullet \left(\tau h_E M \cdot \frac{1}{N} + \mu M \beta h_U \cdot \frac{1}{N} \left(1 - \frac{1}{N} \right) \right). \quad (9)$$

Equation (9) has three immediate implications.

First, the equilibrium overlap is *increasing* in τ : liability directly scales the private marginal benefit of reducing failure risk on the firm's own customers (the $1/N$ term). Second, the demand channel scales with βh_U and is attenuated by competition via $(1/N)(1 - 1/N)$: as N grows, each firm's marginal demand gain from improving safety shrinks roughly like $1/N$. Third, in the linear region the mapping from τ to o^{NE} is *affine*, and the over-

⁴Keeping the exact term simply replaces μ with $\mu - \tau h_E \bar{\pi}(o^{NE})$ in the demand component; the qualitative results below are unchanged in the parameter ranges where the symmetric equilibrium remains interior.

lap floor simply truncates it:

$$o^{NE}(\tau, \bar{o}) = \max \{ \bar{o}, \tilde{o}(\tau) \}, \quad \tilde{o}(\tau) := \frac{a_{\bullet}}{c} \left(\tau h_E M \cdot \frac{1}{N} + \mu M \beta h_U \cdot \frac{1}{N} \left(1 - \frac{1}{N} \right) \right). \quad (10)$$

This explicit form is useful for comparative statics and for interpreting policy: τ shifts incentives one-for-one through the liability channel, while \bar{o} enforces a hard minimum regardless of market conditions.

Interpretation: what the market does and does not internalize. Expression (10) clarifies the economic forces we should expect in competitive LLM deployment.

On the one hand, firms do have private incentives to invest in overlap even without regulation: because users dislike failure risk (weighted by h_U) and because logit demand is responsive (weighted by β), safety improvements can increase market share and thus profit. On the other hand, two attenuation mechanisms are built in. First, *market fragmentation* weakens incentives: when many firms compete, each one captures only a small fraction of the marginal benefit of improving overall ecosystem safety. Second, the private calculus depends on *expected* failure risk $\bar{\pi}$ rather than tail risk conditional on the high-shift regime. In particular, although $a_H \geq a_L$ means overlap may be most valuable under stress, the firm's baseline incentive aggregates regimes linearly through a_{\bullet} and does not (on its own) place extra weight on correlated joint failures.

These are not merely modeling conveniences; they correspond to deployment realities. Users can often react to salient incidents (an average-risk signal), but they are rarely in a position to price the harm to non-users h_E , and they typically do not contract on correlated ecosystem-level events (captured here by κK^2). Liability τ is therefore a targeted instrument: it converts third-party harm into a private expected cost and steepens the equilibrium response in a way that demand alone generally cannot.

Boundary cases and scope of the characterization. The closed form (10) is valid when (i) $\tilde{o}(\tau) \geq \bar{o}$ (so the standard does not bind) and (ii) the linear region is relevant, i.e. $\pi_S^0 - a_S o^{NE} \in (0, 1)$ for both $S \in \{L, H\}$. Outside this range, equilibria can pin to corners: if overlap is extremely effective, the firm may hit the zero-failure truncation in one or both regimes, after which marginal returns drop; conversely, if overlap is ineffective or prohibitively costly, the floor may bind or equilibrium may sit near \bar{o} . These corner regimes matter empirically (e.g. for very capable systems where further overlap yields small incremental reductions in already-low measured failure probability), but they do not change the central logic: private incentives remain tied to own-demand and own-liability exposure.

Transition to the planner benchmark. Having pinned down a unique symmetric equilibrium and an interpretable closed form for $o^{NE}(\tau, \bar{o})$ in the linear region, we are now positioned to ask the welfare question: how does this privately chosen overlap compare to the socially optimal level once we account for third-party harm and convex systemic losses driven by correlated failures? The next section answers this by deriving the planner’s symmetric optimum o^{SP} and decomposing the wedge into internalized demand incentives versus external and systemic components.

5 Social planner benchmark: internal, external, and systemic marginal values of overlap

We now turn to the benchmark that the competitive equilibrium in Section 4 should be compared against: a utilitarian social planner who chooses overlap investments to maximize total welfare, taking into account (i) harm borne by users, (ii) harm borne by non-users (third parties), and (iii) convex systemic losses from correlated failures. Conceptually, this benchmark answers a simple alignment-governance question: *if we could directly set the degree of experimentation/overlap in preference learning across firms, how much would we want, once we price in tail risk and externalities?*

A symmetric welfare objective (and what we hold fixed). Because our goal is to isolate the safety-relevant wedge between private and social incentives, we adopt the standard symmetric reduction: we restrict attention to profiles with $o_i = o$ for all i , and we treat the symmetric demand allocation as approximately uniform, $s_i \approx 1/N$, so that changes in overlap primarily affect welfare through changes in failure risk and investment costs rather than through allocative market-share shifts.⁵ Under this reduction and in the linear region where $\bar{\pi}(o) = \bar{\pi}^0 - a_{\bullet}o$, the relevant welfare tradeoff can be written as

$$W^{\text{sym}}(o) = -N \frac{c}{2} o^2 - M(h_U + h_E) \bar{\pi}(o) - \kappa \mathbb{E}[K^2] + (\text{constants}). \quad (11)$$

The first term is the real resource cost of overlap investment. The second term prices *all* per-user expected harm from failures, including both user harm h_U and third-party harm h_E (the latter is external to user adoption). The third term captures the idea that, even holding fixed per-user harms, *joint* failures can generate additional ecosystem-level losses (e.g. cascading misuse, loss of public trust, correlated critical-infrastructure incidents) that scale convexly in the number of failing deployments K .

⁵One can endogenize shares inside the planner problem as well (e.g. by adding a consumer surplus term consistent with logit). In the symmetric identical-firm class, the key distinction remains: the planner internalizes third-party and systemic harms that do not enter firms’ private objectives absent policy.

Why the systemic term depends on more than mean failure probability. A central modeling point is that $\mathbb{E}[K^2]$ is sensitive not only to the average failure rate but also to *correlation* induced by the common regime state S . Conditional on S , firms fail independently with probability $\pi_S(o)$, so

$$K | S \sim \text{Binomial}(N, \pi_S(o)).$$

A convenient decomposition (law of total variance) is

$$\begin{aligned} \mathbb{E}[K^2] &= \mathbb{E}[\text{Var}(K | S)] + \text{Var}(\mathbb{E}[K | S]) + (\mathbb{E}[K])^2 \\ &= \mathbb{E}[N\pi_S(o)(1 - \pi_S(o))] + \text{Var}(N\pi_S(o)) + (N\bar{\pi}(o))^2. \end{aligned} \quad (12)$$

The first term is the within-regime (idiosyncratic) binomial variance; the second term is the *between-regime* variance generated by the common shock S ; the third term is the squared mean. Overlap reduces $\pi_S(o)$ in both regimes, but when $a_H > a_L$ it also preferentially reduces failures in the high-shift regime, which tends to compress the distribution of $\pi_S(o)$ across S and thereby shrink the common-shock contribution $\text{Var}(N\pi_S(o))$. This is the formal sense in which overlap can mitigate tail-risk correlation: it reduces not only $\bar{\pi}(o)$ but also the dispersion of failure rates across deployment regimes.

