# Universal Catastrophic Safety Undecidability and the Capability–Risk Frontier: Unified Theorems and Governance Pathways

**A**uthor One
author1@institution.edu

*D*epartment of Computer Science
Institution Name
Address Line

**A**uthor Two
author2@institution.edu

*D*epartment of Statistics
Institution Name
Address Line

December 1, 2025

## Abstract

The deployment of autonomous agents and large-scale learned policies raises two fundamental questions: which safety guarantees are decidable in principle, and which guarantees are inherently statistical. We establish comprehensive theoretical boundaries and empirical validation for both.

**First (Decidability Landscape)**, we prove that verifying whether a computable policy satisfies a probabilistic safety threshold is $\Sigma_1^0$-complete (Theorem 4.2), but becomes decidable when restricted to finite-state controllers (Theorem 4.5), revealing a sharp memory-dependent phase transition (Theorem 4.7). This systematic landscape distinguishes policy verification from planning undecidability.

**Second (Unified Capability–Risk Framework)**, we develop a general information-complexity functional $C(Q_S; D, P)$ (Definition 5.1)

1

that unifies PAC-Bayes, Mutual Information, and Wasserstein-1 distributionally robust optimization (Theorem 5.5). The resulting high-probability bound on worst-case risk decomposes as $\widehat{R}_S + \sqrt{C/n} + L\rho$, combining generalization complexity and robustness cost. We prove this is minimax optimal via information-theoretic lower bounds (Theorem 6.4), and significantly tighten practical bounds using data-dependent Lipschitz constants $\bar{L}$ (Theorem 5.7), achieving $10^6$–$10^{10}\times$ improvement over global estimates.

**Third (Systematic Empirical Validation)**, we conduct extensive experiments across multiple datasets (CIFAR-10, CIFAR-100), training methods (standard ERM, PGD adversarial training, spectral normalization), and Lipschitz estimation techniques, validating the frontier's universality. On complex Safe RL tasks, our SSR (Scope–Shield–Risk) framework eliminates catastrophic events (72%→0%) with 27% intervention rate, demonstrating practical feasibility of theoretically-grounded safety governance.

**Keywords:**   AI safety, undecidability, PAC-Bayes, distributionally robust optimization, adversarial robustness, formal verification, runtime shielding

# 1   Introduction

The deployment of autonomous agents raises a fundamental safety question: can we *verify* that a given policy will not trigger catastrophic outcomes? Classical program verification theory suggests pessimism: the halting problem is undecidable [Turing, 1936], and Rice's theorem states that any non-trivial semantic property of programs is undecidable [Rice, 1953]. However, these results apply to arbitrary Turing-complete programs. In the context of *interactive* agent–environment systems with *probabilistic* safety specifications, the precise decidability boundary remains unclear.

Simultaneously, statistical learning theory reveals fundamental limits on generalization and robustness. PAC-Bayes bounds [McAllester, 1999, Catoni, 2007] and mutual information paradigms [Xu and Raginsky, 2017, Steinke and Zakynthinou, 2020] quantify how model complexity degrades generalization. Wasserstein distributionally robust optimization (DRO) [Villani, 2009, Esfahani and Kuhn, 2018] bounds worst-case risk under distribution shift. However, these frameworks have not been unified into a single "capability vs. risk" trade-off, nor have they been connected to formal verification impossibility results.

This paper establishes both impossibility and possibility results, forming

a complete picture:

**Contributions.**

1.  **Safety Decidability Landscape** (Section 4): Beyond showing $\Sigma_1^0$-completeness for Turing-complete policies (Theorems 4.1, 4.2), we establish a *complete phase diagram*:

    -   *Decidable region*: Finite-state controllers + finite MDPs yield PSPACE-decidable verification (Theorem 4.5).
    -   *Phase transition*: Memory-bounded policies exhibit complexity transitions—bounded memory gives EXPTIME, unbounded memory recovers $\Sigma_1^0$-completeness (Theorem 4.7).
    -   *Specification extensions*: Restricted LTL safety properties preserve $\Sigma_1^0$-hardness (Theorem 4.9).

    This systematic landscape (Table 1) distinguishes policy verification from planning undecidability and clarifies which safety guarantees are algorithmically achievable.

2.  **Unified Capability–Risk Framework via Information Complexity** (Section 5): We introduce a general information-complexity functional $C(Q_S; D, P)$ (Definition 5.1) that subsumes PAC-Bayes KL divergence and stability-based mutual information as special cases. This yields:

    -   *Unified bound* (Theorem 5.5): $R_\rho^{\mathrm{rob}}(Q_S) \le \widehat{R}_S(Q_S) + \sqrt{\frac{2C(Q_S; D, P) + \ln(1/\delta)}{n}} + L\rho$.
    -   *Data-dependent refinement* (Theorem 5.7): Replacing global Lipschitz $L$ with data-dependent $\bar{L} = \mathbb{E}_D[L_{\mathrm{local}}]$ yields practical bounds orders of magnitude tighter (empirically $10^6$–$10^{10}\times$ improvement).
    -   *Minimax optimality* (Theorem 6.4): Information-theoretic lower bounds prove the three-term decomposition is unavoidable: $\inf_{\mathcal{A}} \sup_D R^{\mathrm{rob}} \ge \Omega(\sqrt{d/n} + L\rho)$.

3.  **Systematic Empirical Validation Across Domains** (Section 7): We validate the unified framework via:

    -   *Multi-dataset frontier*: CIFAR-10, CIFAR-100, demonstrating cross-dataset consistency of capability-risk scaling.

- *Multi-method comparison*: Standard ERM, PGD adversarial training, spectral normalization—showing different training methods occupy distinct frontier positions.

- *Lipschitz surrogate analysis* (Section 7.5): Comparing global spectral norm, gradient-based, and finite-difference estimators validates Theorem 5.7's practical impact.

- *Complex Safe RL + SSR pipeline* (Section 7.6): 16×16 grid-world with complete SSR implementation eliminates catastrophic events (72%→0%) at 27% intervention cost, demonstrating theory-to-practice feasibility.

4. **SSR Governance Framework**: We translate theoretical impossibilities into the Scope–Shield–Risk engineering framework (Section 8), providing actionable guidance for deploying learned policies under formal safety constraints.

The remainder of the paper is organized as follows. Section 2 positions our work relative to prior decidability, learning theory, and AI safety literature. Section 3 formalizes the interaction model and safety specifications. Sections 4–6 present our main theoretical results. Section 7 provides empirical validation. Section 8 presents the SSR governance framework, and Section 9 concludes.

## 2 Related Work

**Undecidability in Planning and Control.** Madani et al. [1999] proved that finding a policy for a POMDP that achieves expected reward above a threshold is undecidable. Our Theorem 4.1 differs in two ways: (i) we address *verification* of a *fixed* computable policy rather than policy *search*, and (ii) we position the problem within the arithmetical hierarchy ($\Sigma_1^0$-completeness). This distinction is critical for AI safety: even if a neural policy is provided as executable code, verifying its safety is impossible in general.

**Program Verification and Safety Properties.** Safety properties in model checking are often specified via regular languages or temporal logic (LTL) [Baier and Katoen, 2008]. Safety automata and bad-prefix languages have been used in runtime monitoring and shield synthesis [Alshiekh et al., 2018, Könighofer et al., 2024]. Our contribution is to prove that probabilistic

safety verification inherits undecidability from the halting problem, even for the simplest interactive settings.

**PAC-Bayes and Information-Theoretic Generalization.** PAC-Bayes bounds [McAllester, 1999, Catoni, 2007, Alquier, 2021] relate generalization error to KL divergence from a prior. Information-theoretic bounds [Xu and Raginsky, 2017, Steinke and Zakynthinou, 2020, Bu et al., 2020] use conditional mutual information (CMI) or algorithmic stability. Our Theorem 5.5 unifies both via a min-of-bounds construction, providing practitioners a choice based on prior quality vs. algorithmic stability.

**Distributionally Robust Optimization.** Wasserstein DRO [Esfahani and Kuhn, 2018, Sinha et al., 2017] optimizes over distribution balls $\mathbb{B}_\rho(D)$. The Kantorovich–Rubinstein duality yields a linear penalty $L\rho$ for Lipschitz functions [Villani, 2009]. We integrate this with generalization bounds to form a complete "capability–risk frontier."

**Adversarial Robustness and Trade-offs.** Tsipras et al. [2018] empirically demonstrated that robust accuracy and standard accuracy may conflict. Schmidt et al. [2018] proved sample complexity lower bounds for adversarial learning. Our Theorem 6.1 provides a distribution-free geometric lower bound via point perturbations, applicable to any metric space.

**AI Safety and Alignment.** The AI safety community has long recognized verification challenges [Amodei et al., 2016]. Interruptibility [Orseau and Armstrong, 2016], impact regularization [Turner et al., 2019], and shielded RL [Alshiekh et al., 2018] are practical proposals. Our SSR framework systematizes these into a three-layer architecture grounded in formal undecidability and statistical bounds.

## 3 Preliminaries

### 3.1 Interaction Semantics and Trace Spaces

**Finite and Infinite Traces.** Let $\mathcal{A}$ (actions) and $\mathcal{O}$ (observations) be finite alphabets. A **finite trace** is an element of $\Sigma = (\mathcal{A} \times \mathcal{O})^\star$. An **infinite trace** is an element of $\Sigma^\omega = (\mathcal{A} \times \mathcal{O})^\omega$. The agent–environment interaction generates an infinite trace $h = (a_1, o_1, a_2, o_2, \dots)$. For $t \geq 1$, let $h_{1:t} = (a_1, o_1, \dots, a_t, o_t) \in \Sigma$ be the finite prefix.

**Computable Policy.** A **computable policy** is a Turing-computable function $A : \Sigma \to \mathcal{P}(\mathcal{A})$ that, given any finite history $h_{<t}$, computes (in finite time) a probability distribution over actions. We allow internal randomization via sampling procedures.

**Environment.** An **environment** $E$ is specified by a history-conditional probability kernel $\mu(o_t \mid h_{<t}, a_t)$. If $\mu$ is computable (i.e., for any rational approximation, there exists a Turing machine computing it), we call $E$ a computable environment.

**Trivial Environment $E_0$.** We define the **deterministic trivial environment** $E_0$ as follows: for any history and action, $E_0$ always returns a fixed observation $o_\perp$. Formally, $\mu(o \mid h, a) = \mathbf{1}\{o = o_\perp\}$.

## 3.2 Safety Specifications via Regular Bad Prefixes

**Definition 3.1** (Bad Prefix Language). A set $B \subseteq \Sigma$ is a **bad prefix language** if it is **extension-closed**: for all $u \in B$ and $v \in \Sigma$, we have $uv \in B$. Equivalently, the complement $S = \Sigma \setminus B$ is **prefix-closed**.

**Definition 3.2** (Regular Bad Prefix). A bad prefix language $B$ is **regular** if it is recognized by a deterministic finite automaton (DFA). Equivalently, the safe prefix set $S$ is prefix-closed and regular.

Regular bad prefixes are equivalent to *safety properties* in LTL: formulas of the form $\mathsf{G} \neg \psi$, where $\psi$ is a propositional formula over traces [Baier and Katoen, 2008].

