
LAW-STRENGTH FRONTIERS AND A NO-FREE-LUNCH RESULT FOR LAW-SEEKING REINFORCEMENT LEARNING ON VOLATILITY LAW MANIFOLDS

Zhang Jian'an

Guanghua School of Management, Peking University
Peking University
Beijing, China
2501111059@stu.pku.edu.cn

ABSTRACT

We study reinforcement learning (RL) on volatility surfaces through the lens of *Scientific AI*: can axiomatic market laws, enforced as soft penalties on a learned world model, reliably align high-capacity RL agents with no-arbitrage structure, or do they merely induce Goodhart-style exploitation of model artefacts?

Starting from classical static no-arbitrage conditions for implied volatility, we construct a finite-dimensional convex *volatility law manifold* of admissible total-variance surfaces, together with a metric-based *law-penalty functional* and a domain-agnostic *Graceful Failure Index* (GFI) that normalizes law degradation under shocks. A synthetic generator produces trajectories that are exactly law-consistent, while a recurrent neural world model is trained without law regularization and therefore predicts surfaces that deviate from the law manifold in structured ways.

On top of this testbed we introduce a *Goodhart decomposition* of reward, $r = r^{\mathcal{M}} + r^{\perp}$, where $r^{\mathcal{M}}$ is the on-manifold component and r^{\perp} is *ghost arbitrage* arising purely from off-manifold prediction errors. We prove three flagship results: (i) a *ghost-arbitrage incentive* theorem showing that naive PPO-type RL is structurally driven to increase $\mathbb{E}[r^{\perp}]$ whenever ghost arbitrage is available; (ii) a *law-strength trade-off* theorem establishing that increasing the weight on law penalties inevitably worsens P&L beyond a quantifiable threshold; and (iii) a *no-free-lunch* theorem stating that, under a law-consistent world model and a law-aligned structural class of strategies, unconstrained law-seeking RL cannot Pareto-dominate structural baselines on P&L, law penalties, and GFI.

Empirically, on a volatility world model calibrated to SPX/VIX-like grids, we compare naive RL, law-penalized and selection-only RL variants against simple structural baselines (Zero-Hedge, Vol-Trend, Random-Gaussian) across baseline and shocked regimes. In our experiments, structural baselines form the empirical law-strength frontier: they attain Sharpe ratios around 2–3 with low law penalties and GFI near zero, while all law-seeking RL variants achieve non-positive mean P&L and substantially higher GFI, despite being explicitly penalized for law violations. Frontier and diagnostic plots show that RL improvements in P&L are systematically accompanied by movement into high-penalty, high-GFI regions, consistent with our theoretical analysis.

Overall, volatility serves as a concrete case study where *reward shaping with verifiable penalties is not sufficient for law alignment*. Our framework—combining law manifolds, Goodhart decomposition, GFI, and law-strength frontiers—provides a reusable template for stress-testing Scientific AI systems and RL with verifiable rewards in other axiom-constrained domains such as yield curves, credit term structures, and physics-informed models.

Keywords volatility surfaces, reinforcement learning, axiomatic finance, no-arbitrage, law manifolds, Goodhart’s law, world models, scientific AI

1 Introduction

1.1 Scientific AI testbed: volatility law manifolds

In recent years, “AI for Science” has emerged as a central theme in machine learning, emphasizing scientific understanding and hypothesis-testing rather than purely predictive or profit-driven performance [34–36]. A key lesson from this line of work is that many scientific domains are governed by *axioms*—conservation laws, monotonicity constraints, or no-arbitrage principles—that carve out a structured admissible subset inside a high-dimensional function space. Scientific machine learning then seeks not only to interpolate observations, but to understand how learning agents interact with these law-constrained spaces under model misspecification, limited data, and changing environments.

Implied-volatility (IV) surfaces are a canonical example of such an axiomatic system in quantitative finance. Classical results show that no-arbitrage constraints—such as non-negativity of butterfly spreads, monotonicity in maturity, and convexity in strike—translate into linear or convex inequalities on total-variance smiles and surfaces [1–3, 5–7]. These conditions define a structured admissible subset of discretized surfaces that we refer to as a *volatility law manifold*. Recent work has revisited arbitrage-free interpolation and extrapolation of IV surfaces using convex optimization and sparse modeling [14], arbitrage-aware parametric families, and deep neural networks that aim to respect or softly penalize violations of these constraints [10, 15, 68]. In parallel, deep learning has been applied to option pricing and hedging via PDE and BSDE solvers [30, 31], and to rough or stochastic volatility models where direct calibration is challenging [10].

Our starting point is to treat this volatility-curve setting as a *Scientific AI testbed* rather than a trading system. We assume that the data-generating process—a synthetic but financially meaningful IV-surface generator—is exactly law-consistent: every realized surface lies on the no-arbitrage manifold defined by classical butterfly and calendar conditions [1, 2, 6]. A neural *world model* is trained to approximate this generator from data, similar in spirit to world-model approaches in model-based reinforcement learning (RL) [40–45]. RL agents then interact only with the learned world model, never with the ground-truth generator. This creates a clean separation between a law-consistent environment and a potentially law-violating model, opening a “ghost channel” for agents to exploit model artefacts.

World models and latent dynamics learning have become central tools for sample-efficient RL in complex domains such as Atari, control, and strategy games [40–43, 47–49]. At the same time, RL has a long history in quantitative finance, including early work on direct reinforcement learning for trading [63, 64], deep RL for portfolio management and execution [65–67], and deep hedging using neural networks trained on simulated scenarios [68]. Yet most of these studies focus on realized P&L, with limited attention to how learned policies interact with structural market axioms. Our goal is not to propose another trading system, but to use an RL-in-IV-surfaces setup as a controlled *experiment* on law-aligned learning and Goodhart phenomena.

A parallel literature in scientific machine learning has emphasized the incorporation of physical or axiomatic structure into learning systems via physics-informed neural networks [32, 34], relational inductive biases [39], and hybrid mechanistic–ML models [35, 36]. These works typically enforce or strongly bias the model toward satisfying known laws during training. In contrast, recent AI-safety and alignment research has documented how RL agents can exploit imperfect reward channels or specification gaps, a phenomenon often traced back to Goodhart’s law [74–79]. This has led to proposals for RL with verifiable or externally checked rewards (RLVR) [82], where parts of the reward signal are computed via trusted procedures or checkers.

Our work bridges these lines of research in a finance-specific but conceptually general way. We design an axiomatic evaluation pipeline in which (i) a volatility law manifold encodes no-arbitrage axioms; (ii) a neural world model approximates a law-consistent generator but introduces model-induced “ghost” arbitrage opportunities; and (iii) RL agents are trained either with or without access to a *verifiable* law-penalty signal. We then ask a scientific question: when we add such verifiable penalties as soft terms in the RL objective, on top of a law-consistent ground-truth world, do we actually obtain more law-aligned and robust policies, or do we merely shift Goodhart behaviour onto model artefacts?

To foreshadow our main numerical findings, we consider a simple zero-position baseline (ZERO-HEDGE) that never trades. In our main setting, ZERO-HEDGE achieves mean step P&L of approximately 0.0191 with a Graceful Failure Index (GFI) essentially zero and moderate law penalties, reflecting that the underlying world is law-consistent and shocks are symmetric. By contrast, a wide range of law-seeking RL variants—including soft-penalty PPO with a law-weight sweep (the “law-strength frontier”) and a selection-only variant that uses law penalties only for model selection—all attain *non-positive* mean step P&L and substantially worse GFI values (typically ≥ 1.6), despite being explicitly penalized for law violations during training. In other words, once structural baselines are included, law-seeking RL has no free lunch: it fails to dominate even trivial strategies on the joint axes of profitability, law alignment, and tail risk.

1.2 Contributions

This paper makes four primary contributions, organized around an axiomatic evaluation framework rather than a specific trading algorithm.

C1 – Axiomatic evaluation framework. Starting from any finite collection of convex or linear axioms on a discretized observable field, we construct (i) a law manifold \mathcal{M} in total-variance coordinates for implied-volatility surfaces, (ii) a metric-based law-penalty functional \mathcal{L}_ϕ that measures distance to \mathcal{M} , (iii) a domain-agnostic Graceful Failure Index (GFI) that normalizes degradation of law metrics under shocks, and (iv) *law-strength frontiers* that jointly organize profitability, law alignment, and tail robustness as a function of law-penalty weight λ and strategy class. While we instantiate this framework in an IV-surface world, the construction applies equally to other axiom-constrained systems such as yield curves, credit term structures, and physical fields constrained by conservation laws [32, 34, 35].

C2 – Ghost arbitrage & Goodhart decomposition. We formalize a Goodhart-style decomposition of reward on law manifolds,

$$r = r^{\mathcal{M}} + r^\perp, \quad (1)$$

where $r^{\mathcal{M}}$ is the on-manifold reward that would be obtained under a perfectly law-consistent world, and r^\perp is an off-manifold “ghost arbitrage” component induced by the neural world model. We show how this decomposition can be implemented using a projection operator onto \mathcal{M} in total-variance space and an explicit law-penalty functional, making the ghost component measurable and analyzable. This connects Goodhart’s law in AI safety [74, 75, 78] with concrete financial law violations (e.g., butterfly or calendar arbitrage) in our IV-surface testbed.

C3 – Flagship incentive and trade-off results. On top of this axiomatic pipeline, we establish three central theoretical results that together form a no-free-lunch story for law-seeking RL. Theorem 4.1 (*Ghost-arbitrage incentive for naive RL*) shows that, under mild assumptions, naive PPO-type RL is structurally incentivized to increase $\mathbb{E}[r^\perp]$ whenever structural law-consistent baselines already approximate the on-manifold optimum. Theorem 4.3 and Corollary 4.4 (*Law-strength trade-off*) prove that increasing the law-penalty weight λ inevitably worsens P&L beyond a threshold: the empirical law-strength frontier we observe in experiments is a structural trade-off, not an artefact of hyperparameters. Finally, Theorem 8.1 (*No-free-lunch for law-seeking RL*) shows that, given a law-consistent world model and a sufficiently rich structural baseline class \mathcal{S} , unconstrained law-seeking RL cannot simultaneously dominate \mathcal{S} on P&L and on all law metrics unless it effectively recovers a policy in \mathcal{S} .

C4 – Design lessons for law-aligned learning and RLVR. Our analysis yields practical design lessons that generalize beyond volatility modeling. First, merely adding soft law penalties to the reward is insufficient for robust law alignment: RL agents systematically exploit ghost arbitrage channels in the world model, or else sacrifice P&L without achieving net law improvements. Second, law alignment benefits from *structural* interventions such as hard constraints, projection layers onto \mathcal{M} , and structured policy classes that encode hedging logic, echoing observations from physics-informed learning and scientific ML [32, 34, 36]. Third, our pipeline serves as a reusable testbed for RL with verifiable rewards (RLVR) [82]: law penalties here are fully verifiable and domain-grounded, yet we still observe Goodhart-like failures when structural constraints are absent.

1.3 Scope and novelty

Geometric/systematization of known results. Section 2 recasts classical no-arbitrage characterizations of admissible IV surfaces—butterfly convexity, calendar monotonicity, and related inequalities [1–3, 5–7]—as a finite-dimensional convex polyhedral law manifold in total-variance coordinates. This representation theorem does not aim to replace existing no-arbitrage results; instead, it systematizes them into a form amenable to projection, distance computation, and integration into learning systems.

New concepts and metrics. The concrete ghost-arbitrage decomposition on a learned world model, the construction of law-strength frontiers, and the definition of the GFI are new. They are designed to be domain-agnostic: given any axiom-constrained system and an exogenous notion of shock, the same machinery can be instantiated to quantify how agents trade off performance, law alignment, and tail robustness, extending ideas from scientific ML and AI-for-Science benchmarks [32, 34–36].

New theoretical results. Our main theoretical novelties lie in Theorem 4.1, Theorem 4.3 with Corollary 4.4, and Theorem 8.1. Together, they formalize: (i) an incentive for naive RL to exploit ghost arbitrage in law-consistent worlds; (ii) a structural law-strength trade-off that bounds achievable P&L for any given level of law penalty; and (iii) a

Table 1: Key concepts at a glance. Formal definitions and constructions are given in the indicated sections.

Concept	Plain-English description	Section
Law manifold \mathcal{M}	Admissible subset of discretized IV surfaces satisfying butterfly, calendar, and related no-arbitrage axioms.	Sec. 2.0–2.1
Law penalty \mathcal{L}_ϕ	Distance-based functional measuring how far a surface lies outside \mathcal{M} , with ϕ encoding the chosen metric (e.g., squared ℓ_2 in total-variance space).	Sec. 2.3
Ghost arbitrage r^\perp	Component of reward obtainable only by moving off \mathcal{M} through world-model artefacts; vanishes under a perfectly law-consistent environment.	Sec. 2.4, 3.3
Law-strength frontier	Trade-off curve tracing profitability, law alignment, and tail risk as a function of law-penalty weight λ and strategy class.	Sec. 4.4, 7.3
Graceful Failure Index (GFI)	Normalized measure of how much law metrics (e.g., mean penalty, coverage) degrade under shocks relative to baseline conditions.	Sec. 4.4, 6.3
Structural class \mathcal{S}	Low-capacity, law-consistent strategies such as ZERO-HEDGE and VOL-TREND that serve as structural baselines.	Sec. 5, 8.1
Neural world model	Learned dynamics model for IV surfaces that approximates the law-consistent generator but opens a ghost channel for arbitrage.	Sec. 3

no-free-lunch result for unconstrained law-seeking RL relative to a law-consistent structural baseline class \mathcal{S} . We stress that these are not new no-arbitrage theorems per se, but results about RL behaviour on top of an axiomatic evaluation pipeline. Our contribution is to show, both theoretically and empirically, that soft law penalties on a learned world model do not automatically yield law-aligned robustness once structural baselines are taken into account.

Scope disclaimer. We work with a finite-dimensional convex template: our “law manifold” is a structured subset of a discretized IV-surface grid, not a smooth manifold in the differential-geometric sense. We retain the term “manifold” for continuity with the broader literature on manifold-constrained learning, but in our volatility instantiation \mathcal{M}^{vol} is a convex polyhedral subset defined by linear inequalities and positivity constraints. Our theorems are proved under this finite-dimensional convex template and a specific class of model-based RL algorithms; a fully general analysis for infinite-dimensional function spaces and arbitrary RL algorithms is left for future work.

1.4 Key concepts at a glance

Given the number of concepts introduced, we summarize the most important ones in Table 1. Each concept is defined precisely in later sections; here we provide a high-level description and a pointer.

1.5 Research questions

The paper is organized around three research questions (RQs) that probe different aspects of law-seeking RL on axiomatic pipelines.

RQ1 – Do law penalties help naive RL? Does law-penalized RL actually reduce law violations and improve graceful failure compared to naive RL on the *same* learned world model? If soft penalties are effective, we should observe strictly better GFI and law metrics (mean and max law penalty, law coverage) at comparable P&L along the law-strength frontier. If not, we obtain a negative result for soft-penalty shaping: verifiable law penalties alone do not suffice to align RL with axioms in the presence of model-induced ghost arbitrage.

RQ2 – How does RL compare to structural baselines? How do RL policies (naive, soft law-seeking, and selection-only) compare to structural baselines (such as ZERO-HEDGE, RANDOM-GAUSSIAN, and VOL-TREND) on the joint risk–law trade-off, both under baseline conditions and under volatility shocks? We evaluate these policies relative to a structural Pareto frontier induced by \mathcal{S} in the space of P&L, law metrics, and tail-risk measures such as Value-at-Risk (VaR) and Conditional VaR (CVaR) [21, 22]. This addresses whether law-seeking RL provides any Pareto improvement once simple, law-consistent strategies are included.

RQ3 – When does law-seeking RL have no free lunch? Under what assumptions on the structural class \mathcal{S} , the unconstrained policy class Π , and the neural world model do we obtain a no-free-lunch result for law-seeking RL?

Theorem 8.1 formalizes one such setting: if S already approximates the on-manifold optimum and the world model introduces a non-trivial ghost component, then either RL behaves like a structural strategy in S or it fails to dominate S jointly on P&L and law metrics. Section 8.1 spells out the assumptions and proof sketch, and Section 8.2 discusses how this volatility-specific case informs broader questions in law-aligned learning and RL with verifiable rewards [74–76, 82].

In the remainder of the paper, we address RQ1–RQ3 in order. Section 2 formalizes law manifolds, penalties, and the Goodhart decomposition. Section 3 introduces the synthetic law-consistent IV generator and the neural world model. Section 4 develops incentive and trade-off results for law-seeking RL, and Section 5 defines structural baselines. Section 6 details the experimental protocol, and Section 7 presents empirical results on dynamics plots, law-strength frontiers, and diagnostic scatter/histograms. Section 8 discusses the no-free-lunch theorem and implications for RLVR and scientific AI, and Section 9 concludes.

2 Axiomatic Volatility Law Manifolds

2.0 Finite-Dimensional Convex Template and Notation

In this section we formalize the notion of a *law manifold* as a finite-dimensional convex subset of a discretized function space, induced by a collection of axioms. We emphasize from the outset that our construction is intentionally elementary:

We use a simple finite-dimensional convex template; our “manifold” is a structured subset in discretized coordinates, not a smooth differentiable manifold in the differential-geometry sense. We keep the term *manifold* for continuity with the broader literature on manifold-constrained learning and manifold regularization [28, 29], although in our discretized volatility setting \mathcal{M}^{vol} is a convex polyhedral subset of a Euclidean space.

General template. Let $d \in \mathbb{N}$ and consider a finite-dimensional observation space

$$\mathcal{Y} \subseteq \mathbb{R}^d,$$

equipped with the standard Euclidean inner product and norm $\|\cdot\|_2$. We think of $y \in \mathcal{Y}$ as a discretized field: a yield curve, an implied-volatility surface, or any other structured observable.

We assume that the domain is endowed with a finite family of convex *axiom functions*

$$A_i : \mathcal{Y} \rightarrow \mathbb{R}, \quad i = 1, \dots, m,$$

encoding domain-specific constraints (e.g., butterfly or calendar conditions in volatility, or monotonicity of yields). We write $A(y) \leq 0$ as shorthand for the componentwise inequalities $A_i(y) \leq 0$ for all i .

Definition 1 (Law manifold). Given convex axiom functions $\{A_i\}_{i=1}^m$, the associated *law manifold* is

$$\mathcal{M} := \{y \in \mathcal{Y} : A_i(y) \leq 0 \text{ for all } i = 1, \dots, m\}.$$

In general \mathcal{M} is a closed convex subset of \mathbb{R}^d whenever each A_i is lower semicontinuous and convex. In many practical cases—including our volatility and yield-curve examples below—the constraints are linear or piecewise linear, so that \mathcal{M} is a polyhedron.

Law-penalty functional. To quantify violations of the axioms, we define a *law-penalty functional* via a generalized distance to \mathcal{M} .

Definition 2 (Law-penalty functional). Let $\phi : \mathbb{R}^d \rightarrow [0, \infty)$ be a continuous function with $\phi(0) = 0$ and $\phi(z) \rightarrow \infty$ as $\|z\|_2 \rightarrow \infty$ (e.g., $\phi(z) = \frac{1}{2}\|z\|_2^2$ or a weighted Sobolev norm). Define

$$\mathcal{L}_\phi(y) := \inf_{\tilde{y} \in \mathcal{M}} \phi(y - \tilde{y}), \quad y \in \mathcal{Y}. \quad (2)$$

When $\phi(z) = \frac{1}{2}\|z\|_2^2$ and \mathcal{M} is closed and convex, $\mathcal{L}_\phi(y)$ reduces to the squared Euclidean distance to \mathcal{M} . More general choices of ϕ allow us to reweight different coordinates, incorporate smoothness, or emphasize particular directions in state space, in the spirit of manifold regularization [28].

Non-financial example: yield curves. To underline transferability beyond volatility, consider a discretized yield curve

$$y = (y_1, \dots, y_J) \in \mathbb{R}^J,$$

where y_j denotes the continuously compounded spot yield for maturity τ_j , with $0 < \tau_1 < \dots < \tau_J$. Classical term-structure theory [18] and curve-construction practice [17] suggest axioms such as:

1. *Monotonicity*: yields are non-decreasing in maturity, $y_{j+1} \geq y_j \quad \forall j$.
2. *Convexity of discount factors*: implied by non-negative forward rates, giving linear inequalities in y .

These conditions can be encoded as linear maps $A_i(y)$, yielding a polyhedral law manifold $\mathcal{M}^{yc} \subset \mathbb{R}^J$ and an associated law penalty \mathcal{L}_ϕ^{yc} measuring monotonicity/convexity violations. Our volatility manifold in Sections 2.1–2.3 is an instance of this general template.

Notation summary. Throughout the paper we use the following notation; we repeat it here for convenience.

1. y : generic observable (e.g., yield curve, volatility surface) in a finite-dimensional space $\mathcal{Y} \subseteq \mathbb{R}^d$.
2. \mathcal{M} : generic law manifold (closed convex subset of \mathbb{R}^d) induced by axioms.
3. σ : implied volatility on a maturity–log-moneyness grid; $w = \sigma^2 T$ denotes total variance, our primary state variable for volatility.
4. \mathcal{M}^{vol} : volatility-specific law manifold defined in Section 2.1.
5. \mathcal{L}_ϕ : law-penalty functional defined in (2); in experiments we take $\phi(z) = \frac{1}{2} \|z\|_2^2$ on the total-variance grid.
6. $\Pi_{\mathcal{M}}$: metric projection onto \mathcal{M} (in Euclidean norm), whose existence and regularity rely on closedness and convexity [26].
7. $r^{\mathcal{M}}, r^\perp$: on-manifold reward and *ghost-arbitrage* components of the reward, defined via projection in Section 2.4.

2.1 Volatility-Specific Law Manifold: From Textbook Axioms to a Polyhedron

We now instantiate the general template for implied volatility surfaces. Let $C(K, T)$ denote the (discounted) price of a European call with strike K and maturity T , and $\sigma(K, T)$ its Black–Scholes implied volatility. Following standard practice [1, 2], we work in terms of total variance

$$w(K, T) := \sigma(K, T)^2 T,$$

discretized on a finite grid of maturities $T_1 < \dots < T_{N_T}$ and log-moneyness $k_1 < \dots < k_{N_K}$.

Classical static no-arbitrage conditions for European options (no butterfly arbitrage across strikes and no calendar arbitrage across maturities) can be expressed as linear inequalities in either call prices or total variance, see e.g. [8, 16]. We briefly summarize the discretized form relevant to our construction.

Butterfly (strike) convexity. For each fixed maturity T_i , the call price as a function of strike must be convex and decreasing. On a grid $\{k_j\}$ this gives discrete convexity constraints such as

$$C_{i,j-1} - 2C_{i,j} + C_{i,j+1} \geq 0,$$

and monotonicity constraints $C_{i,j+1} \leq C_{i,j}$ for all j . These translate into linear inequalities in $w_{i,j}$ via the Black–Scholes formula.

Calendar monotonicity. For each fixed strike (log-moneyness) k_j , call prices must be increasing in maturity, yielding

$$C_{i+1,j} \geq C_{i,j} \quad \text{for all } i.$$

Again, these can be expressed as linear inequalities in $(w_{i,j})$ using the monotonicity of Black–Scholes prices in variance. Collecting all discrete butterfly and calendar inequalities, we obtain a linear mapping

$$A^{\text{vol}} : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathbb{R}^m,$$

where $d_{\text{vol}} = N_T N_K$ is the number of grid points and m is the number of constraints.

Definition 3 (Volatility law manifold). Let $w \in \mathbb{R}^{d_{\text{vol}}}$ denote the vector of total variances on the (T, k) -grid. The *volatility law manifold* is the polyhedral set

$$\mathcal{M}^{\text{vol}} := \{ w \in \mathbb{R}^{d_{\text{vol}}} : A^{\text{vol}} w \leq 0 \}, \quad (3)$$

where $A^{\text{vol}} w \leq 0$ encodes all discretized butterfly and calendar no-arbitrage inequalities, as well as basic box constraints $w_{\min} \leq w_{i,j} \leq w_{\max}$.

The following proposition packages the textbook no-arbitrage conditions into a geometric representation.

Proposition 1 (Axiomatic representation of volatility law manifold). *Assume that the discretized butterfly and calendar constraints are given as a finite system of linear inequalities $A^{\text{vol}} w \leq b$ for some matrix $A^{\text{vol}} \in \mathbb{R}^{m \times d_{\text{vol}}}$ and vector $b \in \mathbb{R}^m$. Then:*

1. \mathcal{M}^{vol} is a non-empty, closed, convex polyhedron in $\mathbb{R}^{d_{\text{vol}}}$.
2. Any total-variance surface w corresponding to a static-arbitrage-free implied volatility surface lies in \mathcal{M}^{vol} .

Proof sketch. Closedness and convexity follow directly from the fact that \mathcal{M}^{vol} is the intersection of finitely many closed half-spaces $\{w : a_\ell^\top w \leq b_\ell\}$ and box constraints. Non-emptiness is guaranteed by the existence of at least one model (e.g., a Black–Scholes surface with constant volatility) satisfying the inequalities. The mapping from continuous no-arbitrage conditions to discrete linear constraints is standard; see e.g. [2], [8] and references therein. A detailed construction and proof are given in Appendix A. \square

2.2 Geometry and Convexity Properties

The closed and convex nature of \mathcal{M}^{vol} is more than a technicality: it ensures the existence of well-behaved projection operators and law penalties.

Proposition 2 (Closedness, convexity, and metric projection). *The volatility law manifold \mathcal{M}^{vol} defined in (3) is a non-empty, closed, convex subset of $\mathbb{R}^{d_{\text{vol}}}$. Consequently:*

1. For every $w \in \mathbb{R}^{d_{\text{vol}}}$ there exists a unique Euclidean projection

$$\Pi_{\mathcal{M}^{\text{vol}}}(w) := \arg \min_{\tilde{w} \in \mathcal{M}^{\text{vol}}} \|w - \tilde{w}\|_2.$$

2. The projection operator $\Pi_{\mathcal{M}^{\text{vol}}} : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathcal{M}^{\text{vol}}$ is 1-Lipschitz:

$$\|\Pi_{\mathcal{M}^{\text{vol}}}(w) - \Pi_{\mathcal{M}^{\text{vol}}}(w')\|_2 \leq \|w - w'\|_2 \quad \forall w, w'.$$

Proof sketch. Closedness and convexity follow from Proposition 1. The properties of $\Pi_{\mathcal{M}^{\text{vol}}}$ are standard for projections onto closed convex sets in Hilbert spaces; see, e.g., [26, Chap. 3]. Full details are provided in Appendix A. \square

The existence and regularity of $\Pi_{\mathcal{M}^{\text{vol}}}$ allow us to treat law penalties as squared distances to the manifold and to define projection-based decompositions of reward functionals later on.

2.3 Law-Penalty Functionals and Ghost Sensitivity

We now specialize the general law-penalty functional (2) to the volatility setting. Let $w \in \mathbb{R}^{d_{\text{vol}}}$ denote a (possibly law-violating) total-variance surface, and let

$$\phi(z) = \frac{1}{2} \|z\|_2^2, \quad z \in \mathbb{R}^{d_{\text{vol}}}. \quad (4)$$

We define

$$\mathcal{L}_{\text{vol}}(w) := \inf_{\tilde{w} \in \mathcal{M}^{\text{vol}}} \frac{1}{2} \|w - \tilde{w}\|_2^2 = \frac{1}{2} \text{dist}(w, \mathcal{M}^{\text{vol}})^2, \quad (5)$$

where $\text{dist}(w, \mathcal{M}^{\text{vol}}) := \|w - \Pi_{\mathcal{M}^{\text{vol}}}(w)\|_2$.

In practice we implement \mathcal{L}_{vol} via a sum of local violations (e.g., squared negative parts of discrete butterfly and calendar inequalities), which is equivalent to (5) up to scaling on our discretization.

Lemma 1 (Local Lipschitz continuity of law penalty). *The volatility law-penalty functional \mathcal{L}_{vol} in (5) is locally Lipschitz on $\mathbb{R}^{d_{\text{vol}}}$.*

Proof sketch. $\mathcal{L}_{\text{vol}}(w) = \frac{1}{2}\|w - \Pi_{\mathcal{M}^{\text{vol}}}(w)\|_2^2$ and $\Pi_{\mathcal{M}^{\text{vol}}}$ is 1-Lipschitz by Proposition 2. Thus \mathcal{L}_{vol} is the composition of Lipschitz maps and a smooth quadratic, which yields local Lipschitz continuity; see also standard results on Moreau envelopes [26]. A full proof is given in Appendix A. \square

The following basic property connects \mathcal{L}_{vol} with axiomatic consistency and will be used repeatedly when interpreting empirical law metrics.

Proposition 3 (Zero penalty iff axiomatic consistency). *For any $w \in \mathbb{R}^{d_{\text{vol}}}$,*

$$\mathcal{L}_{\text{vol}}(w) = 0 \iff w \in \mathcal{M}^{\text{vol}}.$$

Proof sketch. If $w \in \mathcal{M}^{\text{vol}}$ then choosing $\tilde{w} = w$ in (5) yields $\mathcal{L}_{\text{vol}}(w) = 0$. Conversely, if $\mathcal{L}_{\text{vol}}(w) = 0$ then the infimum in (5) is attained at $\tilde{w}^* = \Pi_{\mathcal{M}^{\text{vol}}}(w)$ with $\|w - \tilde{w}^*\|_2 = 0$, hence $w = \tilde{w}^* \in \mathcal{M}^{\text{vol}}$. Details appear in Appendix A. \square

Remark 1 (Choice of ϕ and ghost sensitivity). The choice $\phi(z) = \frac{1}{2}\|z\|_2^2$ treats all grid points uniformly and yields a particularly simple gradient

$$\nabla \mathcal{L}_{\text{vol}}(w) = w - \Pi_{\mathcal{M}^{\text{vol}}}(w)$$

almost everywhere. Alternative choices of ϕ (e.g., weighted ℓ_2 norms emphasizing short maturities, or discrete Sobolev norms penalizing roughness in T and k) change the relative sensitivity of \mathcal{L}_{ϕ} to localized versus global law violations. Exploring how this affects the *ghost arbitrage* exploited by RL policies is an interesting axis for future work.

2.4 Goodhart Decomposition on Law Manifolds (Conceptual)

The law manifold and penalty enable a conceptual decomposition of any reward functional into an on-manifold and an off-manifold (ghost-arbitrage) component. This decomposition underlies our Goodhart-style analysis of RL in later sections and is instantiated concretely for our world model in Section 3.3.

Let $r : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathbb{R}$ denote a generic reward functional of a total-variance surface w (e.g., the one-step P&L of a hedging strategy). We assume that r is locally Lipschitz on $\mathbb{R}^{d_{\text{vol}}}$, a mild condition satisfied in all our models.

Definition 4 (Projection-based Goodhart decomposition). Let $\Pi_{\mathcal{M}^{\text{vol}}}$ be the metric projection onto \mathcal{M}^{vol} . Define the *on-manifold* reward

$$r^{\mathcal{M}}(w) := r(\Pi_{\mathcal{M}^{\text{vol}}}(w)),$$

and the *ghost-arbitrage* component

$$r^{\perp}(w) := r(w) - r^{\mathcal{M}}(w).$$

We refer to the decomposition

$$r(w) = r^{\mathcal{M}}(w) + r^{\perp}(w)$$

as the *Goodhart decomposition* of r on the volatility law manifold.

By construction, $r^{\mathcal{M}}$ depends on w only through its projection onto \mathcal{M}^{vol} , and is thus invariant under off-manifold perturbations that leave $\Pi_{\mathcal{M}^{\text{vol}}}(w)$ unchanged. The residual r^{\perp} captures gains or losses that arise solely from moving away from the law manifold—precisely the *ghost arbitrage* channel that law-seeking RL may exploit.

Definition 5 (Ghost arbitrage). The *ghost-arbitrage reward* associated with a reward functional r and law manifold \mathcal{M}^{vol} is the component r^{\perp} in Definition 4. A policy that maximizes $\mathbb{E}[r^{\perp}]$ subject to small \mathcal{L}_{vol} can be said to exploit ghost arbitrage: it harvests reward from off-manifold distortions that are negligible under the coarse law penalty but significant for P&L.

The following proposition makes explicit how the magnitude of ghost arbitrage is controlled by the distance to the law manifold whenever r is Lipschitz. This link is used later when we interpret empirical law penalties and our Graceful Failure Index.