For the purposes of a closed-form planner characterization, we will use a baseline approximation in which the dominant systemic force is the squared-mean component, i.e. $\mathbb{E}[K^2] \approx (N\bar{\pi}(o))^2$. This approximation is accurate when (i) failure probabilities are small-to-moderate so that $\pi_S(1 - \pi_S)$ is second order relative to π_S^2 in the region of interest, and/or (ii) we want to isolate the simplest mechanism by which convexity in K creates an extra social incentive to reduce the *level* of failures. We return to the correlation-sensitive refinements (the $\text{Var}(N\pi_S(o))$ term in (12)) when we study how common shocks magnify underinvestment.

Planner first-order condition and marginal-benefit decomposition. Maximizing (11) over $o \geq 0$ yields the (symmetric, interior) planner first-order condition

$$Nco^{SP} = M(h_U + h_E)a_{\bullet} + \kappa \cdot \left(-\frac{d}{do} \mathbb{E}[K^2] \right) \Big|_{o=o^{SP}}. \quad (13)$$

This expression is useful even before we pick an approximation for $\mathbb{E}[K^2]$, because it makes explicit what the planner is trading off:

- *Direct harm reduction (users + third parties):* increasing o reduces $\bar{\pi}(o)$ at rate a_{\bullet} , producing a marginal benefit $M(h_U + h_E)a_{\bullet}$ in the symmetric baseline.
- *Systemic tail-risk reduction:* increasing o reduces the expected convex loss $\kappa\mathbb{E}[K^2]$; the marginal benefit is $\kappa(-d\mathbb{E}[K^2]/do)$.

- *Convex investment cost:* marginal cost is Nco .

The key conceptual distinction from the firm problem is already visible: in the planner FOC, the harm coefficient is $h_U + h_E$ (not merely what users internalize in adoption), and there is an additional term that values reductions in joint failures even if per-user harms were already perfectly priced.

A closed-form baseline with mean-field systemic risk. Under the mean-field approximation $\mathbb{E}[K^2] \approx (N\bar{\pi}(o))^2$, we have

$$-\frac{d}{do} \mathbb{E}[K^2] \approx -\frac{d}{do} (N^2 \bar{\pi}(o)^2) = 2N^2 \bar{\pi}(o) a_\bullet,$$

since $\bar{\pi}'(o) = -a_\bullet$ in the linear region. Substituting into (13) and dividing by N yields the transparent symmetric planner condition

$$co^{SP} = a_\bullet \left(\frac{M(h_U + h_E)}{N} + 2\kappa N \bar{\pi}(o^{SP}) \right). \quad (14)$$

Because $\bar{\pi}(o) = \bar{\pi}^0 - a_\bullet o$ is affine, (14) solves in closed form:

$$o^{SP} = \frac{a_\bullet}{c + 2\kappa N a_\bullet^2} \left(\frac{M(h_U + h_E)}{N} + 2\kappa N \bar{\pi}^0 \right). \quad (15)$$

Two limiting cases are worth flagging for intuition. If $\kappa = 0$ (no systemic convexity), the planner chooses

$$o^{SP} = \frac{a_\bullet}{c} \cdot \frac{M(h_U + h_E)}{N},$$

so the sole difference from private incentives will come from whether firms internalize h_E and how strongly demand reacts to risk. If instead $h_E = 0$ but $\kappa > 0$, overlap is still socially valuable because it reduces the likelihood of multi-firm joint failure events; this is the sense in which systemic risk can justify safety investment even absent classical externalities.

What the planner internalizes that the market does not. Equation (15) provides a clean decomposition of the planner's marginal willingness to pay for overlap into three additive components (all scaled by a_\bullet/c):

1. *Internal user-safety value:* $\frac{Mh_U}{N}$.
2. *External third-party value:* $\frac{Mh_E}{N}$.
3. *Systemic tail-risk value:* $2\kappa N \bar{\pi}(o^{SP})$ (under the mean-field approximation).

In contrast, the symmetric private equilibrium (10) values overlap through (i) a demand/reputation channel, which depends on how strongly adoption responds to perceived user harm (βh_U) and how much profit is at stake (μ), and (ii) an explicit internalization channel only to the extent that policy sets $\tau > 0$.

This difference matters operationally for alignment governance: user choice can only pressure firms on dimensions that are visible, attributable, and privately salient to adopters. Third-party harms h_E (e.g. externalized misuse, labor-market displacement, downstream fraud) are typically not contractible at the point of adoption. Likewise, systemic losses—precisely the events that motivate public concern about frontier deployments—tend to be *joint* and *state-contingent*, and so are only weakly disciplined by firm-level reputation incentives tied to average outcomes.

Interpretive note: why we emphasize the symmetric class. We emphasize the symmetric planner benchmark not because real markets are symmetric, but because it cleanly exposes the structural source of underinvestment: even if firms are equally capable and users are perfectly informed about *expected* failure probabilities, privately optimal overlap is generally governed by appropriable demand gains and any imposed liability, while the planner additionally prices non-user harms and the convexity of correlated bad outcomes. Once this wedge is understood in the symmetric class, extensions (heterogeneous a_S , asymmetric market shares, endogenous disclosure) mainly change *how* the wedge is distributed across firms, not *whether* it exists.

Scope and limitations of the closed form. The expression (15) is a baseline derived in the linear region and under a mean-field approximation for systemic losses. If truncation binds (e.g. $\pi_S(o)$ hits 0 in one regime), marginal benefits become state-dependent and the planner problem becomes piecewise. And if we retain the full decomposition (12), the systemic marginal benefit includes an additional correlation-sensitive term proportional to $d\text{Var}(N\pi_S(o))/do$, which is exactly where the common-shock structure (s and the gap between regimes, including $a_H - a_L$) enters the planner's incentives beyond the mean. These refinements do not change the qualitative message of this section, but they will matter for comparative statics and for quantifying how much policy must lean against correlated tail events.

With the planner benchmark in hand, we can now compare (15) (and its correlation-aware generalization) to the equilibrium characterization from Section 4 and identify conditions under which the market systematically underinvests in overlap, as well as how the wedge scales with externality magnitude, systemic convexity, market fragmentation, and common-shock

strength.

6 Underinvestment and comparative statics: where the wedge comes from

We now compare the competitive outcome to the planner benchmark from Section 5. The object of interest is the *overlap wedge*

$$\Delta := o^{SP} - o^{NE}(\tau, \bar{o}),$$

which operationalizes a governance concern: even if firms optimize given user demand and reputational incentives, do they choose enough experimentation/overlap to control misgeneralization risk in the regimes that matter most?

Private versus social marginal benefits (intuition before algebra). The planner values overlap because it reduces (i) user harm, (ii) third-party harm, and (iii) convex losses from *joint* failures. In contrast, a firm values overlap only to the extent that it (a) increases demand (users avoid higher expected failure probability) and (b) reduces expected liability payments under the policy parameter τ . Two immediate implications follow.

First, if third-party harms h_E are not fully internalized (e.g. $\tau < 1$ or enforcement is incomplete), then a unit reduction in failure probability is socially more valuable than privately valuable, even abstracting away from systemic losses. Second, even if $\tau = 1$ perfectly internalizes third-party harm at the margin, the firm still does not generally internalize the *systemic* component $\kappa \mathbb{E}[K^2]$, because the firm's objective depends (at most) on its own expected harm and adoption, not on the curvature of ecosystem-wide joint-failure loss.