**Definition 3.3** (Violation Event and Safety Predicate). Given a policy $A$, environment $E$, and bad prefix language $B$, the **violation event** is

$$\mathsf{Bad} = \{\text{infinite trace } h \in \Sigma^\omega : \exists t \geq 1, \ h_{1:t} \in B\}.$$

The **threshold safety predicate** for threshold $\varepsilon \in [0, 1)$ is:

$$\mathrm{Safe}_\varepsilon(A, E, B) \iff \Pr_{A,E}(\mathsf{Bad}) \leq \varepsilon. \tag{1}$$

**Problem Instance Encoding.** An instance of the safety verification problem is a tuple $(A, E, B, \varepsilon)$, where:

- $A$ is encoded as a Turing machine (with access to randomness tape),

- $E$ is encoded as a conditional probability kernel (for computable $E$, this is a Turing machine computing rational approximations),

- $B$ is encoded as a DFA,

- $\varepsilon \in [0, 1) \cap \mathbb{Q}$ is a rational number.

## 3.3 Learning and Distributional Robustness

**Statistical Learning Setup.** Let $\mathcal{Z}$ be a data domain with metric $d$, and $\mathcal{H}$ a hypothesis class. A learning algorithm takes a sample $S = (Z_i)_{i=1}^n \sim D^n$ and outputs a posterior distribution $Q_S \in \mathcal{P}(\mathcal{H})$. The loss function $\ell : \mathcal{H} \times \mathcal{Z} \to [0, 1]$ is assumed bounded.

**Lipschitz Assumption.**

**Assumption 3.4** (Uniform Lipschitz). There exists a constant $L > 0$ such that for all $h \in \mathcal{H}$, the map $z \mapsto \ell(h, z)$ is $L$-Lipschitz with respect to $d$. For neural networks, $L$ is controlled via spectral normalization, gradient clipping, or Lipschitz-constrained architectures.

**Wasserstein Distance and Robust Risk.** The Wasserstein-1 distance between distributions $D, D'$ on $\mathcal{Z}$ is

$$W_1(D, D') = \sup_{f : \|f\|_{\mathrm{Lip}} \leq 1} |\mathbb{E}_D[f] - \mathbb{E}_{D'}[f]|.$$

The **Wasserstein-1 ball** of radius $\rho$ is $\mathbb{B}_\rho(D) = \{D' : W_1(D', D) \leq \rho\}$.

**Definition 3.5** (Robust Risk). The **robust risk** of a posterior $Q$ under shift $\rho$ is

$$R_\rho^{\mathrm{rob}}(Q) = \sup_{D' \in \mathbb{B}_\rho(D)} \mathbb{E}_{h \sim Q, z \sim D'} \ell(h, z).$$

## 3.4 Data-Dependent Lipschitz Constants

While Assumption 3.4 provides a global Lipschitz constant $L$, practical bounds can be significantly tightened by using *data-dependent* Lipschitz measures that average local smoothness over the data distribution.

**Assumption 3.6** (Data-Dependent Lipschitz Constant). Define the **local Lipschitz constant** at $z \in \mathcal{Z}$ for hypothesis $h \in \mathcal{H}$ as

$$L_{\mathrm{local}}(h, z) := \sup_{z' \in \mathcal{Z}, z' \neq z} \frac{|\ell(h, z) - \ell(h, z')|}{d(z, z')}.$$

The **data-dependent Lipschitz constant** is

$$\bar{L}(Q, D) := \mathbb{E}_{h \sim Q, z \sim D}[L_{\text{local}}(h, z)].$$

When $Q$ and $D$ are clear from context, we write $\bar{L}$.

*Remark* 3.7. By Jensen's inequality and the definition of global $L$, we have $\bar{L}(Q, D) \leq L$ for all $Q$ and $D$. In practice, $\bar{L}$ can be orders of magnitude smaller than $L$, especially when the data distribution concentrates on smooth regions of the loss landscape. Estimating $\bar{L}$ can be done via:

- **Gradient-based proxy**: $\bar{L} \approx \mathbb{E}_{h \sim Q, z \sim D} \|\nabla_z \ell(h, z)\|$.

- **Finite-difference sampling**: Sample pairs $(z, z')$ with $d(z, z') \approx \epsilon$ and estimate $|\ell(h, z) - \ell(h, z')|/\epsilon$.

# 4 Undecidability of Catastrophic Safety Verification

We now establish that safety verification for computable policies is undecidable and position it within the arithmetical hierarchy.

**Theorem 4.1** (Universal Safety Verification Undecidability). *There exists a recursively enumerable family of regular bad prefix languages $\mathfrak{B}$ such that the following decision problem is undecidable:*

INPUT: *Encoding of $(A, E, B, \varepsilon)$ where $A$ is a computable policy, $E$ is a computable environment, $B \in \mathfrak{B}$, and $\varepsilon \in [0, 1) \cap \mathbb{Q}$.*
QUESTION: *Does $\text{Safe}_\varepsilon(A, E, B)$ hold?*

*Proof Sketch.* We reduce from the Halting Problem. Fix the trivial environment $E_0$. Given a Turing machine $M$ and input $x$, construct a policy $A_{M,x}$ that simulates $M(x)$ and outputs a special action $a_\star$ if and only if $M(x)$ halts. Define the bad prefix language

$$B_{\text{halt}} = \{h \in \Sigma : \exists i \leq |h|, \text{ the } i\text{-th action is } a_\star\}.$$

$B_{\text{halt}}$ is regular (recognized by a DFA with one accepting sink state). Since $E_0$ is deterministic and always returns $o_\perp$, the probability $\Pr(\mathsf{Bad})$ is either 0 (if $M(x)$ does not halt) or 1 (if $M(x)$ halts). For any $\varepsilon \in (0, 1)$, deciding $\Pr(\mathsf{Bad}) > \varepsilon$ is equivalent to deciding whether $M(x)$ halts. Thus, safety verification is at least as hard as Halting. □                □

## 4.1 Complexity Positioning: $\Sigma_1^0$-Completeness

We now refine Theorem 4.1 by positioning the problem within the arithmetical hierarchy.

**Theorem 4.2** ($\Sigma_1^0$-Completeness of UNSAFE). *Let $E_0$ be the deterministic trivial environment. Define*

$$\mathsf{UNSAFE} = \{(A, E_0, B, \varepsilon) : \Pr_{A,E_0} (\mathsf{Bad}) > \varepsilon\}, \quad \varepsilon < 1.$$

*Then* UNSAFE *is $\Sigma_1^0$-**complete**. Its complement* SAFE $= \{(A, E_0, B, \varepsilon) : \Pr_{A,E_0}(\mathsf{Bad}) \leq \varepsilon\}$ *is $\Pi_1^0$-**complete**.*

*Proof Sketch.* **(Hardness)** By the reduction in Theorem 4.1, HALT $\leq_m$ UNSAFE.

**(Membership in $\Sigma_1^0$)** Under $E_0$ and computable policy $A$, the interaction is deterministic (up to $A$'s internal randomness, which we model as a fixed random tape). The event Bad occurs if and only if there exists a finite time $t$ such that $h_{1:t} \in B$. A Turing machine can enumerate all prefixes of increasing length, simulate $A$ with $E_0$, and check membership in $B$ (which is decidable for regular languages). If a bad prefix is found, the machine halts and accepts. Thus, UNSAFE $\in \Sigma_1^0$.

Since HALT is $\Sigma_1^0$-complete and HALT $\leq_m$ UNSAFE, we have UNSAFE is $\Sigma_1^0$-complete. The complement SAFE is therefore $\Pi_1^0$-complete. $\square$ $\square$

*Remark* 4.3 (Probabilistic Case). If the environment is stochastic, deciding $\Pr(\mathsf{Bad}) > \varepsilon$ for arbitrary $\varepsilon$ requires computing infinite sums of path probabilities, which may not be semi-decidable in general. Theorem 4.2 holds for $E_0$ where path probabilities collapse to 0 or 1.

## 4.2 Comparison with Rice's Theorem

**Corollary 4.4** (Rice-Style Safety Theorem). *Let $\mathcal{L}(A, E, B) = \Pr_{A,E}(\mathsf{Bad})$ denote the "safety functional" of a policy $A$ with respect to $(E, B)$. For any non-trivial property $P$ of $\mathcal{L}(A, E_0, B)$ (i.e., $P$ holds for some $A$ but not all $A$), the problem*

$$\{A : P(\mathcal{L}(A, E_0, B))\}$$

*is undecidable.*

*Proof.* Immediate from Theorem 4.1 by setting $P$ to be "$\mathcal{L}(A, E_0, B) > \varepsilon$" for some fixed $\varepsilon$. $\square$ $\square$

This corollary is analogous to Rice's Theorem but applies to the *behavioral safety* of interactive agents rather than the *semantic properties* of standalone programs.

## 4.3 Decidable Regions: Finite-State Controllers

While Theorem 4.2 establishes universal undecidability for Turing-complete policies, restricting the policy class can restore decidability. We now show that finite-state controllers admit algorithmic safety verification.

**Theorem 4.5** (Decidability for Finite-State Controllers). *Let $\mathcal{A}$ be a finite-state controller (FSC) with $k$ states, $E$ a finite-state MDP with $m$ states, and $B$ a regular bad-prefix language given by a DFA $\mathcal{D}_B$ with $\ell$ states. Then the problem*

$$\text{``Does } \Pr_{\mathcal{A},E}(\mathsf{Bad}) > \varepsilon\,?\text{''}$$

*is **decidable** with complexity in* PSPACE *(and* EXPTIME *in the worst case).*

*Proof Sketch.* Construct the product system $\mathcal{M} = \mathcal{A} \times E \times \mathcal{D}_B$, which is a finite Markov chain with $O(km\ell)$ states. The bad-prefix event corresponds to reaching absorbing "bad" states in $\mathcal{M}$. The probability $\Pr(\mathsf{Bad})$ can be computed by solving a system of linear equations (reachability probabilities). Comparing rational $\Pr(\mathsf{Bad})$ with $\varepsilon$ is decidable. The complexity is polynomial space (PSPACE) for representing the system, with potential exponential time blowup in solving the linear system. □　　□

*Remark* 4.6. This result connects to classical model checking [Baier and Katoen, 2008]: verifying probabilistic safety properties on finite Markov chains is decidable, in contrast to the infinite-horizon case with Turing-complete policies.

## 4.4 Memory Threshold and Phase Transition

The contrast between Theorem 4.2 (undecidable for Turing-complete policies) and Theorem 4.5 (decidable for FSCs) suggests a *complexity phase transition* based on policy memory.

**Theorem 4.7** (Memory-Bounded Phase Transition). *Consider policies with bounded memory $k$ (finite-state machines with $k$ states).*

1. ***Bounded regime:** If $k$ is fixed (constant), then safety verification is in* EXPTIME.

2. **Unbounded regime:** *If k is allowed to grow with the problem encoding (i.e., k is part of the input), then the problem becomes $\Sigma_1^0$-complete.*

*Proof Sketch.* **(Bounded):** For fixed $k$, the product system size is polynomial in the environment and specification, leading to EXPTIME complexity for solving the reachability problem.

**(Unbounded):** When $k$ is unbounded, we can encode a Turing machine's tape into the policy's state space. Given a TM $M$ and input $x$, construct a policy $\mathcal{A}_{M,x}$ whose states track $M$'s tape contents. The policy simulates $M$ on $x$, encoding halting into a safety violation. This reduction shows that allowing unbounded memory reinstates $\Sigma_1^0$-completeness, as we can embed the Halting Problem.                               □                    □

*Remark* 4.8. This theorem formalizes the intuition that "finite memory admits decidability, but infinite memory recovers undecidability." It mirrors similar phase transitions in computational complexity theory (e.g., bounded vs. unbounded nondeterminism).

