# Epistemic Constraint Theory

## A Unifying Framework for Inference Limitations Across Bayesian Epistemology, Information Theory, and Decision Theory

Boris Kriger[1]

[1]Institute of Integrative and Interdisciplinary Research,
boriskriger@interdisciplinary-institute.org

### Abstract

We present *Epistemic Constraint Theory*, a conceptual framework that highlights structural parallels among results from Bayesian epistemology, information theory, optimal stopping, and machine learning. We observe that several well-known phenomena—the dominance of inductive bias, the value of query selection, the cost of premature commitment, and the importance of prior calibration—share a common pattern: the structure of the admissible hypothesis space bounds inferential performance. We do not prove these phenomena derive from a single theorem; rather, we provide unified vocabulary and explicit parallels across domains. Our main technical contribution is an information-theoretic reformulation of optimal stopping (Theorem 5.4): under entropy sufficiency conditions (holding exactly for binary hypothesis testing), the optimal stopping rule is characterized by an entropy threshold $\theta^* = \Theta(c/\mu)$. For binary testing, this reformulates the classical Sequential Probability Ratio Test in information-theoretic terms; for symmetric multi-hypothesis problems, it provides new perspective. The theorem does not apply to general asymmetric multi-hypothesis settings. We provide concrete examples and simulations, and derive operational corollaries that restate the main results as behavioral principles. This framework offers pedagogical value for practitioners recognizing similar constraint-type problems across application domains.

**Keywords:** Bayesian inference, hypothesis space constraints, active learning, optimal stopping, inductive bias, information geometry, epistemic limitations

**MSC 2020:** 62C10, 62F15, 94A15, 91B06, 62L15

## 1 Introduction

### 1.1 Motivation and Scope

Researchers across multiple disciplines have independently discovered that inferential performance is often limited not by computational resources or data quantity, but by structural properties of the inference problem itself. In machine learning, the bias-variance tradeoff and no-free-lunch theorems [1] establish that model class expressivity constrains learning. In active learning, query selection strategies dominate passive data accumulation [2, 3]. In decision theory, optimal stopping results characterize when commitment should be delayed [4, 5]. In Bayesian statistics, prior selection fundamentally shapes posterior inference [6, 7].

This paper does not claim to discover new mathematical facts in any of these areas. Rather, we propose that these diverse results share a common structural foundation: *constraints on the admissible hypothesis space bound inferential accuracy independently of—and often more severely than—computational capacity or data volume.*

We call this perspective *Epistemic Constraint Theory*. Our contributions are:

(i) A unified formal framework that reveals structural parallels across disciplines.

(ii) Quantitative bounds for *soft* constraint regimes (extending beyond hard exclusion).

(iii) A novel characterization of optimal commitment timing under residual uncertainty, expressed directly in information-theoretic terms rather than through value functions.

(iv) Concrete examples and numerical simulations illustrating each structural insight.

## 1.2 Related Work

The results we synthesize draw from extensive prior literature:

**Bayesian foundations and prior specification.** The foundational treatment of Bayesian inference by Bernardo and Smith [8] emphasizes that priors are not merely starting points but structural frameworks defining what a system can infer. This perspective is central to our constraint-based formulation. Modern introductions to Bayesian methods [9, 6] continue to stress the fundamental role of prior specification, while asymptotic theory [10] provides rigorous conditions under which posteriors concentrate on true parameters.

**Inductive bias and model selection.** The importance of hypothesis class selection is well-established in statistical learning theory [11, 12]. The bias-variance tradeoff [13] and no-free-lunch theorems [1] formalize constraints on learning. Recent work on double descent [14, 15] and neural network generalization [16, 17] continues this tradition. The Bayesian model selection literature [7, 31] addresses how the choice of model class bounds inferential accuracy.

**Active learning and experimental design.** The value of query selection over passive observation is central to active learning [2, 18], Bayesian experimental design [19, 20], and information-directed sampling [21]. MacKay's work on information-based objective functions [3] and optimal experimental design [22] provides theoretical foundations.

**Optimal stopping and sequential decisions.** The mathematics of when to commit versus when to wait is developed in optimal stopping theory [4, 23, 5]. Classical results include the secretary problem [25] and Gittins indices for multi-armed bandits [26]. Large deviations theory [45] provides tools for analyzing convergence rates. Applications include bandit algorithms [27] and exploration-exploitation tradeoffs in reinforcement learning [28]. Bayesian Model Averaging [29] offers a practical approach to avoiding premature commitment by maintaining multiple hypotheses simultaneously.

**Prior selection, robustness, and cognitive bias.** Bayesian sensitivity analysis [30], objective Bayes methods [31], and the study of prior-data conflict [32] address how prior constraints affect inference. Transfer learning [33] and pretraining [34] operationalize prior modification in deep learning. In cognitive science, confirmation bias [35] and the "need for closure" [36] demonstrate how rigid priors and premature commitment lead to inferential failures in human reasoning. Adversarial robustness [37, 38] studies how inference can fail when the information environment is actively hostile.

Our contribution is not to extend any single line of work, but to reveal their common structure.

## 1.3 Paper Organization

Section 2 establishes the formal framework. Sections 3–6 develop four structural insights, distinguishing our interpretive contribution from established mathematical facts. Section 7 presents the unifying perspective. Section 8 discusses applications. Section 9 concludes. Appendix B provides numerical simulations.

## 2 Formal Framework

### 2.1 Constrained Inference Systems

**Definition 2.1** (Inference System). An *inference system* is a tuple $\mathcal{S} = \langle \Omega, \mathcal{H}, \pi, \ell \rangle$ where:

- $\Omega$ is the set of possible world-states (parameter space);

- $\mathcal{H} \subseteq 2^{\Omega}$ is the hypothesis space;

- $\pi : \mathcal{H} \to [0, 1]$ is a prior distribution over $\mathcal{H}$ with $\sum_{h \in \mathcal{H}} \pi(h) = 1$;

- $\ell : \mathcal{H} \times \mathcal{E} \to \mathbb{R}_+$ is a likelihood function, where $\mathcal{E}$ is the evidence space.

*Remark* 2.2 (Notation Convention). Throughout this paper, $\mathcal{H}$ (calligraphic) denotes the hypothesis space, while $H[\cdot]$ denotes Shannon entropy. This distinguishes the two uses of "H" common in the literature.

**Definition 2.3** (Soft Constraint). A *soft constraint* is a function $w : \mathcal{H} \to [0, 1]$ assigning admissibility weights to hypotheses. The *constrained prior* is:

$$\pi_w(h) = \frac{\pi(h) \cdot w(h)}{Z_w}, \quad \text{where } Z_w = \sum_{h' \in \mathcal{H}} \pi(h') \cdot w(h'). \tag{1}$$

The *hard constraint* case corresponds to $w(h) \in \{0, 1\}$.