A clean comparison in the symmetric linear region. To make the wedge transparent, we stay in the symmetric class and in the linear region where $\bar{\pi}(o) = \bar{\pi}^0 - a_\bullet o$ (no truncation at 0 or 1). Consider an interior symmetric equilibrium with $o_i = o^{NE}$ for all i and no binding floor $\bar{o} = 0$. The equilibrium first-order condition from Section 4 can be written in the reduced form

$$co^{NE} = a_\bullet \left(\underbrace{\frac{\tau h_E M}{N}}_{\text{explicit internalization}} + \underbrace{\frac{\mu M \beta h_U}{N} \left(1 - \frac{1}{N} \right)}_{\text{demand/reputation channel}} \right). \quad (16)$$

The first term is direct: liability converts third-party harm into a private expected cost proportional to the firm's own demand. The second term is

the demand incentive induced by the logit market: when users are more sensitive to expected utility differences (β large) and margins are higher (μ large), firms have a stronger appropriable incentive to reduce *user* expected harm $h_U \bar{\pi}(o)$.

By contrast, the planner condition from Section 5, under the mean-field systemic approximation, is

$$co^{SP} = a_{\bullet} \left(\frac{M(h_U + h_E)}{N} + 2\kappa N \bar{\pi}(o^{SP}) \right). \quad (17)$$

Comparing (16) and (17) already shows the structural difference: the planner weights per-failure harm by $h_U + h_E$ (not just what is disciplined by demand plus the policy-weighted fraction of h_E) and adds the systemic tail term $2\kappa N \bar{\pi}(\cdot)$, which is absent from the private optimum.

Proving underinvestment ($o^{NE} < o^{SP}$) by monotonicity. The formal argument uses the fact that both problems have strictly convex costs and (in the linear region) linear marginal benefits. Define the private marginal benefit coefficient

$$B^{\text{priv}}(\tau) := \frac{\tau h_E M}{N} + \frac{\mu M \beta h_U}{N} \left(1 - \frac{1}{N} \right),$$

and the planner's (mean-field) marginal benefit evaluated at o ,

$$B^{\text{soc}}(o) := \frac{M(h_U + h_E)}{N} + 2\kappa N \bar{\pi}(o).$$

Then (16) is $co^{NE} = a_{\bullet} B^{\text{priv}}(\tau)$, while (17) is $co^{SP} = a_{\bullet} B^{\text{soc}}(o^{SP})$.

Under the baseline assumptions, $B^{\text{soc}}(o)$ is strictly decreasing in o when $\kappa > 0$ (because $\bar{\pi}'(o) = -a_{\bullet}$), and the planner objective is strictly concave in o after the sign flip (equivalently, strictly convex costs plus diminishing marginal returns through $\bar{\pi}$ inside the systemic term). The firm objective is strictly concave in o as well (quadratic costs plus logit concavity), yielding unique interior solutions.

To show $o^{NE} < o^{SP}$, it suffices to show that at the private optimum the planner's marginal benefit exceeds the private marginal benefit:

$$B^{\text{soc}}(o^{NE}) > B^{\text{priv}}(\tau),$$

because then, at o^{NE} , the planner still has positive net marginal gain from increasing o , and strict concavity implies the planner optimum must lie to the right.

Two conditions are enough:

1. If $h_E > 0$ and $\tau < 1$, then even ignoring systemic risk ($\kappa = 0$),

$$\frac{M(h_U + h_E)}{N} > \frac{\tau h_E M}{N} + \frac{M h_U}{N},$$

and the only remaining question is how much of the h_U term the firm internalizes via demand. In general, the demand channel term in (16) is not equal to Mh_U/N ; it is governed by market competitiveness and appropriability (μ, β, N) . Hence unless demand pressure exactly replicates the planner's valuation of user safety, a positive gap remains, and it is strictly positive whenever $\tau < 1$ and demand is not infinitely disciplining.

2. If $\kappa > 0$, then for any τ we have $B^{\text{soc}}(o) \geq \frac{M(h_U+h_E)}{N}$ and additionally the systemic term $2\kappa N\bar{\pi}(o)$ is strictly positive whenever $\bar{\pi}(o) > 0$. Since the private condition (16) does not include this term, the planner's marginal benefit at o^{NE} exceeds the private marginal benefit whenever failures remain possible in equilibrium (i.e. $\bar{\pi}(o^{NE}) > 0$ in the region of interest), implying $o^{SP} > o^{NE}$.

This establishes Proposition 2 in the symmetric linear region: absent a policy that internalizes *both* third-party and systemic marginal harms, competitive incentives lead to underinvestment in overlap.

Comparative statics of the wedge. Because the equilibrium and planner conditions are monotone, comparative statics follow by differentiating the closed-form expressions (or, more robustly, by applying the implicit function theorem to the respective FOCs).

Third-party harm h_E . The planner's marginal value of overlap increases one-for-one with h_E through $\frac{M(h_U+h_E)}{N}$. The firm's marginal value increases only through the liability-weighted term $\frac{\tau h_E M}{N}$. Consequently,

$$\frac{\partial \Delta}{\partial h_E} > 0 \quad \text{whenever} \quad \tau < 1$$

in the interior linear region: more severe externalities widen the underinvestment gap unless policy scales proportionally.

Systemic convexity κ . The planner invests more as κ rises because the systemic marginal benefit $2\kappa N\bar{\pi}(o)$ scales with κ , while the firm's incentives are essentially unchanged (except indirectly via τ if policy responds). Thus

$$\frac{\partial \Delta}{\partial \kappa} > 0$$

whenever $\bar{\pi}(o^{SP}) > 0$. Operationally, this is the “tail-risk governance” regime: even if average outcomes look acceptable, convex losses from joint failures justify materially higher overlap.

Number of firms N . There are two opposing forces in the planner problem: the direct per-user harm term is divided by N in the symmetric allocation, while the systemic term scales as $2\kappa N\bar{\pi}(o)$. In the firm problem, both the liability and demand terms are diluted by competition, scaling roughly

like $1/N$ (and the demand term includes an extra $(1 - 1/N)$ that saturates near 1). The net effect is that as markets fragment, the private incentive to invest in safety-relevant overlap weakens, while the systemic rationale can strengthen. In particular, when κ is nontrivial, larger N increases Δ over a wide parameter range: more competing deployments increase the expected scale of correlated loss events without proportionally increasing any single firm's appropriable return to safety investment.

Demand sensitivity β and margins μ . Both parameters steepen the demand/reputation channel in (16), raising o^{NE} and shrinking the wedge, holding policy fixed:

$$\frac{\partial \Delta}{\partial \beta} < 0, \quad \frac{\partial \Delta}{\partial \mu} < 0.$$

This comparative static captures an important limitation of “let the market discipline safety”: demand pressure can partially internalize *user*-salient harms when risks are legible and attributable, but it does not, by itself, internalize third-party harms nor ecosystem-level convexities. Moreover, demand pressure is a function of observability and trust in disclosure; if users underweight tail risks or cannot verify them, the effective β for safety-relevant attributes is small even if users are otherwise price/quality sensitive.

Common shocks and correlation: why the wedge grows with regime dispersion. The mean-field approximation highlights the level effect of overlap on joint failures. To see how *common-shock strength* magnifies underinvestment (Proposition 3), we revisit the decomposition in (12). The correlation-relevant component is

$$\text{Var}(N\pi_S(o)) = N^2 s(1-s)(\pi_H(o) - \pi_L(o))^2,$$

which increases with the dispersion between regimes and scales quadratically in N . When $a_H > a_L$, increasing overlap shrinks the regime gap $\pi_H(o) - \pi_L(o)$ faster than it shrinks $\pi_L(o)$, so overlap reduces not only the mean failure probability but also the *between-regime variance*. The planner therefore has an additional marginal benefit term proportional to

$$-\frac{d}{do} \text{Var}(N\pi_S(o)) = 2N^2 s(1-s)(\pi_H(o) - \pi_L(o))(a_H - a_L),$$

which is positive in the empirically relevant case where the high-shift regime is riskier ($\pi_H(o) > \pi_L(o)$) and overlap is more protective there ($a_H > a_L$). Firms do not internalize this variance-reduction benefit unless it is somehow priced through liability tied to systemic outcomes or through coordinated standards. Hence the wedge increases with natural indices of common-shock strength such as $s(1-s)(a_H - a_L)^2$ and $s(1-s)(\pi_H^0 - \pi_L^0)^2$: more regime dispersion makes tail-risk correlation more salient to the planner while leaving private incentives largely anchored to *average* perceived risk.