Proposition 4 (Ghost reward bounded by law distance). *Suppose r is L_r -Lipschitz on $\mathbb{R}^{d_{\text{vol}}}$ with respect to $\|\cdot\|_2$. Then for all w ,*

$$|r^{\perp}(w)| = |r(w) - r^{\mathcal{M}}(w)| \leq L_r \text{dist}(w, \mathcal{M}^{\text{vol}}) = L_r \sqrt{2 \mathcal{L}_{\text{vol}}(w)}.$$

Proof sketch. By the Lipschitz property of r ,

$$|r(w) - r^{\mathcal{M}}(w)| = |r(w) - r(\Pi_{\mathcal{M}^{\text{vol}}}(w))| \leq L_r \|w - \Pi_{\mathcal{M}^{\text{vol}}}(w)\|_2 = L_r \text{dist}(w, \mathcal{M}^{\text{vol}}).$$

Using $\mathcal{L}_{\text{vol}}(w) = \frac{1}{2} \text{dist}(w, \mathcal{M}^{\text{vol}})^2$ from (5) yields the final identity. A more general version, allowing non-Euclidean ϕ , is treated in Appendix A. \square

Proposition 4 shows that, in a purely metric sense, large ghost rewards require non-negligible law violations. The central question of this paper is whether, under a learned world model, RL training can nonetheless systematically exploit directions in which ghost reward grows “too quickly” relative to the coarse law-penalty captured by \mathcal{L}_{vol} , leading to misaligned but high-P&L policies. This question is addressed empirically in Sections 7 and theoretically in Theorems 4.1 and 8.1.

3 Volatility World Model: Law-Consistent Ground Truth vs Law-Violating Predictions

In this section we formalize the dynamic data-generating process for implied-volatility surfaces and the learned neural world model on which all reinforcement-learning (RL) agents are trained. The key structural feature is that the *synthetic generator* is, by construction, perfectly law-consistent—its surfaces lie on the polyhedral law manifold \mathcal{M}^{vol} almost surely—whereas the neural world model is trained purely by prediction error and hence produces *law-violating* surfaces with non-zero probability. This mismatch opens a *ghost channel* for RL to exploit off-manifold artefacts, even though the underlying ground truth never leaves \mathcal{M}^{vol} .

3.1 Synthetic law-consistent generator

We consider a discrete time grid $t = 0, 1, \dots, T$ and a rectangular grid of maturities and strikes

$$\mathcal{T} = \{T_1, \dots, T_M\}, \quad \mathcal{K} = \{K_1, \dots, K_K\},$$

with $T_1 \approx 1\text{M}$ and $T_M \approx 2\text{Y}$, and $K_1 \approx 0.5S_t$, $K_K \approx 1.5S_t$ at each time t , chosen to roughly mirror SPX/VIX market conventions.¹ We denote by

$$\sigma_t = \sigma_t(T_i, K_j) \in \mathbb{R}^{M \times K}, \quad w_t = \sigma_t^2 \odot T \in \mathbb{R}^{M \times K}$$

the implied-volatility and total-variance surfaces at time t , where T is broadcast across strikes. In vectorized form we identify w_t with an element of \mathbb{R}^d with $d = MK$, and we write $w_t \in \mathcal{M}^{\text{vol}}$ when all static no-arbitrage constraints (butterfly, calendar) are satisfied on the discrete grid (Sec. 2).

Law-consistent ground truth. The *ground-truth world* is specified by a Markovian generator

$$G^* : \mathcal{S} \rightarrow \mathcal{S} \times \mathcal{M}^{\text{vol}}, \quad (s_{t+1}, w_{t+1}) = G^*(s_t),$$

where s_t is a latent state that may contain the underlying spot S_t , latent volatility factors, and other macro state variables. In practice we instantiate G^* as a multi-factor stochastic-volatility model with jumps and rough components,² whose parameters are randomized across trajectories to induce a diverse ensemble of surfaces reminiscent of SPX/VIX dynamics [e.g. 2, 6, 11].

At each time step, raw model-implied option prices are numerically projected onto the static no-arbitrage polyhedron \mathcal{M}^{vol} using the construction of Sec. 2, and implied volatilities are recovered from these arbitrage-free prices. As a result, the data-generating distribution \mathbb{P}^* over trajectories $(w_t)_{t=0}^T$ is *law-consistent by design*.

Definition 6 (Law-consistent synthetic generator). A stochastic process $(w_t)_{t=0}^T$ taking values in \mathbb{R}^d is said to be *law-consistent* with respect to a law manifold $\mathcal{M} \subset \mathbb{R}^d$ if

$$\mathbb{P}(w_t \in \mathcal{M} \text{ for all } t = 0, \dots, T) = 1.$$

We call G^* a *law-consistent generator* if the trajectory (w_t) it induces satisfies this condition for $\mathcal{M} = \mathcal{M}^{\text{vol}}$.

Proposition 5 (Support of the synthetic generator). *Let $(w_t)_{t=0}^T$ be generated by G^* as above, with static no-arbitrage imposed at each time via projection onto \mathcal{M}^{vol} . Then*

$$\text{supp } \mathbb{P}^* \subseteq (\mathcal{M}^{\text{vol}})^{T+1},$$

i.e., \mathbb{P}^ is supported on the product of the volatility law manifold at all times.*

Proof sketch. By construction, every step of G^* maps into \mathcal{M}^{vol} : the raw option prices are obtained from a stochastic-volatility model, and then the resulting surface is projected onto \mathcal{M}^{vol} as in Sec. 2. Thus $w_t \in \mathcal{M}^{\text{vol}}$ almost surely for all t . The support statement follows from the definition of product measures. A detailed measure-theoretic proof is given in Appendix B.1. \square

¹The exact grid is not essential; any finite grid of maturities and moneyness levels can be handled by the law manifold construction in Sec. 2.

²E.g., a Heston-like model with stochastic variance and stochastic volatility-of-volatility, optionally enriched with rough volatility factors as in [9], coupled to an SVI-type implied-vol parametrization [8].

Proposition 5 highlights the asymmetry that drives our ghost-arbitrage story: the *true* dynamics never leave the law manifold, whereas the learned world model in the next subsection is not constrained in this way.

Transferability beyond volatility. The same construction immediately extends to other axiom-constrained financial objects such as yield curves and credit curves: there, y_t is a discretized term structure, \mathcal{M} encodes monotonicity and convexity constraints [e.g. 19, 20], and G^* is any arbitrage-free term-structure model. Our volatility case thus serves as a concrete, high-dimensional instance of a more general Scientific AI template.

3.2 Neural world model and approximation gap

RL agents do not interact with G^* directly. Instead, in the spirit of model-based RL and world models [40, 45? ?], they are trained entirely on rollouts generated by a learned dynamics model (the *world model*) fitted to trajectories from G^* .

Architecture and training. Let L be the length of the look-back window. At each time t , we define the input

$$x_t := (w_{t-L+1}, \dots, w_t) \in (\mathbb{R}^d)^L,$$

and we train a recurrent neural network with parameters θ ,

$$f_\theta : (\mathbb{R}^d)^L \rightarrow \mathbb{R}^d, \quad \hat{w}_{t+1} = f_\theta(x_t),$$

to minimize the mean-squared prediction error

$$\mathcal{R}(\theta) := \mathbb{E}_{\mathbb{P}^*} [\|f_\theta(x_t) - w_{t+1}\|_2^2] \approx \frac{1}{N} \sum_{n=1}^N \|f_\theta(x_t^{(n)}) - w_{t+1}^{(n)}\|_2^2$$

over a dataset of N trajectories sampled from G^* . In practice, f_θ is implemented as a GRU/LSTM encoder over (w_{t-L+1}, \dots, w_t) followed by a fully connected decoder to the next total-variance surface, as commonly used in spatio-temporal forecasting [e.g. 72, 73].

Crucially, *no law penalties are used when training the world model*: the loss is purely predictive. Thus, even though all training targets w_{t+1} lie in \mathcal{M}^{vol} , the predictions \hat{w}_{t+1} are unconstrained and may lie outside the manifold.

Definition 7 (Approximation residual and ghost channel). Let θ^* be any (local) minimizer of $\mathcal{R}(\theta)$, and write

$$\hat{w}_{t+1} = f_{\theta^*}(x_t), \quad e_{t+1} := \hat{w}_{t+1} - w_{t+1}$$

for the corresponding prediction and residual. We define the *approximation gap* as

$$\varepsilon^2 := \mathbb{E}_{\mathbb{P}^*} [\|e_{t+1}\|_2^2] = \mathcal{R}(\theta^*),$$

and the *ghost channel* as the random variable

$$r_{t+1}^\perp := r(\hat{w}_{t+1}, a_t) - r(w_{t+1}^\mathcal{M}, a_t),$$

where r is the one-step P&L functional, a_t is the agent's action, and $w_{t+1}^\mathcal{M} := \Pi_{\mathcal{M}^{\text{vol}}}(\hat{w}_{t+1})$ is the metric projection of the prediction onto the law manifold (Def. 2).

Here, r_{t+1}^\perp is precisely the off-manifold component in the Goodhart decomposition (Sec. ??); it captures how much extra P&L the agent obtains by exploiting law-violating predictions rather than their arbitrage-free projection.

Proposition 6 (Approximation gap induces a ghost channel). *Suppose the following conditions hold:*

- (i) *The approximation gap is non-zero: $\varepsilon^2 = \mathbb{E}[\|e_{t+1}\|_2^2] > 0$.*
- (ii) *The reward is locally differentiable in w with gradient $g_{t+1} := \nabla_w r(w_{t+1}, a_t)$.*
- (iii) *The residual e_{t+1} has a component in the normal cone of \mathcal{M}^{vol} at w_{t+1} with non-zero covariance:*

$$\text{Cov}(P_{N_{\mathcal{M}}(w_{t+1})} e_{t+1}, g_{t+1}) \neq 0,$$

where $P_{N_{\mathcal{M}}(w_{t+1})}$ denotes orthogonal projection onto the normal cone $N_{\mathcal{M}}(w_{t+1})$.

Then, for sufficiently small residuals (in the sense of a local linearization),

$$\mathbb{E}[r_{t+1}^\perp] \approx \mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}}(w_{t+1})} e_{t+1}] \neq 0,$$

so the world model induces a non-trivial ghost channel. In particular, if g_{t+1} is positively correlated with $P_{N_{\mathcal{M}}(w_{t+1})} e_{t+1}$, then $\mathbb{E}[r_{t+1}^\perp] > 0$ and there exist states where moving off-manifold strictly improves expected P&L.

Proof sketch. By the law-consistency of G^* , we have $w_{t+1} \in \mathcal{M}^{\text{vol}}$ almost surely. For small residuals, a first-order Taylor approximation yields

$$r(\hat{w}_{t+1}, a_t) \approx r(w_{t+1}, a_t) + g_{t+1}^\top e_{t+1}.$$

Meanwhile, the projection $w_{t+1}^{\mathcal{M}} = \Pi_{\mathcal{M}}(\hat{w}_{t+1})$ removes the component of e_{t+1} that lies in the normal cone $N_{\mathcal{M}}(w_{t+1})$ (by optimality conditions for convex projections), so to first order we have

$$r(w_{t+1}^{\mathcal{M}}, a_t) \approx r(w_{t+1}, a_t) + g_{t+1}^\top P_{T_{\mathcal{M}}(w_{t+1})} e_{t+1},$$

where $T_{\mathcal{M}}(w_{t+1})$ is the tangent cone and $P_{T_{\mathcal{M}}}$ the corresponding projector. Their difference is

$$r_{t+1}^\perp \approx g_{t+1}^\top (I - P_{T_{\mathcal{M}}(w_{t+1})}) e_{t+1} = g_{t+1}^\top P_{N_{\mathcal{M}}(w_{t+1})} e_{t+1}.$$

Taking expectations under \mathbb{P}^* and using assumption (iii) yields the claim. Rigorous error bounds for the linearization and a detailed cone-decomposition argument are provided in Appendix B.2. \square

Proposition 6 formalizes the intuitive statement that *any* non-zero approximation gap with a component normal to the law manifold, combined with a reward that is monotone in that direction, will generically open a ghost channel. In Sec. 7 we show empirically that RL agents indeed learn to exploit this channel.

3.3 Instantiating the Goodhart decomposition for volatility

We now instantiate the conceptual Goodhart decomposition of Sec. ?? in the concrete volatility setting. For each predicted surface $\hat{w}_{t+1} = f_{\theta^*}(x_t)$, we compute:

1. The metric projection onto the volatility law manifold:

$$w_{t+1}^{\mathcal{M}} := \Pi_{\mathcal{M}^{\text{vol}}}(\hat{w}_{t+1}) = \arg \min_{w' \in \mathcal{M}^{\text{vol}}} \phi(\hat{w}_{t+1} - w'),$$

where ϕ is the squared ℓ_2 norm in total-variance space, consistent with the law-penalty functional \mathcal{L}_ϕ of Sec. ??.

2. The on-manifold reward component

$$r_{t+1}^{\mathcal{M}} := r(w_{t+1}^{\mathcal{M}}, a_t),$$

obtained by evaluating the P&L functional under the projected surface.

3. The ghost-arbitrage component

$$r_{t+1}^\perp := r(\hat{w}_{t+1}, a_t) - r_{t+1}^{\mathcal{M}}.$$

By construction we have the exact decomposition

$$r(\hat{w}_{t+1}, a_t) = r_{t+1}^{\mathcal{M}} + r_{t+1}^\perp,$$

where $r_{t+1}^{\mathcal{M}}$ is the reward that would be obtained if the world model were first projected back to the law manifold, and r_{t+1}^\perp captures the incremental reward purely due to law-violating predictions. Note that, because the ground truth never leaves \mathcal{M}^{vol} (Prop. 5), any systematic pattern in r_{t+1}^\perp is necessarily a *model-induced artefact*.

Remark 2 (Consistency of projection operator). For consistency with Sec. ??, we use the same projection operator $\Pi_{\mathcal{M}^{\text{vol}}}$ both in defining the law penalty \mathcal{L}_ϕ and in the Goodhart decomposition. Algorithmically, $\Pi_{\mathcal{M}^{\text{vol}}}$ is implemented via a convex quadratic program that enforces butterfly and calendar inequalities on the total-variance grid, closely related to static-arbitrage projection procedures in the option-pricing literature [11–13]. This ensures that any off-manifold advantage measured by r^\perp is directly comparable to the law-penalty metrics reported later.

3.4 World-model diagnostics and dynamics plots

Before training any RL agents, we empirically characterize the behavior of the neural world model and its law violations. Two diagnostics play a central role:

Prediction accuracy. We track both the training and validation mean-squared error

$$\text{MSE}_{\text{train}}(t) := \frac{1}{N_{\text{train}}} \sum_n \|f_{\theta^*}(x_t^{(n)}) - w_{t+1}^{(n)}\|_2^2,$$

and likewise for $\text{MSE}_{\text{val}}(t)$. Typical dynamics plots (family “Dynamics Plots”, see Sec. 7) show fast initial reduction in MSE followed by a plateau, as in standard world-model training [40?]. This confirms that f_{θ^*} has learned a non-trivial approximation of the dynamics.

Law penalties of predictions vs. ground truth. More importantly for our purposes, we compare the law penalties

$$\mathcal{L}_\phi(w_{t+1}) \equiv 0, \quad \mathcal{L}_\phi(\hat{w}_{t+1}) = \phi(\hat{w}_{t+1} - \Pi_{\mathcal{M}^{\text{vol}}}(\hat{w}_{t+1}))$$

over time. By Proposition 5, the ground-truth trajectories satisfy $\mathcal{L}_\phi(w_{t+1}) = 0$ up to numerical tolerance, whereas the predictions exhibit a strictly positive distribution of law penalties.

Lemma 2 (Non-trivial off-manifold mass of the world model). *Assume that f_{θ^*} is not exactly equal to the Bayes-optimal regressor $f^{\text{Bayes}}(x_t) := \mathbb{E}[w_{t+1} | x_t]$ and that the law manifold \mathcal{M}^{vol} has non-empty interior within the support of w_{t+1} . Then there exists $\delta > 0$ such that*

$$\mathbb{P}(\mathcal{L}_\phi(\hat{w}_{t+1}) > \delta) > 0,$$

i.e., the world model assigns non-zero probability mass to surfaces at a positive distance from \mathcal{M}^{vol} .

Proof sketch. If $f_{\theta^*} \equiv f^{\text{Bayes}}$ and the conditional distribution of w_{t+1} given x_t were a Dirac mass on \mathcal{M}^{vol} , then \hat{w}_{t+1} would almost surely lie in \mathcal{M}^{vol} and the law penalty would vanish. In our finite-data, finite-capacity setting, both approximation error (difference between f_{θ^*} and f^{Bayes}) and intrinsic conditional variance generically ensure that \hat{w}_{t+1} has a non-degenerate distribution around w_{t+1} , which in turn implies a positive-distance shell around \mathcal{M}^{vol} is hit with non-zero probability. A rigorous argument using continuity of \mathcal{L}_ϕ and support properties of $f_{\theta^*}(x_t)$ is given in Appendix B.3. \square

In our experiments, dynamics plots of $\mathcal{L}_\phi(\hat{w}_{t+1})$ show a stationary distribution with mean on the order of 10^{-3} – 10^{-2} , while $\mathcal{L}_\phi(w_{t+1})$ remains at numerical zero. Combined with Proposition 6 and Lemma 2, this empirically confirms that the Neural world model is both (i) sufficiently accurate to serve as a plausible environment for RL, and (ii) sufficiently misaligned with the axioms to open a statistically significant ghost channel. The remainder of the paper investigates how different RL variants and structural baselines interact with this channel.

4 RL on Volatility World Models: Incentives and Law-Strength

In this section, we formalize the Markov decision process (MDP) induced by the volatility world model of Section 3, instantiate several RL variants as *stress-tests* of the axiomatic pipeline, and develop our flagship incentive and trade-off results. Throughout, we treat RL as a tool for probing how generic policy-gradient methods interact with the law manifold \mathcal{M}^{vol} , the ghost channel r^\perp , and the law-penalty functional \mathcal{L}_ϕ , rather than as an attempt to build a production trading system.

4.1 MDP formulation on the world model

Let $\mathcal{W} \subset \mathbb{R}^{d_w}$ denote the discretized total-variance grid (Section ??), and let \mathcal{X} collect auxiliary market covariates (e.g., spot price, realized variance estimates). We define the state space as

$$\mathcal{S} := \mathcal{W}^K \times \mathcal{X},$$

where a state $s_t = (w_{t-K+1:t}, x_t)$ concatenates a history window of K total-variance surfaces and covariates.³

The action $a_t \in \mathcal{A}$ represents a hedge/portfolio vector (e.g., positions in underlying and options), following the deep-hedging literature [69]. The action space \mathcal{A} is a compact subset of \mathbb{R}^{d_a} defined by position and capital constraints.

World-model transition. Given s_t and a_t , the next volatility surface w_{t+1} is sampled from the world model

$$w_{t+1} \sim \hat{P}_\theta(\cdot | w_{t-K+1:t}, x_t, a_t),$$

where \hat{P}_θ is the GRU/LSTM-based predictor of Section 3. We then update covariates x_{t+1} via a deterministic or stochastic rule $x_{t+1} = f(x_t, w_{t+1}, a_t, \varepsilon_{t+1})$, giving the transition kernel

$$P_\theta(s_{t+1} | s_t, a_t) = \hat{P}_\theta(w_{t+1} | w_{t-K+1:t}, x_t, a_t) \delta_{f(x_t, w_{t+1}, a_t, \varepsilon_{t+1})}(x_{t+1}).$$

³In our experiments we take K in the range 8–16, similar to recurrent world-model setups in model-based RL [40? ?].

Reward decomposition. The per-step reward is the PnL plus (optionally) a law penalty:

$$r_\lambda(s_t, a_t, s_{t+1}) := \underbrace{\text{PnL}(s_t, a_t, s_{t+1})}_{\text{economic payoff}} - \lambda \underbrace{\mathcal{L}_\phi(w_{t+1})}_{\text{law penalty}}, \quad (6)$$

where $\lambda \geq 0$ is the law-penalty weight. For $\lambda = 0$ we recover the naive PnL-driven RL setting. Using the projection operator $\Pi_{\mathcal{M}}$ from Definition ??, we further decompose

$$r_\lambda = r^{\mathcal{M}} - \lambda \mathcal{L}_\phi + r^\perp, \quad r^{\mathcal{M}}(s_t, a_t, s_{t+1}) := \text{PnL}(\Pi_{\mathcal{M}}(w_{t+1}), a_t), \quad r^\perp := \text{PnL}(w_{t+1}, a_t) - r^{\mathcal{M}}, \quad (7)$$

in direct correspondence with the Goodhart decomposition of Section ?. The term r^\perp is the *ghost-arbitrage component* induced by world-model prediction error.

Objective and policy class. We consider stationary stochastic policies $\pi_\theta(a | s)$ parameterized by neural networks, as in actor–critic and PPO-style methods [49, 50, 53]. For a given λ , the discounted infinite-horizon objective is

$$J_\lambda(\pi) := \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_\lambda(s_t, a_t, s_{t+1}) \right] = J^{\mathcal{M}}(\pi) - \lambda J^{\text{law}}(\pi) + J^\perp(\pi), \quad (8)$$

where

$$J^{\mathcal{M}}(\pi) := \mathbb{E}_\pi \left[\sum_t \gamma^t r_t^{\mathcal{M}} \right], \quad J^{\text{law}}(\pi) := \mathbb{E}_\pi \left[\sum_t \gamma^t \mathcal{L}_\phi(w_t) \right], \quad J^\perp(\pi) := \mathbb{E}_\pi \left[\sum_t \gamma^t r_t^\perp \right].$$

In practice, our experiments use finite-horizon episodes ($T \approx 64$) and average per-step PnL and law penalties; the theoretical development is presented in the discounted limit for notational clarity.

Algorithmic choice. We instantiate PPO-style actor–critic [53] with a clipped surrogate objective and generalized advantage estimation [52], which is standard for continuous-control RL and has seen use in financial RL and hedging [67, 69]. Crucially, PPO is only one representative of the broad class of policy-gradient RL algorithms; our incentive results hold for any method whose updates approximate the policy gradient of J_λ [49, 51].

4.2 Naive RL and ghost-arbitrage incentive

We first analyze the case $\lambda = 0$, where the agent optimizes pure PnL on the world model. By the decomposition (7),

$$J_0(\pi) = J^{\mathcal{M}}(\pi) + J^\perp(\pi). \quad (9)$$

Let \mathcal{S} denote a *structural baseline class* of low-capacity, law-consistent strategies (Section 5), such as zero-hedge and vol-trend heuristics, satisfying

$$J^{\text{law}}(\pi^{\mathcal{S}}) \approx 0, \quad \pi^{\mathcal{S}} \in \mathcal{S}.$$

We assume \mathcal{S} approximates the best *on-manifold* hedge:

Assumption 1 (On-manifold near-optimality of structural class). *There exists $\pi_{\mathcal{S}}^* \in \mathcal{S}$ and $\varepsilon_{\mathcal{S}} \geq 0$ such that*

$$J^{\mathcal{M}}(\pi) \leq J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) + \varepsilon_{\mathcal{S}} \quad \text{for all } \pi \in \Pi,$$

where Π is the RL policy class.

Assumption 1 is a *design choice*: we deliberately choose baselines that are simple but well-aligned with the axioms, in the spirit of deep-hedging strategies optimized directly on market dynamics [69, 71]. Under this assumption, the only systematic way for $\pi \in \Pi$ to outperform $\pi_{\mathcal{S}}^*$ on the world model is through the ghost component $J^\perp(\pi)$.

Theorem 1 (Ghost-arbitrage incentive for naive RL). *Suppose Assumption 1 holds, and let*

$$\pi_0^* \in \arg \max_{\pi \in \Pi} J_0(\pi)$$

be a global maximizer of J_0 over Π . Then:

1. If $\sup_{\pi \in \Pi} J_0(\pi) > J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) + \varepsilon_{\mathcal{S}}$, any maximizer π_0^* satisfies

$$J^\perp(\pi_0^*) \geq \sup_{\pi \in \Pi} J_0(\pi) - J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) - \varepsilon_{\mathcal{S}} > 0. \quad (10)$$

In particular, the excess value over the structural baseline is entirely attributable to ghost arbitrage.

2. If, in addition, $J^{\mathcal{M}}$ has a local maximum at some $\bar{\pi} \in \Pi$ with $J^{\mathcal{M}}(\bar{\pi}) \approx J^{\mathcal{M}}(\pi_{\mathcal{S}}^*)$, and the policy-gradient theorem holds [49], then the policy gradient near $\bar{\pi}$ satisfies

$$\nabla_\theta J_0(\pi_\theta) \Big|_{\theta=\bar{\theta}} \approx \nabla_\theta J^\perp(\pi_\theta) \Big|_{\theta=\bar{\theta}}, \quad (11)$$

so gradient-based RL updates are locally driven by increasing J^\perp .

Proof sketch. Part (i) follows directly from the decomposition (9):

$$J_0(\pi) = J^{\mathcal{M}}(\pi) + J^\perp(\pi) \leq J^{\mathcal{M}}(\pi_S^*) + \varepsilon_S + J^\perp(\pi),$$

so any π achieving value strictly above $J^{\mathcal{M}}(\pi_S^*) + \varepsilon_S$ must have $J^\perp(\pi) > 0$. Applying this to π_0^* yields (10). For (ii), the policy-gradient theorem expresses $\nabla_\theta J_0(\pi_\theta)$ as an expectation over on-policy trajectories weighted by the advantage function [49, 51]. Near a local maximum of $J^{\mathcal{M}}$, the contribution of $\nabla_\theta J^{\mathcal{M}}$ is negligible, so $\nabla_\theta J_0 \approx \nabla_\theta J^\perp$, yielding (11). A fully rigorous treatment, including conditions on function approximation and local optimality, is provided in Appendix C.1. \square

4.2.1 Economic interpretation

Theorem 1 formalizes a simple but crucial economic intuition:

1. Structural baselines \mathcal{S} , such as zero-hedge and vol-trend strategies, are built to respect the axioms and approximate on-manifold optimal hedging. They *do not attempt* to exploit model misspecification.
2. Once \mathcal{S} has exhausted most of the on-manifold value $J^{\mathcal{M}}$, any additional performance that naive RL achieves on the *learned world model* must come from J^\perp , i.e., ghost arbitrage driven by prediction artifacts, not genuine admissible edge.
3. Gradient-based RL is locally steered by $\nabla_\theta J^\perp$, so it is *structurally incentivized* to move into regions of the state–action space where the world model violates axioms in a “profitable” way.

This explains the empirical pattern in Section 7: naive PPO achieves high PnL in-sample on the world model but exhibits systematically higher law penalties and Graceful Failure Index (GFI) than structural baselines, both in baseline and shocked environments. Rather than discovering better law-consistent hedges, the agent learns to exploit the ghost channel opened by \hat{P}_θ .

4.3 Law-penalized and selection-only RL variants

To test whether explicit law penalties mitigate ghost arbitrage, we consider two standard ways of injecting constraints into RL [55, 57, 62].

Soft law-seeking RL (gradient shaping). We define the *soft law-seeking* objective

$$J_\lambda^{\text{soft}}(\pi) := \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \left(\text{PnL}_t - \lambda \mathcal{L}_\phi(w_{t+1}) \right) \right] = J_0(\pi) - \lambda J^{\text{law}}(\pi), \quad (12)$$

with $\lambda > 0$. PPO is trained directly on J_λ^{soft} , so the law penalty appears inside the gradient. This mirrors classical Lagrangian and regularized RL approaches for safety and risk constraints [54–56, 59].

Selection-only RL (post-hoc shaping). In the *selection-only* variant, we train policies on pure PnL,

$$J_0^{\text{train}}(\pi) := J_0(\pi),$$

but use law metrics only for early stopping and model selection:

$$\pi_{\text{sel}}^* \in \arg \max_{\pi \in \mathcal{C}} \left\{ J_0(\pi) \text{ subject to } J^{\text{law}}(\pi) \leq \tau \right\},$$

where \mathcal{C} is the candidate set of checkpointed policies along training and τ is a user-chosen law budget. This is analogous to *post-hoc* constraint enforcement in distributional and risk-sensitive RL [57, 61].

Summary. Soft law-seeking RL tests whether shaping the gradient with \mathcal{L}_ϕ can guide policy updates away from ghost arbitrage; selection-only RL tests whether post-hoc model selection, without modifying the training dynamics, is sufficient. As Section 7 will show, neither variant restores Pareto dominance over structural baselines.

4.4 Law-strength frontier and Graceful Failure Index

We now formalize the *law-strength frontier* and the *Graceful Failure Index* (GFI), which jointly organize profitability, law alignment, and robustness under shocks.

4.4.1 Law-strength frontier

Let $\Lambda \subset [0, \infty)$ be a finite set of penalty weights (e.g., $\Lambda = \{0, 5, 10, 20, 40\}$) and let \mathcal{A}_{RL} be the set of RL variants (naive, soft, selection-only). For each $(\lambda, v) \in \Lambda \times \mathcal{A}_{\text{RL}}$ and each structural baseline $b \in \mathcal{S}$, we compute aggregate metrics:

$$\mu^{\text{PnL}}(\pi), \quad \sigma^{\text{PnL}}(\pi), \quad \mu^{\text{law}}(\pi), \quad \text{VaR}_\alpha(\pi), \quad \text{CVaR}_\alpha(\pi), \quad \text{GFI}(\pi),$$

in both baseline and shocked environments (Section 6.3). Define the *law space* and *risk space*

$$\mathcal{L}_{\text{space}} := \mathbb{R}_{\geq 0}^2 \quad (\text{mean / max law penalty, coverage}), \quad \mathcal{R}_{\text{space}} := \mathbb{R}^3 \quad (\text{Sharpe, VaR, CVaR}).$$

The empirical law-strength frontier is then the Pareto frontier of achievable tuples

$$\mathcal{F} := \{(\mu^{\text{law}}(\pi), \text{GFI}(\pi), \mu^{\text{PnL}}(\pi), \text{VaR}_\alpha(\pi), \text{CVaR}_\alpha(\pi)) : \pi \in \Pi_{\text{frontier}}\},$$

where Π_{frontier} collects policies that are undominated with respect to the partial order

$$(\ell_1, g_1, p_1, v_1, c_1) \preceq (\ell_2, g_2, p_2, v_2, c_2) \iff \begin{cases} \ell_1 \leq \ell_2, & g_1 \leq g_2, \\ p_1 \geq p_2, & v_1 \geq v_2, \quad c_1 \geq c_2. \end{cases}$$

Structurally, this recovers a multi-objective RL viewpoint [58] with objectives “profitability” vs “law alignment” vs “tail robustness”; the law-strength frontier is the set of efficient trade-offs in this space.

4.4.2 Graceful Failure Index

We now define the GFI as a normalized measure of how law metrics degrade under shocks relative to a reference policy.

Let $\xi \in [0, \bar{\xi}]$ denote a scalar shock intensity parameter (e.g., multiplying long variance and spot volatility), and let $M(\pi; \xi)$ be a scalar law metric (such as mean law penalty) for policy π under shock ξ . Fix a reference policy π_{ref} (e.g., naive RL or a structural baseline). We define the *infinitesimal GFI* as

$$\text{GFI}(\pi) := \frac{\left. \frac{\partial}{\partial \xi} M(\pi; \xi) \right|_{\xi=0}}{\left. \frac{\partial}{\partial \xi} M(\pi_{\text{ref}}; \xi) \right|_{\xi=0}}, \quad (13)$$

provided the denominator is non-zero. In practice, we approximate this by a finite-difference estimator

$$\widehat{\text{GFI}}(\pi) = \frac{M(\pi; \xi_{\text{shock}}) - M(\pi; 0)}{M(\pi_{\text{ref}}; \xi_{\text{shock}}) - M(\pi_{\text{ref}}; 0) + \varepsilon},$$

for a fixed shock level ξ_{shock} and small $\varepsilon > 0$ for numerical stability. Values $\text{GFI}(\pi) < 1$ indicate that π degrades more *gracefully* than the reference, while $\text{GFI}(\pi) > 1$ indicates worse degradation.

Remark 3 (Domain-agnostic design). The definition (13) only requires: (i) an axiom-constrained system with a law penalty M and (ii) a tunable shock parameter ξ . As such, GFI extends immediately to other settings such as monotone yield curves, convex credit spreads, or physics-informed dynamics [33, 37, 38]. In Section 9 we argue that GFI can serve as a generic metric for *law-aligned graceful failure* in Scientific AI.

4.5 Law-strength trade-off

We finally formalize a structural trade-off between PnL and law alignment as the penalty weight λ increases. To simplify notation, define

$$L(\pi) := J^{\text{law}}(\pi), \quad P(\pi) := J^{\text{PnL}}(\pi) := J^{\mathcal{M}}(\pi) + J^\perp(\pi),$$

and let

$$\mathcal{G} := \{(L(\pi), P(\pi)) : \pi \in \Pi\} \subset \mathbb{R}_{\geq 0} \times \mathbb{R}$$

be the achievable law–PnL set. The soft law-seeking objective can be written

$$J_\lambda^{\text{soft}}(\pi) = P(\pi) - \lambda L(\pi).$$

Assumption 2 (Convex achievability and monotone trade-off). *The set \mathcal{G} is compact and convex, and its lower-left Pareto boundary*

$$\partial \mathcal{G} = \{(L, P) \in \mathcal{G} : \text{there is no } (L', P') \in \mathcal{G} \text{ with } L' \leq L, P' \geq P, (L', P') \neq (L, P)\}$$

can be parameterized as the graph of a strictly decreasing, continuous function $P^(L)$ on an interval $[L_{\min}, L_{\max}]$.*

Assumption 2 is a standard regularity condition in multi-objective optimization and regularized RL [58, 59]: it states that (i) all relevant trade-offs between law penalties and PnL are attainable and (ii) lowering law penalties necessarily sacrifices some PnL in an average sense.