*Remark* 2.4. The soft constraint formulation generalizes hard exclusion ($w(h^*) = 0$) to penalization ($0 < w(h^*) < 1$). This is essential for non-trivial mathematical results, as we discuss below.

**Definition 2.5** (Constrained Posterior). Given evidence $e \in \mathcal{E}$ and soft constraint $w$, the *constrained posterior* is:

$$\pi_w(h|e) = \frac{\ell(h, e) \cdot \pi(h) \cdot w(h)}{Z_w(e)}, \tag{2}$$

where $Z_w(e) = \sum_{h' \in \mathcal{H}} \ell(h', e) \cdot \pi(h') \cdot w(h')$.

### 2.2 Accuracy Measures

We use log-posterior probability as our accuracy measure, which avoids technical issues with KL-divergence from point masses.

**Definition 2.6** (Log-Posterior Accuracy). For true hypothesis $h^* \in \mathcal{H}$, the *log-posterior accuracy* given evidence $e$ is:

$$A(\mathcal{S}, w, e) = \log \pi_w(h^*|e). \tag{3}$$

*Remark* 2.7. This is equivalent to negative log-loss and relates to KL-divergence via $D_{\mathrm{KL}}(\delta_{h^*} \| \pi_w(\cdot|e)) = -\log \pi_w(h^*|e)$ when the divergence is well-defined. We use log-probability directly to avoid technicalities.

## 3 Constraint Dominance: A Structural Observation

### 3.1 The Core Observation

The following is not a new theorem but a structural observation that unifies existing results on inductive bias:

**Proposition 3.1** (Constraint Bound on Accuracy). *For any inference system $\mathcal{S}$ with soft constraint $w$ and true hypothesis $h^*$:*

$$A(\mathcal{S}, w, e) \leq \log w(h^*) + \log \pi(h^*) + \log \ell(h^*, e) - \log Z_w(e). \tag{4}$$

*In particular, if $w(h^*) = 0$, then $\pi_w(h^*|e) = 0$ for all evidence $e$.*

*Proof.* Direct substitution into the definition of constrained posterior (Equation 2). □

*Remark* 3.2. This is a *definitional consequence*, not a deep theorem. Its value lies in making explicit that log-accuracy decomposes into constraint, prior, likelihood, and normalization terms. The constraint term $\log w(h^*)$ provides a ceiling that no amount of favorable evidence can overcome when $w(h^*)$ is small.

## 3.2 Non-Trivial Result: Soft Constraint Convergence Rate

The interesting mathematics emerges when we consider soft constraints and convergence rates.

**Theorem 3.3** (Convergence Rate Under Soft Constraints). *Let $\mathcal{S}$ be an inference system with true hypothesis $h^*$, soft constraint $w$ with $w(h^*) = \varepsilon \in (0,1)$, and let $\{e_n\}_{n=1}^{\infty}$ be an evidence sequence satisfying:*

*(R1)* ***Independence**: $\{e_n\}$ are independent and identically distributed under the true data-generating process $P_{h^*}$;*

*(R2)* ***Absolute continuity**: For all $h, h' \in \mathcal{H}$ with $w(h), w(h') > 0$, the likelihood ratio $\Lambda(e) = \ell(h^*, e)/\ell(h', e)$ is well-defined (i.e., $\ell(h', e) > 0$ $P_{h^*}$-almost surely);*

*(R3)* ***Finite exponential moments**: $\mathbb{E}_{h^*}[\log \Lambda(e)] = \mu_{h'} > 0$, and there exists $\theta_0 > 0$ such that $\mathbb{E}_{h^*}[e^{\theta \log \Lambda(e)}] < \infty$ for all $|\theta| < \theta_0$ (finite MGF in a neighborhood of zero).*

*Let $\mu = \min_{h' \neq h^*, w(h') > 0} \mu_{h'} > 0$ and $\sigma^2 = \text{Var}_{h^*}[\log \Lambda(e)]$ (which is finite under (R3)). Then:*

*(a)* ***Consistency**: $\pi_w(h^*|e_1, \ldots, e_n) \xrightarrow{a.s.} 1$ as $n \to \infty$.*

*(b)* ***Convergence rate**: For any $\epsilon \in (0, \mu)$, there exists a rate function $I(\epsilon) > 0$ such that*

$$P\left(\pi_w(h^*|e_{1:n}) < 1 - \delta\right) \leq |\mathcal{H}| \exp\left(-nI(\epsilon)\right), \tag{5}$$

*where $I(\epsilon) \geq (\mu - \epsilon)^2/(2\sigma^2)$ for $\epsilon$ near $\mu$ (Gaussian approximation).*

*(c)* ***Sample complexity**: The expected number of observations to reach $\pi_w(h^*|e_{1:n}) \geq 1 - \delta$ satisfies*

$$\mathbb{E}[n^*(\varepsilon, \delta)] = \frac{b}{\mu}, \tag{6}$$

*where $b = \log(1/\varepsilon) + \log((1-\delta)/\delta) + \log(\pi(h^*)/\bar{\pi})$ and $\bar{\pi} = \max_{h' \neq h^*} w(h')\pi(h')$. The variance of the stopping time is $\text{Var}[n^*] = b\sigma^2/\mu^3$.*

*Proof.* Define the log-likelihood ratio process $S_n^{(h')} = \sum_{i=1}^{n} \log \Lambda^{(h')}(e_i)$ for alternative $h' \neq h^*$.
   **Part (a):** Under (R1)–(R3), the finite MGF condition implies finite variance, so the Strong Law of Large Numbers applies: $S_n^{(h')}/n \xrightarrow{a.s.} \mu_{h'} > 0$. Thus $S_n^{(h')} \to +\infty$ a.s. By Bayes' theorem:

$$\frac{\pi_w(h'|e_{1:n})}{\pi_w(h^*|e_{1:n})} = \frac{w(h')\pi(h')}{w(h^*)\pi(h^*)} \cdot e^{-S_n^{(h')}} \xrightarrow{a.s.} 0. \tag{7}$$

Since this holds for all $h' \neq h^*$ with $w(h') > 0$, we have $\pi_w(h^*|e_{1:n}) \to 1$ a.s.