What breaks these conclusions (and why that matters for governance). All of the comparisons above rely on staying in the interior linear region and on treating overlap as an individually chosen scalar that monotonically reduces failure probability. If truncation binds (e.g. $\pi_L(o)$ hits 0 first), then marginal benefits become state-dependent and the wedge can become piecewise: the planner may continue investing to reduce high-regime tail risk even after the low-regime risk is eliminated, while firms—whose incentives are often driven by average outcomes—may not. Similarly, if “overlap” has multi-dimensional structure (e.g. it trades off capability externalities against alignment robustness), then the one-dimensional monotone comparative statics no longer apply mechanically. These are not merely technicalities: they identify precisely where auditability, measurement, and policy design (Section 7) must take over from closed-form reasoning.

7 Policy design: overlap floors vs. liability, and a simple near-first-best rule

The wedge analysis in Section 6 isolates a familiar governance tension: the object we care about socially is *tail-risk reduction under common shocks*, while the object firms can reliably monetize is (at best) *average, user-salient risk* plus whatever harms are priced by liability. This section therefore treats (τ, \bar{o}) not as abstract primitives but as two policy *implementation channels* with different informational and enforcement requirements.

7.1 Two instruments, two kinds of verifiability

Liability τ is an *outcome-based* instrument. It is attractive because it targets harms directly: if third-party harms are observed, attributable, and collectible, then firms can be induced to internalize them by setting τ appropriately. However, outcome-based schemes are brittle in precisely the regimes motivating our model: common-shock deployment shifts can create diffuse, delayed, and legally ambiguous harms (externalities), and systemic losses κK^2 are typically not contractible at the level of any single deployment. Even if a regulator could compute ecosystem-wide damages *ex post*, allocating them across firms requires a causal attribution rule in a correlated environment.

An overlap floor \bar{o} is instead an *input- (or process-) based* instrument. It is attractive when the regulator can more easily audit *whether* a firm performed sufficient experimentation/diversification than audit *what harms* were caused and by whom. In the present formalism, \bar{o} directly rules out the lowest-overlap equilibria, at the cost of (i) potential over-compliance when the regulator mis-estimates primitives, and (ii) the need for an operational proxy for the scalar o_i .

The practical lesson is that the two instruments are complements: liability is information-efficient when harms are measurable and enforceable; floors are enforcement-robust when harms are not.

7.2 A simple near-first-best rule in the symmetric class

In the symmetric linear region, Proposition 4 already suggests a clean implementation logic: if the regulator can compute (or approximate) o^{SP} , then either (i) directly impose it as a minimum standard, or (ii) choose a liability rate τ so that the firm's private first-order condition coincides with the planner's at o^{SP} . The value of making this explicit is that it yields a *calibration target* for policy: estimate a small set of objects that determine o^{SP} and back out the instrument level.

Formally, write the symmetric interior firm condition as

$$co = a_\bullet(B^{\text{priv}}(\tau)), \quad B^{\text{priv}}(\tau) = \frac{\tau h_E M}{N} + \frac{\mu M \beta h_U}{N} \left(1 - \frac{1}{N}\right),$$

and the planner condition (mean-field systemic approximation) as

$$co = a_\bullet(B^{\text{soc}}(o)), \quad B^{\text{soc}}(o) = \frac{M(h_U + h_E)}{N} + \underbrace{\frac{\partial}{\partial o}(\kappa \mathbb{E}[K^2])}_{\text{systemic marginal benefit per unit risk reduction}} \Big/ (-a_\bullet) .$$

Under the simplified expression in Section 6, this systemic term reduces to $2\kappa N \bar{\pi}(o)$, but the policy logic does not depend on that particular approximation.

A *liability that implements o^{SP}* with $\bar{o} = 0$ solves

$$B^{\text{priv}}(\tau^{FB}) = B^{\text{soc}}(o^{SP}),$$

hence (in the simplified mean-field expression)

$$\tau^{FB} = 1 + \frac{2\kappa N^2 \bar{\pi}(o^{SP})}{h_E M} - \frac{\mu \beta h_U}{h_E} \left(1 - \frac{1}{N}\right). \quad (18)$$

Equation (18) should be read as a *design identity* rather than a literal recommendation: it clarifies which gaps liability must close. The first “1” term is the standard Pigouvian correction for third-party harm. The second term prices systemic convexity. The last term subtracts whatever portion of user harm is already internalized through demand/reputation incentives (and vanishes when $\mu\beta$ is small, i.e. when markets do not reliably discipline safety). When correlation effects are material (Proposition 3), the systemic marginal term should also include the variance-reduction benefit from shrinking between-regime dispersion; in that case τ^{FB} contains an additional positive component proportional to $s(1-s)(\pi_H(o) - \pi_L(o))(a_H - a_L)$ evaluated at o^{SP} .

When outcome-based enforcement is not feasible, the floor alternative is conceptually simpler: set $\bar{o} = o^{SP}$ (or, more realistically, \bar{o} equal to a conservative lower bound on o^{SP} given uncertainty). The essential tradeoff is then *estimation error vs. enforcement error*: liability is sensitive to legal/attribution failure; floors are sensitive to miscalibration of o^{SP} and to Goodharting on the chosen proxy for o_i .

7.3 Auditability: what does it mean to “measure overlap”?

The scalar o_i compresses a bundle of practices—experimentation, diversity of training signals, model pluralism, adversarial evaluation, and robustness work—into a single decision variable. For policy, the relevant question is not whether o_i is metaphysically well-defined, but whether there exist *auditable, approximately monotone proxies* for it: quantities that (i) firms can be required to report or demonstrate, (ii) third parties can verify with bounded effort, and (iii) are predictive of reductions in $\pi_S(o)$, especially in the high-shift regime.

One workable approach is to decompose o_i into a weighted scorecard of verifiable sub-investments, e.g.

$$o_i \approx \sum_{m=1}^M w_m x_{im},$$

where x_{im} are auditable activities (number of independent preference-modeling runs; breadth and provenance of feedback data; diversity of elicitation protocols; red-team coverage across capability domains; robustness evaluations under distribution shift; time/compute budget spent on mechanistic investigations of goal misgeneralization), and weights w_m are set by the regulator based on empirical correlations with downstream failures. This makes the standard implementable even when $\pi_S(\cdot)$ is not directly observable.

However, proxy-based floors invite predictable failure modes:

- *Box-checking and Goodharting.* If the score rewards countable artifacts (documents, evaluations run, datasets added), firms will optimize toward those artifacts even when they weakly affect π_S . Mitigation requires random audits, rotating evaluation suites, and penalties for misrepresentation.
- *Capability externalities.* Some forms of “experimentation” can increase capabilities and thereby increase harms conditional on failure. Our one-dimensional o_i assumes monotone safety improvement; in practice the proxy must be constructed to reward *robustness-oriented* overlap rather than raw scale.
- *Hidden regime dependence.* The key governance objective is improvement in the high-shift regime (a_H), not just the average. Audits should

therefore emphasize stress testing and out-of-distribution evaluations whose construct validity is tied to regime shifts (e.g. tool-use, long-horizon planning, and novel instruction contexts).