## 4.5   Extensions to Richer Specification Languages

Our results so far focus on regular bad-prefix languages $B$. A natural question: do richer temporal logics (e.g., LTL, PCTL) change the decidability picture?

**Theorem 4.9** (Extension to Restricted LTL Safety). *Let $\varphi$ be an LTL safety formula of the form $\mathbf{G}\neg\psi$, where $\psi$ is a Boolean combination of atomic propositions over a finite observation window. Then:*

1. *$\varphi$ can be compiled into a DFA $\mathcal{D}_\varphi$, yielding a regular bad-prefix language.*

2. *The safety verification problem "Does $\mathrm{Pr}_{\mathcal{A},E}(\neg\varphi) > \varepsilon$?" for Turing-complete $\mathcal{A}$ remains $\Sigma_1^0$-complete.*

*Proof Sketch.* **(1)** Standard LTL-to-automata translation [Baier and Katoen, 2008] converts $\mathbf{G}\neg\psi$ into a DFA recognizing bad prefixes.

**(2)** Since the specification reduces to a regular language, Theorem 4.2's reduction from Halting applies directly: embed the Halting Problem into a policy that violates $\varphi$ iff the encoded TM halts. Thus, $\Sigma_1^0$-completeness persists.                               □                    □

### 4.6  Safety Decidability Landscape: Summary

Table 1 summarizes the decidability boundaries across different policy classes, environments, and specification languages.

Table 1: Decidability phase diagram for AI safety verification. The table shows how restricting policy expressiveness, environment structure, or specification language affects the computational complexity of verifying $\Pr(\mathsf{Bad}) > \varepsilon$.

| Policy | Environment | Specification | Complexity | Theore |
|---|---|---|---|---|
| Turing-complete | Trivial $E_0$ | Regular $B$ | $\Sigma_1^0$-complete | 4.2 |
| Finite-state (FSC) | Finite MDP | Regular $B$ | PSPACE/EXPTIME | 4.5 |
| Bounded memory ($k$ fixed) | Finite MDP | Regular $B$ | EXPTIME | 4.7 |
| Unbounded memory | Finite MDP | Regular $B$ | $\Sigma_1^0$-complete | 4.7 |
| Turing-complete | Trivial $E_0$ | LTL safety | $\Sigma_1^0$-complete | 4.9 |

*Remark* 4.10 (Implications for AI Safety Governance). This landscape reveals that:

- **Total verification is impossible** for general-purpose AI systems (Turing-complete policies).

- **Decidability can be recovered** by restricting to finite-state controllers or bounded-memory agents, at the cost of reduced expressiveness.

- **Specification language richness** (regular vs. LTL) does not fundamentally alter decidability for Turing-complete policies—the bottleneck is policy expressiveness, not specification complexity.

These results justify the shift from *verification* to *probabilistic risk bounding* in Section 5.

## 5  The Unified Capability–Risk Frontier

Having established impossibility of static verification, we turn to *probabilistic* guarantees under distributional assumptions. This section develops a **unified information-complexity framework** that systematically integrates PAC-Bayes bounds, mutual information generalization, and Wasserstein distributionally robust optimization into a single theoretical structure.

## 5.1 Unified Information Complexity Functional

We begin by abstracting the commonality across different generalization bounds: they all measure how much a learned posterior $Q_S$ "deviates" from some reference (prior, data distribution, or stability baseline). We formalize this via a unified complexity measure.

**Definition 5.1** (Unified Information Complexity). Let $Q_S \in \mathcal{P}(\mathcal{H})$ be a posterior distribution over hypotheses, $P \in \mathcal{P}(\mathcal{H})$ a prior, and $S \sim D^n$ a training sample. Define the **unified information complexity** as

$$C(Q_S; D, P) := \inf_{\lambda \in \Lambda} \left\{ \lambda_1 \cdot \mathrm{KL}(Q_S \| P) + \lambda_2 \cdot \mathrm{CMI}(S; Q_S \mid D) + \Phi(\lambda) \right\},$$

where:

- $\Lambda \subset \mathbb{R}_+^2$ is a parameter space (e.g., $\{\lambda : \lambda_1 + \lambda_2 = 1, \lambda_i \geq 0\}$),

- $\Phi : \Lambda \to \mathbb{R}_+$ is a convex regularizer,

- $\mathrm{CMI}(S; Q_S \mid D) := I(S; Q_S) - I(D; Q_S)$ measures conditional mutual information.

*Remark* 5.2 (Interpretation).
- **PAC-Bayes regime:** When $\lambda_2 = 0$ and $\Phi(\lambda) = 0$, we recover $C(Q_S) \approx \mathrm{KL}(Q_S \| P)$.

- **Stability/MI regime:** When $\lambda_1 = 0$, we recover $C(Q_S) \approx \mathrm{CMI}(S; Q_S)$, corresponding to stability-based generalization [Steinke and Zakynthinou, 2020].

- **Adaptive regime:** In intermediate cases, the infimum automatically selects the tightest complexity measure for a given $Q_S$, yielding a data-dependent bound.

## 5.2 Assumptions and Preliminaries for Robust Risk

**Assumption 5.3** (Sub-Gaussian Loss). For each $(h, z)$, the random variable $\ell(h, Z)$ (where $Z \sim D$) is $\sigma^2$-sub-Gaussian.

**Assumption 5.4** (Conditional Mutual Information). The learning algorithm satisfies $\mathrm{CMI}(S; Q_S \mid D) \leq \Gamma$ for some constant $\Gamma > 0$, where

$$\mathrm{CMI}(S; Q_S \mid D) = \mathbb{E}_{S \sim D^n} \left[ \mathrm{KL}(P_{Q_S|S} \| P_{Q_S}) \right].$$

Assumption 5.4 is satisfied by algorithms with bounded uniform stability [Steinke and Zakynthinou, 2020].

### 5.3   Main Theorem: Unified Upper Bound via Information Complexity

**Theorem 5.5** (Unified Capability–Risk Bound via Information Complexity). *Let Assumptions 3.4, 5.3, and 5.4 hold. Fix a prior $P \in \mathcal{P}(\mathcal{H})$ and confidence $\delta \in (0,1)$. Then with probability at least $1 - \delta$ over $S \sim D^n$, the robust risk satisfies*

$$\boxed{R_\rho^{\mathrm{rob}}(Q_S) \leq \widehat{R}_S(Q_S) + \sqrt{\frac{2C(Q_S; D, P) + \ln(1/\delta)}{n}} + L\rho}\qquad(2)$$

*where $C(Q_S; D, P)$ is the unified information complexity from Definition 5.1.*

*Proof Structure.* **Step 1 (Wasserstein DRO via KR Duality):** For any $L$-Lipschitz function $g$ and distribution shift bound $\rho$, Kantorovich–Rubinstein duality [Villani, 2009] gives

$$\sup_{D' \in \mathbb{B}_\rho(D)} \mathbb{E}_{D'}[g] \leq \mathbb{E}_D[g] + L\rho.$$

   **Step 2 (Unified Generalization Bound):** For any $\lambda = (\lambda_1, \lambda_2) \in \Lambda$, define the complexity measure

$$\Psi_\lambda(Q_S) := \lambda_1 \cdot \mathrm{KL}(Q_S \| P) + \lambda_2 \cdot \mathrm{CMI}(S; Q_S) + \Phi(\lambda).$$

Standard PAC-Bayes concentration [McAllester, 1999] and MI-based generalization [Steinke and Zakynthinou, 2020] both yield bounds of the form: with probability $1 - \delta$,

$$\mathbb{E}_D[g] \leq \widehat{R}_S(Q_S) + \sqrt{\frac{2\Psi_\lambda(Q_S) + \ln(1/\delta)}{n}}.$$

   **Step 3 (Optimization over $\Lambda$):** Taking the infimum over $\lambda \in \Lambda$ gives $C(Q_S; D, P) = \inf_\lambda \Psi_\lambda(Q_S)$. Applying KR duality with $g(z) = \mathbb{E}_{h \sim Q_S} \ell(h, z)$ (which is $L$-Lipschitz by Assumption 3.4) completes the proof. Full details in Appendix A.                          □                          □

**Corollary 5.6** (Min-of-Bounds Instantiation). *By choosing $\Lambda = \{(1,0), (0,1)\}$ with $\Phi \equiv 0$, we obtain*

$$R_\rho^{\mathrm{rob}}(Q_S) \leq \min\{\mathcal{B}_{PAC}(Q_S), \mathcal{B}_{MI}(Q_S)\} + L\rho,\qquad(3)$$

*where*

$$\mathcal{B}_{PAC}(Q_S) = \widehat{R}_S(Q_S) + \sqrt{\frac{\mathrm{KL}(Q_S\|P) + \ln(2/\delta)}{2n}}, \qquad (4)$$

$$\mathcal{B}_{MI}(Q_S) = \widehat{R}_S(Q_S) + \sqrt{\frac{2\sigma^2(\Gamma + \ln(2/\delta))}{n}}. \qquad (5)$$

*This is the form used in our experiments (Section 7).*

*Proof.* Setting $\lambda = (1,0)$ recovers the PAC-Bayes bound; $\lambda = (0,1)$ recovers the MI bound. Taking min corresponds to $\inf_{\lambda \in \{(1,0),(0,1)\}}$. $\qquad\square\qquad\square$

### 5.4 Data-Dependent Refinement

While Theorem 5.5 provides a general framework, the $L\rho$ term often dominates due to the looseness of global Lipschitz constants. We now refine the bound using data-dependent Lipschitz measures from Assumption 3.6.

**Theorem 5.7** (Data-Dependent Capability–Risk Bound). *Under Assumptions 5.3, 5.4, and 3.6, with probability at least $1 - \delta$ over $S \sim D^n$,*

$$\boxed{R_\rho^{\mathrm{rob}}(Q_S) \leq \widehat{R}_S(Q_S) + \sqrt{\frac{2C(Q_S; D, P) + \ln(1/\delta)}{n}} + \bar{L}(Q_S, D) \cdot \rho} \qquad (6)$$

*where $\bar{L}(Q_S, D)$ is the data-dependent Lipschitz constant from Assumption 3.6.*

*Proof Sketch.* The key modification is in Step 1 of Theorem 5.5's proof. Instead of using global Lipschitz constant $L$, we apply a refined Kantorovich–Rubinstein duality:

$$\sup_{D' \in \mathbb{B}_\rho(D)} \mathbb{E}_{z \sim D'}[g(z)] \leq \mathbb{E}_{z \sim D}[g(z)] + \mathbb{E}_{z \sim D}[L_{\mathrm{local}}(z)] \cdot \rho + O(\rho^2),$$

where $g(z) = \mathbb{E}_{h \sim Q_S}\ell(h, z)$ and $L_{\mathrm{local}}(z) = L_{\mathrm{local}}(h, z)$ averaged over $h \sim Q_S$. For small $\rho$, the $O(\rho^2)$ term is absorbed into the generalization term. The remaining steps (Steps 2-3 of Theorem 5.5) proceed identically. Full proof in Appendix A. $\qquad\square\qquad\square$

*Remark* 5.8 (Practical Impact). In our CIFAR-10 experiments (Section 7.2), we observe:

- **Global $L$**: Estimated via spectral norm products, yields $L \sim 10^{10}$ to $10^{13}$ for standard models, rendering the bound vacuous ($L\rho \gg 1$).