**Part (b):** By Cramér's large deviation theorem [45], which requires finite MGF (condition (R3)), for $\alpha < \mu$:

$$P(S_n^{(h')} < n\alpha) \leq \exp(-nI(\alpha)), \tag{8}$$

where $I(\alpha) = \sup_\theta(\theta\alpha - \Lambda^*(\theta))$ is the Cramér rate function and $\Lambda^*(\theta) = \log \mathbb{E}[e^{\theta \log \Lambda}]$ is the cumulant generating function. Note: we do *not* use Hoeffding's inequality, which would require bounded random variables; instead, the large deviations approach handles unbounded log-likelihood ratios under the MGF condition. By the Gaussian approximation near $\mu$, $I(\alpha) \approx (\mu - \alpha)^2/(2\sigma^2)$. Union bound over $|\mathcal{H}|$ alternatives gives the result.

**Part (c):** For $\pi_w(h^*|e_{1:n}) \geq 1 - \delta$, we need $S_n \geq b$ where $b = \log(1/\varepsilon) + \log(\bar{\pi}/\pi(h^*)) + \log((1 - \delta)/\delta)$ and $S_n = \min_{h'} S_n^{(h')}$. By Wald's identity for random walks with drift $\mu$ and variance $\sigma^2$ per step, the hitting time $\tau_b = \inf\{n : S_n \geq b\}$ satisfies:

$$\mathbb{E}[\tau_b] = \frac{b}{\mu}, \qquad \text{Var}[\tau_b] = \frac{b\sigma^2}{\mu^3}. \tag{9}$$

The expected sample complexity is thus exactly $b/\mu$ with no first-order correction term. (An earlier version of this paper incorrectly included an $O(\sigma/\mu)$ correction; this conflated variance of hitting time with bias in expected hitting time.) $\qquad\square$

*Remark* 3.4. Condition (R3) (finite MGF) is stronger than finite variance but standard in large deviations theory [45]. It is satisfied by most common distributions (Gaussian, bounded, exponential family) but excludes heavy-tailed distributions. For heavy-tailed likelihoods, weaker polynomial concentration bounds apply but with slower rates. The simulations in Appendix B focus on the binary Bernoulli case where the MGF condition is trivially satisfied.

**Example 3.5** (Binary Hypothesis Testing)**.** Consider $\mathcal{H} = \{h_0, h_1\}$ with $h_1 = h^*$, uniform prior $\pi(h_0) = \pi(h_1) = 0.5$, and constraint $w(h_0) = 1$, $w(h_1) = \varepsilon$. Evidence consists of i.i.d. Bernoulli trials with $P(e = 1|h_1) = 0.8$ and $P(e = 1|h_0) = 0.2$.

The constrained prior is $\pi_w(h_1) = \varepsilon/(1 + \varepsilon)$. After $n$ observations with $k$ successes:

$$\pi_w(h_1|e_{1:n}) = \frac{\varepsilon \cdot 0.8^k \cdot 0.2^{n-k}}{\varepsilon \cdot 0.8^k \cdot 0.2^{n-k} + 0.2^k \cdot 0.8^{n-k}}. \tag{10}$$

Here $\mu = D_{\text{KL}}(\text{Bern}(0.8)\|\text{Bern}(0.2)) = 0.8 \log 4 + 0.2 \log(1/4) \approx 0.83$ nats. By Equation 6, achieving 95% posterior with $\varepsilon = 0.01$ requires approximately $n^* \approx \log(100)/0.83 \approx 5.5$ additional observations compared to $\varepsilon = 1$.

## 3.3   Connection to Inductive Bias Literature

Proposition 3.1 and Theorem 3.3 formalize the well-known principle that model class selection (inductive bias) fundamentally constrains learning [39]. The soft constraint $w$ can be interpreted as:

- **Model architecture**: Neural network architecture choices exclude certain function classes.

- **Regularization**: $L_2$ regularization corresponds to Gaussian prior, penalizing large weights.

- **Feature selection**: Excluding features reduces hypothesis space dimensionality.

The key insight is not mathematical novelty but *unified language*: the same mathematical structure governs constraints in statistical learning, cognitive psychology, and scientific methodology.

# 4 Query Selection: The Value of Information

## 4.1 Framework

**Definition 4.1** (Query and Information Gain). A *query* $Q : \Omega \to \mathcal{Y}$ maps world-states to observations. The *information gain* of query $Q$ is:

$$\mathrm{IG}(Q) = H[\mathcal{H}] - \mathbb{E}_Q[H[\mathcal{H}|Q]] = I(\mathcal{H}; Q), \tag{11}$$

where $H[\cdot]$ denotes Shannon entropy and $I(\cdot; \cdot)$ denotes mutual information.

**Definition 4.2** (Query Strategy). A *query strategy* $\sigma$ is a policy mapping current beliefs and remaining budget to the next query. The *passive strategy* $\sigma_0$ draws observations without selection. The *active strategy* $\sigma^*$ selects queries to maximize expected information gain.

## 4.2 A Standard Result, Unified Interpretation

The following is well-known in active learning [2, 3]:

**Proposition 4.3** (Query Selection Dominance). *Let $\mathcal{Q} = \{Q_1, \ldots, Q_m\}$ be available queries. Under budget constraint $B$ (number of queries), define:*

- *$\sigma^*$: greedy strategy selecting $Q_i^* = \mathrm{argmax}_Q I(\mathcal{H}; Q|Q_1^*, \ldots, Q_{i-1}^*)$ at step $i$;*

- *$\sigma_0$: passive strategy selecting queries uniformly at random.*

*If queries have non-negative conditional information gain (i.e., $I(\mathcal{H}; Q|past) \geq 0$ for all $Q$ and all histories), then:*

$$\mathbb{E}[IG(\sigma^*, B)] \geq \mathbb{E}[IG(\sigma_0, B)]. \tag{12}$$

*More precisely, if $I(\mathcal{H}; Q_i^*|Q_1^*, \ldots, Q_{i-1}^*) \geq \gamma > 0$ for all $i \leq B$ (queries remain informative), then:*

$$\mathbb{E}[IG(\sigma^*, B)] - \mathbb{E}[IG(\sigma_0, B)] \geq B \cdot \left(\gamma - \frac{1}{m} \sum_{Q \in \mathcal{Q}} I(\mathcal{H}; Q)\right), \tag{13}$$

*which is positive when optimal queries exceed average informativeness.*

*Proof.* The non-negativity of mutual information ensures that total information gain is non-decreasing with budget. The greedy strategy achieves at least $(1 - 1/e)$ of the optimal by submodularity of mutual information [40]. The difference bound follows from comparing per-step gains. $\square$

*Remark* 4.4. The condition that queries remain informative ($I(\mathcal{H}; Q_i^*|\mathrm{past}) \geq \gamma > 0$) is non-trivial. It fails when:

- Early queries exhaust available information (entropy reaches zero);

- Queries are highly correlated (conditioning on past queries makes new queries uninformative).

Without this condition, the sum in the original bound could be negative after some steps. This is not a new result; the value of Proposition 4.3 in our framework is interpretive: query selection is a form of *dynamic constraint modification*.