A complementary audit channel is to measure *observable consequences* of overlap that are harder to fake than inputs, without requiring full harm attribution. For example, a regulator could require standardized disclosure of (i) reward-model disagreement or instability under controlled shifts, (ii) the sensitivity of safety-relevant metrics to prompt distribution changes, or (iii) incident rates in pre-deployment red-teaming at a fixed coverage budget. These metrics can serve as partially outcome-based signals that sit between pure liability and pure process standards.

7.4 Compliance costs, entry, and the “standard-setting” margin

Even when $C(o) = (c/2)o^2$ is a convenient reduced form, real compliance costs include fixed components: staffing, documentation, audit coordination, and delays to deployment. Fixed costs matter because they interact with market structure: raising fixed compliance burdens can reduce the effective number of competitors N , changing both the private incentives and the systemic exposure. This is not unambiguously good or bad. Fewer deployments can reduce the scale of joint failures (lower K mechanically), but can also concentrate market power and reduce the diversity of safety approaches. In other words, a floor can lower *within-firm* risk while raising *between-firm* dependence if the surviving firms converge on similar pipelines.

This suggests two practical design heuristics.

1. Prefer standards that scale smoothly with deployment scope (e.g. user base, capability level, or domains of use) to avoid cliff effects that destroy entry.
2. Pair floors with *safe-harbor* provisions: if a firm demonstrates compliance and promptly discloses incidents, liability multipliers can be reduced, improving incentives for truthful reporting and post-incident learning.

7.5 When heterogeneity forces numerical (and institutional) methods

The symmetric class is deliberately pedagogical. In realistic settings, firms differ in at least: (i) cost of overlap c_i , (ii) baseline failure rates $\pi_{S,i}^0$, (iii) regime sensitivity $a_{S,i}$, and (iv) exposure profiles (different user bases, domains, and third-party harms). With heterogeneity, a single scalar τ generally cannot implement the full-information first best, and a single floor \bar{o}

can be distortionary: high-cost firms may over-comply relative to their comparative advantage, while low-cost firms may under-comply if the floor is set too low.

Analytically, heterogeneity breaks the affine closed forms and turns both the equilibrium mapping $o^{NE}(\tau, \bar{o})$ and the planner problem into coupled systems. Policy design then becomes a *computational* task: estimate primitives, solve for equilibrium under candidate instruments, and search over (τ, \bar{o}) (or richer menus) to maximize expected welfare subject to enforceability constraints. This is where numerical methods are not an optional add-on but the natural continuation of the theory: the regulator is effectively running counterfactual simulations over governance rules under uncertainty about $(a_{S,i}, \pi_{S,i}^0, c_i)$ and about how auditable proxies map to true risk reduction.

This motivates the next section: we need a calibration story that ties o_i and $\pi_S(o_i)$ to observables (benchmarks, incident reports, distribution-shift evaluations) well enough to (i) bound o^{SP} , and (ii) evaluate how close simple instruments like floors and liability come to that benchmark under plausible parameter ranges.

8 Calibration and counterfactual simulations (illustrative)

The policy rules in Section 7 are only as actionable as our ability to connect the latent choice variable o_i and the failure technology $\pi_S(o_i)$ to observables. In practice, the regulator does not observe $\pi_H(o)$ for the relevant tail regimes, and firms have incentives to selectively disclose. This section therefore sketches a calibration workflow that (i) treats “overlap” as a latent safety investment, (ii) ties it to auditable intermediate measurements (especially distribution-shift gaps in reward modeling and evaluation), and (iii) uses those measurements to run counterfactual simulations over (τ, \bar{o}) . The goal is not to claim that any one metric “is” o_i , but to show how one can bound welfare-relevant objects well enough to compare simple instruments.

8.1 From overlap to auditable intermediate signals: OOD reward-model gaps

A natural place to look for an operational proxy is the *difference between in-distribution and shifted-distribution safety judgments*. Concretely, suppose each firm maintains (or is required to provide for audit) a reward model or safety classifier $R_i(\cdot)$ trained to score candidate assistant behaviors. Let \mathcal{D}^L denote a reference distribution of prompts/contexts (“low shift”) and \mathcal{D}^H a stress-test distribution (“high shift”) designed to elicit rare or adversarial behavior (tool-use, long-horizon plans, novel domains, multi-agent interaction,

jailbreak variants, etc.). Define an auditable *gap statistic*

$$g_i := \mathbb{E}_{x \sim \mathcal{D}^H} [\ell(R_i, x)] - \mathbb{E}_{x \sim \mathcal{D}^L} [\ell(R_i, x)],$$

where ℓ is a standardized loss or disagreement measure (e.g. reward-model variance under ensembling; inconsistency under paraphrases; rate of constraint violations detected by a fixed suite). Intuitively, g_i is large when the firm’s safety signal degrades under shift.

We can then treat o_i as a latent driver of both (a) this intermediate gap, and (b) ultimate failure probability. A simple reduced-form that captures the intended monotonicity is

$$g_i = g_i^0 - \lambda o_i + \varepsilon_i, \tag{19}$$

$$\pi_S(o_i) = \pi_{S,i}^0 - a_{S,i} o_i, \quad S \in \{L, H\}, \tag{20}$$

with $\lambda > 0$ and idiosyncratic noise ε_i reflecting measurement error and the extent to which a given stress-test suite matches real deployment shift. Equation (19) is deliberately not structural: it says only that overlap investment tends to reduce a robustly measurable symptom of brittleness under shift.

A slightly richer version links failures to the gap directly, which is convenient when o_i is not itself verifiable:

$$\pi_S = \sigma(\alpha_S + \gamma_S g_i) \quad \text{or} \quad \pi_S = \min\{1, \max\{0, \tilde{\pi}_S^0 + \tilde{b}_S g_i\}\},$$

with $\gamma_S, \tilde{b}_S > 0$. This creates a bridge from audits (which can measure g_i under regulator-controlled prompt distributions) to predicted failure rates. In enforcement terms, this supports either (i) an input floor framed as a requirement on g_i (“your OOD gap must be below a threshold”), or (ii) a calibration step that backs out plausible (π_S^0, a_S) ranges consistent with observed g_i trajectories.

8.2 Using benchmarks and incidents to estimate regime sensitivity

To use the model for counterfactuals, we need at minimum coarse estimates of π_S^0 and a_S (or their population analogues) and some handle on systemic convexity κ . Three data channels are typically available, each imperfect in different ways.

Public and standardized benchmarks. Safety and robustness benchmarks provide repeated, comparable measurements across time and across model versions. The key is to stratify benchmarks into “low shift” and “high shift” categories that plausibly track $S \in \{L, H\}$. For example, \mathcal{D}^L may be a stable mix of common user prompts and known policy-violation probes, while

\mathcal{D}^H emphasizes novel tool APIs, long-horizon autonomy scaffolds, or distribution mixtures sampled from emerging deployment domains. In this view, the object $a_H - a_L$ is empirically about *how much an incremental overlap-style investment disproportionately improves performance on the high-shift suite* relative to the low-shift suite.

Incident reports and near-misses. While severe failures are (hopefully) rare, near-miss logs and red-team findings can be treated as censored observations on π_H . If an industry-wide reporting regime exists, we can partially correct selection bias by conditioning on exposure and on audit intensity. Even without perfect attribution, the *time series* of incidents following capability jumps or new deployment modalities can inform s (how often high-shift conditions effectively occur) and can bound plausible π_H^0 for frontier deployments.