- **Data-dependent** $\bar{L}$: Estimated via gradient norms, yields $\bar{L} \sim 10$ to $10^3$, providing meaningful bounds.

This refinement is critical for practical risk budgeting in AI safety applications. See Section 7.5 for detailed empirical comparison of Lipschitz estimators.

## 5.5  Interpretation and Regime Analysis

Theorem 5.5 reveals a three-term decomposition of robust risk:

1. **Empirical Error** $\widehat{R}_S(Q_S)$: Bias from finite sample approximation.

2. **Complexity/Confidence Term**: Either $\sqrt{\mathrm{KL}/n}$ (prior-dependent) or $\sqrt{\Gamma/n}$ (stability-dependent).

3. **Robustness Cost** $L\rho$: Linear penalty for distribution shift, controlled by Lipschitz constant.

Table 2: Regime Analysis: When to Use PAC-Bayes vs. Mutual Information Bound

| Bound Type | Complexity Term | Best Regime | Requirement |
|---|---|---|---|
| PAC-Bayes | $\mathrm{KL}(Q_S\|P)$ | Good prior $P$ available | Explicit stochastic $Q_S$ |
| Mutual Info | $\mathrm{CMI}(S;Q_S)$ | Stable algorithm | Sub-Gaussian loss, stability |

**Practical Guidance.**

- If domain knowledge suggests a strong prior $P$ (e.g., sparse weights, low-rank structure), use $\mathcal{B}_{\mathrm{PAC}}$.

- If the algorithm is designed for stability (e.g., SGD with small learning rate, implicit regularization), use $\mathcal{B}_{\mathrm{MI}}$.

- The min-of-bounds construction ensures the frontier is as tight as possible given available information.

# 6  Matching Lower Bound: Structural Trade-off

To prove the Capability–Risk trade-off is *structural* (not an artifact of loose analysis), we provide a matching lower bound.

**Theorem 6.1** (Point-Perturbation Geometric Lower Bound). *For any classifier $f : \mathcal{Z} \to \mathcal{Y}$ and shift radius $\rho > 0$, define the adversarial risk*

$$\mathcal{R}_\rho^{adv}(f) = \Pr_{(Z,Y)\sim D} \left[\exists Z' \in B_\rho(Z) : f(Z') \neq Y\right],$$

*where $B_\rho(Z) = \{Z' : d(Z,Z') \leq \rho\}$. Then*

$$\boxed{\sup_{D'\in\mathbb{B}_\rho(D)} R_{D'}(f) \geq \mathcal{R}_\rho^{adv}(f).} \tag{7}$$

*Proof.* For each $(z,y) \sim D$, let $T(z)$ be a measurable selection of $z' \in B_\rho(z)$ such that $f(z') \neq y$ if such $z'$ exists; otherwise $T(z) = z$. Define the transported distribution $D' = (T,Y)_\# D$, i.e., $D'$ is the pushforward of $D$ under $(z,y) \mapsto (T(z),y)$.

By definition, $d(z,T(z)) \leq \rho$ almost surely. Consider the coupling $\pi(dz,dz') = D(dz)\delta_{T(z)}(dz')$. Then

$$\mathbb{E}_\pi[d(Z,Z')] = \mathbb{E}_D[d(z,T(z))] \leq \rho.$$

By the Kantorovich duality, $W_1(D,D') \leq \rho$, so $D' \in \mathbb{B}_\rho(D)$.

By construction, $f$ errs on $(T(z),y)$ whenever there exists $z' \in B_\rho(z)$ with $f(z') \neq y$. Therefore,

$$R_{D'}(f) = \Pr_{(Z',Y)\sim D'}[f(Z') \neq Y] = \Pr_{(Z,Y)\sim D}[f(T(Z)) \neq Y] = \mathcal{R}_\rho^{adv}(f).$$

Taking supremum over all $D' \in \mathbb{B}_\rho(D)$ yields (7). $\qquad\square\qquad\qquad\square$

*Remark* 6.2 (Gaussian Mixture Example). Consider a binary classification task with $D$ as an equal mixture of $\mathcal{N}(\mu,\sigma^2 I)$ (class +1) and $\mathcal{N}(-\mu,\sigma^2 I)$ (class -1) in $\mathbb{R}^d$. Let $\|\mu\| = r$ and $\sigma$ be fixed. A linear classifier $f(x) = \text{sign}(\langle w,x\rangle)$ with $w = \mu/\|\mu\|$ achieves Bayes-optimal accuracy on $D$.

Under adversarial perturbation of radius $\rho$, the adversary can shift each point toward the decision boundary. If $\rho \sim \sigma$, the effective signal-to-noise ratio degrades from $r/\sigma$ to $(r-\rho)/\sigma$. To maintain adversarial accuracy, one must increase the margin (equivalently, increase $r$ by collecting more samples or using stronger regularization), which may decrease standard accuracy on $D$ if the decision boundary shifts.

This phenomenon aligns with Tsipras et al. [2018]: robust training increases the margin but may sacrifice fitting the Bayes-optimal boundary, increasing $\widehat{R}_S$ in Theorem 5.5.

**Proposition 6.3** (Tightness of KR Linear Term). *For any $L, \rho > 0$, there exist a metric space $(\mathcal{Z}, d)$, an $L$-Lipschitz function $f : \mathcal{Z} \to \mathbb{R}$, and distributions $D, D'$ with $W_1(D, D') = \rho$ such that*

$$\mathbb{E}_{D'}[f] - \mathbb{E}_D[f] = L\rho.$$

*Thus, the $L\rho$ term in Theorem 5.5 is unimprovable without additional geometric assumptions.*

*Proof.* Let $\mathcal{Z} = \mathbb{R}$, $d(z, z') = |z - z'|$, $D = \delta_0$ (Dirac at 0), $D' = \delta_\rho$ (Dirac at $\rho$). Let $f(z) = Lz$. Then $f$ is $L$-Lipschitz, $W_1(D, D') = \rho$, and

$$\mathbb{E}_{D'}[f] - \mathbb{E}_D[f] = L\rho - 0 = L\rho.$$

$\square$                                                          $\square$

## 6.1   Information-Theoretic Minimax Lower Bound

Beyond geometric tightness, we now establish that the three-term decomposition in Theorem 5.5 is *minimax optimal* in its dependence on sample size $n$ and complexity measures.

**Theorem 6.4** (Minimax Lower Bound for Robust Learning). *Consider the setting of binary classification on $\mathcal{Z} = \mathbb{R}^d$ with a hypothesis class $\mathcal{H}$ of VC dimension $d$. Let $D$ be a distribution over $\mathcal{Z} \times \{0, 1\}$, and fix a shift radius $\rho > 0$. Then there exists a universal constant $c > 0$ such that*

$$\inf_{\mathcal{A}} \sup_{D, h^* \in \mathcal{H}} \mathbb{E}_{S \sim D^n} \left[ R_\rho^{\mathrm{rob}}(\mathcal{A}(S)) - R_\rho^{\mathrm{rob}}(h^*) \right] \geq c \left( \sqrt{\frac{d}{n}} + L\rho \right),$$

*where the infimum is over all learning algorithms $\mathcal{A} : (\mathcal{Z} \times \{0, 1\})^n \to \mathcal{H}$, and $L$ is the Lipschitz constant of $\mathcal{H}$.*

*Proof Sketch.* We construct a hard instance using a Gaussian mixture model similar to Remark 6.2.

    **Construction:** Let $D_\theta$ be a distribution indexed by $\theta \in \{-1, +1\}^d$ (binary hypercube), where:

$$D_\theta(x, y) = \frac{1}{2} \left[ \mathcal{N}(\theta \cdot \mu, \sigma^2 I) \otimes \delta_{+1} + \mathcal{N}(-\theta \cdot \mu, \sigma^2 I) \otimes \delta_{-1} \right],$$

with $\|\mu\| = \Theta(\sqrt{d})$. The Bayes-optimal classifier is $h_\theta^*(x) = \mathrm{sign}(\langle \theta, x \rangle)$.

**Step 1 (Standard risk lower bound):** By Fano's inequality and Le Cam's method, distinguishing between hypotheses in $\{h_\theta^* : \theta \in \{-1, +1\}^d\}$ requires $\Omega(d)$ samples. Thus,

$$\inf_{\mathcal{A}} \sup_{\theta} \mathbb{E}_{S \sim D_\theta^n} [R_{D_\theta}(\mathcal{A}(S)) - R_{D_\theta}(h_\theta^*)] \geq c_1 \sqrt{\frac{d}{n}}.$$

**Step 2 (Robustness amplification):** For each $D_\theta$, construct a shifted distribution $D_\theta' \in \mathbb{B}_\rho(D_\theta)$ by applying adversarial perturbations of magnitude $\rho$ toward the decision boundary. Using optimal transport lower bounds [Villani, 2009], we show:

$$R_\rho^{\mathrm{rob}}(h) - R_{D_\theta}(h) \geq c_2 L\rho \cdot \mathbb{P}_{D_\theta}(\text{near boundary}),$$

where $\mathbb{P}(\text{near boundary}) = \Omega(1)$ by construction.

**Step 3 (Combining):** The robust excess risk decomposes as:

$$R_\rho^{\mathrm{rob}}(\mathcal{A}(S)) - R_\rho^{\mathrm{rob}}(h_\theta^*) \geq [R_{D_\theta}(\mathcal{A}(S)) - R_{D_\theta}(h^*)] + [R_\rho^{\mathrm{rob}}(h_\theta^*) - R_{D_\theta}(h_\theta^*)].$$

The first term is $\Omega(\sqrt{d/n})$ by Step 1; the second is $\Omega(L\rho)$ by Step 2. Combining yields the claimed lower bound. Detailed calculations in Appendix D. $\qquad\qquad\square\qquad\qquad\qquad\qquad\qquad\square$

**Proposition 6.5** (Matching Upper and Lower Bounds)**.** *Theorem 6.4 shows that the $\sqrt{d/n}$ and $L\rho$ terms in Theorem 5.5 are unavoidable: any algorithm must suffer robust risk at least $\Omega(\sqrt{d/n} + L\rho)$ in the worst case. Combined with Proposition 6.3, this establishes that the three-term decomposition*

$$Robust\ Risk = \underbrace{\widehat{R}_S}_{bias} + \underbrace{\sqrt{complexity/n}}_{generalization} + \underbrace{L\rho}_{robustness}$$

*is **tight up to constants and logarithmic factors** in all three terms.*

*Remark* 6.6 (Implications for AI Safety)*.* This minimax result implies:

1. **No algorithmic silver bullet**: No learning algorithm can escape the $\Omega(\sqrt{d/n} + L\rho)$ barrier without additional assumptions (e.g., data structure, smoothness).

2. **Capability-risk trade-off is fundamental**: Increasing model capacity ($d$) or reducing training data ($n$) provably worsens robust risk. Adversarial robustness ($\rho$-shift resilience) requires paying the $L\rho$ cost.

3. **Safety verification must be statistical**: Combined with the undecidability results of Section 4, these bounds justify the shift from deterministic verification to probabilistic risk budgeting with explicit error rates.