### 4.3 Value of Information Characterization

Following [41, 42]:

**Definition 4.5** (Value of Information). The *value of information* (VoI) of query $Q$ for decision problem with utility $U$ is:

$$\text{VoI}(Q) = \mathbb{E}_Q\left[\max_a \mathbb{E}[U(a,\omega)|Q]\right] - \max_a \mathbb{E}[U(a,\omega)]. \tag{14}$$

**Proposition 4.6** (VoI-Information Gain Relationship). *Under logarithmic utility $U(a,\omega) = \log P(\omega|a)$:*

$$VoI(Q) = I(\mathcal{H}; Q). \tag{15}$$

This well-known result [20] connects information theory to decision theory, providing another perspective on why query selection matters.

**Example 4.7** (Medical Diagnosis). Consider diagnosing a rare disease ($h_1$, prevalence 1%) versus healthy ($h_0$). Two tests are available:

- Test A: Sensitivity 99%, specificity 90%.

- Test B: Sensitivity 80%, specificity 99%.

For a patient with positive Test A, the information gain from Test B (specificity-focused) exceeds that from repeating Test A. Active query selection achieves diagnosis with fewer tests than random testing.

## 5 Uncertainty Tolerance: Optimal Commitment Timing

This section contains our most substantive mathematical contribution: a characterization of optimal commitment timing that, while building on optimal stopping theory [4, 5], provides a novel formulation in terms of residual uncertainty.

### 5.1 Relation to Classical Optimal Stopping

Classical optimal stopping results, such as the secretary problem [25] and Gittins indices [26], express stopping rules in terms of value functions or index policies. Our contribution differs in three ways:

1. We express the stopping threshold directly in *information-theoretic terms* (residual entropy), rather than through value iteration.

2. We characterize when entropy alone is a sufficient statistic for optimal stopping (Lemma 5.3), clarifying the scope of entropy-based rules.

3. We provide explicit bounds relating the threshold to the cost-information tradeoff (Part (c) of Theorem 5.5).

This bridges the gap between decision-theoretic stopping rules and information-theoretic measures of uncertainty.

## 5.2 Sequential Decision Framework

**Definition 5.1** (Sequential Inference Problem). A *sequential inference problem* consists of:

- Time horizon $T \in \mathbb{N} \cup \{\infty\}$;

- Evidence process $\{E_t\}_{t=1}^T$ with $E_t$ revealed at time $t$;

- Posterior process $\{\pi_t\}_{t=0}^T$ with $\pi_t = \pi(\cdot | E_1, \ldots, E_t)$;

- Commitment action $\Gamma$ that selects $\hat{h} = \text{argmax}_h \pi_\tau(h)$ at stopping time $\tau$;

- Accuracy reward $R(\tau) = \mathbf{1}[\hat{h} = h^*]$;

- Waiting cost $c > 0$ per time step.

**Definition 5.2** (Residual Uncertainty). The *residual uncertainty* at time $t$ is:

$$\text{RU}(t) = H[\pi_t] = -\sum_{h \in \mathcal{H}} \pi_t(h) \log \pi_t(h). \tag{16}$$

## 5.3 Main Result

The following lemma clarifies when entropy-based thresholds are sufficient statistics for optimal stopping.

**Lemma 5.3** (Entropy Sufficiency Conditions). *Let $\{\pi_t\}$ be a posterior process on finite hypothesis space $\mathcal{H}$. The residual uncertainty $RU(t) = H[\pi_t]$ is a sufficient statistic for optimal stopping decisions if either:*

*(S1) **Symmetric information**: The conditional distribution of $RU(t+1)$ given $\mathcal{F}_t$ depends on $\pi_t$ only through $RU(t)$; or*

*(S2) **Binary hypothesis space**: $|\mathcal{H}| = 2$; or*

*(S3) **Bounded posterior heterogeneity**: There exists $\eta > 0$ such that for all posteriors $\pi, \pi'$ with $H[\pi] = H[\pi']$:*

$$\max_h |\pi(h) - \pi'(h)| \leq \eta.$$

*This is a geometric condition on entropy level sets in the probability simplex.*

*Proof.* (S1): If the transition kernel for $RU(t+1)$ depends only on $RU(t)$, then the value function inherits this property by backward induction on the Bellman equation.

(S2): For $|\mathcal{H}| = 2$, entropy is a bijection with $\max_h \pi(h)$ on $[0.5, 1]$: $H[\pi] = -p \log p - (1-p) \log(1-p)$ where $p = \max_h \pi(h)$. Thus both the immediate reward $\max_h \pi(h)$ and the continuation value (which depends on the distribution of future posteriors) are determined by $RU(t)$.

(S3): For $|\mathcal{H}| = k$, the entropy level set $\{\pi : H[\pi] = r\}$ is a $(k-2)$-dimensional manifold in the $(k-1)$-simplex. Condition (S3) requires this manifold to have bounded diameter in the $L^\infty$ metric. For $k = 2$, the level set consists of exactly two points (related by symmetry), so $\eta = 0$ trivially. For $k > 2$, entropy level sets can have large diameter (e.g., $(0.5, 0.5, 0, \ldots, 0)$ and $(1/k, \ldots, 1/k)$ can have similar entropy but very different max probability). When (S3) holds with small $\eta$, the immediate reward $\max_h \pi(h)$ varies by at most $\eta$ across the level set, implying the value function varies by at most $\eta T$ over horizon $T$. This yields an $\eta T$-suboptimal threshold policy. $\square$

*Remark* 5.4. Condition (S3) is now stated in verifiable geometric terms rather than in terms of the (unknown) optimal value function $V$. For most multi-hypothesis problems with $|\mathcal{H}| > 2$, condition (S3) fails because entropy level sets have large diameter. We therefore restrict Theorem 5.5 to exact sufficiency conditions (S1) or (S2).