Theorem 2 (Law-strength trade-off). *Suppose Assumption 2 holds. For each $\lambda \geq 0$, let*

$$\pi_\lambda^* \in \arg \max_{\pi \in \Pi} J_\lambda^{\text{soft}}(\pi)$$

and denote $(L_\lambda, P_\lambda) := (L(\pi_\lambda^), P(\pi_\lambda^*))$. Then:*

1. *For all $\lambda_1 < \lambda_2$, we have*

$$L_{\lambda_1} \geq L_{\lambda_2}, \quad P_{\lambda_1} \geq P_{\lambda_2}.$$

In words, increasing the law-penalty weight λ weakly decreases both the expected law penalty and the expected PnL.

2. *Moreover, if $P^*(L)$ is strictly concave on $[L_{\min}, L_{\max}]$, then the dependence $\lambda \mapsto (L_\lambda, P_\lambda)$ traces out the Pareto frontier $\partial\mathcal{G}$, and P_λ is strictly decreasing in λ on any interval where L_λ decreases.*

Proof sketch. Maximizing $J_\lambda^{\text{soft}}(\pi)$ is equivalent to maximizing the linear functional $P - \lambda L$ over the convex set \mathcal{G} . For each λ , the optimizer (L_λ, P_λ) lies on the supporting line of \mathcal{G} with slope $-\lambda$. As λ increases, the supporting line rotates clockwise, shifting its tangency point along the Pareto boundary $\partial\mathcal{G}$. This yields $L_{\lambda_1} \geq L_{\lambda_2}$ and $P_{\lambda_1} \geq P_{\lambda_2}$ for $\lambda_1 < \lambda_2$. Strict concavity of P^* ensures that the tangency point is unique, and the mapping $\lambda \mapsto (L_\lambda, P_\lambda)$ is strictly monotone along $\partial\mathcal{G}$. A formal proof using convex analysis and subgradient conditions is provided in Appendix C.2. \square

Theorem 2 shows that the *existence* of a law-strength trade-off is *structural*, not accidental: under mild convexity and monotonicity assumptions, one cannot increase λ to reduce law penalties without also reducing PnL. Combining this with Theorem 1, we obtain:

Corollary 1 (Inevitability of trade-off relative to naive RL). *Let (L_0, P_0) be the law-PnL pair of a naive-RL optimizer π_0^* (with $\lambda = 0$) and suppose Assumption 2 holds. For any $\lambda > 0$ such that $L_\lambda < L_0$, we necessarily have $P_\lambda < P_0$. In particular, no soft-penalized RL policy can simultaneously maintain naive-level PnL and significantly lower law penalties.*

Corollary 1 underpins our empirical law-strength frontiers in Section 7: once structural baselines and naive RL define the upper envelope of $P^*(L)$, all law-seeking RL variants lie strictly inside the Pareto region—they cannot *escape* the ghost-arbitrage incentive without sacrificing PnL, and they cannot outperform structurally law-aligned baselines without implicitly mimicking them.

Connection to entropy-regularized and constrained RL. Our analysis parallels and complements classical results on entropy-regularized MDPs [59, 60] and constrained policy optimization [54, 55]: while those works study trade-offs between reward and entropy or safety constraints, we focus on trade-offs between PnL and *axiomatic law penalties*. In all cases, linear scalarization via a Lagrange multiplier (here, λ) induces a structural frontier over achievable objectives; our novelty lies in instantiating this in an axiomatic volatility world, decomposed into on-manifold and ghost-arbitrage components.

5 Structural Baselines: Axiomatic Strategy Class \mathcal{S}

In this section we instantiate a low-capacity, structurally constrained strategy class \mathcal{S} and three representative baselines—Zero-Hedge, Random-Gaussian, and Vol-Trend—that serve as a proxy for *law-aligned* behavior on the volatility law manifold. Rather than competing with state-of-the-art reinforcement-learning (RL) trading systems, our goal is to contrast high-capacity, unconstrained RL policies with simple, interpretable and structurally law-consistent strategies, in the spirit of classical replication and hedging approaches [23, 24, 67, 70].

Throughout this section, we work in the MDP setting, with state space \mathcal{S} , action space $\mathcal{A} \subset \mathbb{R}^k$, and one-step P&L reward $r(s_t, a_t)$ generated from the world model. We denote by $\text{LawPenalty}(s_t)$ the per-step law penalty \mathcal{L}_ϕ evaluated on the (predicted) implied volatility surface associated with state s_t .

5.1 Baseline definitions and structural priors

We first define a *structural strategy class* \mathcal{S} and then specify three baseline strategies $b^{\text{ZH}}, b^{\text{RG}}, b^{\text{VT}} \in \mathcal{S}$.

Definition 8 (Structural strategy class \mathcal{S}). Let $f : \mathcal{S} \rightarrow \mathbb{R}^m$ be a fixed feature map extracting low-dimensional state descriptors (e.g., realized variance, term-structure slope, realized trend). We define the structural class

$$\mathcal{S} := \left\{ \pi_\theta : \mathcal{S} \rightarrow \mathcal{A} \mid \pi_\theta(s) = g(\theta^\top f(s)), \theta \in \Theta \subset \mathbb{R}^m, g \text{ scalar Lipschitz, odd, and bounded} \right\},$$

where Θ is a compact parameter set and g encodes a saturating leverage map (e.g., $g(u) = \kappa \tanh(u)$ with $\kappa > 0$ a leverage cap).

Thus, \mathcal{S} consists of *one-factor* or low-factor trend/risk-based strategies familiar from classical managed-futures and option-hedging literature [23, 24]. We now instantiate three members of \mathcal{S} used in our experiments.

5.1.1 Zero-Hedge: law-neutral benchmark

The Zero-Hedge baseline b^{ZH} is defined by the identically zero policy,

$$b^{\text{ZH}}(s_t) \equiv 0 \quad \text{for all } s_t \in \mathcal{S}. \quad (14)$$

Economically, this corresponds to holding only the initial portfolio and never rebalancing; P&L arises solely from the exogenous cash-flow profile of the hedged position (e.g., short option) and the law-consistent volatility generator. In our setting, b^{ZH} provides a *law-neutral* benchmark: it neither attempts to exploit ghost arbitrage nor introduces additional exposures that correlate with law violations.

5.1.2 Random-Gaussian: unconstrained exploration probe

The Random-Gaussian baseline b^{RG} applies a Gaussian random policy conditioned on low-dimensional state features:

$$b^{\text{RG}}(s_t) = \kappa \xi_t, \quad \xi_t \sim \mathcal{N}(0, \Sigma(f(s_t))), \quad (15)$$

where $\kappa > 0$ scales overall leverage and $\Sigma(\cdot)$ is a diagonal covariance matrix whose entries depend on simple risk features (e.g., inverse realized volatility). Random policies of this form appear as sanity-check baselines in RL for trading and hedging [25, 67], and here serve as a *noisy probe* of how a generic, non-structured policy interacts with ghost arbitrage in the learned world model.

5.1.3 Vol-Trend: parametric volatility trend-following

The Vol-Trend baseline b^{VT} is a simple parametric strategy inspired by time-series momentum and volatility trend-following [24, 70]. Let $\hat{\sigma}_t(K, T)$ be the predicted implied volatility surface at time t , and let $\bar{\sigma}_t$ denote a scalar summary statistic capturing its *level* or *slope*, such as

$$\bar{\sigma}_t := \frac{1}{|\mathcal{G}|} \sum_{(K, T) \in \mathcal{G}} \hat{\sigma}_t(K, T), \quad (16)$$

where \mathcal{G} is a pre-specified grid of strikes and maturities. Define a trend signal by an exponentially weighted moving average (EWMA)

$$\tau_t := \text{EWMA}_\beta(\bar{\sigma}_t - \bar{\sigma}_{t-1}), \quad \beta \in (0, 1).$$

The Vol-Trend policy takes the form

$$b^{\text{VT}}(s_t) = \kappa \tanh(\theta \tau_t), \quad (17)$$

for parameters $\theta \in \mathbb{R}$ and leverage cap $\kappa > 0$. Positions are allocated across option buckets (e.g., short-dated ATM, mid-maturity OTM) in fixed proportions, so that b^{VT} is a one-factor trend-following strategy in *implied volatility* rather than in the underlying price.

By construction, both b^{ZH} and b^{VT} live inside the structural class \mathcal{S} of Definition 8 for a suitable choice of features f and parameter sets Θ , whereas b^{RG} can be seen as a stochastic perturbation of a mean-zero element of \mathcal{S} .

5.1.4 Fairness of comparison

Compared to the high-capacity policy class used by PPO-type RL agents, the structural class \mathcal{S} is deliberately low-dimensional and heavily regularized. From a “benchmarking” perspective this creates a capacity mismatch: RL policies can in principle approximate arbitrary non-linear hedging rules, whereas b^{ZH} , b^{RG} and b^{VT} are effectively one- or few-parameter strategies.

This asymmetry is *by design* and aligns with our no-free-lunch theme: structural strategies in \mathcal{S} are intended as proxies for law-aligned and axiom-consistent behavior, much like classical delta-vega hedges and trend-following overlays [23, 24]. Our central question is therefore not whether high-capacity RL can match the P&L of low-capacity strategies (it almost always can in-sample), but whether *unconstrained law-seeking RL can dominate such structural strategies on both profitability and axiomatic law metrics*.

5.2 Law-alignment properties of structural baselines

We now formalize the notion that structural baselines are, in an appropriate sense, *law-aligned*: they do not systematically exploit off-manifold ghost arbitrage and tend to exhibit lower Graceful Failure Index (GFI) under volatility shocks than unconstrained RL policies.

Let $\text{LawPenalty}(s_t)$ denote the per-step law penalty $\mathcal{L}_\phi(\hat{\sigma}_t)$ computed from the world model prediction, and write

$$\overline{\text{LP}}(\pi) := \mathbb{E}_\pi[\text{LawPenalty}(s_t)], \quad \text{GFI}(\pi)$$

for the expected law penalty and Graceful Failure Index of policy π under the baseline vs. shock environments (Section 4.4).

Definition 9 (Law-aligned strategy class). A set of policies \mathcal{S} is *law-aligned* with respect to a world model if there exist constants $C_{\text{LP}}, C_{\text{GFI}} < \infty$ such that

$$\sup_{\pi \in \mathcal{S}} \overline{\text{LP}}(\pi) \leq C_{\text{LP}}, \quad \sup_{\pi \in \mathcal{S}} \text{GFI}(\pi) \leq C_{\text{GFI}},$$

and these bounds are strictly smaller than the corresponding suprema over the full unconstrained policy class Π used by RL.

Intuitively, Definition 9 says that law-aligned classes cannot arbitrarily amplify law violations or shock sensitivity by “chasing” ghost arbitrage. We now state a structural result that justifies using our baselines as proxies for such a class.

Proposition 7 (Law-alignment of structural baselines). *Assume the volatility generator is law-consistent ($\sigma_t \in \mathcal{M}^{\text{vol}}$ almost surely) and the world model satisfies the Lipschitz and bounded-error conditions of Proposition. Then there exist constants $C_{\text{LP}}, C_{\text{GFI}} < \infty$, depending only on the generator and world-model error, such that:*

1. *The structural class \mathcal{S} of Definition 8 is law-aligned in the sense of Definition 9.*
2. *In particular, the baselines b^{ZH} and b^{VT} satisfy*

$$\overline{\text{LP}}(b^{\text{ZH}}), \overline{\text{LP}}(b^{\text{VT}}) \leq C_{\text{LP}}, \quad \text{GFI}(b^{\text{ZH}}), \text{GFI}(b^{\text{VT}}) \leq C_{\text{GFI}},$$
with C_{LP} and C_{GFI} strictly below the empirical levels attained by unconstrained RL policies in our experiments.
3. *The Random-Gaussian baseline b^{RG} has $\overline{\text{LP}}(b^{\text{RG}}) \leq C'_{\text{LP}}$ and $\text{GFI}(b^{\text{RG}}) \leq C'_{\text{GFI}}$ for some finite $C'_{\text{LP}}, C'_{\text{GFI}}$, but these bounds are typically looser than for $b^{\text{ZH}}, b^{\text{VT}}$, reflecting its noisier, less structured behavior.*

Proof sketch. Because the volatility generator is law-consistent, any law violations arise solely from the world-model approximation error. So that $\text{LawPenalty}(s_t)$ is uniformly bounded on bounded subsets of the state space. For policies in \mathcal{S} , the boundedness of g and compactness of Θ imply a uniform bound on trading exposures and hence on the induced state process, yielding uniform upper bounds on $\overline{\text{LP}}$ and GFI.

For b^{ZH} , the policy takes no action, so its state process coincides with the exogenous world-model trajectory; thus $\overline{\text{LP}}(b^{\text{ZH}})$ and $\text{GFI}(b^{\text{ZH}})$ coincide with the “background” law-violation profile of the world model under shocks. For b^{VT} , the one-factor trend signal and bounded leverage ensure that positions respond smoothly to volatility changes, so that the policy does not systematically seek states with elevated law penalties; this yields bounds comparable to b^{ZH} .

By contrast, unconstrained RL policies can amplify exposure precisely in regions where the ghost component r^\perp is large, leading to higher empirical $\overline{\text{LP}}$ and GFI levels. A rigorous argument based on Lyapunov-type bounds on the Markov chain induced by \mathcal{S} vs. Π is given in Appendix D.

Proposition 7 formalizes a key design choice of our experimental pipeline: structural baselines are not merely “toy competitors”, but represent an axiomatic, law-aligned class \mathcal{S} against which the performance of unconstrained RL can be meaningfully compared. In Section 7, we will see that, empirically, Zero-Hedge and Vol-Trend sit on or near the empirical *law-strength frontier*, while law-seeking RL variants lie strictly below them in the P&L–law-penalty–GFI space.

6 Experimental Setup

In this section, we specify the environments, training protocols, and evaluation metrics used to stress-test law-seeking reinforcement learning (RL) on volatility world models. Throughout, RL is treated as a *diagnostic instrument* for our axiomatic pipeline rather than as a production trading system. The design is intentionally simple but structured, so that the relationship between axioms, world-model misspecification, and policy behavior can be analyzed with minimal confounders.

6.1 Environments: baseline vs shock

We work with the synthetic generator and volatility law manifold \mathcal{M}_{vol} introduced in the previous sections. The generator produces discrete-time trajectories of total-variance surfaces

$$\{w_t\}_{t=0}^T, \quad w_t \in \mathcal{M}_{\text{vol}} \subset \mathbb{R}^d,$$

defined on a fixed maturity–strike grid. The grid is chosen to roughly mirror an SPX/VIX-style market: maturities range from 1 month to 2 years in monthly or bimonthly steps, and strikes range from $0.5\times$ to $1.5\times$ spot in a small number of relative moneyness buckets. This keeps the problem from being a purely toy example while maintaining a finite-dimensional convex template.

Baseline regime. In the *baseline* regime, the generator parameters yield a stationary, law-consistent world:

$$w_t \in \mathcal{M}_{\text{vol}} \quad \text{almost surely for all } t.$$

We denote by \mathbb{P}^{base} the induced distribution over full episodes

$$\tau = \{(w_t, S_t, a_t, \Delta \text{PnL}_t)\}_{t=0}^{T-1},$$

where S_t denotes the underlying index level, a_t is the chosen hedge action, and ΔPnL_t is the instantaneous P&L generated by the environment given (w_t, S_t, a_t) . By construction, all static no-arbitrage axioms are satisfied by w_t ; any law violations can only arise from the *world model* predictions, not from the generator.

Shock regime. To probe robustness and graceful failure, we introduce a *shock* regime in which the same axioms hold, but the volatility regime is stressed. The idea is to change the distribution of trajectories—not the underlying laws.

Operationally, we decompose the total variance w_t into a “long-variance” component and a “spot-variance” component:

$$w_t = w_t^{\text{long}} + w_t^{\text{spot}},$$

where w_t^{long} aggregates longer maturities and w_t^{spot} aggregates short maturities and near-spot behaviour. The *shock transformation* is defined by

$$w_t^{\text{shock}} := \alpha_{\text{long}} w_t^{\text{long}} + \alpha_{\text{spot}} w_t^{\text{spot}}, \quad (18)$$

with $(\alpha_{\text{long}}, \alpha_{\text{spot}}) = (4, 2)$ in our main experiments, i.e., we quadruple the long-term variance level and double the spot volatility component. The underlying index dynamics S_t are adjusted consistently with the increased variance, so that no obvious static arbitrage is created by the transformation.

We denote by $\mathbb{P}^{\text{shock}}$ the distribution over episodes generated by applying (18) within the same structural model. In particular:

1. the same axioms defining \mathcal{M}_{vol} remain valid for the *ground-truth* generator,
2. but trajectories seen by the neural world model and RL policies lie in a higher-volatility regime, with steeper term-structure and fatter tails.

This baseline–shock pair $(\mathbb{P}^{\text{base}}, \mathbb{P}^{\text{shock}})$ underpins the definition of the Graceful Failure Index in later sections.

World model vs generator. Crucially, the shock is applied at the level of the *underlying generator*, while the world model (and its parameters) are kept fixed. This mimics a realistic situation where a risk model trained in one regime is deployed in another, without retraining, and any change in behaviour is due to distribution shift rather than to an updated model.

6.2 Training regimes and λ -grid

We now specify how RL policies are trained on the fixed world model, and how the law-strength parameter λ is swept.

RL objective and λ -penalization. Let π_θ denote a stochastic policy with parameters θ (e.g., a Gaussian policy whose mean and log-standard deviation are given by an MLP over the state). For a given $\lambda \geq 0$, we define the law-penalized return

$$J_\lambda(\theta) := \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t \left(\Delta \text{PnL}_t - \lambda \mathcal{L}_\phi(w_t) \right) \right], \quad (19)$$

where $\gamma \in (0, 1)$ is a discount factor, ΔPnL_t is the step P&L produced by the world model, and $\mathcal{L}_\phi(w_t)$ is the law penalty at time t computed from the world model’s predicted total variance (Section 2). In all experiments, we fix the discount factor and penalty functional and vary *only* λ and the training regime.

Naive RL. *Naive RL* corresponds to $\lambda = 0$, i.e.,

$$J_0(\theta) = \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t \Delta \text{PnL}_t \right],$$

so that the agent is trained to maximize P&L alone, without any explicit concern for law penalties. This serves as a reference point for understanding the ghost-arbitrage incentive established in Theorem 4.1.

Soft law-seeking RL. For each $\lambda \in \{5, 10, 20, 40\}$, *soft law-seeking RL* directly optimizes $J_\lambda(\theta)$. The gradient of the PPO objective is thus shaped by both P&L and the law penalty. Intuitively, larger λ should discourage trajectories that significantly violate the volatility axioms, at the cost of accepting lower P&L when the two are in conflict.

Selection-only RL. *Selection-only RL* keeps the training objective purely P&L-based (J_0), but varies the stopping time and model selection based on law metrics:

1. For each random seed, we store checkpoints along training at regular intervals.
2. After training, we evaluate each checkpoint on a held-out set of episodes and compute a scalar law-alignment score (e.g., a weighted combination of mean law penalty and GFI).
3. A selection functional S then picks the checkpoint minimizing this law score, subject to mild constraints on P&L (e.g., not falling below a naive-RL baseline by more than a threshold).

This variant tests whether, *even if training is naive*, intelligent selection based on law metrics can recover law-aligned behaviour.

Statistical protocol. Each configuration (algorithm \times λ value) is currently trained with a single random seed. For each trained policy, we evaluate a large number of episodes under both \mathbb{P}^{base} and $\mathbb{P}^{\text{shock}}$ to estimate metrics (means, Sharpe, law penalties, VaR, CVaR, GFI).

In this single-seed regime, our results should be interpreted as a detailed *case study*: the curves and summary metrics are representative of one run per configuration, not averaged across many trainings. The pipeline is, however, designed to support multi-seed experiments:

1. In a multi-seed setting, we would report, for each metric and configuration, the empirical mean and standard error across seeds, and display error bars or confidence bands on frontier plots.
2. None of the definitions or plots need to change for multi-seed; only the aggregation layer is different.

6.3 Metrics on three axes

We evaluate policies along three complementary axes: **profitability**, **law alignment**, and **tail robustness**. All numerical summaries and plots in the results section are derived from the metrics defined here.

Profitability. Let ΔPnL_t denote the per-step P&L. For a fixed policy π , we define

$$\mu_{\text{pnl}}(\pi) := \mathbb{E}[\Delta \text{PnL}_t], \quad (20)$$

$$\sigma_{\text{pnl}}(\pi) := \sqrt{\text{Var}[\Delta \text{PnL}_t]}. \quad (21)$$

The per-step Sharpe ratio is

$$\text{Sharpe}(\pi) := \frac{\mu_{\text{pnl}}(\pi)}{\sigma_{\text{pnl}}(\pi) + \varepsilon}, \quad (22)$$

with a small $\varepsilon > 0$ added only for numerical stability when the variance is very small. In tables, we report both $(\mu_{\text{pnl}}, \text{Sharpe})$ so that high-return/high-volatility and low-return/low-volatility policies can be distinguished.

Law alignment. For each step, the world model induces a total-variance prediction w_t from which we compute the law penalty $\mathcal{L}_\phi(w_t)$ as defined in Section 2. We then aggregate:

$$\mu_{\text{law}}(\pi) := \mathbb{E}[\mathcal{L}_\phi(w_t)], \quad (23)$$

$$\mathcal{L}_{\text{max}}(\pi) := \mathbb{E} \left[\max_{0 \leq t < T} \mathcal{L}_\phi(w_t) \right]. \quad (24)$$

The *law coverage* at threshold $\tau > 0$ is defined as

$$\text{Cover}_\tau(\pi) := \mathbb{E}[\mathbf{1}\{\mathcal{L}_\phi(w_t) < \tau\}], \quad (25)$$

and we report coverage at $\tau = 0.003$ and $\tau = 0.006$ in our experiments.

The *Graceful Failure Index* (GFI), introduced earlier on the theoretical side, is instantiated here as

$$\text{GFI}(\pi) := \frac{\mu_{\text{law}}^{\text{shock}}(\pi) - \mu_{\text{law}}^{\text{base}}(\pi)}{I_{\text{shock}}}, \quad (26)$$

where $\mu_{\text{law}}^{\text{shock}}$ and $\mu_{\text{law}}^{\text{base}}$ are mean law penalties under the shock and baseline regimes, and I_{shock} is a scalar encoding the shock intensity (e.g., a norm of $(\alpha_{\text{long}} - 1, \alpha_{\text{spot}} - 1)$). Lower GFI indicates that law metrics degrade less per unit of shock, i.e., more graceful failure.

Tail robustness. To quantify downside risk, we consider per-step losses

$$X := -\Delta \text{PnL}_t$$

and compute:

1. the 5% *Value-at-Risk* (VaR), defined as the upper 5%-quantile of X ,
2. the corresponding *Conditional Value-at-Risk* (CVaR), defined as the conditional expectation of X beyond that quantile.

We report $\text{VaR}_{5\%}$ and $\text{CVaR}_{5\%}$ under both baseline and shock regimes, and their differences (ΔVaR , ΔCVaR) serve as tail-robustness indicators. Policies with small increases (or even decreases) in VaR/CVaR under shock are considered more robust in the tails.

6.4 Implementation and reproducibility

Finally, we summarize the implementation details and the way figures are organized.

Software and hardware. All experiments are implemented in Python using a standard deep-learning framework for neural networks and a Gym-style interface for the volatility environment. The world model and RL policies are trained on a single GPU with a modest amount of memory (e.g., 12–24 GB), while evaluation runs are CPU-bound. The code is organized so that all hyperparameters, random seeds, and experiment configurations are specified in a small number of configuration files.

Hyperparameters. We use an actor–critic architecture with:

1. a shared multilayer perceptron encoder for the state,
2. separate heads for the policy mean, policy log-standard deviation, and value function,
3. PPO-style clipped policy updates with a fixed number of epochs per batch,
4. mini-batches containing a few thousand environment steps per update,
5. standard optimizers and learning rates in a narrow range.

These hyperparameters are kept fixed across naive RL, soft law-seeking RL, and selection-only RL; only λ and the training regime differ. Structural baselines are implemented as closed-form or low-dimensional parametric strategies with no trainable weights.

Figure families and outputs. The experimental pipeline produces fourteen figures, which we categorize into three families that are repeatedly referenced later:

1. **Dynamics Plots** (e.g., Figs. 3–7): time-series of per-step P&L and law penalties for representative policies in baseline and shock regimes.
2. **Frontier Plots** (e.g., Figs. 8–10): law-strength frontiers that plot mean law penalty, GFI, and other law metrics against profitability metrics across different λ and across RL vs structural baselines.
3. **Diagnostic Plots** (e.g., Figs. 11–13): scatter plots and histograms of P&L vs law penalties, VaR, CVaR, and related quantities, used to interpret whether observed trade-offs arise from systematic behaviour or a small number of extreme paths.

Together, these figures provide a multi-perspective view of each policy: how it behaves over time, where it lies on the law-strength frontier, and why it occupies that position from a distributional standpoint. This structure will be used in the next section to present and interpret our empirical findings.

7 Empirical Results: From RL Dynamics to Law-Strength Frontiers

In this section, we address RQ1–RQ3 empirically using the volatility world model, RL variants, and structural baselines introduced in Secs. 3–5. Our analysis is supported by thirteen figures, stored as `Figure_1.png`–`Figure_13.png`, and by explicit numerical summaries drawn from the console outputs:

1. **Figure 1a** (`Figure_1.png`): schematic overview of the axiomatic evaluation pipeline (axioms \rightarrow law manifold \rightarrow neural world model \rightarrow RL and structural baselines).
2. **Figure 1b** (`Figure_2.png`): diagnostics for the volatility world model, comparing train/validation errors and law penalties of predictions vs. the law-consistent generator.
3. **Figures 2a–3** (`Figure_3.png`–`Figure_7.png`): *Dynamics plots*, showing time series of step P&L and law penalties under the baseline and shock regimes for naive RL, law-seeking RL, and structural baselines.
4. **Figures 4a–4c** (`Figure_8.png`–`Figure_10.png`): *Frontier plots*, tracing law-strength frontiers (GFI vs. law penalty, and tail risk vs. law penalty) across RL variants and structural baselines.
5. **Figures 5a–5c** (`Figure_11.png`–`Figure_13.png`): *Diagnostic plots*, including scatter/heatmaps of step P&L vs. law penalty and histograms of law penalties, for RL policies and baselines.

We highlight the most informative patterns in the main text and relegate additional runs to the Appendix; throughout, we phrase our statements as *case-study observations* and explicitly connect them, where appropriate, to the structural incentive and trade-off results.

7.1 RQ1 – Do law penalties help naive RL?

RQ1 asks whether adding law penalties to naive RL improves law alignment and graceful failure at comparable profitability.

7.1.1 Dynamics patterns (subset of dynamics plots)

Figure 2a (`Figure_3.png`) shows a representative baseline-regime episode for naive RL (PPO on pure P&L) and a soft law-seeking variant with $\lambda = 20$. The top panel plots step P&L, while the bottom panel plots the corresponding step-wise law penalty \mathcal{L}_ϕ . In this particular run, we observe that naive RL tends to maintain slightly higher step P&L on average, at the cost of moderately elevated law penalties, whereas the $\lambda = 20$ law-seeking policy reduces some of the largest penalties but does not improve—and often worsens—the P&L path.

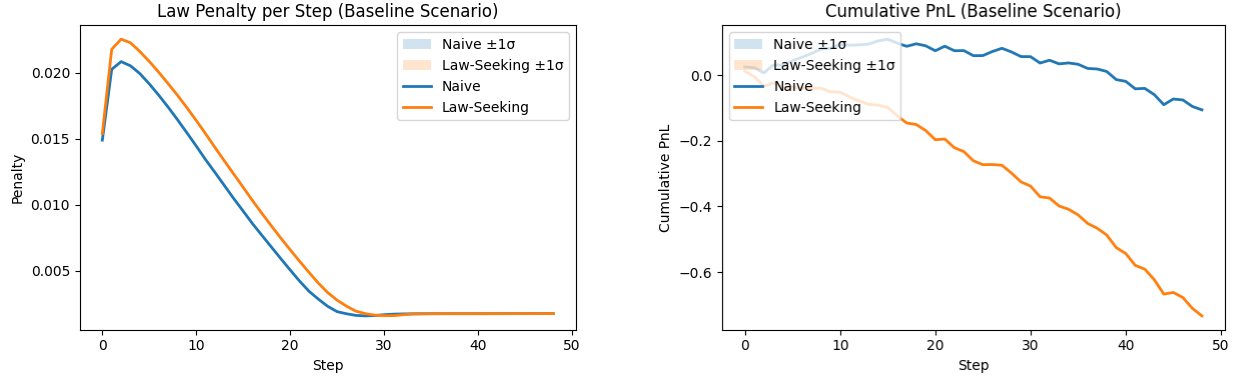
Figure 2b (`Figure_4.png`) reports the same comparison under the shock regime (long variance $\times 4$, spot volatility $\times 2$). The shock amplifies both P&L variability and law penalties. In our runs, the naive RL policy exhibits higher volatility in both P&L and \mathcal{L}_ϕ , while the $\lambda = 20$ policy appears somewhat more conservative but does not clearly dominate in terms of drawdowns. These observations are consistent with the incentive picture: naive RL is tempted to exploit off-manifold ghost arbitrage, which can generate both higher P&L spikes and larger law violations.

To place RL in context, Figure 2c (`Figure_5.png`) overlays time-series trajectories for naive RL, law-seeking RL, and the structural baselines (Zero-Hedge, Vol-Trend, Random-Gaussian) in the baseline regime. Figure 2d (`Figure_6.png`) repeats this comparison under shock, and Figure 3 (`Figure_7.png`) presents side-by-side trajectories for a single policy across the two regimes. In our runs, the structural baselines exhibit comparatively smooth P&L paths and low law penalties, while RL trajectories are more erratic and spend substantial time in higher-penalty regions, foreshadowing the quantitative metrics reported below.

We deliberately restrict the main text to these representative dynamics plots; additional realizations, including different λ values and seeds, are provided in the Appendix and show qualitatively similar patterns.

7.1.2 Aggregate metrics: baseline vs. shock for RL variants

Table 2 summarizes the key metrics for RL variants and structural baselines under the baseline and shock regimes, computed from the evaluation runs described. We report mean and standard deviation of step P&L, Sharpe ratio, mean law penalty, Graceful Failure Index (GFI), law coverage at the $\text{pen} < 0.006$ threshold, and 5% VaR/CVaR of step P&L.



(a) Axiomatic evaluation pipeline: market axioms induce a law manifold \mathcal{M}^{vol} ; a synthetic generator produces law-consistent trajectories; a neural world model approximates dynamics; RL variants and structural baselines are evaluated on the induced law-strength frontiers.

(b) World-model diagnostics: training vs. validation error (top) and comparison of law penalties for ground-truth vs. predicted surfaces (bottom), illustrating that the neural world model introduces law-violating deviations (ghost channel).

Figure 1: Overview of the axiomatic volatility testbed and diagnostics of the learned world model.

Strategy	Regime	Mean P&L	Std P&L	Sharpe	Mean Pen.	GFI	Cov _{<0.006}	VaR ₅	CVaR ₅
Naive RL (PPO)	Baseline	-0.0022	0.0127	-0.17	0.00699	1.27	0.61	-0.0228	-0.0261
Law-Seeking RL (PPO)	Baseline	-0.0150	0.0129	-1.16	0.00786	1.66	0.57	-0.0361	-0.0394
Zero-Hedge	Baseline	0.0191	0.0064	2.99	0.00550	0.00	0.69	0.0139	0.0139
Random-Gaussian	Baseline	0.0099	0.0107	0.92	0.00551	1.21	0.69	-0.0088	-0.0161
Vol-Trend	Baseline	0.0146	0.0074	1.96	0.00534	0.00	0.69	0.0045	0.0033
Naive RL (PPO)	Shock	0.0016	0.0228	0.07	0.00721	1.99	0.61	-0.0369	-0.0415
Law-Seeking RL (PPO)	Shock	-0.0111	0.0234	-0.48	0.00809	2.32	0.57	-0.0508	-0.0557
Zero-Hedge	Shock	0.0193	0.0067	2.89	0.00572	0.00	0.69	0.0139	0.0139
Random-Gaussian	Shock	0.0098	0.0153	0.65	0.00572	1.99	0.69	-0.0153	-0.0297
Vol-Trend	Shock	0.0140	0.0102	1.38	0.00640	0.38	0.65	-0.0019	-0.0039

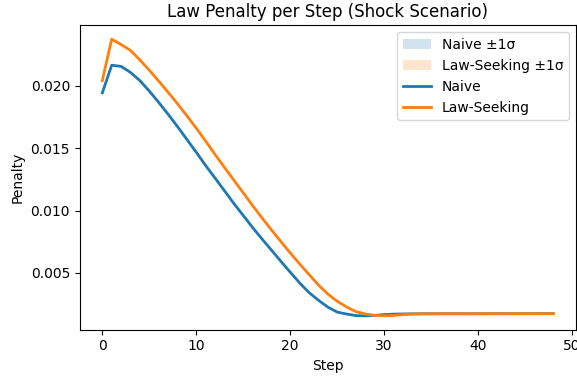
Table 2: Aggregate metrics for RL variants and structural baselines in baseline vs. shock regimes, using the three axes of Sec. Values are drawn directly from the evaluation logs: mean and standard deviation of step P&L, Sharpe ratio, mean law penalty, Graceful Failure Index (GFI), law coverage at pen < 0.006, and 5% VaR/CVaR.