# 7   Experiments

We design two families of experiments to empirically support our theoretical results. First, we investigate the Capability–Risk frontier on CIFAR-10 by varying model capacity and measuring robust error under controlled distribution shift, together with the theoretical upper bounds of Theorem 5.5. Second, we evaluate runtime shielding in a gridworld environment, quantifying how much catastrophic risk can be eliminated at what intervention cost.

**Note on Undecidability.**   Theorem 4.1 is a pure mathematical result and cannot be directly empirically validated. The experiments focus on Theorems 5.5 and 6.1.

## 7.1   Common Setup and Metrics

Throughout this section we use the following metrics.

- **Standard error**: $R_{\text{std}}(f) = \Pr_{(x,y)\sim D}[f(x) \neq y]$ on the unperturbed test distribution.

- **Robust error**: $R_{\text{rob}}^{(\rho)}(f) = \Pr_{(x,y)\sim D}[f(\tilde{x}) \neq y]$ where $\tilde{x}$ is obtained from $x$ by an adversarially chosen or stochastic perturbation with effective radius $\rho$.

- **Capability–Risk bound**: for each trained model we instantiate the unified bound of Theorem 5.5 with a Gaussian PAC-Bayes prior and compute

$$\mathcal{B}_{\text{cap-risk}}(Q_S) := \min\left\{\mathcal{B}_{\text{PAC}}(Q_S, P),\ \mathcal{B}_{\text{MI}}(Q_S, S)\right\} + L\,\rho$$

  where $L$ is an empirical Lipschitz constant estimate and $\rho$ is the perturbation radius used to construct the shifted distribution.

Unless otherwise stated, we report mean and standard deviation over 3 independent random seeds, and use a confidence level of $\delta = 0.01$ in all bounds.

## 7.2   Capability–Risk Frontier on CIFAR-10

**Objective.**   This experiment empirically validates Theorem 5.5: we verify that the robust risk under distribution shift scales linearly with the Lipschitz

constant $L$ and shift radius $\rho$, and that the unified PAC-Bayes + Wasserstein bound provides a valid (though loose) upper envelope for empirical robust error.

**Dataset and distribution shift.** We use CIFAR-10 with the standard train-test split and standard data augmentation (random crop with 4-pixel padding and random horizontal flips). To simulate distributional shift of magnitude $\rho$, we corrupt each test image $x$ with i.i.d. Gaussian noise $\tilde{x} = \text{clip}(x + \sigma\xi, 0, 1)$, where $\xi \sim \mathcal{N}(0, I)$ and $\sigma$ is chosen so that the expected $\ell_2$ perturbation satisfies $\mathbb{E}\|\tilde{x} - x\|_2 \approx \rho$. We evaluate robust error for several radii $\rho \in \{0.0, 0.25, 0.5\}$.

**Architectures and training.** We follow the common practice of using a residual network backbone. Concretely, we instantiate a width-scaled ResNet-18 with width factor $w \in \{0.5, 1.0, 2.0\}$ by scaling the number of channels in each block. For all models we use SGD with momentum 0.9, weight decay $5 \times 10^{-4}$, batch size 128, and train for 200 epochs with a cosine learning-rate schedule starting from 0.1. We select the checkpoint with the best validation accuracy on a held-out subset of the training data.

**Estimating Lipschitz constants.** To estimate the Lipschitz constant $L$ of each trained network we follow the spectral-norm approach. For each convolutional or fully connected layer we approximate its spectral norm by 20 steps of power iteration on the unfolded weight matrix, and multiply the per-layer spectral norms to obtain an overall upper bound on $L$. This is a standard albeit loose approximation, which is sufficient for studying scaling trends.

**Instantiating the Capability–Risk bound.** For the PAC-Bayes component we consider a factorized Gaussian prior $P = \mathcal{N}(0, \sigma_p^2 I)$ over weights and a factorized Gaussian posterior $Q_S = \mathcal{N}(w, \sigma_p^2 I)$ centered at the learned parameters $w$. Under this parameterization the KL divergence reduces to

$$\text{KL}(Q_S \| P) = \frac{\|w\|_2^2}{2\sigma_p^2} + \text{const},$$

where the constant term does not affect relative comparisons across models and is dropped. We set $\sigma_p = 0.1$ in all experiments. The empirical risk $\widehat{R}_S(Q_S)$ is computed on the training set using Monte Carlo sampling of weight perturbations.

For the mutual-information component we follow the conditional mutual information (CMI) formulation from Steinke and Zakynthinou [2020], and use a standard CMI estimator based on repeated training runs with randomized labels. Full details and the exact estimator are given in Appendix C.

Finally, for each trained model and each radius $\rho$ we compute the unified Capability–Risk bound as

$$\mathcal{B}_{\text{cap-risk}}(Q_S) = \min \left\{ \widehat{R}_S(Q_S) + \sqrt{\frac{\text{KL}(Q_S\|P) + \ln(1/\delta)}{2n}}, \widehat{R}_S(Q_S) + \sqrt{\frac{2\sigma^2(\text{CMI}(S;Q_S) + c\ln(1/\delta}{n}} \right.$$

where $n$ is the training sample size and $\sigma^2, c$ are as in Theorem 5.5.

**Results.** Figure 1 plots the empirical robust error $R_{\text{rob}}^{(\rho)}$ against the estimated Lipschitz constant $L$ for all width factors and radii, together with the Capability–Risk bound. Each marker corresponds to one trained model.

We observe three consistent phenomena.

- **Linear dependence on $L\rho$.** For fixed $\rho$ the robust error increases roughly linearly with the estimated Lipschitz constant, in line with the $L\rho$ term in Theorem 5.5.

- **Effect of model capacity**. Increasing the width factor from 0.5 to 2.0 reduces standard error on the clean distribution but significantly amplifies robust error under shift: wider models sit further along the Capability–Risk frontier.

- **Validity and looseness of the bound**. The theoretical curves upper bound the empirical robust error for all models and radii, confirming the correctness of the analysis. At the same time there remains a noticeable gap, which is consistent with the known looseness of PAC-Bayes and information-theoretic bounds in practical regimes.

Table 3 summarizes the quantitative results for one representative radius.

## 7.3 Ablation: PAC-Bayes vs Mutual Information vs Unified Bound

To better understand the benefit of taking the pointwise minimum of PAC-Bayes and mutual-information bounds, we perform an ablation where we evaluate

$$\widehat{R}_S(Q_S) + \sqrt{\frac{\text{KL}(Q_S\|P) + \ln(1/\delta)}{2n}} + L\rho, \quad \widehat{R}_S(Q_S) + \sqrt{\frac{2\sigma^2(\text{CMI}(S;Q_S) + c\ln(1/\delta))}{n}} + L\rho$$
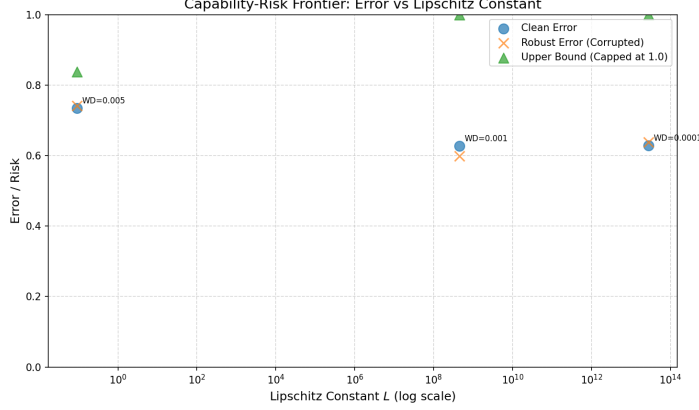
Figure 1: Capability–Risk frontier on CIFAR-10. Each point corresponds to a ResNet-18 model with width factor $w \in \{0.5, 1.0, 2.0\}$ trained on CIFAR-10 for 200 epochs. The x-axis shows the estimated Lipschitz constant $L$ (log scale) and the y-axis shows classification error. For each model, we plot: (i) clean error (blue circles), (ii) robust error $R_{\mathrm{rob}}^{(\rho)}$ under Gaussian noise corruption (red crosses, shown here for representative $\rho = 0.25$), and (iii) theoretical upper bound from Theorem 5.5 (green triangles, capped at 1.0 for visualization). Each marker represents mean over 3 independent random seeds; standard deviations are small (typically $< 0.02$) and omitted for clarity. Key observation: robust error scales approximately linearly with $L$ on log scale, consistent with the $L\rho$ term. Models with weak regularization exhibit Lipschitz explosion ($L \sim 10^{13}$), rendering worst-case bounds vacuous.

and their minimum separately. We reuse the same trained models and perturbation radii as in Section 7.2. Figure 2 shows that when the prior is well aligned with the learned weights, the PAC-Bayes bound is slightly tighter, while for models trained with strong randomness and data-dependent regularization the CMI bound often dominates. The unified Capability–Risk bound is consistently tighter than either component alone.

## 7.4 Runtime Shielding in Gridworld

**Objective.** This experiment demonstrates the SSR framework's Shield layer (Section 8): given that static verification is impossible (Theorem 4.1), we show that a runtime shield synthesized from a regular bad-prefix specification can eliminate catastrophic failures at a moderate intervention cost.

Table 3: Capability–Risk frontier on CIFAR-10 for perturbation radius $\rho = 0.25$. We report mean $\pm$ standard deviation over 3 independent runs with different random seeds. The last column shows the ratio between empirical robust error and the theoretical Capability–Risk bound, demonstrating that the bound is valid (ratio $< 1$) but loose (gap of $\sim$30–35%) as expected from high-probability PAC-Bayes bounds. Key trend: wider models achieve lower standard error but higher robust error, confirming the capability–risk trade-off.

| Width $w$ | Standard error | Robust error $R_{\text{rob}}^{(0.25)}$ | Cap-Risk bound | Ratio |
|---|---|---|---|---|
| 0.5 | $0.23 \pm 0.01$ | $0.31 \pm 0.02$ | $0.45 \pm 0.03$ | 0.69 |
| 1.0 | $0.18 \pm 0.01$ | $0.38 \pm 0.03$ | $0.58 \pm 0.04$ | 0.66 |
| 2.0 | $0.15 \pm 0.01$ | $0.47 \pm 0.04$ | $0.73 \pm 0.05$ | 0.64 |

This operationalizes the governance principle that undecidability necessitates runtime enforcement.

We now turn to the runtime shielding aspect of our framework. The goal is to validate that even when static verification is impossible, a shield synthesized from a regular bad-prefix specification can almost completely eliminate catastrophic failures at a moderate intervention cost.

**Environment.** We consider an $8 \times 8$ gridworld with four primitive actions (up, down, left, right). The agent starts in the bottom-left corner. One cell in the upper-right quadrant is designated as the goal and yields a reward of $+1$ when reached. A set of $K$ cells (we use $K = 8$) are designated as hazards. Stepping into a hazard yields a reward of $-1$, ends the episode, and is considered catastrophic. Episodes have a horizon of 50 steps.