**Theorem 5.5** (Entropy Reformulation of Binary Optimal Stopping). *Consider a sequential inference problem with:*

(A1) **Informative evidence**: $\mathbb{E}[RU(t+1)|\mathcal{F}_t] < RU(t)$ *whenever* $RU(t) > 0$;

(A2) **Bounded information rate**: *There exist* $0 < \underline{\mu} \le \bar{\mu} < \infty$ *such that*

$$\underline{\mu} \le \mathbb{E}[RU(t) - RU(t+1)|\mathcal{F}_t] \le \bar{\mu}$$

*whenever* $RU(t) > 0$;

(A3) **Positive waiting cost**: $c > 0$;

(A4) **Binary hypothesis space**: $|\mathcal{H}| = 2$;

(A5) **Symmetric structure**: $\pi(h_0) = \pi(h_1) = 1/2$ *and the accuracy reward* $R(\tau) = \mathbf{1}[\hat{h} = h^*]$ *treats both hypotheses symmetrically.*

*Then:*

(a) *The optimal stopping time* $\tau^*$ *exists and is finite almost surely.*

(b) $\tau^*$ *is characterized by a threshold policy: there exists* $\theta^* > 0$ *such that*

$$\tau^* = \inf\{t \ge 0 : RU(t) \le \theta^*\}. \tag{17}$$

(c) *The threshold satisfies* $\theta^* = \Theta(c/\mu)$ *when* $c/\mu$ *is in a moderate range (bounded away from 0 and* $\log 2$*). More precisely:*

$$\theta^* \in \left[ \frac{c}{\bar{\mu} \cdot \max_\theta |f'(\theta)|}, \frac{c}{\underline{\mu} \cdot \min_{\theta \in [\delta, \log 2 - \delta]} |f'(\theta)|} \right] \tag{18}$$

*for any* $\delta > 0$, *where* $f'(\theta) = dp/dH$ *is the derivative of the entropy-to-probability map.*

(d) **(Qualitative)** *Premature commitment at* $\tau < \tau^*$ *results in suboptimal expected reward:* $\mathbb{E}[R(\tau)] < \mathbb{E}[R(\tau^*)]$. *The gap* $V_{cont}(RU(\tau)) - V_{stop}(RU(\tau)) > 0$ *depends on* $RU(\tau) - \theta^*$ *but no closed-form bound is available without additional distributional assumptions.*

*Remark* 5.6 (On Assumption (A5)). Assumption (A5) (symmetric priors and costs) is essential for the symmetry-based proof. Without it, the value function $V(\pi)$ depends on the full posterior $\pi$, not just entropy or $|L|$. For asymmetric priors, the optimal stopping boundary in $(p, 1 - p)$ space is not symmetric about $p = 1/2$. The theorem can be extended to asymmetric binary testing, but the threshold would be in terms of log-likelihood ratio rather than entropy.

*Remark* 5.7 (Scope Limitation). This theorem applies only under the stated assumptions. In practice:

- **Symmetric binary testing**: The theorem holds exactly. This covers symmetric clinical trials and A/B testing with equal priors.

- **Asymmetric binary testing**: The symmetry argument fails; optimal stopping depends on the log-likelihood ratio, not entropy. SPRT [24] provides the appropriate framework.

- **Multi-hypothesis** ($|\mathcal{H}| > 2$): The theorem does not apply. Entropy is not a sufficient statistic.

*Proof.* **Part (a):** Under (A1), RU($t$) is a supermartingale bounded below by 0. By the martingale convergence theorem, RU($t$) → RU$_\infty$ a.s. for some RU$_\infty \geq 0$. Combined with (A2), if RU$_\infty > 0$, then $\mathbb{E}[\text{RU}(t) - \text{RU}(t+1)] \geq \underline{\mu} > 0$ infinitely often, contradicting convergence. Thus RU($t$) → 0 a.s., and $\tau_0 = \inf\{t : \text{RU}(t) = 0\}$ satisfies $\mathbb{E}[\tau_0] \leq \text{RU}(0)/\underline{\mu} < \infty$.

Under (A3), the expected cost of waiting until $\tau_0$ is $c \cdot \mathbb{E}[\tau_0] < \infty$. The optimal stopping problem is thus well-posed with finite optimal stopping time [5, 23].

**Part (b):** Define the value function $V(t, \pi_t)$ as the expected net reward (accuracy minus waiting cost) from state $\pi_t$ at time $t$ under optimal policy. The Bellman equation is:

$$V(t, \pi_t) = \max\left\{ \max_h \pi_t(h),\ -c + \mathbb{E}[V(t+1, \pi_{t+1})|\mathcal{F}_t] \right\}. \tag{19}$$

Under (A4) with exact sufficiency, by Lemma 5.3, both the immediate reward $\max_h \pi_t(h)$ and the continuation value depend on $\pi_t$ only through RU($t$). Write $V(t, \pi_t) = \tilde{V}(t, \text{RU}(t))$.

Define $\tilde{V}_{\text{stop}}(r) = f(r)$ where $f$ maps entropy to $\max_h \pi(h)$ (this is a bijection under (S2)). Define $\tilde{V}_{\text{cont}}(r) = -c + \mathbb{E}[\tilde{V}(t+1, \text{RU}(t+1))|\text{RU}(t) = r]$.

Key property: $\tilde{V}_{\text{stop}}(r) = f(r)$ is strictly decreasing in $r$ (lower entropy ⇒ higher maximum probability, since $f$ is the inverse of binary entropy restricted to [0.5, 1]).

**Establishing threshold optimality (binary case):** For $|\mathcal{H}| = 2$, we establish threshold optimality via the log-likelihood ratio representation and symmetry, *without* claiming that $\tilde{V}_{\text{cont}}$ is monotonic in entropy.

Let $L_t = \log(\pi_t(h_1)/\pi_t(h_0))$ be the log-posterior ratio. For binary testing, $L_t$ evolves as a random walk with drift toward $+\infty$ or $-\infty$ depending on the true hypothesis. Crucially, entropy $\text{RU}_t = H(\pi_t)$ is a deterministic function of $|L_t|$:

$$\text{RU} = H(\sigma(L)) = -\sigma(L)\log\sigma(L) - \sigma(-L)\log\sigma(-L), \tag{20}$$

where $\sigma(L) = 1/(1 + e^{-L})$ is the sigmoid function. This function is symmetric in $L$: $H(\sigma(L)) = H(\sigma(-L))$, and strictly decreasing in $|L|$.

**Symmetry-based argument:** The binary testing problem is symmetric under the exchange $h_0 \leftrightarrow h_1$, which corresponds to $L \leftrightarrow -L$. Under this symmetry:

- The immediate reward $\max_h \pi(h) = \max(\sigma(L), \sigma(-L)) = \sigma(|L|)$ depends only on $|L|$;

- The transition dynamics are symmetric in $L$;

- Therefore, the optimal value function depends on $\pi$ only through $|L|$, i.e., $V(\pi) = \tilde{V}(|L|)$.

For the reduced one-dimensional problem in $|L|$, at $|L| = 0$ (maximum entropy), continuation is preferred for small $c$; at $|L| = \infty$ (certainty), stopping is preferred. By continuity, there exists a threshold $L^* > 0$ such that the optimal stopping region is $\{|L| \geq L^*\}$.

Since entropy $\text{RU} = g(|L|)$ where $g$ is strictly decreasing, the stopping region $|L| \geq L^*$ corresponds exactly to $\text{RU} \leq \theta^*$ where $\theta^* = g(L^*)$. This establishes threshold optimality in entropy space.