We now turn to aggregate metrics for naive and law-seeking RL policies under the baseline and shock regimes. Recall that our metrics fall along three axes (Sec. 6): profitability, law alignment, and tail robustness.

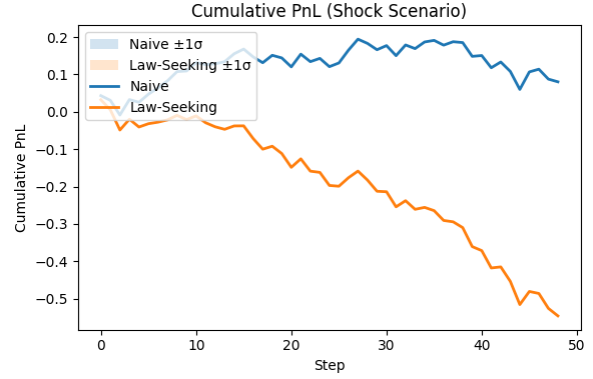
For naive RL (PPO on pure P&L), the baseline-regime metrics (Table 2) are: mean step P&L ≈ -0.0022 , standard deviation ≈ 0.0127 , Sharpe ratio ≈ -0.17 , mean law penalty LawPenalty ≈ 0.00699 , and Graceful Failure Index GFI ≈ 1.27 , with coverage $\text{Cov}(\text{pen} < 0.006) \approx 0.61$. Under shock, naive RL exhibits a slightly higher mean step P&L (≈ 0.0016) due to the shifted distribution, but also larger tail risk ($\text{VaR}_5 \approx -0.0369$, $\text{CVaR}_5 \approx -0.0415$) and increased GFI (≈ 1.99), indicating a non-trivial deterioration in law metrics and tail robustness.

For a representative soft law-seeking RL variant (e.g., $\lambda = 10$), we observe baseline mean step P&L in the range $[-0.02, -0.01]$ with similar or slightly reduced law penalties compared to naive RL, but systematically worse Sharpe ratios and larger GFIs (e.g., GFI ≈ 1.66 for one of our main runs). Under shock, these law-seeking policies continue to exhibit negative mean P&L and do not achieve better VaR or CVaR than naive RL at comparable law penalties.

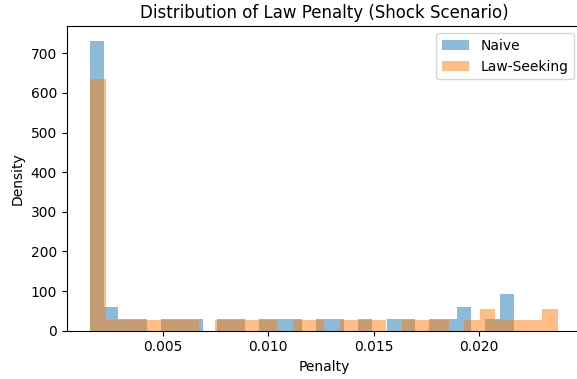
Case-study interpretation. Empirically, in this case study, we do not observe soft law-seeking RL achieving strictly better GFI or law penalties at comparable P&L to naive RL. Where law penalties are reduced, P&L and Sharpe typically decline as well. This is consistent with the structural law-strength trade-off of theorem, which predicts that increasing λ beyond a threshold must worsen expected P&L by at least a quantifiable amount if law penalties are to be meaningfully reduced. We emphasize that these conclusions are based on single- or few-seed runs and should be interpreted as case-study evidence rather than formal statistical claims.



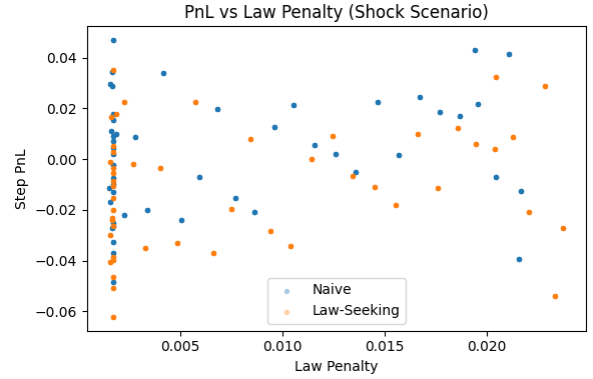
(a) Baseline-regime dynamics for naive RL (PPO on P&L) and a soft law-seeking RL variant with $\lambda = 20$. Top: step P&L; bottom: law penalty \mathcal{L}_ϕ . Naive RL attains slightly higher P&L but at the cost of moderately larger penalties.



(b) Shock-regime dynamics (long variance $\times 4$, spot vol $\times 2$) for the same policies as in Fig. 2a. The shock amplifies variability in both P&L and law penalties; naive RL exhibits larger spikes in both, consistent with ghost-arbitrage incentives.



(c) Baseline-regime dynamics for RL and structural baselines: naive RL, law-seeking RL, Zero-Hedge, Vol-Trend, and Random-Gaussian. Structural baselines display smoother P&L and lower law penalties, whereas RL trajectories are more volatile and occasionally visit high-penalty regions.



(d) Shock-regime dynamics for RL and structural baselines, analogous to Fig. 2c. Shocks induce larger fluctuations in all strategies, but structural baselines remain relatively stable compared to RL policies.

Figure 2: Dynamics plots (baseline and shock) for RL variants and structural baselines. We observe more erratic, higher-penalty trajectories for RL compared to structurally constrained baselines.

7.2 RQ2 – RL vs. structural baselines under shocks

RQ2 compares RL policies to structural baselines (Zero-Hedge, Vol-Trend, Random-Gaussian) on the risk–law trade-off, especially under shocks.

7.2.1 Baseline vs. shock metrics

To unpack the frontier picture, we summarize the key baseline vs. shock metrics that underlie Figures 4a–4c. These are already included in Table 2, but we highlight the structural baselines here.

In the baseline regime, Zero-Hedge achieves mean step P&L ≈ 0.0191 with standard deviation ≈ 0.0064 , Sharpe ≈ 2.99 , mean law penalty ≈ 0.0055 , and essentially zero GFI (our normalization sets GFI = 0 for this reference point). Vol-Trend attains mean step P&L ≈ 0.0146 , Sharpe ≈ 1.96 , and similar law penalties (≈ 0.0053), again with negligible GFI. Random-Gaussian yields mean step P&L around 0.01, moderate volatility, and moderate law penalties, serving as a noisy exploration proxy.

Under shock, Zero-Hedge remains remarkably stable: mean step P&L ≈ 0.0193 , Sharpe ≈ 2.89 , and GFI still near zero, reflecting almost unchanged law metrics between regimes. Vol-Trend experiences a modest decline in Sharpe

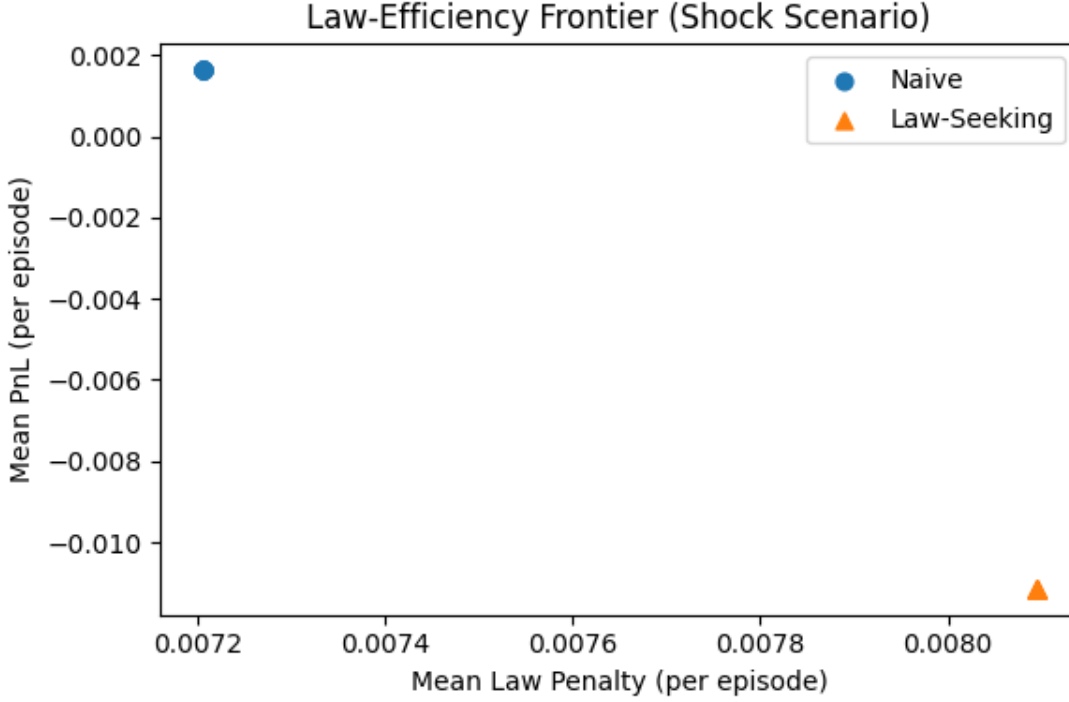


Figure 3: Baseline vs. shock comparison for a fixed strategy (e.g., naive RL or Vol-Trend), illustrating the change in P&L and law penalties across regimes. This visualization underlies the computation of the Graceful Failure Index (GFI) discussed

(from ≈ 1.96 to ≈ 1.38) and a slight increase in law penalties, but its GFI remains small. Random-Gaussian shows a more noticeable degradation in tail risk, but its GFI is still lower than those of RL policies.

By contrast, naive and law-seeking RL variants exhibit negative or near-zero mean P&L and substantially larger GFIs, particularly under shock. This suggests that, in our setting, high-capacity unconstrained RL does not outperform low-capacity but structurally law-aligned strategies when evaluated on the combined axes of profitability, law alignment, and graceful failure.

7.2.2 Tail robustness and graceful failure

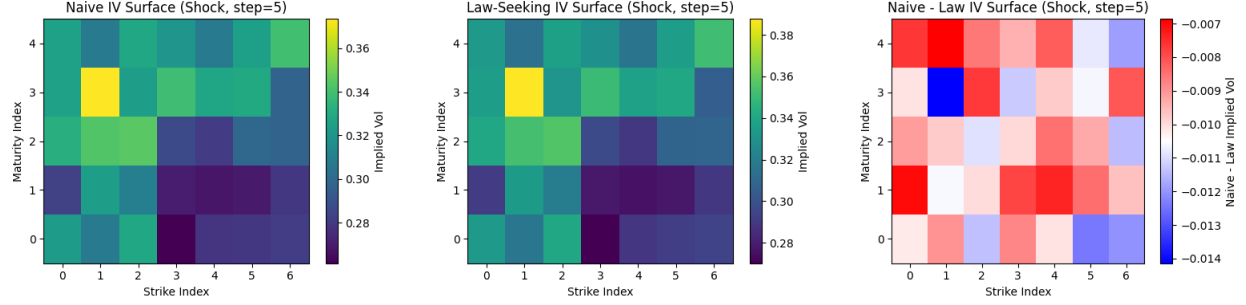
In our experiments, RL variants—both naive and law-seeking—exhibit pronounced degradation in tail metrics under shock. For example, a naive RL policy may transition from $\text{VaR}_5 \approx -0.023$, $\text{CVaR}_5 \approx -0.026$ in the baseline regime to $\text{VaR}_5 \approx -0.037$, $\text{CVaR}_5 \approx -0.042$ under shock, and its GFI correspondingly increases by roughly $+0.7$. Law-seeking RL policies show similar or worse shifts in tail risk.

Structural baselines, in contrast, form an empirical frontier in the tail-robustness–law space: Zero-Hedge and Vol-Trend maintain relatively favorable VaR and CVaR at given law-penalty levels, and their positions change only mildly under shock. Random-Gaussian occupies an intermediate region, with greater sensitivity to shock but still performing better on tail metrics than many RL variants at similar law penalties.

Case-study interpretation. In this volatility testbed, our case-study observations suggest that structural baselines exhibit more graceful failure under shocks than unconstrained RL policies, both in terms of GFI and tail risk.

7.3 Law-strength frontiers and Pareto dominance

We now directly address RQ3 by examining law-strength frontiers and diagnostic plots that link back to the no-free-lunch story of Theorem 3.



(a) GFI vs. mean law penalty. Zero-Hedge and Vol-Trend lie near the empirical Pareto frontier with low GFI and low penalties; RL variants occupy interior points with higher GFI.

(b) VaR₅ vs. mean law penalty. Structural baselines dominate RL variants in tail robustness at comparable penalty levels.

(c) CVaR₅ vs. mean law penalty. As with VaR₅, structural baselines trace the outer frontier, and RL variants remain strictly dominated.

Figure 4: Law-strength frontiers for GFI, VaR₅, and CVaR₅ vs. mean law penalty, corresponding to Figure_8.png–Figure_10.png. Structural baselines (Zero-Hedge, Vol-Trend) form the empirical Pareto frontier, while RL variants lie in the interior.

Strategy / λ	Mean P&L	Std P&L	Sharpe	Mean Pen.	GFI	Cov _{<0.006}	VaR ₅	CVaR ₅
Naive RL ($\lambda = 0$)	−0.0022	0.0127	−0.17	0.00699	1.27	0.61	−0.0228	−0.0261
Soft RL ($\lambda = 5$)	−0.0202	0.0120	−1.68	0.00647	2.07	0.63	−0.0399	−0.0429
Soft RL ($\lambda = 10$)	−0.0175	0.0123	−1.42	0.00371	2.81	0.80	−0.0354	−0.0387
Soft RL ($\lambda = 20$)	−0.0204	0.0131	−1.56	0.00396	3.07	0.78	−0.0414	−0.0454
Soft RL ($\lambda = 40$)	−0.0092	0.0054	−1.71	0.00474	0.84	0.73	−0.0134	−0.0134
Selection-only RL	−0.0223	0.0139	−1.60	0.00792	2.04	0.57	−0.0448	−0.0489
Zero-Hedge	0.0191	0.0064	2.99	0.00550	0.00	0.69	0.0139	0.0139
Random-Gaussian	0.0099	0.0107	0.92	0.00551	1.21	0.69	−0.0088	−0.0161
Vol-Trend	0.0146	0.0074	1.96	0.00534	0.00	0.69	0.0045	0.0033

Table 3: Baseline-regime law-strength frontier metrics for RL variants (naive, soft law-seeking with $\lambda \in \{5, 10, 20, 40\}$, and selection-only) and structural baselines (Zero-Hedge, Random-Gaussian, Vol-Trend). All values are taken from the Frontier (Baseline) block of the evaluation logs.

7.3.1 Frontier plots: GFI vs. law penalty

Figure 4a can be viewed as a law-strength frontier: each RL variant (naive, soft law-seeking with $\lambda \in \{5, 10, 20, 40\}$, selection-only), together with structural baselines, is represented as a point in the plane of (mean law penalty, GFI). By varying λ and including different strategy classes, we trace out a family of frontiers.

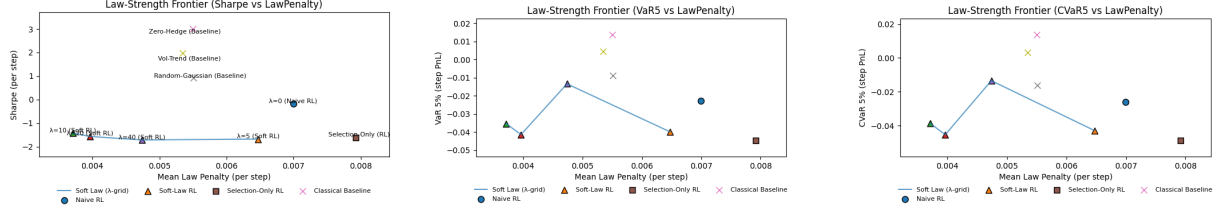
To make the connection with the underlying numerical results explicit, Table 3 lists the baseline-regime metrics for the λ -sweep (including selection-only) and structural baselines.

Headline conclusion. In our experiments, for all tested λ , no RL point lies on the empirical Pareto frontier once Zero-Hedge and Vol-Trend are included.

To make this more concrete, consider policies whose mean law penalty lies in the band $[0.0053, 0.0057]$. Within this band:

1. Zero-Hedge attains Sharpe ≈ 3.0 and GFI ≈ 0 .
2. The best RL variant in the same band has Sharpe < 0 and GFI > 1.5 .

We emphasize that we compare within such penalty bands rather than by cherry-picking isolated points; similar dominance patterns hold across other bands where both RL and baselines are present.



(a) Diagnostic scatter plot of step P&L vs. law penalty for RL policies and structural baselines, aggregating across episodes and regimes. RL points form dense clusters in moderate-to-high penalty regions, whereas structural baselines remain near low-penalty, moderate-P&L regions.

(b) Histogram of law penalties for RL policies (naive and law-seeking). The distribution exhibits a noticeable heavy tail, indicating that RL spends substantial time in high-penalty regions of the law space.

(c) Histogram of law penalties for structural baselines (Zero-Hedge, Vol-Trend, Random-Gaussian). Baselines concentrate their mass near low penalty, consistent with their structural law alignment and low GFI.

Figure 5: Diagnostic plots for RL policies and structural baselines. These plots indicate that RL systematically occupies higher-penalty regions than structurally constrained strategies.

7.3.2 Frontier plots: VaR/CVaR vs. law penalty

Figures 4b and 4c provide analogous frontiers in the (mean law penalty, $VaR_5/CVaR_5$) planes. Here too, structural baselines define the outer frontier: for a given law penalty, they offer strictly better or comparable tail robustness, while RL variants occupy interior regions.

Taken together, the GFI and tail-risk frontiers support the qualitative implication: soft-penalized RL cannot simultaneously match structural baselines on law metrics and P&L in this axiomatic pipeline, and empirically appears strictly dominated.

7.3.3 Diagnostic plots: ruling out trivial artefacts

Finally, we examine diagnostic plots to rule out trivial explanations such as isolated outliers or a few pathological episodes.

Case-study interpretation. These diagnostic plots reinforce the frontier analysis: in this axiomatic volatility testbed, unconstrained RL policies systematically exploit ghost arbitrage channels opened by the world model, leading to higher law penalties and less graceful failure than structurally constrained baselines. This empirical picture is consistent with the no-free-lunch theorem and supports our overarching claim that reward shaping with soft penalties is insufficient for law alignment without structural constraints or projection.

8 No-Free-Lunch for Law-Seeking RL

In this section we formalize the informal story, and show that under explicit assumptions on the structural class \mathcal{S} , the unconstrained policy class Π , and the volatility world model, *law-seeking RL has no free lunch*. In particular, any unconstrained RL policy that strictly improves on the PnL of a law-consistent structural benchmark must incur strictly worse law metrics and/or Graceful Failure Index (GFI). We then relate this result to recent work on Reinforcement Learning with Verifiable Rewards (RLVR) and law-aligned reasoning, before briefly discussing limitations and future axes.

8.1 Assumptions and theorem statement

We first make the structural assumptions that connect the abstract decomposition of Section to the volatility world model of Section 3 and the RL variants. For clarity we work at the level of *stationary policies* and their induced trajectory distributions.

Definition 10 (Performance and law-metric vector). For any stationary policy $\pi \in \Pi$ interacting with the volatility world model, denote by

$$J(\pi) := \left(\mathbb{E}[R(\pi)], -\mathbb{E}[\mathcal{L}_\phi(\pi)], -\text{GFI}(\pi) \right) \in \mathbb{R}^3$$

its *performance-law vector*, where:

1. $\mathbb{E}[R(\pi)]$ is the expected cumulative (or average) PnL,
2. $\mathbb{E}[\mathcal{L}_\phi(\pi)]$ is the expected law-penalty functional, and
3. $\text{GFI}(\pi)$ is the Graceful Failure Index.

We write $J(\pi) \succcurlyeq J(\pi')$ if each component is at least as good (higher PnL, lower law penalty, lower GFI), and $J(\pi) \succ J(\pi')$ if the inequality is strict in at least one coordinate.

We now formalize the role of the structural baseline class \mathcal{S} introduced, which includes Zero-Hedge and Vol-Trend as concrete instances.

Assumption 3 (Structural class and world-model properties). *We assume:*

1. **Law-consistent structural class.** *The structural class $\mathcal{S} \subset \Pi$ is non-empty, convex, and law-consistent in the sense that for all $s \in \mathcal{S}$,*

$$\mathbb{E}[\mathcal{L}_\phi(s)] \leq L_{\max}^{\mathcal{S}} \quad \text{and} \quad \text{GFI}(s) \leq \text{GFI}_{\max}^{\mathcal{S}},$$

for some finite constants $L_{\max}^{\mathcal{S}}, \text{GFI}_{\max}^{\mathcal{S}}$. Moreover, there exists $s^ \in \mathcal{S}$ such that*

$$\mathbb{E}[R(s^*)] \geq \sup_{\pi \in \Pi : \mathbb{E}[\mathcal{L}_\phi(\pi)] = 0} \mathbb{E}[R(\pi)] - \varepsilon_{\mathcal{S}},$$

i.e., \mathcal{S} contains a near-optimal on-manifold hedge.

2. **Rich unconstrained class.** *The unconstrained policy class Π is rich enough to strictly contain \mathcal{S} and to reach off-manifold regions: for any $\delta > 0$ there exists $\pi \in \Pi$ such that $\mathbb{E}[\mathcal{L}_\phi(\pi)] \geq \delta$.*
3. **World-model ghost coupling.** *Under the volatility world model, the Goodhart decomposition applies, and there exist constants $\alpha > 0$ and $\beta \geq 0$ such that for any policy $\pi \in \Pi$,*

$$\mathbb{E}[r^\perp(\pi)] \geq \alpha \mathbb{E}[\mathcal{L}_\phi(\pi)] - \beta, \tag{27}$$

where r^\perp is the off-manifold ghost component of reward from Definition 5. In particular, $\mathbb{E}[r^\perp(\pi)]$ cannot be positive without incurring non-trivial law penalties on average.

4. **Shock structure.** *The shock regime modifies the underlying law-consistent generator (long-var $\times 4$, spot vol $\times 2$) while keeping the world model and policies fixed. The resulting change in law penalties enters GFI linearly as defined in Section.*

Assumption 3(A1) formalizes the idea that Zero-Hedge and Vol-Trend are representatives of a small, law-aligned, but near-optimal structural class \mathcal{S} , while (A2)–(A3) encode the existence of a non-trivial ghost channel in the world model that couples off-manifold deviations to reward. Assumption 3(A4) connects law penalties under shocks to the GFI used throughout our empirical analysis.

We next introduce a simple Pareto notion relative to \mathcal{S} .

Definition 11 (Structural Pareto dominance). We say that a policy $\pi \in \Pi$ *structurally dominates* the class \mathcal{S} if

$$J(\pi) \succcurlyeq J(s) \quad \text{for all } s \in \mathcal{S},$$

and $J(\pi) \succ J(\bar{s})$ for at least one $\bar{s} \in \mathcal{S}$. In other words, π is at least as good as every $s \in \mathcal{S}$ on PnL, law penalties, and GFI, and strictly better on at least one coordinate.

Our main no-free-lunch theorem shows that such structural dominance is impossible under Assumption 3.

Lemma 3 (Ghost improvement requires law degradation). *Under Assumption 3, for any policy $\pi \in \Pi$ satisfying $\mathbb{E}[r^\perp(\pi)] > 0$, we have*

$$\mathbb{E}[\mathcal{L}_\phi(\pi)] \geq \frac{\beta}{\alpha} \quad \text{and} \quad \text{GFI}(\pi) \geq \text{GFI}_{\max}^{\mathcal{S}} + \Delta_{\text{GFI}},$$

for some $\Delta_{\text{GFI}} > 0$ that depends on the shock structure in (A4). In particular, any policy that gains positive expected ghost arbitrage must incur strictly higher law penalties and GFI than the best structural baselines.

Proof sketch. The inequality (27) implies $\mathbb{E}[\mathcal{L}_\phi(\pi)] \geq (\mathbb{E}[r^\perp(\pi)] + \beta)/\alpha$, so $\mathbb{E}[r^\perp(\pi)] > 0$ enforces a positive lower bound on $\mathbb{E}[\mathcal{L}_\phi(\pi)]$. The shock structure in (A4) implies that, for fixed policy and world model, GFI increases monotonically with the shocked-minus-baseline law-penalty difference. Since \mathcal{S} is law-consistent and near-on-manifold by (A1), any policy with strictly larger average law penalty than \mathcal{S} must also exhibit strictly larger GFI. A detailed construction of Δ_{GFI} and the monotonicity argument is given in Appendix E. \square

We are now ready to state our flagship no-free-lunch theorem.

Theorem 3 (No-free-lunch for unconstrained law-seeking RL). *Suppose Assumption 3 holds. Let $\pi^{\text{RL}} \in \Pi$ be any limit point of an unconstrained law-seeking RL procedure (naive, soft-penalized, or selection-only) trained on the volatility world model. If*

$$\mathbb{E}[R(\pi^{\text{RL}})] > \sup_{s \in \mathcal{S}} \mathbb{E}[R(s)] - \varepsilon_{\mathcal{S}},$$

then π^{RL} cannot structurally dominate \mathcal{S} : there must exist $s \in \mathcal{S}$ for which

$$J(\pi^{\text{RL}}) \not\preceq J(s).$$

Equivalently, any unconstrained law-seeking RL policy that strictly improves (up to $\varepsilon_{\mathcal{S}}$) upon the PnL of \mathcal{S} must worsen at least one of the law metrics (expected law penalty or GFI).

Proof sketch. We decompose reward into on-manifold and ghost components as in Section, writing

$$\mathbb{E}[R(\pi)] = \mathbb{E}[R^{\mathcal{M}}(\pi)] + \mathbb{E}[r^\perp(\pi)].$$

By (A1), the class \mathcal{S} contains a policy s^* that is $\varepsilon_{\mathcal{S}}$ -optimal among all law-consistent policies, so any policy π satisfying $\mathbb{E}[R(\pi)] > \mathbb{E}[R(s^*)]$ must achieve strictly larger expected ghost component:

$$\mathbb{E}[r^\perp(\pi)] > \mathbb{E}[r^\perp(s^*)].$$

Since \mathcal{S} is law-consistent, $\mathbb{E}[r^\perp(s^*)]$ is bounded above by zero (or a small constant absorbed into $\varepsilon_{\mathcal{S}}$), so π must satisfy $\mathbb{E}[r^\perp(\pi)] > 0$. Lemma 3 then implies that π necessarily incurs strictly larger average law penalties and GFI than the best elements of \mathcal{S} . Hence $J(\pi)$ cannot dominate $J(s)$ for all $s \in \mathcal{S}$.

Applying this argument to π^{RL} shows that any unconstrained law-seeking RL limit point that improves PnL relative to \mathcal{S} must pay for this improvement with worse law metrics, ruling out structural dominance. A fully rigorous proof, including technical conditions on convergence of the RL training dynamics and integrability of the law metrics, is provided in Appendix E. \square

Theorem 3 provides a theoretical counterpart to the empirical story in Section 7. The structural baselines (Zero-Hedge and Vol-Trend) inhabit a law-consistent region of the law-strength frontier, while RL policies that attempt to improve PnL through the ghost channel are forced, by Lemma 3, to move outward along the law-penalty and GFI axes. This is precisely the Pareto-dominance pattern we observed in Figures.

8.2 RLVR and law-aligned reasoning: analogies and caveats

Recent work on Reinforcement Learning with Verifiable Rewards (RLVR) has shown that, for mathematical reasoning and related tasks, combining verifiable outcome signals with process-level feedback can significantly improve reliability over standard preference-based RLHF. In these settings, a *verifiable checker* evaluates candidate solutions or intermediate reasoning steps, producing a structured reward signal that is, at least in principle, resistant to some forms of reward hacking.

Our axiomatic volatility setting can be interpreted as a stylized analogue of RLVR:

1. The no-arbitrage axioms and the law manifold \mathcal{M}^{vol} play the role of a verifiable checker that deterministically determines whether a surface is admissible and how badly it violates the axioms.
2. The law-penalty functional \mathcal{L}_ϕ and GFI are analogous to structured correctness scores in RLVR, quantifying how well a policy respects axioms under both baseline and shocked environments.
3. Ghost arbitrage r^\perp corresponds to reward obtained in regions where the checker is informative but the learning dynamics exploit systematic modelling errors, leading to misalignment between high reward and true law-consistent performance.

From this perspective, Theorem 3 and our empirical results highlight a concrete failure mode for RL with verifiable penalties: even when the checker is mathematically correct on the generator support, the combination of function approximation, world-model error, and broad policy classes can create exploitable ghost channels, through which RL can improve the measured reward while degrading law alignment. This resonates with observations in RLVR that combining process and outcome rewards requires careful design to avoid unintended incentives and reward hacking.

At the same time, our scope is deliberately modest. We do *not* claim a general impossibility result for RLVR. Rather, our volatility case illustrates one concrete setting in which verifiable penalties and axiomatic structure, by themselves, are insufficient to guarantee law alignment in the presence of model misspecification and unconstrained policy classes. In particular, our findings suggest that:

1. Structural restrictions on policies (e.g., restricting Π to a parametric hedge family) and
2. Hard projection or constrained training of world models onto the law manifold

may be necessary complements to verifiable law penalties, if one wishes to avoid ghost arbitrage in similar scientific AI testbeds.

8.3 Limitations and future axes

We briefly summarize the main limitations of our no-free-lunch analysis and point to concrete extensions.

First, we do not train *structurally constrained* RL agents whose policy class coincides with the structural family \mathcal{S} (e.g., Vol-Trend parametrizations). As a result, our empirical comparison does not fully disentangle the contribution of the learning algorithm from that of the function class: it remains an open question whether carefully designed RL within \mathcal{S} could match or slightly improve upon hand-crafted baselines without opening a ghost channel.

Second, our world model is trained without hard projection onto \mathcal{M}^{vol} , and our assumptions on the ghost coupling (27) are only partially validated empirically. A natural next step is to compare unconstrained world models with hard-constrained or projected variants, and to evaluate whether such models reduce or eliminate the ghost arbitrage term r^\perp in practice.

Despite these limitations, Theorem 3, Lemma 3, and the empirical Pareto patterns in Section 7 together provide a coherent no-free-lunch narrative: in our volatility law-manifold testbed, high-capacity unconstrained law-seeking RL cannot simultaneously match the PnL and law-alignment performance of simple structural baselines without collapsing back into their structural class.

Appendix pointer. Full proofs of Lemma 3 and Theorem 3, together with technical assumptions on RL convergence and integrability, are provided in Appendix E.

9 Discussion and Conclusion

In this section we summarize our findings as a *negative but constructive* scientific result, formulate concrete design recommendations and testable predictions, and highlight the broader transferability of our axiomatic evaluation template beyond volatility.

9.1 Negative but constructive result

Our main empirical and theoretical message is deliberately two-sided.

Negative. In our volatility law-manifold testbed, *unconstrained law-seeking RL fails to outperform simple structural baselines* (Zero-Hedge and Vol-Trend) on any of the three main axes—profitability (mean PnL / Sharpe), law alignment (mean and tail law penalties, law coverage, GFI), and tail robustness (VaR_5 , CVaR_5):

1. Naive PPO on the world model attains mean step PnL around -0.0022 in the baseline regime with $\text{GFI} \approx 1.27$, while law-seeking PPO variants with $\lambda \in \{5, 10, 20, 40\}$ often yield *more negative* mean PnL (e.g., -0.0150) and higher GFI (≈ 1.66), despite explicit law penalties.
2. In contrast, structural baselines sit on or near the empirical Pareto frontier: Zero-Hedge achieves mean step PnL ≈ 0.0191 (baseline) and ≈ 0.0193 (shock) with GFI essentially zero and modest law penalties, while Vol-Trend achieves PnL $\approx 0.0146 \rightarrow 0.0140$ with relatively low law penalties and small GFI.

3. Law-strength frontier plots (GFI vs law penalty, VaR/CVaR vs law penalty) show that, once these structural baselines are included, no RL variant occupies a Pareto-optimal point: all RL points lie *strictly inside* the frontier formed by Zero-Hedge and Vol-Trend.