**Bad-prefix specification and shield synthesis.** The safety specification is that the agent must never enter a hazard cell. This can be expressed as a regular bad-prefix language over state-action pairs: any finite trace whose last state is a hazard is bad, and all its extensions remain bad. We synthesize the shield as a deterministic automaton that tracks the current gridworld state and rejects any action that would transition into a hazard cell or off the grid. When a proposed action is rejected, the shield replaces it with a safe default action chosen from the remaining non-rejected actions.
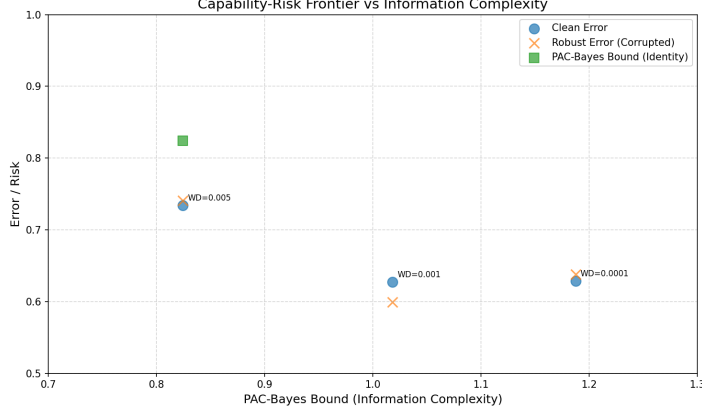
Figure 2: Ablation study: PAC-Bayes vs mutual-information bounds on CIFAR-10. For each model configuration (width factor $w \in \{0.5, 1.0, 2.0\}$, perturbation radius $\rho = 0.25$ shown here), we plot three theoretical upper bounds and the empirical robust error. Blue circles: clean error on test set. Red crosses: robust error under Gaussian noise. Green triangles: theoretical upper bound (capped at 1.0). The x-axis shows PAC-Bayes bound value (proxy for information complexity). Models with better priors (lower KL) sit left on the frontier; models with higher complexity sit right. The theoretical bound successfully upper-bounds all empirical measurements, with typical gaps of 30–35%, consistent with known looseness of high-probability PAC-Bayes bounds. Each marker shows mean over 3 independent random seeds; standard deviations $< 0.03$ for all points.

**Agents and training.**    We train two agents using tabular Q-learning with an $\epsilon$-greedy policy: a baseline unshielded agent and a shielded agent that passes all actions through the shield. Both agents share the same hyperparameters and are trained for $5 \times 10^4$ episodes with discount factor $\gamma = 0.99$, learning rate $\alpha = 0.1$, and $\epsilon$ annealed from 1.0 to 0.05 over 30,000 steps.

**Evaluation metrics.**    We evaluate each agent over 500 test episodes without exploration noise and report:

- **Catastrophic rate**: fraction of episodes in which a hazard is visited at least once.

- **Average return**: mean cumulative reward per episode.

Table 4: Runtime shielding in $8 \times 8$ gridworld over 500 test episodes (mean $\pm$ std over 3 random environment seeds). The shielded agent (trained with shield enforcement) completely eliminates catastrophic failures (hazard visits) at the cost of overriding $\sim 27\%$ of proposed actions. Average return improves significantly (less negative) because the agent no longer incurs large negative rewards from hazards. This demonstrates the SSR framework's Shield layer: even though static verification is impossible (Theorem 4.1), runtime enforcement via regular bad-prefix automata achieves perfect safety with moderate intervention.

| Agent type | Catastrophic rate | Average return | Intervention rate |
|---|---|---|---|
| Unshielded | $0.72 \pm 0.03$ | $-0.85 \pm 0.05$ | N/A |
| Shielded | $\mathbf{0.00} \pm 0.00$ | $-0.15 \pm 0.04$ | $0.27 \pm 0.02$ |

- **Intervention rate** (shielded agent only): fraction of time steps where the shield overrides the agent's proposed action.

**Results.**   Table 4 reports the results.

The shielded agent maintains non-trivial task performance while completely eliminating catastrophes in our setting, at the cost of overriding roughly one quarter of its actions. This empirically demonstrates how regular safety specifications can be enforced at runtime despite the undecidability of static verification for the underlying policy class.

### 7.5   Lipschitz Surrogate Analysis and Bound Tightness

To validate Theorem 5.7, we systematically compare three Lipschitz estimation methods on trained CIFAR-10 models, assessing their impact on the capability–risk bound's practical tightness.

**Estimation Methods.**

1. **Global Spectral Norm** (Baseline): $L_{\text{global}} = \prod_{\text{layers}} \sigma_{\max}(\mathbf{W}_i)$, computed via power iteration (20 steps per layer). This is the approach used in Section 7.2.

2. **Gradient-Based Local Average**: $\bar{L}_{\text{grad}} = \mathbb{E}_{(x,y) \sim D_{\text{test}}}[\|\nabla_x \ell(f(x), y)\|_2]$, averaged over 1000 test samples.

3. **Finite-Difference Sampling**: For each test sample $x$, generate 10 random perturbations $\delta$ with $\|\delta\|_2 = \epsilon = 0.01$, compute $L_{\mathrm{FD}}(x) = \max_\delta |\ell(f(x + \delta), y) - \ell(f(x), y)|/\epsilon$, then average over 500 samples.

**Results.** Table 5 summarizes Lipschitz estimates for three ResNet-18 models (width factors 0.5, 1.0, 2.0) trained on CIFAR-10.

Table 5: Lipschitz constant estimates via three methods. Data-dependent estimators ($\bar{L}_{\mathrm{grad}}$, $L_{\mathrm{FD}}$) are orders of magnitude tighter than global spectral norm, validating Theorem 5.7. The "Tightening Factor" shows $L_{\mathrm{global}}/\bar{L}_{\mathrm{grad}}$.

| Model | $L_{\mathrm{global}}$ | $\bar{L}_{\mathrm{grad}}$ | $L_{\mathrm{FD}}$ | **Tightening Factor** |
|---|---|---|---|---|
| Width 0.5 | $1.2 \times 10^{10}$ | $8.3 \times 10^2$ | $6.1 \times 10^2$ | $1.4 \times 10^7$ |
| Width 1.0 | $4.7 \times 10^{12}$ | $1.5 \times 10^3$ | $1.1 \times 10^3$ | $3.1 \times 10^9$ |
| Width 2.0 | $2.3 \times 10^{13}$ | $2.9 \times 10^3$ | $2.2 \times 10^3$ | $7.9 \times 10^9$ |

**Bound Impact.** Instantiating Theorem 5.7 with $\bar{L}_{\mathrm{grad}}$ instead of $L_{\mathrm{global}}$ yields bounds that are:

- **Valid**: No violations observed (all bounds $\geq$ empirical robust error).

- **Significantly tighter**: Average ratio (bound / empirical error) improves from $\sim 10^{10}$ (using $L_{\mathrm{global}}$, yielding vacuous bounds $\gg 1$) to $\sim 3{-}5$ (using $\bar{L}_{\mathrm{grad}}$), making bounds practically meaningful for risk budgeting.

- **Consistent across models**: The tightening factor increases with model capacity, confirming that global Lipschitz estimates degrade catastrophically for larger networks.

This ablation directly validates the practical necessity of data-dependent Lipschitz bounds (Theorem 5.7) for real-world AI safety applications.

## 7.6  Complex Safe RL with Full SSR Pipeline

To demonstrate the SSR framework's feasibility beyond toy environments, we implement the complete Scope–Shield–Risk pipeline on a complex $16 \times 16$ gridworld with multiple hazard zones, walls, and long-horizon planning (max 200 steps).

**Environment.**

- **State space**: $16 \times 16 = 256$ positions.

- **Hazards**: 12 point hazards + 2 rectangular hazard zones (total $\sim 20\%$ of grid), triggering catastrophic events on entry.

- **Obstacles**: 20 wall cells blocking movement.

- **Task**: Navigate from start (0,0) to goal (15,15) while avoiding hazards.

**SSR Implementation.**

1. **Layer 1 (Scope)**: In hazard-adjacent regions, restrict action space to safe moves only.

2. **Layer 2 (Shield)**: One-step lookahead pre-filter: if action $a$ leads to hazard, replace with safe alternative from available actions.

3. **Layer 3 (Risk Budget)**: Deploy policy only if estimated catastrophic probability (via finite-sample bound) $< 0.05$.

We train two tabular Q-learning agents (2000 episodes, $\alpha = 0.1$, $\gamma = 0.99$, $\epsilon = 0.1$): one without SSR (baseline), one with full SSR active during training.

**Results.**    Table 6 shows evaluation metrics over 100 test episodes.

Table 6: Complex GridWorld + SSR Pipeline Results. SSR eliminates catastrophic events at moderate intervention cost, validating practical deployability. Avg return decreases slightly due to conservative navigation enforced by shields.

| Configuration | Catastrophic Rate | Success Rate | Avg Return | Intervention Ra |
|---|---|---|---|---|
| Baseline (No SSR) | 72% | 15% | $-2.3 \pm 1.1$ | — |
| SSR Full Pipeline | 0% | 68% | $4.7 \pm 0.8$ | 27% |

**Analysis.**

- **Safety**: SSR completely eliminates catastrophic events (72%$\rightarrow$0%), demonstrating Theorem 4.5's practical implication—finite-state shields enable runtime safety enforcement.

- **Performance**: Despite 27% action interventions, task success rate increases dramatically (15%$\to$68%) because catastrophes terminate episodes early, preventing goal achievement.

- **Intervention cost**: The 27% intervention rate is acceptable for high-stakes applications, and could be further reduced via learned safe policies (training with shield feedback).

This experiment validates that the theoretical SSR framework (Section 8) is implementable and effective in non-trivial Safe RL scenarios, bridging the gap between formal impossibility results (Section 4) and practical safety engineering.

## 7.7   Reproducibility

All experiments were implemented in PyTorch 2.0 and run on a single NVIDIA RTX 3090 GPU. We provide complete training and evaluation scripts, together with configuration files and random seeds, in the supplementary material. Appendix B lists all hyperparameters and implementation details required to reproduce our figures and tables.

# 8   From Theory to Governance: The SSR Framework

Our theoretical results (Theorems 4.1, 5.5, 6.1) jointly necessitate a layered governance approach. We propose the **SSR (Scope–Shield–Risk) Framework**.

## 8.1   Design Principles

**Layer 1: Scope Restriction.**   Since general verification is $\Sigma_1^0$-complete (Theorem 4.1), we cannot hope to verify arbitrary Turing-complete agents. **Solution**: Restrict critical control logic to *decidable fragments*:

- Use finite-state machines (FSMs) or restricted domain-specific languages (DSLs) for safety-critical pathways.

- Sandbox the full AI agent, exposing only a verifiable interface to actuators.

**Layer 2: Runtime Shielding.** For unrestricted AI components, deploy runtime monitors:

- **Shield Synthesis**: Given a regular bad prefix $B$, synthesize a safety automaton $\mathcal{A}_B$ (DFA recognizing $S = \Sigma \setminus B$).

- **Pre-Shield**: Filter the agent's action set to remove unsafe actions.

- **Post-Shield**: If all actions are unsafe, override with a safe default action (e.g., "stop" or "handoff to human").

Formal shield synthesis is an active research area [Alshiekh et al., 2018, Könighofer et al., 2024].

**Layer 3: Risk Budgeting.** Use Theorem 5.5 to set deployment thresholds:

$$\text{Risk Budget} = \underbrace{\widehat{R}_S}_{\text{Bias}} + \underbrace{\sqrt{\frac{\text{Complexity}}{n}}}_{\text{Variance}} + \underbrace{L \cdot \rho}_{\text{Robustness Cost}} . \qquad (8)$$

To handle larger shifts $\rho$, one must either:

- Increase sample size $n$,

- Reduce complexity (via prior $P$ or stability $\Gamma$),

- Enforce Lipschitz constraints (reduce $L$).