**What we do not claim:** We do not claim $\tilde{V}_{\text{cont}}(r)$ is monotonic in entropy $r$. The threshold structure follows from symmetry and the one-dimensional reduction, not from monotonicity of the continuation value.

**Threshold existence:** The threshold $\theta^* > 0$ exists because:

- At $\text{RU} = 0$ (certainty): stopping is optimal since no information gain is possible.

- At $\text{RU} = \log 2$ (maximum uncertainty): continuation is optimal for small $c$ since expected information gain exceeds cost.

The indifference point $\theta^* \in (0, \log 2)$ exists by continuity.

**Part (c):** At the threshold $\theta^*$, indifference holds: $\tilde{V}_{\text{stop}}(\theta^*) = \tilde{V}_{\text{cont}}(\theta^*)$. For binary $|\mathcal{H}| = 2$, the entropy-probability relationship is explicit: $p = \max_h \pi(h)$ satisfies $H(p) = -p \log p - (1 - p) \log(1 - p) = \theta^*$. The indifference condition implies that the expected gain from one more observation equals the cost $c$.

The derivative $f'(\theta) = dp/dH = -1/\log(p/(1-p))$ satisfies $|f'(\theta)| \to \infty$ as $\theta \to 0$ (certainty) and as $\theta \to \log 2$ (maximum uncertainty). However, $|f'|$ is bounded below on any compact subinterval $[\delta, \log 2 - \delta]$, with minimum approximately 1.13 at $p \approx 0.77$.

Under (A2), expected entropy reduction is between $\underline{\mu}$ and $\bar{\mu}$. The indifference condition gives:

$$c \approx |f'(\theta^*)| \cdot \mu, \tag{21}$$

so $\theta^* = \Theta(c/\mu)$ holds when $\theta^*$ lies in a region where $|f'|$ is bounded (i.e., when $c/\mu$ is in a moderate range). For very small $c/\mu$, the threshold $\theta^* \to 0$; for very large $c/\mu$, $\theta^* \to \log 2$.

**Part (d):** At $\tau < \tau^*$ with $\text{RU}(\tau) > \theta^*$, the optimal policy prescribes continuation, so $\tilde{V}_{\text{cont}}(\text{RU}(\tau)) > \tilde{V}_{\text{stop}}(\text{RU}(\tau))$. The suboptimality of premature stopping is exactly this gap. $\square$

*Remark* 5.8 (On the Proof Strategy). An earlier version of this proof attempted to establish monotonicity of $\tilde{V}_{\text{cont}}$ via convexity of the value function. This was erroneous: the function $f(r) = \max_h \pi(h)$ mapping entropy to max probability is *concave*, not convex, so that argument fails. The corrected proof uses symmetry of the binary testing problem to reduce to a one-dimensional problem in $|L|$. Critically, we do not claim $\tilde{V}_{\text{cont}}$ is monotonic—the threshold structure arises from symmetry, not monotonicity.

*Remark* 5.9 (Relation to SPRT and Novelty). For binary hypothesis testing, Theorem 5.5 is a reformulation, not a new result. The SPRT of Wald [24] characterizes optimal sequential testing under Type I/Type II error constraints; our theorem uses a cost-per-observation model with accuracy reward. These are distinct optimization problems, but both admit threshold solutions due to the shared symmetry of binary testing.

Our contribution is expressing the threshold in information-theoretic terms: stop when entropy falls below $\theta^* = \Theta(c/\mu)$. This is a change of variables from likelihood-ratio thresholds to entropy thresholds, providing alternative vocabulary rather than new optimality results. For genuine novelty beyond binary testing, one would need to prove threshold optimality for $|\mathcal{H}| > 2$ with explicit conditions and bounds.

*Remark* 5.10. Assumption (A4) restricts scope to binary ($|\mathcal{H}| = 2$) or symmetric multi-hypothesis problems. For general asymmetric $|\mathcal{H}| > 2$, the theorem does not apply because the symmetry argument fails and entropy is not a sufficient statistic. We view the contribution as providing unified vocabulary connecting optimal stopping to information theory, not as proving fundamentally new stopping rules.

*Remark* 5.11. Assumption (A1) (informative evidence) is crucial and non-trivial. It fails when:

- **Adversarial information environments**: If an adversary controls evidence to maintain uncertainty, as studied in adversarial robustness [37, 38], then $\mathbb{E}[\text{RU}(t+1)|\mathcal{F}_t]$ may equal or exceed $\text{RU}(t)$.

- **Non-identifiable problems**: If multiple hypotheses generate identical observation distributions, no evidence can separate them.

- **Extremely noisy channels**: If the channel capacity is below the entropy rate of the prior, information accumulation may be negligible.

Without (A1), early commitment can be optimal, and our result does not apply. This connects to the adversarial robustness literature's finding that standard inference can fail catastrophically under distribution shift.

**Example 5.12** (Clinical Trial Stopping). Consider a trial comparing treatment $(h_1)$ versus placebo $(h_0)$. Patient outcomes $E_t \sim \text{Bernoulli}(p_{h^*})$ with $p_{h_1} = 0.7$, $p_{h_0} = 0.5$. Daily cost $c = 0.01$ (normalized).

This is a binary hypothesis problem, so (S2) holds exactly. By Theorem 5.5, the optimal stopping threshold satisfies $\theta^* \in [c/\bar{\mu}, c/\underline{\mu}]$. Here $\mu = D_{\text{KL}}(\text{Bern}(0.7) \| \text{Bern}(0.5)) \approx 0.082$ nats, giving $\theta^* \approx 0.01/0.082 \approx 0.12$ nats.

This corresponds to posterior confidence $\max_h \pi(h) \approx 0.88$ before stopping. See Appendix B for simulation confirming this threshold.

## 5.4 Connection to Exploration-Exploitation

Theorem 5.5 relates to the exploration-exploitation tradeoff in reinforcement learning [28, 27]. "Exploration" corresponds to waiting (gathering information), while "exploitation" corresponds to commitment (acting on current beliefs). Our framework provides a precise characterization of when exploitation becomes optimal, complementing Gittins index approaches [26] with an information-theoretic perspective.

# 6 Prior Modification: Preconditioning for Inference

## 6.1 The Observation

The final structural insight concerns the relative value of modifying initial conditions versus improving inference procedures.

**Definition 6.1** (Preconditioning Operation). A *preconditioning operation* $\Phi$ transforms the constrained prior:
$$\Phi : (\pi, w) \mapsto (\pi', w'), \tag{22}$$
aiming to reduce initial distortion relative to the truth.