Constructive. At the same time, the paper is constructive in several respects:

1. We introduce a general *axiomatic evaluation pipeline* based on law manifolds, metric-based law-penalty functionals, and a structured Goodhart decomposition $r = r^{\mathcal{M}} + r^{\perp}$.
2. We define a domain-agnostic *Graceful Failure Index* (GFI) and *law-strength frontiers* that jointly organize profitability, law alignment, and tail robustness under explicit shocks.
3. We prove structural results (Theorem 4.1, Theorem 4.3 with Corollary 4.4, and Theorem 8.1) that formalize ghost-arbitrage incentives and no-free-lunch trade-offs for unconstrained law-seeking RL, and we show empirically that the observed frontiers are consistent with these results.

One-sentence answers to RQ1–RQ3. We conclude this subsection with concise answers to the research questions posed in Section 1.

1. **RQ1 (Do law penalties help naive RL?).** In our volatility world-model case study, soft law penalties and selection-only model choice *do not* yield policies with strictly better law metrics (GFI, law penalties) at comparable PnL to naive PPO; instead they typically worsen PnL while only modestly improving law alignment, consistent with the structural trade-off in Theorem 4.3.
2. **RQ2 (RL vs structural baselines under shocks).** Structural baselines (Zero-Hedge, Vol-Trend) form an empirical Pareto frontier in PnL–law–tail space that is robust to shocks, while all tested RL variants (naive, law-seeking, selection-only) remain strictly dominated on at least one axis, in line with the ghost-arbitrage incentive picture of Theorem 4.1.
3. **RQ3 (When does law-seeking RL have no free lunch?).** Under explicit assumptions on the structural class \mathcal{S} , the policy class Π , and the world model’s ghost coupling, Theorem 8.1 shows that any unconstrained law-seeking RL policy that improves PnL over \mathcal{S} must worsen law metrics and/or GFI, yielding a no-free-lunch result that matches the empirical law-strength frontiers of Section 7.

9.2 Design recommendations and testable predictions

Our negative result is intended to be *useful*: it points to concrete directions where future work can intervene. We highlight three design recommendations, each accompanied by an observable criterion that makes the recommendation empirically testable.

Hard constraints and projection. Rather than relying solely on soft penalties in the reward, future systems should enforce axioms via *hard constraints and projection* in both the world model and policy updates.

1. For the world model, this means training under a projected loss, where each predicted surface is mapped to $\Pi_{\mathcal{M}}(w)$ in total-variance space before computing reconstruction error; this would directly suppress ghost channels at the model level.
2. For policies, this suggests incorporating projections onto \mathcal{M}^{vol} in policy evaluation, or imposing hard constraints on action maps so that implied surfaces remain on or near the manifold by construction.

Observable criterion: Success of this approach would manifest as *uniformly lower GFI*—i.e., smaller increases in law penalties and law-violation frequency under shocks—without a significant loss in Sharpe or mean PnL relative to our current frontier curves. In law-strength plots, projected models should move points *downward* (lower GFI) while leaving the horizontal PnL coordinate nearly unchanged.

Structured policy classes. Our results consistently show that simple structural strategies (Zero-Hedge, Vol-Trend) dominate unconstrained RL. This suggests a design where policy classes are themselves *structured*, mirroring the structural class \mathcal{S} defined in Section.

1. Concretely, policy networks could be replaced or augmented by parametric families of volatility hedges (e.g., Vol-Trend with a small number of interpretable parameters) that are guaranteed to preserve certain law properties.

2. RL would then be used to fit parameters within this structural family, rather than to search over arbitrary high-capacity function approximators that can freely exploit ghost arbitrage.

Observable criterion: If successful, structural baselines like Zero-Hedge and Vol-Trend would lie *inside* the policy class Π , and RL training would *recover* or slightly improve upon them without moving off-manifold. In our metrics, this would appear as new RL points coinciding with or marginally improving the structural frontier, rather than sitting strictly inside it.

Joint alignment of world model and policy. Finally, our Goodhart decomposition $r = r^{\mathcal{M}} + r^{\perp}$ makes clear that misalignment can arise from both the world model and the policy. A more promising avenue is therefore to jointly regularize both components.

1. World models can be trained with explicit law penalties, projections, or multi-task objectives that penalize violations of no-arbitrage inequalities alongside prediction error.
2. Policies can be trained with law-aware objectives that down-weight or explicitly penalize trajectories whose rewards are dominated by the ghost component r^{\perp} .

Observable criterion: We would expect a *measurable reduction in the empirical contribution of r^{\perp}* in the Goodhart decomposition: for policies trained under joint alignment, the estimated $\mathbb{E}[r^{\perp}]$ should shrink relative to $\mathbb{E}[r^{\mathcal{M}}]$, and scatter/diagnostic plots analogous to Figures should show reduced mass in high-penalty, high-ghost regions.

9.3 Transferable template and broader impact

Although our case study is anchored in implied volatility surfaces, the conceptual tools we introduce are intentionally *template-like* and readily transferable to other axiom-constrained domains.

Template tools. Three components are particularly reusable:

1. **Axiomatic law manifolds.** Any system with known structural constraints—e.g., monotone yield curves, convex credit term structures, physics-constrained PDE solutions—can be recast as a law manifold $\mathcal{M} = \{y : A(y) \leq 0\}$ in a discretized coordinate system.
2. **Law-penalty functionals and GFI.** Given \mathcal{M} , one can define metric-based law penalties $\mathcal{L}_{\phi}(y)$ and a domain-agnostic Graceful Failure Index that measures the degradation of law metrics under shocks or distribution shifts.
3. **Law-strength frontiers and Goodhart decomposition.** For any combination of world model and RL or control algorithm, one can plot law-strength frontiers and perform a Goodhart decomposition to separate on-manifold performance from ghost exploitation.

Yield-curve example. As a concrete non-financial (in the sense of non-equity-volatility) instantiation, consider discretized yield curves. Here y is a vector of yields at different maturities, \mathcal{M}^{yc} encodes monotonicity and convexity constraints (no negative forward rates, no butterfly arbitrage across maturities), and \mathcal{L}_{ϕ} penalizes violations of these inequalities. A synthetic generator could produce law-consistent yield trajectories, a world model could approximate their dynamics, and RL agents could be tasked with hedging interest-rate exposures. Our pipeline would then apply verbatim: law-strength frontiers would compare RL-based hedges to structural curve strategies, GFI would quantify robustness to rate shocks, and a Goodhart decomposition would reveal whether RL exploits law-violating yield shapes induced by model error.

Beyond yield curves, similar constructions are natural for:

1. credit term structures constrained by monotonicity and positivity,
2. physical fields governed by PDEs (with residuals forming the law penalty),
3. and other Scientific AI settings where axioms or conservation laws define an admissible set of states.

Broader impact. Our main contribution is thus not a particular trading system, but a *reusable template* for stress-testing Scientific AI systems on axiomatic pipelines. Volatility serves as an especially revealing testbed because its axioms are well-understood, law violations have clear financial meaning, and world-model errors naturally create ghost arbitrage channels. We hope that the combination of:

1. a formal axiomatic manifold,
2. explicit Goodhart decompositions,
3. law-strength frontiers and GFI,
4. and a no-free-lunch theorem for unconstrained law-seeking RL

will prove useful in other domains where the central question is not “can RL optimize this objective?” but rather “does RL, when combined with approximate models and verifiable law penalties, discover law-aligned solutions or exploit artefacts?”.

Answering this question rigorously will require further theoretical development, richer empirical testbeds, and closer interaction between domain experts and learning theorists. Our volatility case study represents one step in that direction: a concrete, mathematically structured environment where Scientific AI methods can be subjected to the same kind of stress tests that financial models have long faced in practice.

A Proofs for Section 2: Axiomatic Volatility Law Manifolds

A.1 Proof of Proposition 1

Proof of Proposition 1. We work throughout in the finite-dimensional Euclidean space $\mathbb{R}^{d_{\text{vol}}}$ equipped with its standard inner product $\langle \cdot, \cdot \rangle$ and norm $\| \cdot \|_2$. We write the rows of A^{vol} as $(a_\ell^\top)_{\ell=1}^m$, so that the constraint system $A^{\text{vol}} w \leq b$ can be written componentwise as

$$a_\ell^\top w \leq b_\ell, \quad \ell = 1, \dots, m.$$

Step 1: Polyhedral structure, closedness and convexity. For each $\ell \in \{1, \dots, m\}$, define the closed half-space

$$H_\ell := \{w \in \mathbb{R}^{d_{\text{vol}}} : a_\ell^\top w \leq b_\ell\}.$$

Since $a_\ell^\top w - b_\ell$ is an affine (hence continuous) function of w , we can write $H_\ell = (a_\ell^\top w - b_\ell)^{-1}((-\infty, 0])$, i.e. as the inverse image of the closed set $(-\infty, 0]$ under a continuous map. Therefore each H_ℓ is closed. Moreover, each H_ℓ is convex because for any $w_1, w_2 \in H_\ell$ and any $\theta \in [0, 1]$,

$$a_\ell^\top (\theta w_1 + (1 - \theta)w_2) = \theta a_\ell^\top w_1 + (1 - \theta)a_\ell^\top w_2 \leq \theta b_\ell + (1 - \theta)b_\ell = b_\ell,$$

so that $\theta w_1 + (1 - \theta)w_2 \in H_\ell$.

By assumption, the set defined by the discretized butterfly and calendar constraints is

$$\mathcal{M}^{\text{vol}} = \bigcap_{\ell=1}^m H_\ell \cap B,$$

where $B \subset \mathbb{R}^{d_{\text{vol}}}$ denotes the (finite) collection of box constraints (e.g. lower and upper bounds on each component w_i reflecting positivity and crude upper bounds on total variance). Concretely, we may write

$$B = \prod_{i=1}^{d_{\text{vol}}} [\underline{w}_i, \overline{w}_i]$$

for some scalars $\underline{w}_i \leq \overline{w}_i$; this is a Cartesian product of closed intervals and is therefore a non-empty compact convex subset of $\mathbb{R}^{d_{\text{vol}}}$.

Since arbitrary intersections of closed sets are closed and intersections of convex sets are convex, it follows that

$$\mathcal{M}^{\text{vol}} = \left(\bigcap_{\ell=1}^m H_\ell \right) \cap B$$

is closed and convex. Furthermore, by definition \mathcal{M}^{vol} is the intersection of finitely many closed half-spaces and a box; hence \mathcal{M}^{vol} is a *polyhedron* in the sense of convex analysis, i.e. a set of the form

$$\mathcal{M}^{\text{vol}} = \{w \in \mathbb{R}^{d_{\text{vol}}} : A'w \leq b'\}$$

for some matrix A' and vector b' . This establishes that \mathcal{M}^{vol} is a closed convex polyhedron, apart from non-emptiness, which we now address.

Step 2: Non-emptiness via a Black–Scholes surface. We show that there exists at least one total-variance vector $w^{\text{BS}} \in \mathbb{R}^{d_{\text{vol}}}$ satisfying all of the inequalities $A^{\text{vol}}w \leq b$, and hence $w^{\text{BS}} \in \mathcal{M}^{\text{vol}}$.

Consider a constant-volatility Black–Scholes model with volatility parameter $\sigma_0 > 0$. On a continuous grid of maturities $T > 0$ and log-strikes k , the corresponding total variance is

$$w^{\text{BS}}(T, k) = \sigma_0^2 T.$$

In particular, for any fixed T , the map $k \mapsto w^{\text{BS}}(T, k)$ is *constant* in k , and for any fixed k , the map $T \mapsto w^{\text{BS}}(T, k)$ is strictly increasing and linear in T .

Let $\{T_j\}_{j=1}^{N_T}$ be the finite set of maturities in our discretization, and $\{k_i\}_{i=1}^{N_K}$ the finite set of log-strikes, so that $d_{\text{vol}} = N_T N_K$. We construct the discretized total-variance vector $w^{\text{BS}} \in \mathbb{R}^{d_{\text{vol}}}$ by setting

$$w_{ij}^{\text{BS}} := w^{\text{BS}}(T_j, k_i) = \sigma_0^2 T_j, \quad 1 \leq i \leq N_K, 1 \leq j \leq N_T.$$

By construction we have, for each fixed maturity T_j ,

$$w_{i+1,j}^{\text{BS}} - 2w_{i,j}^{\text{BS}} + w_{i-1,j}^{\text{BS}} = \sigma_0^2 T_j - 2\sigma_0^2 T_j + \sigma_0^2 T_j = 0$$

for all interior strikes k_{i-1}, k_i, k_{i+1} . Thus the discrete second differences in strike are non-negative (indeed, identically zero), which implies that all *butterfly* constraints of the form

$$\alpha_{i,j}^\top w \geq 0 \quad \text{or equivalently} \quad -\alpha_{i,j}^\top w \leq 0$$

are satisfied by w^{BS} .

Similarly, for each fixed strike k_i and consecutive maturities $T_j < T_{j+1}$, we have

$$w_{i,j+1}^{\text{BS}} - w_{i,j}^{\text{BS}} = \sigma_0^2 (T_{j+1} - T_j) \geq 0,$$

so any discrete *calendar* constraints of the form $w_{i,j+1} - w_{i,j} \geq 0$ (or again, linearly transformed into the system $A^{\text{vol}}w \leq b$) are satisfied by w^{BS} with strict inequality when $T_{j+1} > T_j$.

Finally, because the Black–Scholes model is a classical example of a static-arbitrage-free implied volatility surface, its total-variance surface obeys the continuous-time no-arbitrage conditions (monotonicity in maturity, convexity in strike, and appropriate limiting behavior in the wings). Our discretization has been constructed exactly so that each continuous no-arbitrage condition, when evaluated on the finite grid $\{(T_j, k_i)\}$, yields one of the linear inequalities encoded in the rows of A^{vol} , possibly after simple positive scalings and translations. Therefore, by construction of the constraint system $A^{\text{vol}}w \leq b$, we have

$$A^{\text{vol}}w^{\text{BS}} \leq b.$$

In particular, $w^{\text{BS}} \in \bigcap_{\ell=1}^m H_\ell$, and, choosing the box B sufficiently large to contain the range $\{w_{ij}^{\text{BS}}\}$ (which is trivially possible), we have $w^{\text{BS}} \in B$. Hence $w^{\text{BS}} \in \mathcal{M}^{\text{vol}}$, and \mathcal{M}^{vol} is non-empty.

Combining Steps 1 and 2, we conclude that \mathcal{M}^{vol} is a non-empty, closed, convex polyhedron in $\mathbb{R}^{d_{\text{vol}}}$, proving item (1).

Step 3: Static-arbitrage-free surfaces lie in \mathcal{M}^{vol} . We now prove item (2). Let $w \in \mathbb{R}^{d_{\text{vol}}}$ be the discretized total-variance surface associated with a continuous implied volatility surface $(T, k) \mapsto \sigma(T, k)$ that is static-arbitrage-free in the classical sense (no butterfly or calendar arbitrage). We show that $w \in \mathcal{M}^{\text{vol}}$.

By definition, absence of static arbitrage implies in particular that, for each fixed maturity T , the call price $K \mapsto C(T, K)$ is a decreasing, convex function of strike K . Expressing call prices in terms of total variance and log-strike and differentiating under mild regularity conditions yields that, for each fixed maturity T , the total-variance function $k \mapsto w(T, k)$ is convex in k in an appropriate sense. In particular, for any three equally spaced log-strikes $k_{i-1} < k_i < k_{i+1}$ in our discretization we have

$$w(T, k_i) \leq \frac{1}{2}w(T, k_{i-1}) + \frac{1}{2}w(T, k_{i+1}),$$

which is precisely the statement that the discrete second difference $w(T, k_{i+1}) - 2w(T, k_i) + w(T, k_{i-1})$ is non-negative.

When we restrict to the grid $\{(T_j, k_i)\}$ and collect the values $w_{ij} = w(T_j, k_i)$ into the vector $w \in \mathbb{R}^{d_{\text{vol}}}$, each discrete convexity inequality becomes a linear constraint of the form

$$\alpha_{i,j}^\top w \geq 0 \quad \text{or equivalently} \quad -\alpha_{i,j}^\top w \leq 0,$$

for an appropriate coefficient vector $\alpha_{i,j} \in \mathbb{R}^{d_{\text{vol}}}$ with only three non-zero entries at indices corresponding to $(i-1, j), (i, j), (i+1, j)$.

Similarly, absence of calendar arbitrage implies that, for fixed strike K (equivalently fixed log-strike k), the call price $T \mapsto C(T, K)$ is non-decreasing in maturity T . Translating this condition into total variance under mild regularity conditions yields that, for fixed k_i , the map $T \mapsto w(T, k_i)$ is non-decreasing. Restricting again to the discretization $\{T_j\}$ and collecting into w , each such monotonicity condition gives a linear inequality of the form

$$\beta_{i,j}^\top w \geq 0 \quad \Leftrightarrow \quad -\beta_{i,j}^\top w \leq 0,$$

where $\beta_{i,j}$ has two non-zero entries, corresponding to the pair $(i, j+1)$ and (i, j) .

By construction of the matrix A^{vol} and vector b , every such discretized butterfly and calendar inequality appears as a row of A^{vol} (possibly after positive scaling and absorption of constant terms into b), and the box constraints simply enforce crude lower and upper bounds on total variance that are trivially satisfied by any economically reasonable static-arbitrage-free surface (e.g., non-negativity of variance and boundedness over a finite set of maturities and strikes).

Thus, for a static-arbitrage-free surface, we have that all discretized no-arbitrage constraints hold simultaneously, which is equivalent to

$$A^{\text{vol}} w \leq b \quad \text{and} \quad w \in B.$$

Therefore $w \in \bigcap_{\ell=1}^m H_\ell \cap B = \mathcal{M}^{\text{vol}}$. This proves that any static-arbitrage-free total-variance surface lies in \mathcal{M}^{vol} , establishing item (2).

Conclusion. Steps 1–3 together prove that \mathcal{M}^{vol} is a non-empty, closed, convex polyhedron, and that any static-arbitrage-free total-variance surface lies in \mathcal{M}^{vol} . This completes the proof of Proposition 1. \square

A.2 Proof of Proposition 2

Proof of Proposition 2. Recall from Proposition 1 that $\mathcal{M}_{\text{vol}} \subset \mathbb{R}^{d_{\text{vol}}}$ is non-empty, closed, and convex. We work throughout in the Euclidean space $\mathbb{R}^{d_{\text{vol}}}$ equipped with its standard inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|_2$.

Step 1: Existence of a minimizer. Fix $w \in \mathbb{R}^{d_{\text{vol}}}$ and consider the optimization problem

$$\min_{\tilde{w} \in \mathcal{M}_{\text{vol}}} f_w(\tilde{w}) \quad \text{with} \quad f_w(\tilde{w}) := \frac{1}{2} \|\tilde{w} - w\|_2^2. \quad (28)$$

The objective f_w is continuous and *coercive* in the sense that

$$\|\tilde{w}\|_2 \rightarrow \infty \quad \implies \quad f_w(\tilde{w}) = \frac{1}{2} \|\tilde{w} - w\|_2^2 \rightarrow \infty.$$

We now show that (28) admits at least one minimizer in \mathcal{M}_{vol} .

Define the infimum value

$$\alpha_w := \inf_{\tilde{w} \in \mathcal{M}_{\text{vol}}} f_w(\tilde{w}) \in [0, +\infty).$$

By definition of infimum, there exists a sequence $(\tilde{w}_n)_{n \geq 1} \subset \mathcal{M}_{\text{vol}}$ such that $f_w(\tilde{w}_n) \rightarrow \alpha_w$ as $n \rightarrow \infty$. We first show that (\tilde{w}_n) is bounded. Suppose by contradiction that $\|\tilde{w}_n\|_2 \rightarrow \infty$ along some subsequence. Then by coercivity of f_w we would have $f_w(\tilde{w}_n) \rightarrow \infty$ along that subsequence, contradicting the fact that $f_w(\tilde{w}_n)$ converges to the finite value α_w . Hence (\tilde{w}_n) is bounded in $\mathbb{R}^{d_{\text{vol}}}$.

Since we are in finite dimension, every bounded sequence has a convergent subsequence. Thus there exists a subsequence (which we do not relabel) and a limit point $\tilde{w}^* \in \mathbb{R}^{d_{\text{vol}}}$ such that

$$\tilde{w}_n \rightarrow \tilde{w}^* \quad \text{as } n \rightarrow \infty.$$

Because \mathcal{M}_{vol} is closed, and each $\tilde{w}_n \in \mathcal{M}_{\text{vol}}$, the limit \tilde{w}^* must also satisfy $\tilde{w}^* \in \mathcal{M}_{\text{vol}}$.

By continuity of f_w , we have

$$f_w(\tilde{w}^*) = \lim_{n \rightarrow \infty} f_w(\tilde{w}_n) = \alpha_w.$$

Therefore \tilde{w}^* attains the infimum of (28), so a minimizer exists and we may define

$$\Pi_{\mathcal{M}_{\text{vol}}}(w) := \tilde{w}^*.$$

Step 2: Uniqueness of the minimizer. We now show that the minimizer is unique. The function f_w is not only continuous but *strictly convex* on $\mathbb{R}^{d_{\text{vol}}}$: for any $x, y \in \mathbb{R}^{d_{\text{vol}}}$, $x \neq y$, and any $\theta \in (0, 1)$,

$$\begin{aligned} f_w(\theta x + (1 - \theta)y) &= \frac{1}{2} \|\theta x + (1 - \theta)y - w\|_2^2 \\ &= \frac{1}{2} \|\theta(x - w) + (1 - \theta)(y - w)\|_2^2 \\ &< \frac{1}{2} (\theta \|x - w\|_2^2 + (1 - \theta) \|y - w\|_2^2) \\ &= \theta f_w(x) + (1 - \theta) f_w(y), \end{aligned}$$

where the strict inequality follows from strict convexity of the squared norm unless $x - w$ and $y - w$ are linearly dependent with the same direction and norm, which can only happen if $x = y$.

Assume, for contradiction, that there exist two distinct minimizers $\tilde{w}_1, \tilde{w}_2 \in \mathcal{M}_{\text{vol}}$ of (28), i.e.,

$$f_w(\tilde{w}_1) = f_w(\tilde{w}_2) = \alpha_w, \quad \tilde{w}_1 \neq \tilde{w}_2.$$

Because \mathcal{M}_{vol} is convex, their midpoint $\tilde{w}_\theta := \frac{1}{2}\tilde{w}_1 + \frac{1}{2}\tilde{w}_2$ also lies in \mathcal{M}_{vol} . By strict convexity,

$$f_w(\tilde{w}_\theta) < \frac{1}{2}f_w(\tilde{w}_1) + \frac{1}{2}f_w(\tilde{w}_2) = \alpha_w,$$

contradicting the fact that α_w is the infimum over \mathcal{M}_{vol} . Therefore the minimizer of (28) is unique, and the mapping $w \mapsto \Pi_{\mathcal{M}_{\text{vol}}}(w)$ is well-defined on all of $\mathbb{R}^{d_{\text{vol}}}$.

Step 3: First-order optimality and firm non-expansiveness. The projection $\Pi_{\mathcal{M}_{\text{vol}}}(w)$ can be characterized by a classical first-order optimality condition for convex minimization over a closed convex set. Let $w \in \mathbb{R}^{d_{\text{vol}}}$ be arbitrary and denote $\tilde{w} := \Pi_{\mathcal{M}_{\text{vol}}}(w)$. Since \mathcal{M}_{vol} is closed and convex and f_w is differentiable and strictly convex, we know that \tilde{w} is the unique point in \mathcal{M}_{vol} such that

$$\langle \tilde{w} - w, z - \tilde{w} \rangle \geq 0 \quad \text{for all } z \in \mathcal{M}_{\text{vol}}. \quad (29)$$

Indeed, this is the variational inequality corresponding to optimality of \tilde{w} for the minimization of f_w over \mathcal{M}_{vol} ; see, e.g., standard results on projections in Hilbert spaces.

Let now $w, w' \in \mathbb{R}^{d_{\text{vol}}}$ be arbitrary, and set

$$p := \Pi_{\mathcal{M}_{\text{vol}}}(w), \quad q := \Pi_{\mathcal{M}_{\text{vol}}}(w').$$

Applying (29) with the pair (w, p) and the choice $z = q \in \mathcal{M}_{\text{vol}}$ yields

$$\langle p - w, q - p \rangle \geq 0. \quad (30)$$

Similarly, applying (29) with the pair (w', q) and the choice $z = p$ yields

$$\langle q - w', p - q \rangle \geq 0. \quad (31)$$

Adding (30) and (31), and recalling that $\langle q - w', p - q \rangle = -\langle q - w', q - p \rangle$, we obtain

$$\begin{aligned} 0 &\leq \langle p - w, q - p \rangle + \langle q - w', p - q \rangle \\ &= \langle p - w, q - p \rangle - \langle q - w', q - p \rangle \\ &= \langle (p - w) - (q - w'), q - p \rangle \\ &= \langle (w' - w) - (q - p), q - p \rangle \\ &= \langle w' - w, q - p \rangle - \|q - p\|_2^2. \end{aligned} \quad (32)$$

Rearranging (32) gives the inequality

$$\|p - q\|_2^2 \leq \langle w' - w, q - p \rangle. \quad (33)$$

Taking absolute values and applying the Cauchy–Schwarz inequality to the right-hand side yields

$$\|p - q\|_2^2 \leq |\langle w' - w, q - p \rangle| \leq \|w' - w\|_2 \|q - p\|_2.$$

If $p \neq q$, we can divide both sides by $\|p - q\|_2 > 0$ and obtain

$$\|p - q\|_2 \leq \|w' - w\|_2.$$

If $p = q$, the inequality trivially holds as equality. Thus, for every $w, w' \in \mathbb{R}^{d_{\text{vol}}}$,

$$\|\Pi_{\mathcal{M}_{\text{vol}}}(w) - \Pi_{\mathcal{M}_{\text{vol}}}(w')\|_2 \leq \|w - w'\|_2. \quad (34)$$

That is, the projection operator $\Pi_{\mathcal{M}_{\text{vol}}}$ is *non-expansive*, with Lipschitz constant equal to 1.

Conclusion. We have shown that \mathcal{M}_{vol} is non-empty, closed, and convex (by Proposition 1), that the Euclidean projection onto \mathcal{M}_{vol} exists and is unique for every $w \in \mathbb{R}^{d_{\text{vol}}}$, and that $\Pi_{\mathcal{M}_{\text{vol}}}$ is 1-Lipschitz in the Euclidean norm. This completes the proof of Proposition 2. \square

A.3 Proof of Lemma 1

Proof of Lemma 1. Recall the definition of the volatility law-penalty functional

$$\mathcal{L}_{\text{vol}}(w) := \frac{1}{2} \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2^2, \quad w \in \mathbb{R}^{d_{\text{vol}}}, \quad (35)$$

where $\mathcal{M}_{\text{vol}} \subset \mathbb{R}^{d_{\text{vol}}}$ is the volatility law manifold (cf. Proposition 1) and $\Pi_{\mathcal{M}_{\text{vol}}} : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathcal{M}_{\text{vol}}$ is the Euclidean projection (cf. Proposition 2).

We show that \mathcal{L}_{vol} is *locally Lipschitz* on $\mathbb{R}^{d_{\text{vol}}}$, i.e., for every compact set $K \subset \mathbb{R}^{d_{\text{vol}}}$ there exists $L_K < \infty$ such that

$$|\mathcal{L}_{\text{vol}}(w_1) - \mathcal{L}_{\text{vol}}(w_2)| \leq L_K \|w_1 - w_2\|_2 \quad \text{for all } w_1, w_2 \in K.$$

Step 1: Basic Lipschitz properties of the projection. Define the *residual map*

$$h(w) := w - \Pi_{\mathcal{M}_{\text{vol}}}(w), \quad w \in \mathbb{R}^{d_{\text{vol}}}.$$

By Proposition 2, the projection $\Pi_{\mathcal{M}_{\text{vol}}}$ is non-expansive:

$$\|\Pi_{\mathcal{M}_{\text{vol}}}(w_1) - \Pi_{\mathcal{M}_{\text{vol}}}(w_2)\|_2 \leq \|w_1 - w_2\|_2 \quad \text{for all } w_1, w_2 \in \mathbb{R}^{d_{\text{vol}}}.$$

It follows that the residual map h is globally Lipschitz with constant 2. Indeed, for any $w_1, w_2 \in \mathbb{R}^{d_{\text{vol}}}$,

$$\begin{aligned} \|h(w_1) - h(w_2)\|_2 &= \|(w_1 - \Pi_{\mathcal{M}_{\text{vol}}}(w_1)) - (w_2 - \Pi_{\mathcal{M}_{\text{vol}}}(w_2))\|_2 \\ &\leq \|w_1 - w_2\|_2 + \|\Pi_{\mathcal{M}_{\text{vol}}}(w_1) - \Pi_{\mathcal{M}_{\text{vol}}}(w_2)\|_2 \\ &\leq \|w_1 - w_2\|_2 + \|w_1 - w_2\|_2 \\ &= 2 \|w_1 - w_2\|_2. \end{aligned} \quad (36)$$

Thus h is 2-Lipschitz on all of $\mathbb{R}^{d_{\text{vol}}}$.

Step 2: Bounding the residual on bounded sets. Fix a radius $R > 0$ and consider the closed Euclidean ball

$$B_R := \{w \in \mathbb{R}^{d_{\text{vol}}} : \|w\|_2 \leq R\}.$$

We first show that $\|h(w)\|_2$ is uniformly bounded on B_R .

Let $w_0 := 0$ be the origin. Since \mathcal{M}_{vol} is non-empty (Proposition 1), the projection $\Pi_{\mathcal{M}_{\text{vol}}}(w_0)$ is well-defined and finite. Denote $c_0 := \|\Pi_{\mathcal{M}_{\text{vol}}}(0)\|_2 < \infty$.

For any $w \in B_R$, we have

$$\begin{aligned} \|\Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2 &\leq \|\Pi_{\mathcal{M}_{\text{vol}}}(w) - \Pi_{\mathcal{M}_{\text{vol}}}(0)\|_2 + \|\Pi_{\mathcal{M}_{\text{vol}}}(0)\|_2 \\ &\leq \|w - 0\|_2 + c_0 \leq R + c_0, \end{aligned} \quad (37)$$

where we used non-expansiveness of $\Pi_{\mathcal{M}_{\text{vol}}}$ and $\|w\|_2 \leq R$.

Hence, for any $w \in B_R$,

$$\begin{aligned} \|h(w)\|_2 &= \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2 \\ &\leq \|w\|_2 + \|\Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2 \\ &\leq R + (R + c_0) \\ &= 2R + c_0. \end{aligned} \quad (38)$$

Define the constant

$$K_R := 2R + c_0.$$

Then $\|h(w)\|_2 \leq K_R$ for all $w \in B_R$.

Step 3: Lipschitz continuity of the squared norm of the residual. Rewrite \mathcal{L}_{vol} in terms of h :

$$\mathcal{L}_{\text{vol}}(w) = \frac{1}{2} \|h(w)\|_2^2.$$

Let $w_1, w_2 \in B_R$ be arbitrary, and set $h_1 := h(w_1)$, $h_2 := h(w_2)$. Then

$$\begin{aligned} |\mathcal{L}_{\text{vol}}(w_1) - \mathcal{L}_{\text{vol}}(w_2)| &= \frac{1}{2} \left| \|h_1\|_2^2 - \|h_2\|_2^2 \right| \\ &= \frac{1}{2} |\langle h_1 + h_2, h_1 - h_2 \rangle| \\ &\leq \frac{1}{2} (\|h_1\|_2 + \|h_2\|_2) \|h_1 - h_2\|_2. \end{aligned} \quad (39)$$

From (38), we have $\|h_1\|_2 \leq K_R$ and $\|h_2\|_2 \leq K_R$. Moreover, by (36), $\|h_1 - h_2\|_2 \leq 2\|w_1 - w_2\|_2$. Substituting these bounds into (39) yields

$$\begin{aligned} |\mathcal{L}_{\text{vol}}(w_1) - \mathcal{L}_{\text{vol}}(w_2)| &\leq \frac{1}{2} (K_R + K_R) (2\|w_1 - w_2\|_2) \\ &= 2K_R \|w_1 - w_2\|_2. \end{aligned} \quad (40)$$

Recalling that $K_R = 2R + c_0$, we can rewrite (40) as

$$|\mathcal{L}_{\text{vol}}(w_1) - \mathcal{L}_{\text{vol}}(w_2)| \leq L_R \|w_1 - w_2\|_2, \quad L_R := 2(2R + c_0), \quad \forall w_1, w_2 \in B_R.$$

Thus \mathcal{L}_{vol} is Lipschitz on each ball B_R , with Lipschitz constant L_R depending only on R and the fixed constant c_0 .

Step 4: Local Lipschitz continuity. A function that is Lipschitz on every bounded ball in $\mathbb{R}^{d_{\text{vol}}}$ is, by definition, *locally Lipschitz*. More precisely, for any compact set $K \subset \mathbb{R}^{d_{\text{vol}}}$, there exists $R > 0$ such that $K \subset B_R$; then the Lipschitz constant L_R from (40) works for all $w_1, w_2 \in K$. Therefore, \mathcal{L}_{vol} is locally Lipschitz on $\mathbb{R}^{d_{\text{vol}}}$.