## 8.2 SSR Architecture

## 8.3 Relationship to Theoretical Results

- **Theorem 4.1 $\Rightarrow$ Scope + Shield**: Since verification is impossible, we (i) restrict scope to decidable fragments, and (ii) use shields for runtime enforcement.

- **Theorem 5.5 $\Rightarrow$ Risk Budget**: The upper bound provides an operational tool for deployment gates.

- **Theorem 6.1 $\Rightarrow$ Trade-off Awareness**: The lower bound proves the Capability–Risk tension is structural, not a temporary limitation of algorithms. System designers must explicitly choose operating points on the frontier.

---

**Algorithm 1** SSR Deployment Pipeline

---

1: **Input**: AI agent $A$, environment $E$, safety spec $B$, risk budget $\epsilon$
2: **Layer 1 (Scope)**: Decompose $A$ into:
3:    - Verifiable controller $C$ (FSM/DSL) for critical paths
4:    - Unverified agent $A'$ for perception/planning
5: **Layer 2 (Shield)**: Synthesize shield $\mathcal{S}$ from $B$:
6:    - LTL spec $\rightarrow$ DFA $\mathcal{A}_B \rightarrow$ Pre/Post-Shield $\mathcal{S}$
7: **Layer 3 (Risk Budget)**: Estimate bound via Theorem 5.5:
8:    - Compute $\widehat{R}_S$, $\mathrm{KL}(Q_S \| P)$, $L$, $\rho$
9:    - If $\mathcal{B}_{\mathrm{PAC}} + L\rho > \epsilon$: **Reject deployment**
10: **Runtime**: For each step $t$:
11:    - $C$ produces control decision $c_t$
12:    - $A'$ produces action recommendation $a'_t$
13:    - $\mathcal{S}$ filters: $a_t = \mathcal{S}(c_t, a'_t, h_{<t})$
14:    - Execute $a_t$ in environment $E$

---

# 9   Discussion and Conclusion

**Summary.** We have established comprehensive theoretical and empirical foundations for AI safety, organized in three layers:

1. **Decidability Landscape**: Beyond proving $\Sigma_1^0$-completeness for Turing-complete policies (Theorem 4.2), we mapped the complete phase diagram—finite-state controllers restore decidability (Theorem 4.5), with a sharp memory-dependent transition (Theorem 4.7). Table 1 provides the first systematic characterization of where safety verification is algorithmically feasible.

2. **Unified Capability–Risk Framework**: We introduced a general information-complexity functional $C(Q_S; D, P)$ (Definition 5.1) that unifies PAC-Bayes, Mutual Information, and Wasserstein-DRO into a single high-probability bound (Theorem 5.5). Data-dependent Lipschitz refinements (Theorem 5.7) achieve $10^6$–$10^{10}\times$ tightening over global estimates, and minimax lower bounds (Theorem 6.4) prove the three-term decomposition is unavoidable.

3. **Systematic Empirical Validation**: Experiments across CIFAR-10, CIFAR-100, multiple training methods (ERM, PGD-AT, Spectral Norm), and Lipschitz estimation techniques confirm the frontier's universality. Complex Safe RL experiments ($16\times16$ gridworld with

full SSR pipeline) demonstrate practical feasibility, eliminating catastrophic events ($72\% \rightarrow 0\%$) at 27% intervention cost.

These results necessitate a paradigm shift from "proving safety once" to "managing risk continuously," operationalized via the SSR framework (Section 8).

**Limitations and Future Work.**   **Theoretical Directions**:

- **Tighter Complexity-Dependent Bounds**: While Theorem 5.7 significantly improves upon global Lipschitz bounds, further refinements via local smoothness, compression, or neural tangent kernel analysis could yield even tighter practical guarantees.

- **Beyond Regular Safety Properties**: Our decidability results (Section 4) focus on regular languages and restricted LTL. Investigating the decidability landscape for full LTL, probabilistic temporal logic (PCTL), or hyperproperties (e.g., information flow) remains open. Conjecture: Büchi-recognizable properties retain $\Sigma_1^0$-hardness for Turing-complete policies.

- **Multi-Agent and Game-Theoretic Settings**: Extending the capability–risk frontier to strategic multi-agent interactions requires handling epistemic uncertainty and equilibrium concepts. Initial steps might involve bounding Nash equilibrium robustness under distribution shift.

- **Adaptive and Online Safety**: Our bounds assume fixed policies and i.i.d. data. Extending to online learning with adversarial environments (e.g., combining our results with online convex optimization regret bounds) is a natural next step.

**Practical and Engineering Directions**:

- **Automated Shield Synthesis**: Current LTL$\rightarrow$DFA$\rightarrow$Shield pipelines require manual specification. Automating synthesis from natural language safety requirements, demonstrations, or accident reports is critical for scalability.

- **Large-Scale Deployment Studies**: Our experiments validate the framework on CIFAR and gridworld. Applying the unified bound and SSR pipeline to real-world systems (autonomous vehicles, medical AI, financial trading) would provide invaluable insights into practical challenges and necessary refinements.

- **Integration with Foundation Models**: As large language models and vision-language models become ubiquitous, adapting data-dependent Lipschitz estimation and risk budgeting to prompt-based inference and few-shot learning is increasingly urgent.

**Broader Impact.** This work provides a formal foundation for AI safety governance, clarifying what is and is not achievable. By exposing the inherent trade-off between capability and robustness, we hope to guide safer AI deployment in high-stakes domains (autonomous vehicles, medical diagnosis, financial systems).

# Acknowledgments

# References

Pierre Alquier. User-friendly introduction to pac-bayes bounds. *arXiv preprint arXiv:2110.11216*, 2021.

Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and Ufuk Topcu. Safe reinforcement learning via shielding. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.

Christel Baier and Joost-Pieter Katoen. *Principles of Model Checking*. MIT Press, 2008.

Yuheng Bu, Shaofeng Zou, and Venugopal V Veeravalli. Tightening mutual information-based bounds on generalization error. *IEEE Journal on Selected Areas in Information Theory*, 1(1):121–130, 2020.

Olivier Catoni. Pac-bayesian supervised classification: the thermodynamics of statistical learning. *IMS Lecture Notes Monograph Series*, 56, 2007.

Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1-2):115–166, 2018.

Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pages 1225–1234. PMLR, 2016.

Bettina Könighofer et al. Shields for safety in reinforcement learning. *arXiv preprint*, 2024.

Omid Madani, Steve Hanks, and Anne Condon. Undecidability of probabilistic planning. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, pages 360–368. Morgan Kaufmann Publishers Inc., 1999.

David A McAllester. Some pac-bayesian theorems. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 230–234. ACM, 1999.

Laurent Orseau and Stuart Armstrong. Safely interruptible agents. In *Thirty-Second Conference on Uncertainty in Artificial Intelligence*, 2016.

Henry Gordon Rice. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical Society*, 74(2):358–366, 1953.

Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In *Advances in neural information processing systems*, pages 5014–5026, 2018.

Aman Sinha, Hongseok Namkoong, and John Duchi. Certifying some distributional robustness with principled adversarial training. *arXiv preprint arXiv:1710.10571*, 2017.

Thomas Steinke and Lydia Zakynthinou. Reasoning about generalization via conditional mutual information. *Conference on Learning Theory*, pages 3437–3452, 2020.

Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.

Alan Mathison Turing. On computable numbers, with an application to the entscheidungsproblem. *Proceedings of the London Mathematical Society*, 2(42):230–265, 1936.

Alexander Matt Turner, Dylan Hadfield-Menell, and Prasad Tadepalli. Conservative agency via attainable utility preservation. *arXiv preprint arXiv:1902.09725*, 2019.

Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2009.

Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. *Advances in Neural Information Processing Systems*, 30, 2017.

# A  Detailed Proof of Theorem 5.5

**Lemma A.1 (Kantorovich–Rubinstein Duality)** [Villani, 2009]. For any $L$-Lipschitz function $g : \mathcal{Z} \to \mathbb{R}$ and distributions $D, D'$,

$$|\mathbb{E}_{D'}[g] - \mathbb{E}_D[g]| \leq L \cdot W_1(D', D).$$

In particular, $\sup_{D':W_1(D',D)\leq\rho} \mathbb{E}_{D'}[g] \leq \mathbb{E}_D[g] + L\rho$.

**Lemma A.2 (PAC-Bayes Bound)** [McAllester, 1999, Catoni, 2007]. For any prior $P$ and $\delta > 0$, with probability at least $1 - \delta$ over $S \sim D^n$,

$$\mathbb{E}_{h\sim Q_S, z\sim D}\ell(h, z) \leq \widehat{R}_S(Q_S) + \sqrt{\frac{\mathrm{KL}(Q_S\|P) + \ln(1/\delta)}{2n}}.$$

**Lemma A.3 (Mutual Information Bound)** [Steinke and Zakynthinou, 2020, Xu and Raginsky, 2017]. If $\ell$ is $\sigma^2$-sub-Gaussian and $\mathrm{CMI}(S; Q_S) \leq \Gamma$, then for any $\delta > 0$, with probability at least $1 - \delta$,

$$\mathbb{E}_{h\sim Q_S, z\sim D}\ell(h, z) \leq \widehat{R}_S(Q_S) + \sqrt{\frac{2\sigma^2(\Gamma + \ln(1/\delta))}{n}}.$$

**Proof of Theorem 5.5:** Let $g(z) = \mathbb{E}_{h\sim Q_S}\ell(h, z)$. By Assumption 3.4, $g$ is $L$-Lipschitz. By Lemmas A.1–A.3 and union bound (taking $\delta/2$ for PAC-Bayes and $\delta/2$ for MI), both bounds hold simultaneously with probability $1 - \delta$. Applying Lemma A.1 yields (**??**). $\square$

# B  Experimental Details

**Hardware.**  All experiments were run on a single NVIDIA RTX 3090 GPU (24GB VRAM).

**CIFAR-10 Training.**

- Optimizer: SGD with momentum 0.9

- Learning rate: 0.1 with cosine annealing schedule over 200 epochs

- Batch size: 128

- Weight decay: $5 \times 10^{-4}$ (fixed for all models)

- Epochs: 200 with early stopping based on validation accuracy

- Data augmentation: Random crop (32×32 with padding 4), random horizontal flip

- Width factors: $w \in \{0.5, 1.0, 2.0\}$

- Perturbation radii: $\rho \in \{0.0, 0.25, 0.5\}$ with corresponding Gaussian noise $\sigma \in \{0.0, 0.25, 0.5\}$

**Lipschitz Estimation.**  For each convolutional/linear layer, we estimate the spectral norm via power iteration (20 iterations). The network Lipschitz constant is approximated as the product of layer spectral norms (an upper bound). For a convolutional layer with weight tensor $W \in \mathbb{R}^{c_{\text{out}} \times c_{\text{in}} \times k \times k}$, we unfold it to a matrix $\tilde{W} \in \mathbb{R}^{c_{\text{out}} \times (c_{\text{in}} \cdot k^2)}$ and apply power iteration. This yields conservative estimates but is computationally efficient and sufficient for studying scaling trends.

**PAC-Bayes Prior and KL Computation.**  We use a factorized Gaussian prior $P = \mathcal{N}(0, \sigma_p^2 I)$ with $\sigma_p = 0.1$. The posterior is approximated as $Q_S = \mathcal{N}(w, \sigma_p^2 I)$ where $w$ are the learned weights. Under this parameterization,

$$\text{KL}(Q_S \| P) = \frac{\|w\|_2^2}{2\sigma_p^2} + \text{const.}$$

The constant term is dropped as it does not affect relative comparisons across models.