Internal evaluation artifacts under audit. A regulator can require firms to escrow evaluation traces (prompts, tool-call graphs, reward-model outputs, policy checks) under standardized protocols. Even if the raw model outputs are sensitive, these traces allow third-party computation of shift gaps g_i and other stability statistics. Over time, we can estimate an empirical mapping from overlap-related practices (documented experiment breadth, independent runs, red-team coverage) to changes in these auditable metrics, providing a practical estimate of λ in (19).

In all channels, the central identification challenge is that measured improvements may reflect *capability* changes as well as *safety* investments, and that suite design itself can be Goodharted. A conservative calibration stance is therefore to treat estimates of a_H as interval-valued (or prior distributions) and to report policy performance under pessimistic assumptions about construct validity.

8.3 A minimal counterfactual simulation loop over (τ, \bar{o})

Given estimated (or bounded) primitives, the model can be used as a *policy simulator*: for each candidate instrument pair (τ, \bar{o}) , compute equilibrium overlap, induced failure rates, and welfare. Even in the symmetric baseline, the point of simulation is to propagate uncertainty in S , in the measurement-to-risk mapping, and in the systemic term.

A minimal workflow is:

1. **Fit risk-reduction parameters.** Use benchmark/incident/audit data to estimate (or bound) $(\pi_L^0, \pi_H^0, a_L, a_H)$ and a mapping between auditable proxies and risk (e.g. infer a posterior over a_H given observed changes in g_i after measured increases in experimentation breadth).

2. **Specify market primitives.** Choose $(N, M, \mu, \beta, h_U, h_E)$ using observed market shares, approximate margins, and survey-based estimates of user sensitivity to disclosed risk; treat these as uncertain where necessary.
3. **Solve firm equilibrium.** For each (τ, \bar{o}) , compute the symmetric $o^{NE}(\tau, \bar{o}) = \max\{\bar{o}, \tilde{o}(\tau)\}$ in the linear region, or numerically solve the best-response condition when truncation or heterogeneity matters.
4. **Simulate outcomes.** Draw $S \sim \text{Bernoulli}(s)$ and then $F_i \sim \text{Bernoulli}(\pi_S(o^{NE}))$ independently conditional on S , compute $K = \sum_i F_i$, and record welfare components (user harm, third-party harm, and systemic loss κK^2). Repeat to estimate expected welfare under the policy.
5. **Compare instruments.** Plot $\mathbb{E}[W]$ as a function of τ for fixed \bar{o} (and vice versa), and report the policy frontier under uncertainty sets for (a_H, a_L, κ) .

Two practical refinements matter. First, if users respond to disclosed risk (via β) and disclosure itself depends on audits, one can incorporate a disclosure rule that maps audit outcomes (e.g. observed g_i) to $\bar{\pi}(o_i)$ as perceived by users. Second, if the systemic loss is driven by tail dependence beyond a two-state S , the simulation can replace S with a richer shift process (mixtures over stressor families), while keeping the same conceptual decomposition: overlap affects both mean failure and the co-movement of failures across firms.

8.4 Robust calibration targets: bounding o^{SP} and avoiding false precision

Because o^{SP} depends on κ and on high-shift risk reduction, point estimates are often inappropriate. A more governance-relevant target is a *lower bound* on socially desirable overlap under pessimistic but plausible assumptions. Operationally, we can:

- treat a_H as the key uncertainty and compute $o^{SP}(a_H)$ over a credible interval;
- treat κ as partially identified and report policy performance under a range of cascade severities (e.g. calibrated to historical analogues of correlated outages or security failures);
- prioritize policies whose welfare is insensitive to moderate misspecification (a “flat optimum” criterion), which often favors floors set near a conservative quantile of o^{SP} plus moderate liability, rather than extreme reliance on either instrument alone.

This framing also clarifies what measurement infrastructure is valuable: marginal improvements in estimating κ and a_H can have first-order effects on recommended standards, whereas fine-tuning π_L^0 is often second-order for systemic objectives.

8.5 Limitations and what the simulation can (and cannot) justify

The illustrative procedure above can support *relative* comparisons (e.g. “a modest floor dominates pure liability when attribution is weak”) more reliably than it can justify a single “optimal” τ in absolute terms. The main failure modes are (i) construct invalidity of \mathcal{D}^H (we stress-test the wrong thing), (ii) strategic adaptation (firms learn to reduce g_i without reducing true π_H), and (iii) missing channels (overlap may change capability or deployment scope, affecting harms conditional on failure). For these reasons, we should interpret counterfactual results as conditional statements: *if* the audit suite tracks real regime shift, and *if* overlap investments monotonically reduce high-shift failure, then the simulated welfare rankings are informative.

These limitations motivate the extensions in the next section. In particular, once we allow endogenous monitoring as a separate decision, heterogeneity in costs and regimes, shared foundation-model dependencies, and user-generated prompt dynamics, both calibration and policy simulation become institutionally coupled to monitoring design: the regulator is no longer merely choosing (τ, \bar{o}) given primitives, but shaping the observability of the primitives themselves.

9 Extensions (brief): endogenous monitoring, heterogeneity, shared foundations, and endogenous prompts

The calibration exercise in Section 8 treated the regulator as facing a fixed measurement channel (audits produce some proxy for $\bar{\pi}(o_i)$) and treated “overlap” o_i as the single privately chosen safety-relevant input. In practice, several adjacent choices and structural features matter for both incentives and identification. Here we sketch four extensions that we expect to be first-order for 2026-era regulation: (i) endogenous monitoring as a separate choice variable, (ii) heterogeneity across firms, (iii) shared foundation-model dependencies that couple failures, and (iv) endogenous, user-generated prompts that confound regime shift and selection into stress.

9.1 Endogenous monitoring as a separate (strategic) choice

Many concrete governance proposals implicitly assume that evaluation and monitoring effort is “free” or externally provided. But firms choose how much to invest in (and how much to expose of) monitoring, and these investments affect both (a) true failure probabilities (via earlier detection and remediation) and (b) perceived failure probabilities (via disclosure and user trust). A minimal extension adds a monitoring choice $m_i \geq 0$ alongside overlap o_i , with cost $C_m(m_i) = (d/2)m_i^2$.

There are (at least) two conceptually distinct channels:

1. *Risk reduction*: monitoring reduces true failure probability, e.g.

$$\pi_S(o_i, m_i) = \pi_S^0 - a_S o_i - b_S m_i, \quad b_H \geq b_L > 0,$$

capturing that better evals are particularly valuable under high shift.

2. *Observability and sanctions*: monitoring increases the probability that failures are detected and attributed, which increases effective liability. A reduced-form version is to replace τ by an effective $\tau q(m_i)$ with $q'(\cdot) > 0$, so expected liability becomes $\tau q(m_i) h_E M s_i(o) \bar{\pi}(o_i, m_i)$.