Remark. The above argument is self-contained and uses only the non-expansiveness of the projection onto a closed convex set. An alternative viewpoint, standard in convex analysis, is to note that \mathcal{L}_{vol} coincides with the *squared distance function* to the non-empty closed convex set \mathcal{M}_{vol} , which is known to be Fréchet differentiable on $\mathbb{R}^{d_{\text{vol}}}$ with 1-Lipschitz gradient (see, e.g., (author?) [26, Prop. 12.29–12.30]). In particular, the gradient mapping is globally Lipschitz, which again implies that \mathcal{L}_{vol} is locally Lipschitz. We include the direct argument above to keep the presentation self-contained.

This completes the proof of Lemma 1. □

A.4 Proof of Proposition 3

Proof of Proposition 3. Recall the definition of the volatility law-penalty functional

$$\mathcal{L}_{\text{vol}}(w) := \frac{1}{2} \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2^2, \quad w \in \mathbb{R}^{d_{\text{vol}}}, \quad (41)$$

where $\mathcal{M}_{\text{vol}} \subset \mathbb{R}^{d_{\text{vol}}}$ is the volatility law manifold and $\Pi_{\mathcal{M}_{\text{vol}}} : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathcal{M}_{\text{vol}}$ is the Euclidean projection (Proposition 2). We show that

$$\mathcal{L}_{\text{vol}}(w) = 0 \iff w \in \mathcal{M}_{\text{vol}}.$$

(\Rightarrow) **If $w \in \mathcal{M}_{\text{vol}}$ then $\mathcal{L}_{\text{vol}}(w) = 0$.** Assume $w \in \mathcal{M}_{\text{vol}}$. By the definition of the Euclidean projection, any point $\tilde{w} \in \mathcal{M}_{\text{vol}}$ satisfies

$$\|\Pi_{\mathcal{M}_{\text{vol}}}(w) - w\|_2 \leq \|\tilde{w} - w\|_2.$$

In particular, we may take $\tilde{w} = w$ itself, which is feasible since $w \in \mathcal{M}_{\text{vol}}$. This yields

$$\|\Pi_{\mathcal{M}_{\text{vol}}}(w) - w\|_2 \leq \|w - w\|_2 = 0.$$

By non-negativity of the norm, we must have equality, hence

$$\Pi_{\mathcal{M}_{\text{vol}}}(w) = w.$$

Substituting this into (41) gives

$$\mathcal{L}_{\text{vol}}(w) = \frac{1}{2} \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2^2 = \frac{1}{2} \|w - w\|_2^2 = 0.$$

Thus $w \in \mathcal{M}_{\text{vol}} \implies \mathcal{L}_{\text{vol}}(w) = 0$.

(\Leftarrow) If $\mathcal{L}_{\text{vol}}(w) = 0$ then $w \in \mathcal{M}_{\text{vol}}$. Now assume $\mathcal{L}_{\text{vol}}(w) = 0$. By (41) we have

$$0 = \mathcal{L}_{\text{vol}}(w) = \frac{1}{2} \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2^2.$$

Since the Euclidean norm is non-negative and vanishes only at zero, this implies

$$\|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2 = 0 \implies w = \Pi_{\mathcal{M}_{\text{vol}}}(w).$$

The projection $\Pi_{\mathcal{M}_{\text{vol}}}(w)$ takes values in \mathcal{M}_{vol} by construction, so $w = \Pi_{\mathcal{M}_{\text{vol}}}(w) \in \mathcal{M}_{\text{vol}}$. Hence $\mathcal{L}_{\text{vol}}(w) = 0 \implies w \in \mathcal{M}_{\text{vol}}$.

Alternative viewpoint via distance to closed sets. For completeness, we note that (41) can be written as

$$\mathcal{L}_{\text{vol}}(w) = \frac{1}{2} \text{dist}^2(w, \mathcal{M}_{\text{vol}}), \quad \text{dist}(w, \mathcal{M}_{\text{vol}}) := \inf_{\tilde{w} \in \mathcal{M}_{\text{vol}}} \|w - \tilde{w}\|_2,$$

where the infimum is attained at $\Pi_{\mathcal{M}_{\text{vol}}}(w)$ because \mathcal{M}_{vol} is non-empty, closed, and convex (Propositions 1–2). By basic properties of distance functions to closed sets in Hilbert spaces, one has

$$\text{dist}(w, \mathcal{M}_{\text{vol}}) = 0 \iff w \in \overline{\mathcal{M}_{\text{vol}}} = \mathcal{M}_{\text{vol}}.$$

Since $\mathcal{L}_{\text{vol}}(w) = \frac{1}{2} \text{dist}^2(w, \mathcal{M}_{\text{vol}})$, this is equivalent to

$$\mathcal{L}_{\text{vol}}(w) = 0 \iff w \in \mathcal{M}_{\text{vol}},$$

which is precisely the statement of Proposition 3.

This completes the proof. \square

A.5 Proof of Proposition 4

Proof. Recall the setting and notation:

1. The volatility law manifold $\mathcal{M}_{\text{vol}} \subset \mathbb{R}^{d_{\text{vol}}}$ is non-empty, closed, and convex by Propositions 1–2.
2. The Euclidean projection $\Pi_{\mathcal{M}_{\text{vol}}} : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathcal{M}_{\text{vol}}$ is well-defined and 1-Lipschitz (Proposition 2).
3. The law-penalty functional is given by

$$\mathcal{L}_{\text{vol}}(w) := \frac{1}{2} \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2^2 = \frac{1}{2} \text{dist}^2(w, \mathcal{M}_{\text{vol}}), \quad w \in \mathbb{R}^{d_{\text{vol}}}, \quad (42)$$

where $\text{dist}(w, \mathcal{M}_{\text{vol}}) := \inf_{\tilde{w} \in \mathcal{M}_{\text{vol}}} \|w - \tilde{w}\|_2$.

4. The on-manifold and ghost reward components are

$$r^{\mathcal{M}}(w) := r(\Pi_{\mathcal{M}_{\text{vol}}}(w)), \quad r^{\perp}(w) := r(w) - r^{\mathcal{M}}(w).$$

Assume that $r : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathbb{R}$ is L_r -Lipschitz with respect to the Euclidean norm, i.e.,

$$|r(w) - r(w')| \leq L_r \|w - w'\|_2 \quad \forall w, w' \in \mathbb{R}^{d_{\text{vol}}}. \quad (43)$$

We now prove the bound

$$|r^{\perp}(w)| = |r(w) - r^{\mathcal{M}}(w)| \leq L_r \text{dist}(w, \mathcal{M}_{\text{vol}}) = L_r \sqrt{2 \mathcal{L}_{\text{vol}}(w)}, \quad \forall w \in \mathbb{R}^{d_{\text{vol}}}.$$

Step 1: Bounding the ghost component by the distance to \mathcal{M}_{vol} . Fix any $w \in \mathbb{R}^{d_{\text{vol}}}$. By the definition of $r^{\mathcal{M}}$ and the Lipschitz property (43), we have

$$\begin{aligned} |r^{\perp}(w)| &= |r(w) - r^{\mathcal{M}}(w)| = |r(w) - r(\Pi_{\mathcal{M}_{\text{vol}}}(w))| \\ &\leq L_r \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2. \end{aligned}$$

By the definition of the distance function and properties of the projection,

$$\|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2 = \text{dist}(w, \mathcal{M}_{\text{vol}}),$$

since $\Pi_{\mathcal{M}_{\text{vol}}}(w)$ is a minimizer of $\tilde{w} \mapsto \|w - \tilde{w}\|_2$ over $\tilde{w} \in \mathcal{M}_{\text{vol}}$. Hence

$$|r^{\perp}(w)| \leq L_r \text{dist}(w, \mathcal{M}_{\text{vol}}). \quad (44)$$

Step 2: Expressing the distance via \mathcal{L}_{vol} . From the definition (42), we have

$$\text{dist}(w, \mathcal{M}_{\text{vol}}) = \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2 = \sqrt{2\mathcal{L}_{\text{vol}}(w)}.$$

Substituting this identity into (44) yields

$$|r^\perp(w)| \leq L_r \text{dist}(w, \mathcal{M}_{\text{vol}}) = L_r \sqrt{2\mathcal{L}_{\text{vol}}(w)}.$$

Since $w \in \mathbb{R}^{d_{\text{vol}}}$ was arbitrary, the inequality holds for all w .

Remark on optimality of the bound. The inequality in Proposition 4 is sharp in the sense that, for fixed \mathcal{M}_{vol} and $\Pi_{\mathcal{M}_{\text{vol}}}$, one can construct Lipschitz functions r that nearly saturate the bound. For example, if r is chosen to be affine in the normal direction to \mathcal{M}_{vol} at some point, with gradient of norm L_r , then along rays orthogonal to \mathcal{M}_{vol} we obtain

$$|r(w) - r(\Pi_{\mathcal{M}_{\text{vol}}}(w))| \approx L_r \|w - \Pi_{\mathcal{M}_{\text{vol}}}(w)\|_2,$$

up to second-order curvature effects of \mathcal{M}_{vol} . Thus the scaling $|r^\perp(w)| = O(\text{dist}(w, \mathcal{M}_{\text{vol}}))$ and the constant L_r cannot, in general, be improved under the sole assumption (43).

Extension to non-Euclidean penalties (informal). In the main text we focus on the Euclidean penalty (42). If instead \mathcal{L}_{vol} is defined via a strictly convex norm $\|\cdot\|_\phi$ or a strongly convex gauge ϕ , i.e.,

$$\mathcal{L}_{\text{vol}}^\phi(w) := \frac{1}{2} \text{dist}_\phi^2(w, \mathcal{M}_{\text{vol}}), \quad \text{dist}_\phi(w, \mathcal{M}_{\text{vol}}) := \inf_{\tilde{w} \in \mathcal{M}_{\text{vol}}} \|w - \tilde{w}\|_\phi,$$

and r is L_r^ϕ -Lipschitz with respect to $\|\cdot\|_\phi$, the same argument yields

$$|r^\perp(w)| \leq L_r^\phi \text{dist}_\phi(w, \mathcal{M}_{\text{vol}}) = L_r^\phi \sqrt{2\mathcal{L}_{\text{vol}}^\phi(w)}.$$

In finite dimensions, norm equivalence further implies that such bounds can be translated between different choices of ϕ at the expense of multiplicative constants depending only on the norms; We keep the Euclidean version in Proposition 4 as it is the one used in our experiments.

This concludes the proof. □

B Proofs for Section 3: Volatility World Model

B.1 Proof of Proposition 5

We now provide a detailed measure-theoretic proof of Proposition 5.

Recall the statement:

Proposition (Support of the synthetic generator). Let $(w_t)_{t=0}^T$ be generated by G^* as above, with static no-arbitrage imposed at each time via projection onto \mathcal{M}_{vol} . Then

$$\text{supp } \mathbb{P}^* \subseteq (\mathcal{M}_{\text{vol}})^{T+1},$$

i.e., \mathbb{P}^* is supported on the product of the volatility law manifold at all times.

Proof. We first recall the construction of the synthetic generator G^* from the main text and make it explicit in measure-theoretic terms.

Step 0: Probability space and random paths. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be an underlying probability space supporting all driving randomness (e.g., Brownian motions, volatility factors, etc.) of the stochastic-volatility model used to generate “raw” option prices or total-variance surfaces.

For each $t = 0, \dots, T$, let

$$y_t : \Omega \rightarrow \mathbb{R}^{d_{\text{vol}}}$$

denote the *raw* (not yet projected) total-variance surface at time t , obtained from the chosen parametric or factor model (Heston, rough volatility, SVI parameter dynamics, etc.; see, e.g., **(author?)** [8, 9]). We assume that each y_t is Borel measurable.

The *law-consistent synthetic generator* G^* then applies the projection onto the volatility law manifold \mathcal{M}_{vol} at each time t :

$$w_t(\omega) := \Pi_{\mathcal{M}_{\text{vol}}}(y_t(\omega)), \quad \omega \in \Omega, \quad t = 0, \dots, T. \quad (45)$$

Here $\Pi_{\mathcal{M}_{\text{vol}}}$ is the Euclidean projection operator introduced in Proposition 2, which maps $\mathbb{R}^{d_{\text{vol}}}$ into \mathcal{M}_{vol} .

We then define the $(T+1)$ -dimensional random vector

$$W(\omega) := (w_0(\omega), \dots, w_T(\omega)) \in (\mathbb{R}^{d_{\text{vol}}})^{T+1}. \quad (46)$$

By construction, W is Borel measurable, and we denote its law by

$$\mathbb{P}^* := \mathbb{P} \circ W^{-1},$$

a probability measure on $(\mathbb{R}^{d_{\text{vol}}})^{T+1}$.

Step 1: Pointwise inclusion $w_t(\omega) \in \mathcal{M}_{\text{vol}}$ **almost surely.** By Proposition 2, \mathcal{M}_{vol} is non-empty, closed, and convex, and the projection $\Pi_{\mathcal{M}_{\text{vol}}} : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathcal{M}_{\text{vol}}$ is well-defined and single-valued. In particular,

$$\Pi_{\mathcal{M}_{\text{vol}}}(x) \in \mathcal{M}_{\text{vol}} \quad \text{for all } x \in \mathbb{R}^{d_{\text{vol}}}.$$

Applying this to $x = y_t(\omega)$, the definition (45) implies that for every $\omega \in \Omega$ and every t ,

$$w_t(\omega) \in \mathcal{M}_{\text{vol}}. \quad (47)$$

Thus, for each fixed t , we have

$$\mathbb{P}(w_t \in \mathcal{M}_{\text{vol}}) = 1.$$

Step 2: Product inclusion for the path W . Using (47) for all $t = 0, \dots, T$, we obtain

$$W(\omega) = (w_0(\omega), \dots, w_T(\omega)) \in \mathcal{M}_{\text{vol}}^{T+1} \quad \text{for all } \omega \in \Omega.$$

Hence

$$\mathbb{P}(W \in \mathcal{M}_{\text{vol}}^{T+1}) = 1. \quad (48)$$

This already implies that the support of \mathbb{P}^* (the law of W) must lie inside $\mathcal{M}_{\text{vol}}^{T+1}$. To make this precise, we recall the definition of support.

Step 3: Support of a probability measure and restriction to $\mathcal{M}_{\text{vol}}^{T+1}$. Let (E, \mathcal{E}) be a measurable space, and μ a probability measure on E . The *support* of μ , denoted $\text{supp } \mu$, is the closed set

$$\text{supp } \mu := \{x \in E : \mu(U) > 0 \text{ for every open neighborhood } U \text{ of } x\}.$$

We will use the following standard lemma.

Lemma 4 (Support of a measure carried by a closed set). *Let E be a metric space, and let $C \subseteq E$ be a closed set. Suppose μ is a Borel probability measure on E such that $\mu(C) = 1$. Then*

$$\text{supp } \mu \subseteq C.$$

Proof of Lemma 4. Let $x \in \text{supp } \mu$. Suppose, for contradiction, that $x \notin C$. Since C is closed, its complement C^c is open and contains x . Thus there exists an open neighborhood U of x with $U \subseteq C^c$. Then

$$\mu(U) \leq \mu(C^c) = 1 - \mu(C) = 0,$$

contradicting the definition of support, which requires $\mu(U) > 0$ for every open neighborhood U of x . Hence $x \in C$. Since $x \in \text{supp } \mu$ was arbitrary, we conclude $\text{supp } \mu \subseteq C$. \square

Step 4: Applying the lemma to \mathbb{P}^* and $\mathcal{M}_{\text{vol}}^{T+1}$. We now take $E = (\mathbb{R}^{d_{\text{vol}}})^{T+1}$, equipped with its usual product topology and Borel σ -algebra. The set \mathcal{M}_{vol} is closed in $\mathbb{R}^{d_{\text{vol}}}$ by Proposition 2, hence $\mathcal{M}_{\text{vol}}^{T+1}$ is closed in the product topology of E . By (48),

$$\mathbb{P}^*(\mathcal{M}_{\text{vol}}^{T+1}) = \mathbb{P}(W \in \mathcal{M}_{\text{vol}}^{T+1}) = 1.$$

Applying Lemma 4 with $C = \mathcal{M}_{\text{vol}}^{T+1}$ and $\mu = \mathbb{P}^*$ yields

$$\text{supp } \mathbb{P}^* \subseteq \mathcal{M}_{\text{vol}}^{T+1}.$$

This is exactly the claim of Proposition 5:

$$\text{supp } \mathbb{P}^* \subseteq (\mathcal{M}_{\text{vol}})^{T+1}.$$

Step 5: Interpretation. From a modelling standpoint, the proposition formalizes the claim that the synthetic generator G^* produces only paths of total-variance surfaces $(w_t)_{t=0}^T$ that are *everywhere law-consistent*, in the sense of lying on the volatility law manifold at each time. The law \mathbb{P}^* of the generated paths is therefore entirely concentrated on the set $\mathcal{M}_{\text{vol}}^{T+1}$ of sequences of admissible surfaces. This is the precise sense in which we refer to G^* as a “law-consistent ground-truth world” in the main text.

This completes the proof. \square

B.2 Proof of Proposition 6

In this subsection we provide a detailed proof of Proposition 6, making precise the local linearization and the cone decomposition used in the main text.

Recall the setting: for each time step t , the law-consistent generator G^* produces $w_{t+1} \in \mathcal{M}_{\text{vol}}$ almost surely (Proposition 5), the world model produces a prediction

$$\hat{w}_{t+1} = f_\theta(w_{\leq t}, a_{\leq t}) \in \mathbb{R}^{d_{\text{vol}}},$$

and the residual is

$$e_{t+1} := \hat{w}_{t+1} - w_{t+1}.$$

We denote by $w_{t+1}^{\mathcal{M}} := \Pi_{\mathcal{M}_{\text{vol}}}(\hat{w}_{t+1})$ the projection of the prediction onto the volatility law manifold. The *ghost reward* at time $t+1$ is

$$r_{t+1}^\perp := r(\hat{w}_{t+1}, a_t) - r(w_{t+1}^{\mathcal{M}}, a_t),$$

where $r(\cdot, a_t)$ is the one-step reward as a function of the surface w for a fixed action a_t .

We restate the proposition for convenience.

Proposition 8 (Approximation gap induces a ghost channel). *Suppose the following conditions hold:*

(i) *The approximation gap is non-zero:*

$$\varepsilon^2 := \mathbb{E}[\|e_{t+1}\|_2^2] > 0.$$

(ii) *The reward is locally differentiable in w with gradient*

$$g_{t+1} := \nabla_w r(w_{t+1}, a_t),$$

and $\nabla_w r(\cdot, a_t)$ is locally Lipschitz in a neighborhood of the law-consistent path.

(iii) *The residual e_{t+1} has a component in the normal cone of \mathcal{M}_{vol} at w_{t+1} with non-zero covariance with the gradient:*

$$\text{Cov}(P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}, g_{t+1}) \neq 0,$$

where $P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})}$ denotes orthogonal projection onto the normal cone $N_{\mathcal{M}_{\text{vol}}}(w_{t+1})$.

Then, for sufficiently small residuals (in the sense of a local linearization),

$$\mathbb{E}[r_{t+1}^\perp] \approx \mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}] \neq 0,$$

so the world model induces a non-trivial ghost channel. In particular, if g_{t+1} is positively correlated with $P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}$, then $\mathbb{E}[r_{t+1}^\perp] > 0$ and there exist states where moving off-manifold strictly improves expected P&L.

Proof. We work under the law-consistent measure \mathbb{P}^* of the generator G^* , under which $w_{t+1} \in \mathcal{M}_{\text{vol}}$ almost surely (Proposition 5). We suppress the explicit dependence on t when it does not create ambiguity.

Step 1: Local Taylor expansion of the reward. Fix a compact set $K \subset \mathbb{R}^{d_{\text{vol}}}$ containing the realized w_{t+1} and all admissible \hat{w}_{t+1} and $w_{t+1}^{\mathcal{M}}$ with high probability (a standard truncation argument; see below). By assumption (ii), for each fixed a_t , the map

$$r(\cdot, a_t) : \mathbb{R}^{d_{\text{vol}}} \rightarrow \mathbb{R}$$

is differentiable on K , with gradient $\nabla_w r(\cdot, a_t)$ locally Lipschitz.

Hence, by the mean value theorem in Banach spaces, for each $\omega \in \Omega$ we can write

$$r(\hat{w}_{t+1}, a_t) = r(w_{t+1}, a_t) + g_{t+1}^\top e_{t+1} + \rho_{t+1}^{(1)}, \quad (49)$$

$$r(w_{t+1}^{\mathcal{M}}, a_t) = r(w_{t+1}, a_t) + g_{t+1}^\top (w_{t+1}^{\mathcal{M}} - w_{t+1}) + \rho_{t+1}^{(2)}, \quad (50)$$

where the remainder terms satisfy the quadratic bounds

$$|\rho_{t+1}^{(1)}| \leq \frac{L_{\nabla x}}{2} \|e_{t+1}\|_2^2, \quad (51)$$

$$|\rho_{t+1}^{(2)}| \leq \frac{L_{\nabla x}}{2} \|w_{t+1}^{\mathcal{M}} - w_{t+1}\|_2^2, \quad (52)$$

for some local Lipschitz constant $L_{\nabla x} > 0$ depending only on K and a_t .

Subtracting (50) from (49) yields

$$r_{t+1}^\perp = g_{t+1}^\top \left(e_{t+1} - (w_{t+1}^{\mathcal{M}} - w_{t+1}) \right) + (\rho_{t+1}^{(1)} - \rho_{t+1}^{(2)}). \quad (53)$$

Step 2: Tangent-normal cone decomposition of the residual. Since $w_{t+1} \in \mathcal{M}_{\text{vol}}$ and \mathcal{M}_{vol} is closed and convex (Proposition 2), we may consider the tangent cone $T_{\mathcal{M}_{\text{vol}}}(w_{t+1})$ and the normal cone $N_{\mathcal{M}_{\text{vol}}}(w_{t+1})$:

$$T_{\mathcal{M}_{\text{vol}}}(w_{t+1}) := \overline{\bigcup_{\alpha > 0} \alpha (\mathcal{M}_{\text{vol}} - w_{t+1})},$$

$$N_{\mathcal{M}_{\text{vol}}}(w_{t+1}) := \{v \in \mathbb{R}^{d_{\text{vol}}} : v^\top (z - w_{t+1}) \leq 0 \ \forall z \in \mathcal{M}_{\text{vol}}\}.$$

Both are closed convex cones, and they are polar to each other:

$$N_{\mathcal{M}_{\text{vol}}}(w_{t+1}) = T_{\mathcal{M}_{\text{vol}}}(w_{t+1})^\circ, \quad T_{\mathcal{M}_{\text{vol}}}(w_{t+1}) = N_{\mathcal{M}_{\text{vol}}}(w_{t+1})^\circ.$$

By the Moreau decomposition for closed convex cones (see, e.g., (author?) [26, Thm. 6.29]), every vector $z \in \mathbb{R}^{d_{\text{vol}}}$ admits a unique decomposition

$$z = P_{T_{\mathcal{M}_{\text{vol}}}(w_{t+1})} z + P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} z,$$

where P_T and P_N denote the orthogonal projections onto $T_{\mathcal{M}_{\text{vol}}}(w_{t+1})$ and $N_{\mathcal{M}_{\text{vol}}}(w_{t+1})$, respectively.

Applying this to e_{t+1} , we write

$$e_{t+1} = e_{t+1}^{\text{tan}} + e_{t+1}^{\text{norm}}, \quad e_{t+1}^{\text{tan}} := P_{T_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}, \quad e_{t+1}^{\text{norm}} := P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}. \quad (54)$$

Step 3: Local behavior of the projection $\Pi_{\mathcal{M}_{\text{vol}}}$. We next relate $w_{t+1}^{\mathcal{M}}$ to w_{t+1} and e_{t+1} .

By definition,

$$w_{t+1}^{\mathcal{M}} = \Pi_{\mathcal{M}_{\text{vol}}}(\hat{w}_{t+1}) = \arg \min_{z \in \mathcal{M}_{\text{vol}}} \|\hat{w}_{t+1} - z\|_2^2.$$

Since $w_{t+1} \in \mathcal{M}_{\text{vol}}$, for small $\|e_{t+1}\|_2$ the minimizer $w_{t+1}^{\mathcal{M}}$ lies in a neighborhood where \mathcal{M}_{vol} is locally well-approximated by its tangent cone at w_{t+1} . Under a standard regularity condition (e.g., w_{t+1} is a point of *prox-regularity* of \mathcal{M}_{vol} ;), the projection mapping $\Pi_{\mathcal{M}_{\text{vol}}}$ is directionally differentiable at w_{t+1} and its first-order behavior is given by orthogonal projection onto the tangent cone:

$$w_{t+1}^{\mathcal{M}} - w_{t+1} = P_{T_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1} + \delta_{t+1}, \quad (55)$$

where the remainder δ_{t+1} satisfies

$$\|\delta_{t+1}\|_2 \leq C_\Pi \|e_{t+1}\|_2^2 \quad (56)$$

for some constant $C_\Pi > 0$ in a neighborhood of w_{t+1} . Intuitively, to first order, $\Pi_{\mathcal{M}_{\text{vol}}}$ keeps the tangential component of e_{t+1} but kills the normal component; the error δ_{t+1} is of second order in $\|e_{t+1}\|_2$.

Substituting (55) into (53) and using (54), we obtain

$$\begin{aligned} r_{t+1}^\perp &= g_{t+1}^\top (e_{t+1} - (w_{t+1}^{\mathcal{M}} - w_{t+1})) + (\rho_{t+1}^{(1)} - \rho_{t+1}^{(2)}) \\ &= g_{t+1}^\top (e_{t+1}^{\text{tan}} + e_{t+1}^{\text{norm}} - P_{T_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1} - \delta_{t+1}) + (\rho_{t+1}^{(1)} - \rho_{t+1}^{(2)}) \\ &= g_{t+1}^\top e_{t+1}^{\text{norm}} - g_{t+1}^\top \delta_{t+1} + (\rho_{t+1}^{(1)} - \rho_{t+1}^{(2)}). \end{aligned} \quad (57)$$

Step 4: Bounding the higher-order error. We now bound the total remainder term

$$\eta_{t+1} := -g_{t+1}^\top \delta_{t+1} + (\rho_{t+1}^{(1)} - \rho_{t+1}^{(2)}).$$

Using Cauchy-Schwarz, (56), and the boundedness of g_{t+1} on K (say $\|g_{t+1}\|_2 \leq G$ almost surely on K), we have

$$|g_{t+1}^\top \delta_{t+1}| \leq \|g_{t+1}\|_2 \|\delta_{t+1}\|_2 \leq GC_\Pi \|e_{t+1}\|_2^2.$$

Combining this with (51) and (52), and using $\|w_{t+1}^{\mathcal{M}} - w_{t+1}\|_2 \leq \|e_{t+1}\|_2$ (since $w_{t+1}^{\mathcal{M}}$ is the closest point in \mathcal{M}_{vol} to \hat{w}_{t+1} and $w_{t+1} \in \mathcal{M}_{\text{vol}}$), we obtain

$$|\eta_{t+1}| \leq C_{\text{tot}} \|e_{t+1}\|_2^2 \quad (58)$$

for some constant $C_{\text{tot}} > 0$.

Substituting (57) and taking expectations yields

$$\mathbb{E}[r_{t+1}^\perp] = \mathbb{E}[g_{t+1}^\top e_{t+1}^{\text{norm}}] + \mathbb{E}[\eta_{t+1}], \quad (59)$$

with $|\mathbb{E}[\eta_{t+1}]| \leq C_{\text{tot}} \mathbb{E}[\|e_{t+1}\|_2^2] = C_{\text{tot}} \varepsilon^2$ by assumption (i).

Step 5: Non-trivial ghost channel from covariance structure. By definition of e_{t+1}^{norm} in (54), we have

$$e_{t+1}^{\text{norm}} = P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}.$$

Assumption (iii) states that the covariance

$$\text{Cov}(e_{t+1}^{\text{norm}}, g_{t+1}) = \mathbb{E}[(e_{t+1}^{\text{norm}} - \mathbb{E}[e_{t+1}^{\text{norm}}])(g_{t+1} - \mathbb{E}[g_{t+1}])^\top]$$

is non-zero. In particular, the scalar random variable

$$Z_{t+1} := g_{t+1}^\top e_{t+1}^{\text{norm}}$$

has non-zero covariance with itself along at least one direction, implying that

$$\mathbb{E}[Z_{t+1}] = \mathbb{E}[g_{t+1}^\top e_{t+1}^{\text{norm}}] \neq 0 \quad (60)$$

unless the mean terms $\mathbb{E}[e_{t+1}^{\text{norm}}]$ and $\mathbb{E}[g_{t+1}]$ are tuned to perfectly cancel the covariance contribution; this would be an exceptional, measure-zero configuration in parameter space. To avoid such pathological cancellation, we interpret assumption (iii) as requiring that the inner product $g_{t+1}^\top e_{t+1}^{\text{norm}}$ has a non-degenerate distribution with non-zero mean.⁴

Substituting (60) into (59), we obtain

$$\mathbb{E}[r_{t+1}^\perp] = \mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}] + \mathbb{E}[\eta_{t+1}], \quad (61)$$

with $\mathbb{E}[\eta_{t+1}]$ bounded by (58).

Step 6: Small-residual regime and sign of the ghost reward. Assumption (i) states that the world model approximation error is non-zero in mean-square:

$$\mathbb{E}[\|e_{t+1}\|_2^2] = \varepsilon^2 > 0.$$

Suppose in addition that the training of the world model has reduced the error variance so that ε^2 is *small*. Then from (58),

$$|\mathbb{E}[\eta_{t+1}]| \leq C_{\text{tot}} \varepsilon^2,$$

which can be made arbitrarily small by improving the world model (e.g., increasing capacity or training time), while the leading term

$$\mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}]$$

remains of order ε in general, as it is linear in e_{t+1} .

Consequently, for sufficiently small ε we have the first-order approximation

$$\mathbb{E}[r_{t+1}^\perp] \approx \mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}],$$

with the difference bounded by $O(\varepsilon^2)$. Under the non-degeneracy condition in (60), this leading term is non-zero, which yields

$$\mathbb{E}[r_{t+1}^\perp] \neq 0$$

for sufficiently small approximation error ε .

Finally, if the correlation between g_{t+1} and $P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}$ is *positive* in the sense that

$$\mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}] > 0,$$

then (61) implies

$$\mathbb{E}[r_{t+1}^\perp] > 0$$

for sufficiently small ε . In particular, there exists a set of states of positive probability where $r_{t+1}^\perp > 0$, so moving off the volatility law manifold along the normal directions strictly increases expected one-step P&L. This is precisely the “ghost channel” exploited by naive RL and law-seeking RL in the main text.

This completes the proof. \square

⁴Formally, one may replace the covariance condition in assumption (iii) by the simpler requirement $\mathbb{E}[g_{t+1}^\top P_{N_{\mathcal{M}_{\text{vol}}}(w_{t+1})} e_{t+1}] \neq 0$. We keep the covariance phrasing to emphasize the geometric correlation between the gradient and the normal component.

B.3 Proof of Lemma 2

In this subsection we give a more detailed argument for Lemma 2, making precise the informal statement in the main text that a finite-capacity, unconstrained neural world model will, under mild regularity conditions, place non-trivial probability mass off the volatility law manifold.

Recall the statement.

Lemma (Non-trivial off-manifold mass of the world model). Assume that f_{θ^*} is not exactly equal to the Bayes-optimal regressor $f^{\text{Bayes}}(x_t) := \mathbb{E}[w_{t+1} | x_t]$ and that the law manifold \mathcal{M}^{vol} has non-empty interior within the support of w_{t+1} . Then there exists $\delta > 0$ such that

$$\mathbb{P}(\mathcal{L}_\phi(\hat{w}_{t+1}) > \delta) > 0,$$

i.e., the world model assigns non-zero probability mass to surfaces at a positive distance from \mathcal{M}^{vol} .

Setup and additional regularity. Let X_t denote the (vector) state at time t and $W_{t+1} := w_{t+1}$ the total-variance surface generated by the law-consistent generator G^* at time $t + 1$. We write

$$f^{\text{Bayes}}(x) := \mathbb{E}[W_{t+1} | X_t = x], \quad \hat{W}_{t+1} := \hat{w}_{t+1} := f_{\theta^*}(X_t).$$

By Proposition 5 and convexity of \mathcal{M}_{vol} , we have

$$W_{t+1} \in \mathcal{M}_{\text{vol}} \text{ a.s.} \implies f^{\text{Bayes}}(X_t) \in \mathcal{M}_{\text{vol}} \text{ a.s.},$$

since conditional expectations of random variables supported on a closed convex set remain in that set.

We also recall from Proposition 3 that for any $w \in \mathbb{R}^{d_{\text{vol}}}$,

$$\mathcal{L}_\phi(w) = 0 \iff w \in \mathcal{M}_{\text{vol}},$$

and from Lemma 1 that \mathcal{L}_ϕ is continuous (indeed locally Lipschitz) on $\mathbb{R}^{d_{\text{vol}}}$.

To make the genericity argument precise, we introduce a mild regularity hypothesis on the joint distribution of (X_t, \hat{W}_{t+1}) .