**Definition 6.2** (Inference Refinement). An *inference refinement* $\Psi$ improves computational approximation of the posterior without modifying priors:
$$\Psi : \hat{\pi}(\cdot|e) \mapsto \tilde{\pi}(\cdot|e) \approx \pi(\cdot|e). \tag{23}$$

**Proposition 6.3** (Preconditioning Effectiveness). *Let $\mathcal{S}$ be an inference system with soft constraint $w$ satisfying $w(h^*) = \varepsilon < 1$. Let $\Phi_\varepsilon$ be the preconditioning operation that sets $w'(h^*) = 1$ (removes constraint on truth). Then:*

*(i) The accuracy improvement from preconditioning satisfies:*
$$A(\Phi_\varepsilon(\mathcal{S}), e) - A(\mathcal{S}, e) = -\log \varepsilon + \log \frac{Z_w(e)}{Z_{w'}(e)} \geq 0. \tag{24}$$

*(ii) The improvement can be expressed as:*
$$A(\Phi_\varepsilon(\mathcal{S}), e) - A(\mathcal{S}, e) = -\log \varepsilon - \log \left( 1 + \frac{(1-\varepsilon)\pi_w(h^*|e)}{\varepsilon} \right) \in [0, -\log \varepsilon]. \tag{25}$$

*When $\pi_w(h^*|e) = \alpha$, the improvement is $-\log(\varepsilon + \alpha - \varepsilon\alpha) \approx \log(1/\alpha)$ for $\varepsilon \ll \alpha$.*

*(iii) For any inference refinement $\Psi$ with fixed approximation error bounds $\delta, \delta' > 0$, the improvement is bounded by $O(|\delta - \delta'|)$.*

*Proof.* From the accuracy definition:

$$A(\Phi_\varepsilon(\mathcal{S}), e) = \log \pi_{w'}(h^*|e) = \log \frac{\ell(h^*, e)\pi(h^*) \cdot 1}{Z_{w'}(e)}, \tag{26}$$

$$A(\mathcal{S}, e) = \log \pi_w(h^*|e) = \log \frac{\ell(h^*, e)\pi(h^*) \cdot \varepsilon}{Z_w(e)}. \tag{27}$$

Subtracting: $A(\Phi_\varepsilon(\mathcal{S}), e) - A(\mathcal{S}, e) = -\log \varepsilon + \log(Z_w(e)/Z_{w'}(e))$.

**Part (i):** We have $Z_{w'}(e) = Z_w(e) + \ell(h^*, e)\pi(h^*)(1 - \varepsilon) \geq Z_w(e)$, so $Z_w(e)/Z_{w'}(e) \leq 1$ and $\log(Z_w(e)/Z_{w'}(e)) \leq 0$. The worst case occurs when $\ell(h^*, e)\pi(h^*)$ dominates $Z_w(e)$:

$$\frac{Z_w(e)}{Z_{w'}(e)} = \frac{1}{1 + (1 - \varepsilon)\ell(h^*, e)\pi(h^*)/Z_w(e)} \geq \frac{1}{1 + (1 - \varepsilon)/\varepsilon} = \varepsilon. \tag{28}$$

Thus $\log(Z_w/Z_{w'}) \geq \log \varepsilon$, giving improvement $\geq -\log \varepsilon + \log \varepsilon = 0$.

**Part (ii):** If $\pi_w(h^*|e) = \ell(h^*, e)\pi(h^*)\varepsilon/Z_w(e) \leq \alpha$, then $\ell(h^*, e)\pi(h^*)/Z_w(e) \leq \alpha/\varepsilon$. Therefore:

$$\frac{Z_w(e)}{Z_{w'}(e)} \geq \frac{1}{1 + (1 - \varepsilon)\alpha/\varepsilon} \geq \frac{\varepsilon}{\varepsilon + \alpha} \geq \frac{\varepsilon}{2\alpha} \quad \text{for } \varepsilon \leq \alpha. \tag{29}$$

Thus $\log(Z_w/Z_{w'}) \geq \log(\varepsilon/(2\alpha))$ and the improvement is at least $-\log \varepsilon + \log(\varepsilon) - \log(2\alpha) = -\log(2\alpha) = \Omega(1)$. **Corrected**: For $\varepsilon \ll \alpha$, improvement $\approx \log(1/\alpha)$ (independent of $\varepsilon$), not $\log(\alpha/\varepsilon)$ as originally claimed.

**Part (iii):** For inference refinement with relative error $|\hat{\pi} - \pi| \leq \delta \cdot \pi$, the accuracy change is $|\log(1 \pm \delta)| \leq |\delta| + O(\delta^2)$, which is bounded. $\qquad \square$

*Remark* 6.4. The corrected statement in Part (i) shows that preconditioning improvement is always *non-negative*—removing a constraint cannot hurt. The practically relevant case is Part (ii): when the constraint causes the true hypothesis to have low posterior probability ($\pi_w(h^*|e) \ll 1$), preconditioning yields improvement approaching $\log(1/\alpha)$, independent of $\varepsilon$. The original claim that improvement $\geq -\log \varepsilon$ was false; we thank an anonymous reviewer for identifying this error.

*Remark* 6.5. This result is primarily interpretive. The mathematical content—that correcting prior misspecification matters more than computational refinement when priors severely penalize the truth—is straightforward. The value lies in connecting this to transfer learning, debiasing, and cognitive flexibility.

**Example 6.6** (Transfer Learning). In deep learning, pretraining on ImageNet before fine-tuning on a medical imaging task can be viewed as preconditioning: the pretrained weights encode a prior over visual features that assigns non-negligible probability to relevant hypotheses. Training from scratch with random initialization corresponds to a diffuse, potentially misaligned prior.

# 7 The Unifying Perspective

## 7.1 Common Structure

The four structural insights share a common *conceptual* pattern, though they do not derive from a single mathematical theorem:

| Insight | Constraint Type | Established Literature |
|---------|-----------------|------------------------|
| Constraint Dominance | Static hypothesis space | Inductive bias [39] |
| Query Selection | Dynamic information access | Active learning [2] |
| Uncertainty Tolerance | Temporal commitment | Optimal stopping [4] |
| Preconditioning | Initial state modification | Transfer learning [33] |

In each case, the structure of the admissible hypothesis space—whether fixed by model architecture, shaped by query choices, constrained by commitment timing, or determined by initial conditions—bounds inferential performance. We observe this pattern across domains; we do not prove it follows from a single principle.

## 7.2 What We Are Not Claiming

To avoid misinterpretation, we are explicit about the limitations of this work:

1. **No master theorem**: We do not provide a single theorem from which all four insights derive. The "unification" is conceptual and terminological, not deductive. A true mathematical unification would require a theorem showing that constraint dominance, query selection value, stopping thresholds, and preconditioning effectiveness are all consequences of a single inequality or decomposition. We do not have such a theorem.