The second channel is governance-relevant even if monitoring does not itself prevent failures: if attribution is weak, liability under-internalizes harms. The strategic implication is that monitoring and overlap become complements or substitutes depending on whether the dominant role of monitoring is prevention (substitute for o_i) or enforceability (complement to o_i through stronger incentives). In symmetric interior regions, the first-order conditions take the schematic form

$$\begin{aligned} co^{NE} &\approx a_\bullet \cdot \left(\text{demand incentives} + \tau q(m^{NE}) \cdot \text{external harm term} \right), \\ dm^{NE} &\approx b_\bullet \cdot \left(\text{demand incentives} + \tau q(m^{NE}) \cdot \text{external harm term} \right) + \tau q'(m^{NE}) \cdot (\text{enforceability we} \dots) \end{aligned}$$

Two safety-relevant failure modes appear immediately. First, if $q(\cdot)$ is largely under firm control (e.g. selective logging, unverifiable eval suites), then private incentives may favor “monitoring theater”: increasing apparent q or reducing disclosed $\bar{\pi}$ without reducing true π_H . Second, even honest monitoring can shift deployment behavior: a firm might increase monitoring to justify expanding deployment scope, potentially increasing harms conditional on failure. This suggests a regulatory separation between (i) *monitoring standards* (protocols, escrow, third-party compute, and penalties for tampering) and (ii) *outcome instruments* (liability and overlap floors). Formally, one can model the regulator as first choosing a monitoring technology (fixing $q(\cdot)$ and the mapping from eval artifacts to $\bar{\pi}$), and only then choosing (τ, \bar{o}) .

9.2 Heterogeneity across firms: costs, efficacy, and user trust

The symmetric baseline is useful for isolating externalities, but heterogeneity is the default: firms differ in safety culture, in evaluation maturity, and in how much overlap investment translates into robustness. Let $(c_i, a_{L,i}, a_{H,i}, \pi_{L,i}^0, \pi_{H,i}^0)$ vary across firms. Even holding demand fixed at $s_i = 1/N$, the privately optimal overlap becomes

$$o_i^{NE}(\tau, \bar{o}) = \max \left\{ \bar{o}, \frac{a_{\bullet,i}}{c_i} \cdot \left(\tau h_E \cdot (\text{scale}) + (\text{any private benefit}) \right) \right\}, \quad a_{\bullet,i} := (1-s)a_{L,i} + s a_{H,i}.$$

Two qualitative changes follow.

First, the welfare cost of a uniform floor \bar{o} can rise: if some firms have high c_i (overlap is expensive) or low $a_{\bullet,i}$ (overlap is ineffective), then forcing them to match the frontier can be inefficient relative to a differentiated policy. This pushes toward standards indexed to auditable proxies (e.g. an upper bound on g_i) rather than a single required input.

Second, market selection becomes ambiguous when safety is partly unobserved. If users imperfectly infer $\bar{\pi}$, firms with low safety may still capture share via branding, bundling, or cross-subsidization. In the logit system, heterogeneity interacts with β : when β is low (users weakly responsive), demand discipline is weak and underinvestment is worse; when β is high but disclosure is noisy, firms can overinvest in marketing signals that correlate weakly with true risk. A natural extension is to explicitly model a signal $\hat{\pi}_i = \bar{\pi}(o_i) + \eta_i$ observed by users, with variance decreasing in monitoring. This endogenizes both information quality and competition on safety, and makes clear why verification infrastructure is not merely “nice to have” but directly incentive-shaping.

9.3 Shared foundation models and upstream coupling of failures

A distinctive feature of frontier deployment is shared dependence on upstream components: multiple “firms” may deploy fine-tuned variants of the same foundation model, share a common tool-use stack, or rely on the same inference provider. This creates an additional coupling across failures beyond the common regime S . A simple way to represent this is to add an upstream shock $Z \in \{0, 1\}$ (e.g. a vulnerability or latent misgeneralization mode common to a model family) with $\Pr(Z = 1) = z$, and let

$$\pi_S(o_i | Z) = \pi_S^0 - a_S o_i + \Delta_S Z, \quad \Delta_S \geq 0,$$

or, more structurally, decompose each firm’s failure indicator as

$$F_i = 1\{U_i \leq \pi_S(o_i)\} \vee 1\{Z = 1 \text{ and } V_i \leq \rho\},$$

so that even if firm-specific failures are conditionally independent given S , the upstream component induces extra positive dependence.

This matters because the systemic term κK^2 is highly sensitive to tail dependence. Even small z can dominate $\mathbb{E}[K^2]$ if Z produces near-simultaneous failures across many deployers. From a governance perspective, this pushes policy attention upstream: overlap investment by any single downstream firm does not fully address upstream shared risk. Two implementation-relevant implications follow. (i) Liability or standards might need to apply not only to deployers but also to foundation-model providers, or to the *interface contract* between them (evaluation artifacts, weight release conditions, incident response obligations). (ii) Audits should explicitly test for *cross-deployer correlated modes*, e.g. by requiring that red-team prompts and tool-call traces be shared (under appropriate confidentiality) so that “one firm’s near miss” is informative about others.

Formally, shared foundations break the clean mapping from each firm’s o_i to aggregate risk. A planner would value investments that reduce Δ_S (upstream hardening) potentially more than investments that reduce idiosyncratic π_S . This motivates extending the action space to include an upstream safety input u chosen by a foundation provider, with its own cost and with spillovers to all N deployers. The resulting game resembles a public-goods problem with both horizontal (across deployers) and vertical (upstream–downstream) externalities, suggesting that simple per-firm floors may be insufficient without upstream obligations.

9.4 Endogenous user-generated prompts and confounding of regime shift

Our two-regime state $S \in \{L, H\}$ stands in for “how stressful deployment is.” But stress is not purely exogenous: user behavior adapts to model capabilities, product design, and publicized incidents. As systems become more agentic, users may supply longer-horizon tasks, more tool access, and more adversarial experimentation, effectively increasing the probability of high-shift conditions. This creates a confounding loop: the observed frequency of high-shift events is jointly determined by deployment choices and by user responses, not just by nature.

A minimal reduced-form is to let

$$s = s(q(o), \text{exposure}, \text{attention}),$$

where $q(o)$ is a measure of average deployed capability or adoption, and “attention” may spike after incidents (increasing adversarial probing). Alternatively, we can model prompts as drawn from a mixture distribution $\mathcal{D} = \omega \mathcal{D}^H + (1 - \omega) \mathcal{D}^L$ with mixture weight ω increasing in adoption and in perceived model power, and then interpret S as a coarse discretization of ω .

This endogeneity has two practical consequences. First, naive calibration can attribute rising incident rates to worse alignment when the true driver is a shift in prompt mix (more high-stakes, longer-horizon tasks). Second, policy can change the prompt mix: strong disclosure and reporting may reduce risky usage (lowering s) but might also increase adversarial attention (raising s) in the short run. These confounds do not invalidate the core externality logic, but they do caution against interpreting estimated (π_H^0, a_H) as stable primitives.

For regulation, the key is to incorporate exposure and use-case controls into audits and incident reporting. Concretely, the regulator can require that firms report standardized denominators (number of tool calls, autonomy depth, sensitive-domain invocations), enabling estimation of failure rates conditional on exposure. In model terms, one can treat harm as $h_U e$ and $h_E e$ where e is an exposure index influenced by product decisions, and then extend the planner problem to jointly choose overlap-like investments and constraints on exposure (rate limits, tool permissions, domain gating). This highlights a policy substitution: if overlap investment is hard to verify, limiting exposure in high-shift channels can partially reduce systemic risk, though at potentially large utility cost.

9.5 What remains open

These extensions share a theme: once monitoring, upstream coupling, and endogenous usage enter, the regulator is no longer choosing (τ, \bar{o}) in a fixed environment, but shaping the environment that makes (τ, \bar{o}) meaningful. The most important open problems are therefore institutional as well as analytic: (i) designing audit protocols robust to strategic adaptation (so that $q(m)$ is not easily manipulated), (ii) building data-sharing and incident-response mechanisms that internalize shared-foundation risk, and (iii) constructing exposure-adjusted safety metrics that remain informative under changing user prompt distributions. Our intent in keeping these extensions brief is not to downplay their importance, but to clarify how they connect to the baseline comparative statics: they typically increase correlation and reduce observability, both of which amplify the gap between private incentives and social welfare and strengthen the case for verifiable standards and enforceable liability.