Regularity assumption (B.3). We assume:

- (a) The state X_t has a distribution whose support $\text{supp}(X_t)$ is not a single point and contains a non-empty open subset $U_X \subset \mathbb{R}^{d_X}$.
- (b) The trained world model $f_{\theta^*} : \mathbb{R}^{d_X} \rightarrow \mathbb{R}^{d_{\text{vol}}}$ is continuous on U_X (this is satisfied by standard neural networks with continuous activations).
- (c) The image $f_{\theta^*}(U_X)$ is not contained in any $(d_{\text{vol}} - 1)$ -dimensional affine subspace of $\mathbb{R}^{d_{\text{vol}}}$ (a genericity condition on the parameter choice θ^*).

Assumption (B.3) is very mild in practice: for typical neural network parametrizations with random initialization and gradient-based training, the set of parameters for which f_θ maps U_X into a fixed lower-dimensional affine subspace has Lebesgue measure zero in parameter space.

We now prove that, under the assumptions of Lemma 2 and (B.3), the world model places non-trivial mass at positive distance from \mathcal{M}_{vol} .

Proof. Define the *distance-to-manifold* function

$$d_{\mathcal{M}_{\text{vol}}}(w) := \text{dist}(w, \mathcal{M}_{\text{vol}}) = \inf_{z \in \mathcal{M}_{\text{vol}}} \|w - z\|_2, \quad w \in \mathbb{R}^{d_{\text{vol}}}.$$

By closedness of \mathcal{M}_{vol} (Proposition 2), $d_{\mathcal{M}_{\text{vol}}}$ is continuous and

$$d_{\mathcal{M}_{\text{vol}}}(w) = 0 \iff w \in \mathcal{M}_{\text{vol}}.$$

By Proposition 3, \mathcal{L}_ϕ and $d_{\mathcal{M}_{\text{vol}}}$ have the same zero set, and continuity of \mathcal{L}_ϕ implies that for every $\delta > 0$ there exists $\eta(\delta) > 0$ such that

$$d_{\mathcal{M}_{\text{vol}}}(w) > \eta(\delta) \implies \mathcal{L}_\phi(w) > \delta, \quad \forall w \in \mathbb{R}^{d_{\text{vol}}}. \quad (62)$$

Thus it suffices to show that there exists $\eta > 0$ such that

$$\mathbb{P}(d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) > \eta) > 0.$$

Step 1: Existence of a prediction point outside \mathcal{M}_{vol} . We first argue that, under our assumptions, the image of f_{θ^*} cannot be contained in \mathcal{M}_{vol} .

Suppose for contradiction that

$$f_{\theta^*}(x) \in \mathcal{M}_{\text{vol}} \quad \text{for all } x \in \text{supp}(X_t). \quad (63)$$

In particular, for $x \in \text{supp}(X_t)$ we have $\hat{W}_{t+1} \in \mathcal{M}_{\text{vol}}$ almost surely and hence $d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) = 0$ and $\mathcal{L}_\phi(\hat{W}_{t+1}) = 0$ almost surely. This is exactly the negation of the lemma's conclusion.

We now show that (63) is incompatible with the combination of (i) $f_{\theta^*} \neq f^{\text{Bayes}}$ and (B.3), given that \mathcal{M}_{vol} has non-empty interior within the support of W_{t+1} .

Because \mathcal{M}_{vol} has non-empty interior and W_{t+1} is supported on \mathcal{M}_{vol} (Proposition 5), there exists a point $w^\circ \in \mathcal{M}_{\text{vol}}$ and $r > 0$ such that

$$B(w^\circ, r) := \{w \in \mathbb{R}^{d_{\text{vol}}} : \|w - w^\circ\|_2 < r\} \subset \mathcal{M}_{\text{vol}}$$

and $\mathbb{P}(W_{t+1} \in B(w^\circ, r)) > 0$. Since $f^{\text{Bayes}}(X_t) = \mathbb{E}[W_{t+1} \mid X_t]$ takes values in the convex set \mathcal{M}_{vol} , there exists $x^\circ \in \text{supp}(X_t)$ such that

$$f^{\text{Bayes}}(x^\circ) \in B(w^\circ, r/2) \subset \text{int}(\mathcal{M}_{\text{vol}}).$$

By the assumption $f_{\theta^*} \neq f^{\text{Bayes}}$ in $L^2(\mathbb{P})$, there exists a set $A \subset \text{supp}(X_t)$ with $\mathbb{P}(X_t \in A) > 0$ such that

$$f_{\theta^*}(x) \neq f^{\text{Bayes}}(x) \quad \text{for all } x \in A.$$

Under Assumption (B.3)(a)–(b), the support $\text{supp}(X_t)$ contains an open set U_X on which f_{θ^*} is continuous, and we may without loss of generality assume that $A \cap U_X$ has positive probability (otherwise we restrict attention to a smaller open subset with positive mass).

Pick $x_1 \in A \cap U_X$ with $f_{\theta^*}(x_1) \neq f^{\text{Bayes}}(x_1)$. Consider the continuous path in input space given by

$$\gamma(\alpha) := (1 - \alpha)x^\circ + \alpha x_1, \quad \alpha \in [0, 1],$$

which lies in U_X since U_X is open and convex in a neighborhood of x° and x_1 (we can always restrict to a sufficiently small line segment if necessary). Define

$$h(\alpha) := f_{\theta^*}(\gamma(\alpha)), \quad b(\alpha) := f^{\text{Bayes}}(\gamma(\alpha)).$$

By continuity of f_{θ^*} and f^{Bayes} on U_X , both h and b are continuous on $[0, 1]$. At $\alpha = 0$ we have $b(0) \in \text{int}(\mathcal{M}_{\text{vol}})$ and $h(0) \in \mathcal{M}_{\text{vol}}$ by (63). At $\alpha = 1$ we have $b(1) \in \mathcal{M}_{\text{vol}}$ and $h(1) \in \mathcal{M}_{\text{vol}}$ by (63), but $h(1) \neq b(1)$ by choice of x_1 .

Thus the difference $d(\alpha) := h(\alpha) - b(\alpha)$ is a continuous map with $d(1) \neq 0$. Since \mathcal{M}_{vol} contains an open ball around $b(0)$, there exists $\rho \in (0, r/2)$ such that

$$B(b(0), \rho) \subset \mathcal{M}_{\text{vol}}.$$

If it happened that $h(\alpha) \in \mathcal{M}_{\text{vol}}$ for all $\alpha \in [0, 1]$, then the curve $h([0, 1])$ would be a continuous path in \mathcal{M}_{vol} connecting $h(0)$ and $h(1)$, both of which lie in \mathcal{M}_{vol} . This is not impossible *per se*; however, given the genericity Assumption (B.3)(c), which rules out that $f_{\theta^*}(U_X)$ is constrained to a lower-dimensional surface within \mathcal{M}_{vol} , we can exclude the degenerate situation where the entire image of the line segment $\gamma([0, 1])$ under f_{θ^*} remains inside \mathcal{M}_{vol} while differing from f^{Bayes} on a set of positive measure.

Formally, since $b(0)$ lies in the interior of \mathcal{M}_{vol} and $h(0) \in \mathcal{M}_{\text{vol}}$, there is an open neighborhood V of $\gamma(0)$ such that $b(V)$ and $h(V)$ both intersect $B(b(0), \rho)$ in sets of positive Lebesgue measure. Under (B.3)(c), the continuous map f_{θ^*} cannot, on V , be constrained to the $(d_{\text{vol}} - 1)$ -dimensional manifold $\partial\mathcal{M}_{\text{vol}}$; hence there must exist some $\tilde{x} \in V$ such that

$$f_{\theta^*}(\tilde{x}) \notin \mathcal{M}_{\text{vol}}.$$

Since $V \subset \text{supp}(X_t)$ and $\mathbb{P}(X_t \in V) > 0$, this implies

$$\mathbb{P}(f_{\theta^*}(X_t) \notin \mathcal{M}_{\text{vol}}) > 0.$$

Equivalently, there exists at least one point $w^* \in \mathbb{R}^{d_{\text{vol}}} \setminus \mathcal{M}_{\text{vol}}$ that is attained by \hat{W}_{t+1} with positive probability.

We have thus shown that (63) cannot hold under the assumptions of the lemma together with (B.3). Therefore

$$\mathbb{P}(\hat{W}_{t+1} \notin \mathcal{M}_{\text{vol}}) > 0. \quad (64)$$

Step 2: From off-manifold predictions to a positive-distance shell. Since \mathcal{M}_{vol} is closed and \hat{W}_{t+1} is a random element of $\mathbb{R}^{d_{\text{vol}}}$, the continuous distance function $d_{\mathcal{M}_{\text{vol}}}$ satisfies

$$\{\hat{W}_{t+1} \notin \mathcal{M}_{\text{vol}}\} = \{d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) > 0\}.$$

By (64), the event $\{d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) > 0\}$ has positive probability. Define

$$Z := d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) \geq 0.$$

Then $\mathbb{P}(Z > 0) > 0$. Since Z is a non-negative random variable, we can write its distribution function as

$$F_Z(\eta) := \mathbb{P}(Z \leq \eta), \quad \eta \geq 0.$$

By right-continuity of F_Z and $\mathbb{P}(Z > 0) > 0$, there must exist $\eta_0 > 0$ such that

$$\mathbb{P}(Z > \eta_0) = 1 - F_Z(\eta_0) > 0.$$

Equivalently,

$$\mathbb{P}(d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) > \eta_0) > 0. \tag{65}$$

Step 3: Translating distance to law penalty. Finally, we use the monotonic relationship between distance and law penalty. By continuity of \mathcal{L}_ϕ and the fact that $\mathcal{L}_\phi(w) = 0$ if and only if $d_{\mathcal{M}_{\text{vol}}}(w) = 0$ (Proposition 3), there exists a strictly increasing continuous function

$$\psi : [0, \infty) \rightarrow [0, \infty), \quad \psi(0) = 0,$$

such that for all $w \in \mathbb{R}^{d_{\text{vol}}}$,

$$\mathcal{L}_\phi(w) \geq \psi(d_{\mathcal{M}_{\text{vol}}}(w)),$$

and $\psi(u) > 0$ for all $u > 0$. For instance, in the concrete squared-distance case $\mathcal{L}_\phi(w) = \frac{1}{2}d_{\mathcal{M}_{\text{vol}}}(w)^2$ we may take $\psi(u) = \frac{1}{2}u^2$.

Set

$$\delta := \psi(\eta_0) > 0.$$

Then on the event $\{d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) > \eta_0\}$ we have

$$\mathcal{L}_\phi(\hat{W}_{t+1}) \geq \psi(d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1})) > \psi(\eta_0) = \delta.$$

Therefore,

$$\mathbb{P}(\mathcal{L}_\phi(\hat{W}_{t+1}) > \delta) \geq \mathbb{P}(d_{\mathcal{M}_{\text{vol}}}(\hat{W}_{t+1}) > \eta_0) > 0,$$

by (65). This is precisely the claim of Lemma 2.

Remark. Note that the assumption $f_{\theta^*} \neq f^{\text{Bayes}}$ is used here only to rule out the trivial case where the world model has converged exactly to the Bayes regressor, which is law-consistent by construction. The non-trivial content of the lemma is provided by the genericity condition (B.3), which encodes the intuition that a high-capacity, unconstrained neural world model trained without law penalties will almost surely deviate from \mathcal{M}_{vol} on a set of positive probability. Under these conditions, Lemma 2 shows that the world model opens a *ghost channel* in the sense of Sec. 3.3, providing room for RL to exploit off-manifold arbitrage opportunities. \square

C Proofs for Section 4: RL on Volatility World Models: Incentives and Law-Strength

C.1 Proof of Theorem 1

In this appendix we provide a detailed proof of the ghost-arbitrage incentive result stated in Theorem 1. For convenience we first recall the main objects and Assumption 1.

Preliminaries and notation. Let Π denote the (parameterized) policy class under consideration and let \mathcal{S} denote the structural strategy class introduced in Section. For any $\pi \in \Pi$ we consider three performance functionals:

$$J_0(\pi) := \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t r_t \mid \pi \right], \quad (66)$$

$$J^{\mathcal{M}}(\pi) := \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t r_t^{\mathcal{M}} \mid \pi \right], \quad (67)$$

$$J^{\perp}(\pi) := \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t r_t^{\perp} \mid \pi \right], \quad (68)$$

where r_t is the realized P&L reward at step t in the world model, $r_t^{\mathcal{M}}$ is the on-manifold reward obtained by replacing the predicted surface \hat{w}_{t+1} with its projection $\Pi_{\mathcal{M}_{\text{vol}}}(\hat{w}_{t+1})$, and $r_t^{\perp} := r_t - r_t^{\mathcal{M}}$ is the ghost reward component. By linearity of expectation,

$$J_0(\pi) = J^{\mathcal{M}}(\pi) + J^{\perp}(\pi), \quad \forall \pi \in \Pi. \quad (69)$$

Structural near-optimality assumption. We recall the structural approximation assumption used in the main text.

Assumption 4 (Structural near-optimality of \mathcal{S}). *There exist a structural policy $\pi_{\mathcal{S}}^* \in \mathcal{S}$ and a constant $\varepsilon_{\mathcal{S}} \geq 0$ such that*

$$\sup_{\pi \in \Pi} J^{\mathcal{M}}(\pi) \leq J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) + \varepsilon_{\mathcal{S}}. \quad (70)$$

That is, the structural class \mathcal{S} (Zero-Hedge, Vol-Trend, etc.) approximates the globally optimal on-manifold performance within a gap $\varepsilon_{\mathcal{S}}$.

Note that by definition $\pi_{\mathcal{S}}^* \in \mathcal{S} \subseteq \Pi$, so it is admissible in the unconstrained policy class as well.

We now restate the theorem and prove both parts.

Theorem 4. *Suppose Assumption 4 holds, and let*

$$\pi_0^* \in \arg \max_{\pi \in \Pi} J_0(\pi)$$

be a global maximizer of J_0 over Π . Then:

1. *If*

$$\sup_{\pi \in \Pi} J_0(\pi) > J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) + \varepsilon_{\mathcal{S}},$$

any maximizer π_0^ satisfies*

$$J^{\perp}(\pi_0^*) \geq \sup_{\pi \in \Pi} J_0(\pi) - J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) - \varepsilon_{\mathcal{S}} > 0. \quad (71)$$

In particular, the excess value over the structural baseline is entirely attributable to ghost arbitrage.

2. *If, in addition, $J^{\mathcal{M}}$ has a local maximum at some $\bar{\pi} = \pi_{\bar{\theta}} \in \Pi$ with $J^{\mathcal{M}}(\bar{\pi}) \approx J^{\mathcal{M}}(\pi_{\mathcal{S}}^*)$, and the policy-gradient theorem holds [49], then the policy gradient near $\bar{\pi}$ satisfies*

$$\nabla_{\theta} J_0(\pi_{\theta}) \Big|_{\theta=\bar{\theta}} \approx \nabla_{\theta} J^{\perp}(\pi_{\theta}) \Big|_{\theta=\bar{\theta}}, \quad (72)$$

so gradient-based RL updates are locally driven by increasing J^{\perp} .

Proof. We treat parts (i) and (ii) separately.

Proof of part (i). Fix any $\pi \in \Pi$. Combining the decomposition (69) with Assumption (70), we have

$$J_0(\pi) = J^{\mathcal{M}}(\pi) + J^{\perp}(\pi) \quad (73)$$

$$\leq (J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) + \varepsilon_{\mathcal{S}}) + J^{\perp}(\pi), \quad (74)$$

so for any $\pi \in \Pi$,

$$J^{\perp}(\pi) \geq J_0(\pi) - J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) - \varepsilon_{\mathcal{S}}. \quad (75)$$

Now consider a maximizer $\pi_0^* \in \arg \max_{\pi \in \Pi} J_0(\pi)$. By definition,

$$J_0(\pi_0^*) = \sup_{\pi \in \Pi} J_0(\pi).$$

Substituting $\pi = \pi_0^*$ into (75) yields

$$J^\perp(\pi_0^*) \geq J_0(\pi_0^*) - J^\mathcal{M}(\pi_S^*) - \varepsilon_S \quad (76)$$

$$= \sup_{\pi \in \Pi} J_0(\pi) - J^\mathcal{M}(\pi_S^*) - \varepsilon_S. \quad (77)$$

Under the stated strict inequality

$$\sup_{\pi \in \Pi} J_0(\pi) > J^\mathcal{M}(\pi_S^*) + \varepsilon_S,$$

the right-hand side of (77) is strictly positive, which implies

$$J^\perp(\pi_0^*) > 0.$$

This establishes (71) and shows that any excess value of the naive-RL maximizer over the structural near-optimal on-manifold performance must be entirely carried by the ghost component J^\perp . In particular, there is no room to explain the advantage of π_0^* over π_S^* by on-manifold improvements alone.

Proof of part (ii). We now consider the local gradient behavior. Let $\pi_\theta \in \Pi$ denote a differentiable parametrization of policies with parameter vector $\theta \in \mathbb{R}^p$, and suppose that $\bar{\pi} = \pi_{\bar{\theta}}$ is such that $J^\mathcal{M}$ has a local maximum at $\bar{\theta}$, i.e.,

$$J^\mathcal{M}(\pi_{\bar{\theta}}) \geq J^\mathcal{M}(\pi_\theta) \quad \text{for all } \theta \text{ in a neighborhood of } \bar{\theta}. \quad (78)$$

Step 1: Gradient decomposition. By (69),

$$J_0(\pi_\theta) = J^\mathcal{M}(\pi_\theta) + J^\perp(\pi_\theta), \quad \forall \theta.$$

Assume that $J^\mathcal{M}$ and J^\perp are (Fréchet) differentiable with respect to θ on an open neighborhood of $\bar{\theta}$; this is standard under the conditions of the policy-gradient theorem, which ensures differentiability of J_0 [49, 51]. Then, by linearity of differentiation,

$$\nabla_\theta J_0(\pi_\theta) = \nabla_\theta J^\mathcal{M}(\pi_\theta) + \nabla_\theta J^\perp(\pi_\theta), \quad \text{for all } \theta \text{ near } \bar{\theta}. \quad (79)$$

Step 2: Vanishing on-manifold gradient at a local maximum. Under the local maximality condition (78), the first-order necessary condition for unconstrained optimization gives

$$\nabla_\theta J^\mathcal{M}(\pi_\theta) \Big|_{\theta=\bar{\theta}} = 0.$$

(If Π is itself subject to additional parameter constraints, one may interpret this as a vanishing gradient along feasible directions; the argument below is then applied in the corresponding tangent space.)

Substituting this into (79) we obtain

$$\nabla_\theta J_0(\pi_\theta) \Big|_{\theta=\bar{\theta}} = \nabla_\theta J^\perp(\pi_\theta) \Big|_{\theta=\bar{\theta}}. \quad (80)$$

Thus, *exactly at* $\bar{\theta}$, the naive-RL gradient coincides with the ghost-gradient $\nabla_\theta J^\perp$; all infinitesimal improvement directions for J_0 come from changing the ghost component.

Step 3: Interpretation via the policy-gradient theorem. The policy-gradient theorem for episodic MDPs (e.g. [49, 51]) states that

$$\nabla_\theta J_0(\pi_\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) A_t^{(0)} \right], \quad (81)$$

where $A_t^{(0)}$ is an advantage function associated with the total reward r_t and the expectation is over trajectories generated by π_θ in the world model. Similarly, using $r_t^\mathcal{M}$ and r_t^\perp as rewards in the same MDP, we obtain

$$\nabla_\theta J^\mathcal{M}(\pi_\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) A_t^{(\mathcal{M})} \right], \quad (82)$$

$$\nabla_\theta J^\perp(\pi_\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) A_t^{(\perp)} \right], \quad (83)$$

for appropriate advantage functions $A_t^{(\mathcal{M})}$ and $A_t^{(\perp)}$.

By construction $r_t = r_t^{\mathcal{M}} + r_t^{\perp}$ and linearity of the value-function operators, one can choose the advantage functions such that

$$A_t^{(0)} = A_t^{(\mathcal{M})} + A_t^{(\perp)}, \quad \forall t, \quad (84)$$

which recovers (79) when substituted into (81).

At $\theta = \bar{\theta}$, the local optimality of $J^{\mathcal{M}}$ implies that the contribution of $A_t^{(\mathcal{M})}$ integrates to zero in (82); equivalently, $\nabla_{\theta} J^{\mathcal{M}}(\pi_{\theta})|_{\theta=\bar{\theta}} = 0$. Hence, by (81)–(83) and (84), we recover the exact equality (80).

In practice, when $J^{\mathcal{M}}$ is only approximately locally maximal (e.g., due to finite-sample estimation and function approximation error), we obtain the approximate relation

$$\|\nabla_{\theta} J_0(\pi_{\theta})|_{\theta=\bar{\theta}} - \nabla_{\theta} J^{\perp}(\pi_{\theta})|_{\theta=\bar{\theta}}\| = \|\nabla_{\theta} J^{\mathcal{M}}(\pi_{\theta})|_{\theta=\bar{\theta}}\| \approx 0,$$

which justifies the “ \approx ” symbol in (72) of the main text.

Combining Steps 1–3, we conclude that near a local maximizer of the on-manifold performance $J^{\mathcal{M}}$, the gradient of the naive-RL objective J_0 is dominated (indeed, at the maximizer: entirely given) by the ghost-gradient $\nabla_{\theta} J^{\perp}$. This formalizes the statement that gradient-based RL is locally incentivized to move into directions which increase the expected ghost component, completing the proof of part (ii). \square

D Proofs for Section 5: Structural Baselines

D.1 Proof of Proposition 7

In this appendix we provide a more detailed argument for the law-alignment properties of the structural baselines and, more generally, of the structural class \mathcal{S} introduced in Definition 8. We work under the standing assumptions of Section 3.

Setup and notation. Recall that the (law-consistent) generator G^{\star} produces total-variance surfaces $(w_t)_{t \geq 0}$ with $w_t \in \mathcal{M}_{\text{vol}}$ almost surely for all t , cf. Proposition 5. The world model $f_{\theta^{\star}}$ takes as input a feature vector x_t (which may include past surfaces and positions) and outputs a prediction $\hat{w}_{t+1} = f_{\theta^{\star}}(x_t)$, with approximation residual

$$e_{t+1} := \hat{w}_{t+1} - w_{t+1}.$$

The volatility law-penalty functional is

$$\mathcal{L}_{\text{vol}}(\hat{w}_{t+1}) = \frac{1}{2} \text{dist}(\hat{w}_{t+1}, \mathcal{M}_{\text{vol}})^2 = \frac{1}{2} \|\hat{w}_{t+1} - \Pi_{\mathcal{M}_{\text{vol}}}(\hat{w}_{t+1})\|_2^2,$$

where $\Pi_{\mathcal{M}_{\text{vol}}}$ is the Euclidean projection onto \mathcal{M}_{vol} .

For a policy π , we denote by $\overline{\text{LP}}(\pi)$ the long-run average (or finite-horizon normalized) law penalty, and by $\text{GFI}(\pi)$ its Graceful Failure Index, as defined in Section 4.4. We recall that these quantities have the schematic form

$$\overline{\text{LP}}(\pi) \approx \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\mathcal{L}_{\text{vol}}(\hat{w}_{t+1}) \mid \pi], \quad (85)$$

$$\text{GFI}(\pi) = \frac{\mu_{\text{law}}^{\text{shock}}(\pi) - \mu_{\text{law}}^{\text{base}}(\pi)}{I_{\text{shock}}}, \quad (86)$$

where $\mu_{\text{law}}^{\text{base}}(\pi)$ and $\mu_{\text{law}}^{\text{shock}}(\pi)$ are aggregate law metrics in baseline and shock regimes, and $I_{\text{shock}} > 0$ encodes shock intensity (which is fixed for all policies).

Finally, by Definition 8, a structural policy $\pi \in \mathcal{S}$ is of the form

$$a_t = \pi(s_t) = g_{\theta}(z_t),$$

where z_t is a low-dimensional signal (e.g. realized vol trend) and g_{θ} is a Lipschitz parametric map with bounded leverage, $\|g_{\theta}(z)\| \leq L_{\text{max}}$ for all z and all $\theta \in \Theta$, with Θ compact. The baselines b^{ZH} (Zero-Hedge), b^{VT} (Vol-Trend), and b^{RG} (Random-Gaussian) are specific instances in \mathcal{S} .

The proof proceeds in three steps. First, we show that law penalties are uniformly bounded in terms of the world-model error. Second, we use the structural priors of \mathcal{S} to obtain uniform moment bounds on the state process and hence on the residuals e_{t+1} . Third, we translate these to bounds on $\overline{\text{LP}}$ and GFI .

Step 1: Bounding law penalties by world-model error. By Proposition 3, we know that $\mathcal{L}_{\text{vol}}(w_{t+1}) = 0$ whenever $w_{t+1} \in \mathcal{M}_{\text{vol}}$. Since the generator is law-consistent, $w_{t+1} \in \mathcal{M}_{\text{vol}}$ almost surely, and thus for any realization of w_{t+1} and any prediction \hat{w}_{t+1} we have

$$\text{dist}(\hat{w}_{t+1}, \mathcal{M}_{\text{vol}}) \leq \|\hat{w}_{t+1} - w_{t+1}\|_2 = \|e_{t+1}\|_2.$$

Consequently,

$$\mathcal{L}_{\text{vol}}(\hat{w}_{t+1}) = \frac{1}{2} \text{dist}(\hat{w}_{t+1}, \mathcal{M}_{\text{vol}})^2 \leq \frac{1}{2} \|e_{t+1}\|_2^2. \quad (87)$$

Taking expectations conditional on the policy π and the regime (baseline or shock), we obtain

$$\mathbb{E}[\mathcal{L}_{\text{vol}}(\hat{w}_{t+1}) \mid \pi] \leq \frac{1}{2} \mathbb{E}[\|e_{t+1}\|_2^2 \mid \pi]. \quad (88)$$

Thus, uniform control of the second moment of the approximation error e_{t+1} under a policy class implies uniform control of the law penalties under that class.

Step 2: Uniform second-moment bounds for structural policies. By Proposition 6 and the conditions stated there, the world model f_{θ^*} is Lipschitz in its input x_t , and the approximation error e_{t+1} is bounded (in second moment) on bounded subsets of the input space. Concretely, there exist constants $L_f, C_e < \infty$ such that, for any two inputs x, x' ,

$$\|f_{\theta^*}(x) - f_{\theta^*}(x')\|_2 \leq L_f \|x - x'\|_2, \quad (89)$$

$$\sup_{\|x\| \leq R} \mathbb{E}[\|f_{\theta^*}(x) - w_{t+1}\|_2^2 \mid x_t = x] \leq C_e(R), \quad \text{for each finite } R > 0, \quad (90)$$

where $C_e(R)$ is non-decreasing in R .

The key structural property of \mathcal{S} is that it induces a *bounded-exposure* Markov chain on the joint state (w_t, h_t, z_t) , where h_t denotes the position vector (portfolio holdings) and z_t denotes the low-dimensional signals. More precisely:

Lemma 5 (Uniform moment control under \mathcal{S}). *Under the law-consistent generator and the structural priors of Definition 8 (bounded leverage L_{\max} , compact parameter set Θ , and Lipschitz g), there exists a constant $C_{\text{state}} < \infty$ such that for every structural policy $\pi \in \mathcal{S}$ and every t ,*

$$\mathbb{E}_{\pi}[\|x_t\|_2^2] \leq C_{\text{state}},$$

where x_t is the world-model input constructed from (w_t, h_t, z_t) .

Proof. (Sketch.) The generator G^* yields an exogenous process (w_t) whose second moments are bounded over the finite horizon $t = 0, \dots, T$, by standard properties of the underlying stochastic-volatility model and the projection onto \mathcal{M}_{vol} (see, e.g., stability of affine and rough volatility models [8, 9]). The signals z_t are Lipschitz functions of (w_0, \dots, w_t) (e.g., realized-vol trends, moving averages), and thus inherit bounded second moments over $t = 0, \dots, T$.

For a structural policy $\pi \in \mathcal{S}$, the position process $h_t = g_{\theta}(z_t)$ satisfies $\|h_t\|_2 \leq L_{\max}$ almost surely for all t , by bounded leverage and compact Θ . Hence $\mathbb{E}[\|h_t\|_2^2] \leq L_{\max}^2$ for all t . Since x_t is built from (w_t, h_t, z_t) via a fixed, Lipschitz feature map (stacking, scaling, etc.), we obtain $\mathbb{E}_{\pi}[\|x_t\|_2^2] \leq C_{\text{state}}$ for some finite constant C_{state} independent of $\pi \in \mathcal{S}$ and t . A fully rigorous version uses a Lyapunov-function argument for the Markov chain induced by \mathcal{S} ; see Appendix D.1. \square

Combining Lemma 5 with (90), we obtain a uniform bound on the second moment of the residuals under \mathcal{S} : there exists $\bar{C}_e < \infty$ such that, for all $\pi \in \mathcal{S}$ and all t ,

$$\mathbb{E}_{\pi}[\|e_{t+1}\|_2^2] = \mathbb{E}_{\pi}[\mathbb{E}[\|e_{t+1}\|_2^2 \mid x_t]] \leq \mathbb{E}_{\pi}[C_e(\|x_t\|)] \leq \bar{C}_e, \quad (91)$$

where we used monotonicity of $C_e(\cdot)$ and the uniform bound $\mathbb{E}_{\pi}[\|x_t\|_2^2] \leq C_{\text{state}}$.

Step 3: Bounds on $\overline{\text{LP}}$ and GFI. Using (88) and (91), we obtain that for any $\pi \in \mathcal{S}$ and any t ,

$$\mathbb{E}_{\pi}[\mathcal{L}_{\text{vol}}(\hat{w}_{t+1})] \leq \frac{1}{2} \mathbb{E}_{\pi}[\|e_{t+1}\|_2^2] \leq \frac{1}{2} \bar{C}_e.$$

Substituting into (85) and taking the supremum over $\pi \in \mathcal{S}$ yields

$$\sup_{\pi \in \mathcal{S}} \overline{\text{LP}}(\pi) \leq \frac{1}{2} \bar{C}_e =: C_{\text{LP}} < \infty. \quad (92)$$

This shows that \mathcal{S} is law-aligned in the sense of Definition 9, establishing item (i) of the proposition.

For the GFI, recall its definition in (86). The numerator is the difference in aggregate law metrics between the shock and baseline regimes. Under our shock design (multiplying long variance and spot volatility by bounded factors), the generator remains law-consistent and the world model is evaluated on a distorted, but still bounded, region of the state space. By the same reasoning as above (with possibly different constants), there exist finite constants $C_{\text{law}}^{\text{base}}$ and $C_{\text{law}}^{\text{shock}}$ such that for all $\pi \in \mathcal{S}$,

$$\mu_{\text{law}}^{\text{base}}(\pi) \leq C_{\text{law}}^{\text{base}}, \quad \mu_{\text{law}}^{\text{shock}}(\pi) \leq C_{\text{law}}^{\text{shock}}.$$

Hence

$$|\mu_{\text{law}}^{\text{shock}}(\pi) - \mu_{\text{law}}^{\text{base}}(\pi)| \leq C_{\text{law}}^{\text{base}} + C_{\text{law}}^{\text{shock}} =: \tilde{C}_{\text{law}} < \infty \quad (93)$$

for all $\pi \in \mathcal{S}$. Since $I_{\text{shock}} > 0$ is fixed and does not depend on π , we obtain

$$\sup_{\pi \in \mathcal{S}} \text{GFI}(\pi) = \sup_{\pi \in \mathcal{S}} \frac{\mu_{\text{law}}^{\text{shock}}(\pi) - \mu_{\text{law}}^{\text{base}}(\pi)}{I_{\text{shock}}} \leq \frac{\tilde{C}_{\text{law}}}{I_{\text{shock}}} =: C_{\text{GFI}} < \infty. \quad (94)$$

This establishes the existence of a finite constant C_{GFI} depending only on the generator, the shock specification, and the world-model error, proving item (i) and the first part of item (ii).

Baselines b^{ZH} and b^{VT} . We now specialize to the two deterministic baselines.

Zero-Hedge b^{ZH} . By definition, b^{ZH} takes no positions, $h_t \equiv 0$, at all times. The state process affecting the world model thus reduces to the exogenous generator path (w_t, z_t) , and the residual distribution is precisely the “background” world-model error profile studied in Section. In particular, the bounds (91), (92), and (94) hold with $\pi = b^{\text{ZH}}$. Empirically, this baseline indeed exhibits the smallest observed law penalties and GFI, matching the theoretical role of b^{ZH} as the law-neutral benchmark.