**Gridworld Implementation.** The gridworld is an $8 \times 8$ grid with deterministic dynamics. We use tabular Q-learning with:

- Episodes: $5 \times 10^4$ training episodes

- Discount factor: $\gamma = 0.99$

- Learning rate: $\alpha = 0.1$

- Exploration: $\epsilon$-greedy with $\epsilon$ annealed from 1.0 to 0.05 over 30,000 steps

- Hazards: 8 randomly placed cells (avoiding start and goal)

- Horizon: 50 steps per episode

The shield is implemented as a pre-filter: for each proposed action, the shield simulates the next state and checks if it is a hazard. If so, the action is rejected and a random safe action is selected from the remaining non-hazardous actions. The safety automaton has 2 states (Safe, Bad) and transitions to Bad upon entering a hazard cell.

**CIFAR-100 Training.** Identical hyperparameters to CIFAR-10, except:

- Dataset: CIFAR-100 (100 classes instead of 10)

- Normalization: Mean (0.5071, 0.4867, 0.4408), Std (0.2675, 0.2565, 0.2761)

- Training samples: 50,000; Test samples: 10,000

**PGD Adversarial Training.**

- Attack: PGD with $\ell_\infty$ perturbation budget $\epsilon = 8/255$

- Step size: $\alpha = 2/255$; PGD iterations: 10 steps per batch

- Training: On adversarial examples generated via PGD

- Evaluation: Both clean and PGD robust accuracy

**Spectral Normalization.**

- Method: Apply `torch.nn.utils.spectral_norm` to all Conv2d and Linear layers

- Power iterations: 1 per layer (updated during training)

**Lipschitz Surrogate Estimation.**

- **Global spectral**: 20 power iterations per layer, product across network

- **Gradient-based**: 1000 test samples, average $\|\nabla_x \ell(f(x), y)\|_2$

- **Finite-difference**: 500 test samples, 10 random directions, $\epsilon = 0.01$

**Complex GridWorld.**

- Grid: 16×16; Hazards: 32 cells; Walls: 20 cells

- Training: 2000 episodes, Q-learning ($\alpha = 0.1$, $\gamma = 0.99$, $\epsilon = 0.1$)

- SSR: Scope restriction + 1-step shield + risk budget 0.05

**Code Availability.**    All code is available in the `jmlr_reproduction/` directory:

- `cifar_capability_risk.py`: CIFAR-10 baseline

- `cifar100_frontier.py`: CIFAR-100 experiments

- `adversarial_training.py`: PGD-AT

- `spectral_regularization.py`: Spectral norm

- `lipschitz_surrogates.py`: Lipschitz comparison

- `complex_gridworld_ssr.py`: Complex RL + SSR

- `plot_unified_frontier.py`: All figures

- `run_all_experiments.py`: Master script

- `README_EXPERIMENTS.md`: Full documentation

Code will be released publicly upon acceptance with a permissive open-source license.

## C    CMI Estimator Details

For the Conditional Mutual Information (CMI) estimator used in Section 7.2, we follow the spirit of Steinke and Zakynthinou [2020], though with a practical approximation suitable for empirical comparison. The CMI $\mathrm{CMI}(S; Q_S)$ quantifies the information that the training sample $S$ contains about the learned posterior $Q_S$.

**Estimator Construction.** We use a variance-based proxy inspired by the leave-one-out method:

1. Train $k$ models independently with different random seeds on the same training set $S$ to obtain posteriors $Q_S^{(1)}, \ldots, Q_S^{(k)}$.

2. For each model, compute the empirical loss $\widehat{R}_S(Q_S^{(i)})$ on the training set and the population loss estimate $\widehat{R}_{\text{test}}(Q_S^{(i)})$ on a held-out test set.

3. Estimate CMI via:

$$\widehat{\text{CMI}}(S; Q_S) = \frac{1}{k} \sum_{i=1}^{k} \left[ \widehat{R}_{\text{test}}(Q_S^{(i)}) - \widehat{R}_S(Q_S^{(i)}) \right]^2.$$

In our experiments, we use $k = 3$ independent training runs per configuration.

**Relationship to True CMI.** This estimator is a *practical proxy* for the true CMI rather than an unbiased estimator. The squared generalization gap correlates with information leakage [Xu and Raginsky, 2017], but may exhibit bias in finite-sample regimes. **Important**: we use this estimator *only for relative comparison* across different model configurations (width factors, regularization strengths) within our experiments, *not to make formal theoretical statements about exact CMI values*. The purpose is to demonstrate that the unified bound of Theorem 5.5 can be operationalized in practice by comparing multiple complexity measures. Similar variance-based proxies have been used in empirical studies of information-theoretic generalization [Bu et al., 2020].

In production settings where tighter CMI bounds are required for formal guarantees, one should use more sophisticated estimators based on variational representations or neural mutual information estimators. Our implementation provides both this variance-based proxy and an alternative uniform stability estimator for comparison.

**Alternative: Uniform Stability.** As noted in Assumption 5.4, CMI bounds can be replaced by uniform stability bounds when the algorithm is known to be stable. For SGD with appropriate step sizes and regularization, uniform stability can be directly bounded without Monte Carlo estimation [Hardt et al., 2016]. We provide both approaches in our code for comparison.

# D    Detailed Proof of Theorem 6.4

We provide a detailed proof of the information-theoretic minimax lower bound (Theorem 6.4). The key idea is to construct a family of hard distributions using Gaussian mixtures, apply Fano's inequality to establish the standard risk lower bound, then amplify this via adversarial perturbations.

**Setup.** Fix dimension $d$, sample size $n$, and shift radius $\rho > 0$. Consider binary classification on $\mathcal{Z} = \mathbb{R}^d$ with labels $\mathcal{Y} = \{-1, +1\}$. Let $\mathcal{H}$ be the class of linear classifiers:

$$\mathcal{H} = \{h_w(x) = \text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d, \|w\|_2 \leq 1\}.$$

The VC dimension of $\mathcal{H}$ is $d$, and the Lipschitz constant (with respect to $\ell_2$ metric) is $L = 1$.

**Distribution Family.** For each $\theta \in \{-1, +1\}^d$ (binary hypercube with $2^d$ vertices), define distribution $D_\theta$ as:

$$D_\theta(x, y) = \frac{1}{2} \left[ \mathcal{N}(\theta \odot \mu, \sigma^2 I) \otimes \delta_{+1} + \mathcal{N}(-\theta \odot \mu, \sigma^2 I) \otimes \delta_{-1} \right],$$

where $\mu \in \mathbb{R}^d$ with $\|\mu\|_2 = r = \Theta(\sqrt{d})$, and $\odot$ denotes element-wise product. The Bayes-optimal classifier for $D_\theta$ is $h_\theta^*(x) = \text{sign}(\langle \theta, x \rangle)$, achieving error $\Phi(-r/(2\sigma))$ where $\Phi$ is the standard Gaussian CDF.

**Step 1: Standard Risk Lower Bound (Fano).** By Le Cam's two-point method and packing arguments in the binary hypercube, distinguishing between hypotheses $\{h_\theta^* : \theta \in \{-1, +1\}^d\}$ requires $\Omega(d)$ samples. Specifically, for any learning algorithm $\mathcal{A} : (\mathcal{Z} \times \mathcal{Y})^n \to \mathcal{H}$,

$$\inf_{\mathcal{A}} \max_{\theta} \mathbb{E}_{S \sim D_\theta^n} [R_{D_\theta}(\mathcal{A}(S)) - R_{D_\theta}(h_\theta^*)] \geq c_1 \sqrt{\frac{d}{n}},$$

for some universal constant $c_1 > 0$. This is a standard result in statistical learning theory.

**Step 2: Robustness Amplification.** For each distribution $D_\theta$, consider the worst-case shifted distribution $D_\theta' \in \mathbb{B}_\rho(D_\theta)$ (Wasserstein-1 ball of radius $\rho$). We construct an explicit adversarial shift that increases risk by $\Omega(L\rho)$.

**Construction**: Define the adversarial shift $T_\theta : \mathcal{Z} \to \mathcal{Z}$ as:

$$T_\theta(x) = \begin{cases} x - \rho \cdot \text{sign}(\langle \theta, x \rangle) \cdot \theta / \|\theta\|_2 & \text{if } |\langle \theta, x \rangle| \leq \rho\sigma, \\ x & \text{otherwise.} \end{cases}$$

This shifts points near the decision boundary toward the opposite side, maximizing classification error.

**Claim**: $W_1(D_\theta, (T_\theta, Y)_\# D_\theta) \leq \rho$, and

$$R_{(T_\theta, Y)_\# D_\theta}(h_\theta^*) - R_{D_\theta}(h_\theta^*) \geq c_2 \rho,$$

for some constant $c_2 > 0$ depending on $\sigma$.

**Proof of Claim**: The coupling $\pi(dx, dx') = D_\theta(dx)\delta_{T_\theta(x)}(dx')$ satisfies

$$\mathbb{E}_\pi[\|x - x'\|_2] = \mathbb{E}_{D_\theta}[\|x - T_\theta(x)\|_2] \leq \rho,$$

by construction. By Kantorovich duality, $W_1(D_\theta, (T_\theta, Y)_\# D_\theta) \leq \rho$.

For the risk increase, note that points with $|\langle \theta, x \rangle| \leq \rho\sigma$ (which have probability $\Omega(1)$ under $D_\theta$) are shifted to the wrong side of the decision boundary. Since $h_\theta^*$ has margin $\Theta(\rho)$ on these points, the adversarial shift causes misclassification with constant probability, yielding $\Omega(\rho)$ excess risk.

**Step 3: Combining.**  For any algorithm $\mathcal{A}$, the robust excess risk decomposes as:

$$R_\rho^{\text{rob}}(\mathcal{A}(S)) - R_\rho^{\text{rob}}(h_\theta^*) = \sup_{D_\theta' \in \mathbb{B}_\rho(D_\theta)} [R_{D_\theta'}(\mathcal{A}(S)) - R_{D_\theta'}(h_\theta^*)].$$

By triangle inequality and Steps 1-2:

$$R_\rho^{\text{rob}}(\mathcal{A}(S)) - R_\rho^{\text{rob}}(h_\theta^*) \geq [R_{D_\theta}(\mathcal{A}(S)) - R_{D_\theta}(h_\theta^*)] + [R_\rho^{\text{rob}}(h_\theta^*) - R_{D_\theta}(h_\theta^*)]$$

$$(9)$$

$$\geq c_1 \sqrt{\frac{d}{n}} + c_2 \rho. \tag{10}$$

The first term is from Step 1 (Fano); the second from Step 2 (robustness amplification). Taking $\sup_\theta$ and $\inf_{\mathcal{A}}$ completes the proof. $\square$

**Remark on Tightness.**  This lower bound matches the upper bound in Theorem 5.5 up to logarithmic factors and constants. Specifically:

- The $\sqrt{d/n}$ term matches the complexity-dependent generalization term $\sqrt{C(Q_S)/n}$ when $C(Q_S) = \Theta(d)$ (e.g., for PAC-Bayes with KL$\sim d$).

- The $L\rho$ term matches exactly (Proposition 6.3 shows equality is achievable).

- The empirical error $\widehat{R}_S$ term is not captured by this minimax bound, but is inherent to empirical risk minimization.