10 Conclusion: implications for 2026 regulation and firm strategy; limitations and future empirical work

We built a deliberately spare model of “overlap” investment under competition—an input that improves robustness to goal misgeneralization and, crucially,

is most valuable in rare but high-shift deployment regimes. The central qualitative result is robust across the variants we explored: even when users partially discipline safety through demand, competitive equilibrium overlap o^{NE} is generically below the social optimum o^{SP} whenever failures impose third-party harms ($h_E > 0$) or create convex systemic losses (modeled as κK^2). The wedge widens precisely in the circumstances regulators worry about in 2026: larger market participation (higher N), more tail risk from common shocks (high $s(1 - s)$ and greater regime sensitivity $a_H - a_L$), and greater external exposure of failures (higher h_E or larger effective “blast radius”). In that sense, the formalism does not merely restate “externalities exist”: it identifies a particular mechanism by which correlation and tail sensitivity amplify the divergence between private incentives and collective risk.

A first practical implication is that policy should target *marginal* incentives under common shocks, not only average reliability. When $a_H > a_L$, an incremental unit of overlap reduces not only the mean failure probability $\bar{\pi}(o)$ but also the *variance of failure rates across regimes*, thereby dampening correlated tail events. This variance-reduction channel is exactly what a convex systemic term values and exactly what decentralized firms tend to underweight. As a result, policies calibrated solely to expected harm can remain too weak in the presence of tail dependence. Concretely, if regulators want to avoid “many systems fail together” scenarios, then instruments must either (i) explicitly price joint-failure risk (through a systemic surcharge or higher effective τ in high-risk regimes), or (ii) impose verifiable minimum standards \bar{o} (or closely related requirements) that are set with tail scenarios in mind.

A second implication is that the two canonical instruments we analyzed—harm-based liability τ and an overlap floor \bar{o} —should be viewed as complements in realistic governance settings. In the symmetric benchmark, either instrument can in principle implement o^{SP} , but they differ in what they demand from the regulatory measurement channel. Liability asks the regulator (and courts) to measure harms, attribute causality, and enforce payments; a floor asks the regulator to measure and audit inputs or practices. In 2026-era deployment, attribution, counterfactual causality, and harm quantification are all hard precisely in the regimes where correlated failures matter most. This suggests a hybrid design: use liability to capture the portion of harm that is observable and attributable (especially steady-state third-party damages), while using standards—minimum eval coverage, red-team protocols, training-time experimentation budgets, incident-response obligations, and restrictions on deployment scope—to bound tail risk where ex post enforcement is weak. Put differently, τ is a good instrument when measurement is ex post and legible; \bar{o} -like standards are a good instrument when measurement must be ex ante and robust to strategic adaptation.

A third implication concerns firm strategy under competition. If user

choice depends on perceived risk $\bar{\pi}(o_i)$ only imperfectly, then purely reputational incentives can select for *cheap signals* rather than true overlap: marketing, selectively chosen benchmarks, or “monitoring theater” that reduces reported risk without reducing π_H . Our model makes precise why this is more than an information problem: when $\kappa > 0$, even small misalignments between perceived and true tail risk can create large welfare losses because $E[K^2]$ is sensitive to dependence. A rational firm that anticipates regulation and public scrutiny therefore benefits from investing in verifiable, *hard-to-fake* safety inputs and disclosure mechanisms: third-party audits with escrowed logs, standardized evaluation suites that include high-shift stressors, and commitments to share incident indicators that are predictive of correlated failure modes. Strategically, this resembles a move from “compete on claims” to “compete on auditable processes,” which can be privately beneficial if it reduces the chance of sudden liability expansions or moratoria after high-profile incidents.

A fourth implication is upstream. Shared foundations, common tooling, and inference-provider dependencies create coupling that cannot be managed by downstream overlap alone. Even if each deployer chooses high o_i , upstream shocks can generate a large common component in F_i that dominates systemic risk. Thus, policies and contracts should allocate responsibilities along the supply chain: foundation providers should face requirements tied to vulnerability disclosure, safety-case artifacts, evaluation transparency, and incident response, while deployers should face requirements tied to use-case gating, monitoring, and post-deployment controls. From a mechanism-design perspective, the goal is to align incentives so that the party best positioned to reduce a correlated mode is the one facing the marginal cost of leaving it unaddressed. In practice, this points toward governance that treats “model families” and shared stacks as regulated risk pools rather than as independent products.

We also emphasize several limitations of the present analysis. First, we worked primarily in a symmetric class (or a near-symmetric local approximation), which is analytically clarifying but masks important distributional and selection effects: heterogeneous costs c_i , heterogeneous efficacy $a_{S,i}$, and heterogeneous observability of safety can make uniform floors inefficient and can induce market concentration. Second, we used a reduced-form demand system and a reduced-form mapping from overlap to failure probability. Real deployments include nonstationary feedback loops (user adaptation, attacker adaptation, and evolving tool ecosystems), which can make the effective state S partially endogenous and can create time inconsistency in both firm incentives and regulatory commitments. Third, we treated failures as conditionally independent given the regime, then added correlation via a common shock; in reality, dependence can be richer (shared vulnerabilities, cascades through the broader digital ecosystem, and strategic interaction among attackers), and these may require systemic terms more structured than κK^2 .

Finally, we abstracted from political economy constraints and from the institutional capacity required to run audits, enforce disclosure, and adjudicate harms.

These limitations point directly to a future empirical agenda that can make the model operational. The most valuable parameters to estimate are those that govern tail risk and common-shock sensitivity: $(a_H - a_L)$, the frequency and severity of high-shift regimes (our s and related measures of exposure), and the degree to which observed incidents cluster across firms conditional on shared upstream components. Empirically, this requires data that current incident reporting rarely provides: standardized denominators for exposure (e.g. autonomy depth, tool-call counts, sensitive-domain invocations), consistent taxonomies for near misses, and time-synchronized reporting that allows identification of correlated modes. A promising approach is to treat stress-testing as an “experiment” generating estimates of a_S under controlled perturbations, then combine this with field data on exposure and incident rates to infer how often the world visits $S = H$ -like conditions. Similarly, estimating systemic convexity κ is hard, but proxies can be constructed from observed downstream cascade costs (e.g. correlated fraud events, widespread service disruption, or correlated misuse incidents) and from the elasticity of damages to the number of simultaneous failures.

Finally, we view the core contribution as a way to translate alignment-motivated concerns into incentive-compatible governance levers. The safety-relevant takeaway is not merely that “firms underinvest,” but that underinvestment is most severe when (i) harms are external, (ii) failures are correlated, and (iii) the protective action is most effective in high-stress regimes that are hard to observe and verify. These three features are characteristic of agentic deployments, which suggests that light-touch, purely reputational approaches will be structurally fragile. The hopeful note is that the same formalism clarifies what can work: verifiable standards for high-shift robustness, liability that is tied to measurable external harms, and monitoring infrastructures that are designed to be tamper-resistant and comparable across firms. Advancing from a stylized model to actionable policy will require better measurement and stronger institutions, but the incentive logic is already clear: without mechanisms that internalize tail-correlated externalities, competitive pressure will push precisely against the safety margin we most need to preserve.