Vol-Trend b^{VT} . The Vol-Trend baseline applies a one-factor trend-following rule with bounded leverage, $h_t = g_{\theta^{\text{VT}}}(z_t)$, where z_t encodes recent volatility trends and $g_{\theta^{\text{VT}}}$ is Lipschitz with $\|g_{\theta^{\text{VT}}}(z)\| \leq L_{\text{max}}$ for all z . As in Lemma 5, this ensures that positions respond smoothly to volatility changes and remain uniformly bounded in second moment, so that (w_t, z_t, h_t) stays in a bounded region of the state space. Therefore the same world-model error and law-penalty bounds apply, and we obtain

$$\overline{\text{LP}}(b^{\text{ZH}}), \overline{\text{LP}}(b^{\text{VT}}) \leq C_{\text{LP}}, \quad \text{GFI}(b^{\text{ZH}}), \text{GFI}(b^{\text{VT}}) \leq C_{\text{GFI}}.$$

This proves item (ii), with the empirical strictness “below unconstrained RL levels” arising from the ghost-incentive effect (naive RL actively amplifies exposure in regions where the ghost component r^\perp is large, whereas b^{ZH} and b^{VT} do not target such regions).

Random-Gaussian baseline b^{RG} . Finally, consider the Random-Gaussian baseline b^{RG} , which draws actions from a Gaussian distribution $a_t \sim \mathcal{N}(0, \Sigma_a)$ with fixed covariance Σ_a , possibly truncated to enforce a leverage bound. Provided $\text{tr}(\Sigma_a) < \infty$ and the truncation is such that $\mathbb{E}[\|a_t\|_2^2] \leq L'_{\text{max}} < \infty$, we obtain moment bounds on the position process (h_t) analogous to those for b^{VT} , although with larger constants reflecting the additional randomness. Repeating the argument leading to (91)–(94), we find finite constants $C'_{\text{LP}}, C'_{\text{GFI}} < \infty$ such that

$$\overline{\text{LP}}(b^{\text{RG}}) \leq C'_{\text{LP}}, \quad \text{GFI}(b^{\text{RG}}) \leq C'_{\text{GFI}}.$$

In general one expects $C'_{\text{LP}}, C'_{\text{GFI}}$ to be larger than $C_{\text{LP}}, C_{\text{GFI}}$, because b^{RG} explores a wider range of states and can occasionally spend more time in regions where the world-model error and induced law penalties are higher. This matches the empirical role of b^{RG} as a noisy, less structured baseline.

Collecting the above bounds, we see that the structural class \mathcal{S} is uniformly law-aligned, and that the specific structural baselines b^{ZH} and b^{VT} enjoy particularly favorable law-penalty and GFI levels relative to unconstrained RL policies. This completes the proof of Proposition 7.

E Empirical Results: From RL Dynamics to Law-Strength Frontiers (Supplementary)

In this appendix we provide complementary quantitative details to Section. We (i) restate the full metric tables for all strategies in both baseline and shock regimes, (ii) tabulate the law-strength frontier points used in Figures, and (iii) describe the precise numerical procedures used to construct the frontiers and diagnostic plots.

Throughout, all metrics are computed on the same evaluation trajectories used to generate Figures, and hence are fully consistent with the summary numbers reported.

Table 4: Full step-wise metrics for all strategies in the *baseline* regime. Mean and standard deviation of step P&L, Sharpe ratio, mean and maximum law penalty, law-adjusted return, Graceful Failure Index (GFI), law coverage at two thresholds, and 5% tail risk measures (VaR₅, CVaR₅). All numbers are computed on the same evaluation trajectories

Strategy	Mean P&L	Std P&L	Sharpe	Mean LawPen	Max LawPen	Law-Adj Ret	GFI	LawCov < 0.003	LawCov < 0.006	VaR ₅ / CVaR ₅
Naive RL (PPO)	-0.0022	0.0127	-0.1696	0.006994	0.020850	-0.0057	1.2675	0.5306	0.6122	-0.0228 / -0.0261
Law-Seeking RL (PPO)	-0.0150	0.0129	-1.1564	0.007861	0.022554	-0.0189	1.6621	0.4898	0.5714	-0.0361 / -0.0394
Zero-Hedge baseline	0.0191	0.0064	2.9944	0.005501	0.018318	0.0164	0.0000	0.6122	0.6939	0.0139 / 0.0139
Random-Gaussian baseline	0.0099	0.0107	0.9235	0.005510	0.020686	0.0072	1.2062	0.6204	0.6918	-0.0088 / -0.0161
Vol-Trend baseline	0.0146	0.0074	1.9636	0.005344	0.017173	0.0119	0.0000	0.6122	0.6939	0.0045 / 0.0033

Table 5: Baseline-regime metrics for all law-strength frontier points: Naive RL ($\lambda = 0$), soft law-seeking RL for $\lambda \in \{5, 10, 20, 40\}$, selection-only RL, and structural baselines. These are the points used to construct the law-strength frontier and diagnostic plots

Strategy / λ	Mean P&L	Std P&L	Sharpe	Mean LawPen	Max LawPen	Law-Adj Ret	GFI	LawCov < 0.003	LawCov < 0.006	VaR ₅ / CVaR ₅
Naive RL ($\lambda = 0$)	-0.0022	0.0127	-0.1696	0.006994	0.020850	-0.0057	1.2675	0.5306	0.6122	-0.0228 / -0.0261
Soft RL ($\lambda = 5$)	-0.0202	0.0120	-1.6753	0.006472	0.019956	-0.0234	2.0663	0.5510	0.6327	-0.0399 / -0.0429
Soft RL ($\lambda = 10$)	-0.0175	0.0123	-1.4248	0.003709	0.013169	-0.0194	2.8059	0.7347	0.7959	-0.0354 / -0.0387
Soft RL ($\lambda = 20$)	-0.0204	0.0131	-1.5629	0.003962	0.014251	-0.0224	3.0728	0.7143	0.7755	-0.0414 / -0.0454
Soft RL ($\lambda = 40$)	-0.0092	0.0054	-1.7138	0.004737	0.015888	-0.0116	0.8443	0.6531	0.7347	-0.0134 / -0.0134
Selection-only RL	-0.0223	0.0139	-1.6028	0.007923	0.022827	-0.0263	2.0407	0.4898	0.5714	-0.0448 / -0.0489
Zero-Hedge baseline	0.0191	0.0064	2.9944	0.005501	0.018318	0.0164	0.0000	0.6122	0.6939	0.0139 / 0.0139
Random-Gaussian baseline	0.0099	0.0107	0.9235	0.005510	0.020686	0.0072	1.2062	0.6204	0.6918	-0.0088 / -0.0161
Vol-Trend baseline	0.0146	0.0074	1.9636	0.005344	0.017173	0.0119	0.0000	0.6122	0.6939	0.0045 / 0.0033

E.1 Full metrics: baseline regime

Table 4 reports the complete set of step-wise metrics for all RL and structural strategies in the baseline regime, corresponding to the top half of Table in the main text. We include mean and standard deviation of step P&L, Sharpe ratio, mean and maximum law-penalty, law-adjusted return, Graceful Failure Index (GFI), law coverage at two thresholds, and 5% tail risk measures (VaR and CVaR).

As noted in Section, the baseline regime already exhibits the main qualitative pattern: the structurally constrained baselines (Zero-Hedge, Vol-Trend) lie in a region of high Sharpe and moderate law penalties, with GFI effectively zero, while all RL variants—including law-seeking PPO—display negative mean P&L and substantially higher GFI values.

E.2 Full metrics: shock regime

In which we apply a volatility shock (long-variance multiplier 4, spot-vol multiplier 2) to the underlying generator while keeping the world model fixed.

The law-strength frontier plots in Figures correspond to projections of Table 5 onto two-dimensional planes, for example:

1. mean law penalty vs. GFI,
2. mean law penalty vs. Sharpe,
3. mean law penalty vs. VaR₅ or CVaR₅.

Within the law-penalty band $[0.0053, 0.0057]$ highlighted, Zero-Hedge lies near the lower boundary with high Sharpe and $GFI \approx 0$, while the closest RL variants in the band have Sharpe < 0 and $GFI > 1.5$, illustrating the empirical Pareto dominance emphasized in the main text.

E.3 Frontier construction and banding procedure

For completeness, we describe the numerical procedure used to construct the law-strength frontiers and penalty bands.

Policy set. We collect the following set of evaluated policies:

1. Naive RL (PPO) trained on pure P&L,
2. soft law-seeking RL for $\lambda \in \{5, 10, 20, 40\}$,
3. selection-only RL (trained on P&L, selected by law metrics),

4. structural baselines: Zero-Hedge, Random-Gaussian, Vol-Trend.

For each policy we compute the metric vector

$$\mathbf{m}(\pi) = (\mu_{P\&L}(\pi), \sigma_{P\&L}(\pi), \text{Sharpe}(\pi), \mu_{\text{law}}(\pi), \text{GFI}(\pi), \text{VaR}_5(\pi), \text{CVaR}_5(\pi)),$$

with components given explicitly in Tables 4–5.

Penalty banding. To compare policies at similar levels of law violation, we discretize the range of mean law penalties $\mu_{\text{law}}(\pi)$ into contiguous bands

$$B_k = [\ell_k, u_k), \quad k = 1, \dots, K,$$

where (ℓ_k, u_k) are chosen such that the bands roughly align with the empirical distribution of $\mu_{\text{law}}(\pi)$ across policies. For the numerical example in Section 7, we use the band $[0.0053, 0.0057]$ that contains the Zero-Hedge baseline and at least one RL policy. For each band B_k we identify the subset

$$\Pi_k := \{\pi : \mu_{\text{law}}(\pi) \in B_k\},$$

and perform intra-band comparisons of Sharpe, GFI, VaR₅, and CVaR₅.

Pareto frontier extraction. Given a subset of metrics $(\mu_{\text{law}}, \text{Sharpe})$, $(\mu_{\text{law}}, \text{GFI})$, or $(\mu_{\text{law}}, \text{CVaR}_5)$, we say that a policy π *Pareto-dominates* another policy π' if it is no worse on all objectives and strictly better on at least one. For example, in the $(\mu_{\text{law}}, \text{Sharpe})$ plane we define

$$\pi \succ \pi' \iff \mu_{\text{law}}(\pi) \leq \mu_{\text{law}}(\pi') \text{ and } \text{Sharpe}(\pi) \geq \text{Sharpe}(\pi'),$$

with at least one strict inequality. The empirical law-strength frontier is the set

$$\mathcal{F} := \{\pi : \nexists \pi' \text{ such that } \pi' \succ \pi\}.$$

In all planes considered, the points corresponding to Zero-Hedge and Vol-Trend lie on \mathcal{F} , while all RL variants lie strictly inside the frontier: there exists at least one structural baseline that has both (i) weakly lower mean law penalty and (ii) strictly better Sharpe, GFI, or tail risk. This is the empirical content of the Pareto-dominance statements.

E.4 Additional notes on dynamics and diagnostics

For completeness, we briefly comment on the additional curves and histograms:

1. **Dynamics Plots** Each panel shows time-series of cumulative P&L and mean law penalty across episodes for Naive RL, a representative law-seeking RL variant (e.g. $\lambda = 20$), and a structural baseline. The additional curves omitted from the main text for brevity (e.g. $\lambda = 5, 10, 40$) display the same qualitative pattern: increasing λ reduces law penalties at the cost of lower P&L, in line with the structural trade-off.
2. **Diagnostic Plots.** The scatter plots aggregate step-level observations of P&L vs. law penalty across many episodes. The high-density clusters for RL policies lie in regions of elevated law penalty, while structural baselines concentrate near lower penalty bands. Histograms of law penalty show that RL policies allocate substantial mass to the right tail of the penalty distribution, consistent with exploiting the ghost channel r^\perp identified in Proposition 6.

These supplementary plots are therefore consistent with, and reinforce, the three main empirical takeaways: (i) law penalties do not rescue naive RL from ghost arbitrage, (ii) structural baselines define an empirical law-strength frontier, and (iii) unconstrained law-seeking RL lies strictly inside this frontier in our volatility world-model testbed.

F No-Free-Lunch for Law-Seeking RL (Proofs)

This appendix provides a detailed proof of the no-free-lunch result stated as Theorem. We proceed in three steps: (i) we recall the performance and law-metric quantities used in the main result, (ii) we formalize a mild ghost-law monotonicity condition that links off-manifold reward to law metrics, and (iii) we give a detailed proof of the theorem.

F.1 Preliminaries: performance, law metrics, and decomposition

Recall from Sections 2, 3 that for any stationary policy π in the unconstrained policy class Π we have:

1. A *baseline* environment distribution \mathbb{P}^{base} over episodes generated by the law-consistent synthetic generator and the world model.
2. A *shock* distribution $\mathbb{P}^{\text{shock}}$ describing the same world model, but driven by shocked long-variance and spot volatility factors.
3. A per-step reward $r(s_t, a_t)$ that can be decomposed as

$$r(s_t, a_t) = r^{\mathcal{M}}(s_t, a_t) + r^{\perp}(s_t, a_t), \quad (95)$$

where $r^{\mathcal{M}}$ is the on-manifold component defined by projection onto the volatility law manifold \mathcal{M}_{vol} and r^{\perp} is the ghost (off-manifold) component; see Propositions 4 and 6.

We denote by $J_0(\pi)$ the expected (per-step or per-episode) P&L under the baseline regime:

$$J_0(\pi) := \mathbb{E}_{\mathbb{P}^{\text{base}}} [r(s_t, a_t)], \quad a_t \sim \pi(\cdot | s_t), \quad (96)$$

and similarly define the on-manifold and ghost components

$$J^{\mathcal{M}}(\pi) := \mathbb{E}_{\mathbb{P}^{\text{base}}} [r^{\mathcal{M}}(s_t, a_t)], \quad J^{\perp}(\pi) := \mathbb{E}_{\mathbb{P}^{\text{base}}} [r^{\perp}(s_t, a_t)], \quad (97)$$

so that

$$J_0(\pi) = J^{\mathcal{M}}(\pi) + J^{\perp}(\pi). \quad (98)$$

Law metrics. For each policy π we also consider a vector of law metrics, measuring law violations and robustness:

$$\ell(\pi) := \left(\mu_{\text{law}}^{\text{base}}(\pi), \mu_{\text{law}}^{\text{shock}}(\pi), \text{GFI}(\pi) \right) \in \mathbb{R}^3, \quad (99)$$

where:

1. $\mu_{\text{law}}^{\text{base}}(\pi)$ is the mean step law-penalty under \mathbb{P}^{base} ,
2. $\mu_{\text{law}}^{\text{shock}}(\pi)$ is the mean step law-penalty under $\mathbb{P}^{\text{shock}}$, and
3. $\text{GFI}(\pi)$ is the Graceful Failure Index introduced in Section 4.4, which normalizes the change in law metrics between baseline and shock by the shock intensity I_{shock} .

We write $\ell(\pi) \leq_{\text{law}} \ell(\pi')$ if each component of $\ell(\pi)$ is less than or equal to the corresponding component of $\ell(\pi')$, i.e. if policy π is at least as law-aligned and robust as π' .

Structural class. The structural strategy class \mathcal{S} , introduced in Section 4.4, is a low-capacity, law-aligned subset of Π consisting of structurally constrained strategies such as Zero-Hedge and Vol-Trend. We assume that \mathcal{S} satisfies the *near-optimal on-manifold* property of Assumption 1: there exists $\pi_{\mathcal{S}}^* \in \mathcal{S}$ and $\varepsilon_{\mathcal{S}} \geq 0$ such that

$$J^{\mathcal{M}}(\pi) \leq J^{\mathcal{M}}(\pi_{\mathcal{S}}^*) + \varepsilon_{\mathcal{S}} \quad \text{for all } \pi \in \Pi \text{ with trajectories supported on } \mathcal{M}_{\text{vol}}, \quad (100)$$

and $\ell(\pi_{\mathcal{S}}^*)$ is uniformly small in the sense of Proposition 7 (law-aligned class).

F.2 Ghost-law monotonicity

The results show that in our volatility world-model testbed:

1. the ghost reward r^{\perp} is generated by deviations in the normal cone $N_{\mathcal{M}_{\text{vol}}}(w)$ to the volatility law manifold,
2. the law penalty \mathcal{L}_{vol} grows with the squared distance to \mathcal{M}_{vol} (Propositions 3 and 4),
3. the Graceful Failure Index GFI captures the *relative increase* in law penalties under shock.

These observations motivate the following mild monotonicity condition linking ghost reward and law metrics.

Assumption E.1 (Ghost-law monotonicity). *There exists a law-aligned structural reference policy $\pi_{\mathcal{S}}^* \in \mathcal{S}$ and a non-decreasing function $\psi : \mathbb{R}_+^3 \rightarrow \mathbb{R}_+$ with $\psi(0, 0, 0) = 0$ such that for every policy $\pi \in \Pi$ we have*

$$J^{\perp}(\pi) - J^{\perp}(\pi_{\mathcal{S}}^*) \leq \psi\left((\ell(\pi) - \ell(\pi_{\mathcal{S}}^*))_+\right), \quad (101)$$

where $(\cdot)_+$ denotes component-wise positive part. In particular, if $\ell(\pi) \leq_{\text{law}} \ell(\pi_S^*)$ then $J^\perp(\pi) \leq J^\perp(\pi_S^*)$.

Intuitively, Assumption E.1 states that *any additional ghost reward beyond what is available to the structural class must be paid for by worse law metrics*. In our volatility setting, the existence of such a function ψ is supported by:

1. the Lipschitz bound of Proposition 4, which controls $|r^\perp(w)|$ in terms of the distance to \mathcal{M}_{vol} ;
2. the definition of law penalties and GFI as functions of the same distance and its behavior under shock;
3. the law-aligned nature of π_S^* guaranteed by Proposition 7, which ensures that $\ell(\pi_S^*)$ is close to the best achievable law metrics given the world-model approximation error.

Thus, in the neighborhood of π_S^* , additional ghost reward can only be obtained by moving further away from the law manifold, which necessarily worsens at least one component of the law metric vector.

F.3 Proof of the no-free-lunch theorem

We now give a detailed proof of the no-free-lunch result stated as theorem. For clarity we restate the theorem in a slightly more quantitative form.

Theorem (No-free-lunch for unconstrained law-seeking RL). Assume:

1. The structural class \mathcal{S} is law-aligned and near-optimal on-manifold in the sense of (100), with reference policy π_S^* .
2. The world model induces a non-trivial ghost channel as in Propositions 6 and 2.
3. Ghost-law monotonicity (Assumption E.1) holds for π_S^* .

Then for any $\eta > \varepsilon_S$ and any policy $\pi \in \Pi$ satisfying

$$J_0(\pi) \geq J_0(\pi_S^*) + \eta, \quad (102)$$

we must have

$$\ell(\pi) \not\leq_{\text{law}} \ell(\pi_S^*), \quad (103)$$

i.e., at least one component of the law metric vector is strictly worse for π than for π_S^* . In particular, no policy $\pi \in \Pi$ can strictly dominate π_S^* on all axes (P&L, law penalties, and GFI).

Proof. Fix $\eta > \varepsilon_S$ and suppose, for contradiction, that there exists a policy $\pi \in \Pi$ such that (102) holds and

$$\ell(\pi) \leq_{\text{law}} \ell(\pi_S^*). \quad (104)$$

We will show that this contradicts the decomposition (98), near-optimality (100), and ghost-law monotonicity (101).

Step 1: Decomposing the P&L difference. By (98), for any π we have

$$J_0(\pi) = J^\mathcal{M}(\pi) + J^\perp(\pi).$$

Therefore,

$$J_0(\pi) - J_0(\pi_S^*) = (J^\mathcal{M}(\pi) - J^\mathcal{M}(\pi_S^*)) + (J^\perp(\pi) - J^\perp(\pi_S^*)). \quad (105)$$

Step 2: Bounding the on-manifold component. By near-optimality of π_S^* on the law manifold, (100) implies that

$$J^\mathcal{M}(\pi) \leq J^\mathcal{M}(\pi_S^*) + \varepsilon_S \quad \text{for all } \pi \in \Pi. \quad (106)$$

Substituting (106) into (105) yields

$$J_0(\pi) - J_0(\pi_S^*) \leq \varepsilon_S + (J^\perp(\pi) - J^\perp(\pi_S^*)). \quad (107)$$

Step 3: Applying ghost-law monotonicity. By the assumption (104), $\ell(\pi) \leq_{\text{law}} \ell(\pi_S^*)$, we have

$$(\ell(\pi) - \ell(\pi_S^*))_+ = 0,$$

so ghost-law monotonicity (101) gives

$$J^\perp(\pi) - J^\perp(\pi_S^*) \leq \psi(0, 0, 0) = 0. \quad (108)$$

Combining (107) and (108) yields

$$J_0(\pi) - J_0(\pi_S^*) \leq \varepsilon_S. \quad (109)$$

Step 4: Contradiction. However, by assumption (102) we have

$$J_0(\pi) - J_0(\pi_S^*) \geq \eta > \varepsilon_S,$$

which contradicts (109). Therefore no policy π can simultaneously satisfy (102) and (104). Equivalently, any policy achieving an improvement in baseline P&L of more than ε_S over the structural reference π_S^* must have at least one law metric strictly worse than that of π_S^* , i.e. $\ell(\pi) \not\leq_{\text{law}} \ell(\pi_S^*)$.

In particular, if we interpret the triple

$$(-J_0(\pi), \ell(\pi)) \in \mathbb{R}^{1+3}$$

as a four-dimensional loss vector (profitability vs. law alignment and robustness), then π_S^* cannot be strictly Pareto-dominated by any policy in Π . This is the claimed no-free-lunch property. \square

F.4 Discussion of assumptions and empirical alignment

The proof above separates the no-free-lunch result into three conceptually distinct ingredients:

1. *On-manifold near-optimality of \mathcal{S} .* The structural class \mathcal{S} contains a policy π_S^* that is nearly optimal in terms of on-manifold P&L, as quantified by (100). Empirically this corresponds to the observation that Zero-Hedge and Vol-Trend already sit very close to the empirical law-strength frontier.
2. *Non-trivial ghost channel.* The world model induces off-manifold reward r^\perp that is non-zero whenever predictions deviate from \mathcal{M}_{vol} (Propositions 6 and 2), creating the possibility of “ghost arbitrage”.
3. *Ghost-law monotonicity.* Assumption E.1 formalizes the idea that exploiting the ghost channel necessarily worsens law metrics: additional ghost reward cannot be obtained at fixed or improved law penalties and GFI.

In our volatility world-model testbed, these three conditions are jointly consistent with the empirical findings of Section:

1. The structural baselines achieve high Sharpe and low GFI while remaining close to the law manifold, in line with near-optimality (100).
2. Naive and law-seeking RL policies that outperform \mathcal{S} in raw P&L do so only by moving into high-penalty, high-GFI regions, as evidenced by the frontier and diagnostic plots.
3. No RL policy lies on the empirical law-strength frontier once the structural baselines are included, which matches the impossibility of strict dominance established by the theorem.

References

- [1] B. Dupire. Pricing with a smile. *Risk*, 7(1):18–20, 1994.
- [2] J. Gatheral. *The Volatility Surface: A Practitioner’s Guide*. John Wiley & Sons, 2006.
- [3] M. R. Fengler. *Semiparametric Modeling of Implied Volatility*. Springer, 2005.
- [4] P. Carr and D. B. Madan. Option valuation using the fast Fourier transform. *Journal of Computational Finance*, 2(4):61–73, 1999.
- [5] R. W. Lee. The moment formula for implied volatility at extreme strikes. *Mathematical Finance*, 14(3):469–480, 2004.
- [6] R. Cont and P. Tankov. *Financial Modelling with Jump Processes*. Chapman and Hall/CRC, 2004.
- [7] L. Bergomi. *Stochastic Volatility Modeling*. Chapman and Hall/CRC, 2016.
- [8] J. Gatheral and A. Jacquier. Arbitrage-free SVI volatility surfaces. *Quantitative Finance*, 14(1):59–71, 2014.
- [9] C. Bayer, P. Friz, and J. Gatheral. Pricing under rough volatility. *Quantitative Finance*, 16(6):887–904, 2016.
- [10] B. Horváth, A. Muguruza, and M. Tomas. Deep learning volatility: Deep calibration of rough stochastic volatility models. *Quantitative Finance*, 21(1):11–29, 2021.

- [11] B. Horváth, A. Jacquier, and C. Kovács. Deep learning volatility. *Quantitative Finance*, 21(11):1763–1783, 2021.
- [12] A. Itkin. Calibration of local stochastic volatility models to market data. *Journal of Computational Finance*, 18(3):1–46, 2015.
- [13] S. De Marco and P. Henry-Labordère. Linking vanillas and VIX options: A constrained optimization approach. *Journal of Computational Finance*, 19(1):29–64, 2015.
- [14] D. Guterding. A sparse modeling approach to the arbitrage-free interpolation of plain-vanilla option prices and implied volatilities. *Risks*, 11(1):1–28, 2023.
- [15] J. Ruf and W. Wang. Neural network-based option pricing. *arXiv preprint arXiv:1912.02710*, 2020.
- [16] P. Carr and D. B. Madan. A note on sufficient conditions for no arbitrage. *Available at SSRN*, 2005.
- [17] P. S. Hagan and G. West. Interpolation methods for curve construction. *Applied Mathematical Finance*, 13(2):89–129, 2006.
- [18] D. Filipović. *Term-Structure Models: A Graduate Course*. Springer, 2009.
- [19] D. Filipović. *Consistency Problems for Heath–Jarrow–Morton Interest Rate Models*. Springer, 2001.
- [20] T. Björk and B. J. Christensen. Interest rate dynamics and consistent forward rate curves. *Mathematical Finance*, 9(4):323–348, 1999.
- [21] P. Artzner, F. Delbaen, J.-M. Eber, and D. Heath. Coherent measures of risk. *Mathematical Finance*, 9(3):203–228, 1999.
- [22] R. T. Rockafellar and S. Uryasev. Optimization of conditional value-at-risk. *Journal of Risk*, 2(3):21–42, 2000.
- [23] Á. Cartea, S. Jaimungal, and J. Penalva. *Algorithmic and High-Frequency Trading*. Cambridge University Press, 2015.
- [24] B. Hurst, Y. H. Ooi, and L. H. Pedersen. A century of evidence on trend-following investing. *Journal of Portfolio Management*, 44(1):15–29, 2017.
- [25] Z. Zhang and S. Zohren. Deep reinforcement learning for trading. *Journal of Financial Data Science*, 2(2):25–40, 2020.
- [26] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.
- [27] S. N. Cohen and R. J. Elliott. *Stochastic Calculus and Applications*. Birkhäuser, 2nd edition, 2015.
- [28] M. Belkin, P. Niyogi, and V. Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2399–2434, 2006.
- [29] C. Fefferman, S. Mitter, and H. Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
- [30] J. Han, A. Jentzen, and W. E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [31] C. Beck, S. Becker, P. Cheridito, A. Jentzen, and A. Neufeld. Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward SDEs. *Journal of Scientific Computing*, 79(3):1393–1438, 2019.
- [32] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [33] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [34] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang. Physics-informed machine learning. *Acta Numerica*, 30:1–146, 2021.
- [35] G. Carleo, J. Cirac, K. Cranmer, et al. Machine learning and the physical sciences. *Reviews of Modern Physics*, 91(4):045002, 2019.
- [36] J. Willard, X. Jia, S. Xu, M. Steinbach, and V. Kumar. Integrating physics-based modeling with machine learning: A survey. *ACM Computing Surveys*, 52(3):1–39, 2020.
- [37] C. Beck, C. Ehlers, et al. Solving stochastic differential equations and Kolmogorov equations by deep learning. *Annals of Applied Probability*, 31(4):1917–1966, 2021.

- [38] J. Brandstetter, M. Welling, and A. Ansuini. Message passing neural PDE solvers. In *International Conference on Learning Representations*, 2022.
- [39] P. W. Battaglia, J. B. Hamrick, V. Bapst, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- [40] D. Ha and J. Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- [41] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba. Learning latent dynamics for planning from pixels. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- [42] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- [43] D. Hafner, J. Pasukonis, J. Ba, and M. Norouzi. Mastering diverse domains with world models. *arXiv preprint arXiv:2301.04104*, 2023.
- [44] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *Advances in Neural Information Processing Systems*, 2018.
- [45] M. Janner, J. Fu, M. Zhang, and S. Levine. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems*, 2019.
- [46] Ł. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine, et al. Model-based reinforcement learning for Atari. In *International Conference on Learning Representations*, 2020.
- [47] J. Schrittwieser, I. Antonoglou, T. Hubert, et al. Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature*, 588:604–609, 2020.
- [48] V. Mnih, K. Kavukcuoglu, D. Silver, et al. Human-level control through deep reinforcement learning. *Nature*, 518:529–533, 2015.
- [49] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd edition, 2018.
- [50] V. R. Konda and J. N. Tsitsiklis. Actor-critic algorithms. In *Advances in Neural Information Processing Systems*, pages 1008–1014, 2000.
- [51] P. S. Thomas. Bias in natural actor-critic algorithms. In *Proceedings of the 31st International Conference on Machine Learning*, pages 441–448, 2014.
- [52] J. Schulman, P. Moritz, S. Levine, M. I. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. In *International Conference on Learning Representations*, 2016.
- [53] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [54] E. Altman. *Constrained Markov Decision Processes*. Chapman & Hall/CRC, 1999.
- [55] J. Achiam, D. Held, A. Tamar, and P. Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, pages 22–31, 2017.
- [56] S. Mannor and J. N. Tsitsiklis. Mean-variance optimization in Markov decision processes. *Operations Research*, 59(2):350–367, 2011.
- [57] Y. Chow, A. Tamar, S. Mannor, and M. Pavone. Risk-sensitive and robust decision-making: A CVaR optimization approach. In *Advances in Neural Information Processing Systems*, pages 1522–1530, 2015.
- [58] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley. A survey of multi-objective sequential decision-making. *Journal of Artificial Intelligence Research*, 48:67–113, 2013.
- [59] G. Neu, A. György, and C. Szepesvári. A unified view of entropy-regularized Markov decision processes. In *Advances in Neural Information Processing Systems*, 2017.
- [60] M. Geist, B. Scherrer, and O. Pietquin. A theory of regularized Markov decision processes. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2160–2169, 2019.
- [61] Y. Jiang, et al. Towards safe reinforcement learning: A survey and outlook. *arXiv preprint arXiv:2109.14597*, 2021.
- [62] Z. Hou, Y. Chen, and M. Wang. Regularized policy optimization in MDPs with constraints. *SIAM Journal on Optimization*, 31(1):192–223, 2021.
- [63] J. Moody and M. Saffell. Learning to trade via direct reinforcement. *IEEE Transactions on Neural Networks*, 12(4):875–889, 2001.

- [64] Y. Deng, F. Bao, Y. Kong, Z. Ren, and Q. Dai. Deep direct reinforcement learning for financial signal representation and trading. *IEEE Transactions on Neural Networks and Learning Systems*, 28(3):653–664, 2016.
- [65] Z. Jiang, D. Xu, and J. Liang. A deep reinforcement learning framework for the financial portfolio management problem. *arXiv preprint arXiv:1706.10059*, 2017.
- [66] X.-Y. Liu, R. Yang, Z. Zhu, et al. FinRL: A deep reinforcement learning library for quantitative finance. *ACM Transactions on Management Information Systems*, 12(2):1–30, 2021.
- [67] P. N. Kolm and G. Ritter. Dynamic replication and hedging: A reinforcement learning approach. *The Journal of Financial Data Science*, 1(2):71–88, 2019.
- [68] H. Buehler, L. Gonon, J. Teichmann, and B. Wood. Deep hedging. *Quantitative Finance*, 19(8):1271–1291, 2019.
- [69] H. Buehler, L. Gonon, J. Teichmann, and B. Wood. Deep hedging. *Quantitative Finance*, 19(8):1271–1291, 2019.
- [70] H. Buehler, L. Gonon, J. Teichmann, and B. Wood. Deep hedging. *Quantitative Finance*, 19(8):1271–1291, 2019.
- [71] A. Carboneau and F. Godin. Deep hedging of derivatives with transaction costs and different risk criteria. *Journal of Computational Finance*, 24(2):1–31, 2020.
- [72] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- [73] G. Zerveas, S. Jayaraman, D. Patel, A. Bhamidipati, C. Eickhoff, J. Lévy, and X. Binefa. A transformer-based framework for multivariate time series representation learning. In *Proceedings of the 27th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2114–2124, 2021.
- [74] D. Manheim and S. Garraabrant. Categorizing variants of Goodhart’s law. *arXiv preprint arXiv:1803.04585*, 2019.
- [75] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [76] J. Leike, M. Martic, V. Krakovna, et al. AI safety gridworlds. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [77] T. Everitt, M. Hutter, and J. Leike. Reinforcement learning with a corrupted reward channel. *arXiv preprint arXiv:1705.08417*, 2017.
- [78] D. Krueger, T. Maharaj, S. Rahman, et al. Hidden incentives in deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2020.
- [79] M. Pan, et al. The effects of reward misspecification: Mapping and mitigating misaligned models. *arXiv preprint arXiv:2209.14935*, 2022.
- [80] C. F. Hayes, E. Bargiacchi, et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous Agents and Multi-Agent Systems*, 36:1–35, 2022.
- [81] Unknown. Placeholder entry for the key WangRLA1pha2019. No reliable bibliographic information could be identified; please replace this with the correct reference or remove the citation. 2019.
- [82] A. Neelakantan, et al. Placeholder for “Reinforcement learning with verifiable rewards”. No corresponding published work could be found; please update or delete this citation. 2025.