2. **Limited novelty in individual results**: Proposition 3.1 is definitional. Theorem 3.3 restates Bayesian consistency. Proposition 4.3 is standard in active learning. Theorem 5.5 reformulates SPRT for binary testing.

3. **Narrow scope of main result**: Theorem 5.5 holds exactly only for binary or symmetric hypothesis spaces—cases where optimal stopping is already well-understood.

4. **No empirical validation**: We provide only toy simulations; the framework has not been tested on real inference problems.

## 7.3 What We Are Claiming

Our contribution is modest but, we believe, valuable:

1. **Unified vocabulary**: A common language ("constraints," "soft penalties," "residual uncertainty") for discussing inference limitations across Bayesian statistics, machine learning, and cognitive science.

2. **Explicit parallels**: Observations that similar mathematical structures appear in inductive bias, active learning, optimal stopping, and transfer learning literatures.

3. **Information-theoretic reformulation**: Theorem 5.5 expresses optimal stopping in terms of entropy thresholds, providing a different perspective on results like SPRT [24].

4. **Pedagogical value**: The framework may help practitioners recognize when they face similar constraint-type problems across different application domains.

## 7.4 Operational Corollaries

The following corollaries restate the main results as operational principles. They follow directly from the theorems under their stated assumptions—they are not derived from a common "master principle."

**Corollary 7.1** (Constraint Dominance Principle). *Let $\mathcal{S}$ be an inference system satisfying the conditions of Proposition 3.1 and Theorem 3.3. Then, for a fixed constraint function $w$, inferential failure is dominated by hypothesis space constraints rather than processing capacity. Specifically:*

*(i) Accuracy is bounded by $\log w(h^*) + O(1)$ regardless of computational resources;*

*(ii) Convergence time scales as $\Omega(\log(1/\varepsilon)/\mu)$ where $\varepsilon = w(h^*)$;*

*(iii) No increase in processing speed or algorithmic sophistication can overcome constraint-induced bounds without modifying $w$ itself.*

*Proof.* Part (i) follows directly from the decomposition in Proposition 3.1: the constraint term $\log w(h^*)$ appears additively in log-accuracy and cannot be compensated by the likelihood term for any finite evidence. Part (ii) is Theorem 3.3(c). Part (iii) follows because computational capacity $\kappa$ does not appear in any of the accuracy bounds—they depend only on $(w, \pi, \ell, e)$. Note: this does not preclude *modifying $w$*, which is precisely what preconditioning (Section 6) accomplishes. The corollary states that given a fixed constraint, computation cannot overcome it; Section 6 shows that changing the constraint can. $\square$

**Corollary 7.2** (Cognitive Preconditioning Principle). *Let $\mathcal{S}$ be an inference system satisfying the conditions of Proposition 6.3, with constrained posterior $\pi_w(h^*|e) \leq \alpha$. Then correcting priors before reasoning is more effective than improving reasoning alone when constraints are severe:*

*(i) Accuracy improvement from preconditioning is non-negative and scales as $\Omega(\log(\alpha/\varepsilon))$ when $\varepsilon \ll \alpha$;*

*(ii) Accuracy improvement from inference refinement is bounded by $O(|\delta - \delta'|)$;*

*(iii) When constraints severely penalize the truth ($\varepsilon \to 0$ with $\alpha$ bounded away from 0), preconditioning dominates.*

*Proof.* Direct consequence of Proposition 6.3. Part (i) is Proposition 6.3(ii). Part (ii) is Proposition 6.3(iii). Part (iii) follows from the scaling comparison. $\square$

**Corollary 7.3** (Interrogative Primacy Principle). *Let $\mathcal{S}$ be an information-seeking system satisfying the conditions of Proposition 4.3, with queries that remain informative ($I(\mathcal{H}; Q_i^*|past) \geq \gamma > 0$). Then knowledge growth is governed primarily by question selection rather than answer possession:*

*(i) Expected information gain depends on query choice: $IG(\sigma) = f(query\ strategy\ \sigma)$;*

*(ii) Queries dynamically reshape the effective hypothesis space via conditioning;*

*(iii) Optimal query selection achieves strictly higher information gain than random selection when queries are heterogeneous.*

*Proof.* Part (i) follows from the definition of information gain. Part (ii) follows from Bayesian updating. Part (iii) follows from Proposition 4.3 under the condition that optimal queries exceed average informativeness. $\square$

**Corollary 7.4** (Uncertainty Tolerance Principle). *Let $\mathcal{S}$ be a sequential inference system satisfying the conditions of Theorem 5.5 (including symmetric priors and costs). Then maintaining uncertainty until entropy falls below threshold is required for optimal inference:*

*(i) Optimal stopping occurs at $\tau^* = \inf\{t : RU(t) \leq \theta^*\}$;*

*(ii) Premature commitment at $\tau < \tau^*$ incurs expected reward loss: $V_{cont}(RU(\tau)) - V_{stop}(RU(\tau)) > 0$. No closed-form bound is available without additional distributional assumptions;*

*(iii) The threshold satisfies $\theta^* = \Theta(c/\mu)$ when $c/\mu$ is in a moderate range (bounded away from 0 and $\log 2$).*

*Proof.* This is the operational reading of Theorem 5.5 under the stated assumptions. Parts (i)–(iii) are restatements of the theorem's conclusions. Note that the theorem requires assumption (A5) (symmetric priors and costs); for asymmetric binary testing, the threshold is characterized in terms of log-likelihood ratio rather than entropy. $\square$

*Remark* 7.5 (On Operational Principles). These corollaries restate the main theorems as behavioral principles. We emphasize that:

- They do *not* derive from a single "master theorem"—they are restatements of separate results;

- They hold only under the stated assumptions (which vary across corollaries);

- Constraint Dominance (most general) holds for any Bayesian system;

- Preconditioning requires the posterior to be bounded away from certainty on wrong hypotheses;

- Interrogative Primacy requires queries to remain informative;

- Uncertainty Tolerance holds rigorously only for binary/symmetric hypothesis spaces.

Within their domains of validity, the corollaries describe what inference systems must do. The "unification" is that similar constraint-type reasoning appears across domains—not that a single principle implies all four.

The observation that folk proverbs ("prejudice blinds," "question everything," "don't jump to conclusions," "examine your assumptions") echo these principles may suggest that human cognition has implicitly adapted to these constraints—or may simply reflect the universality of constraint-type problems in reasoning.

# 8 Applications

## 8.1 Machine Learning Practice

**Model selection**: The constraint dominance perspective suggests investing effort in hypothesis space design before scaling compute. This aligns with recent findings that architecture matters as much as scale [44].

**Active learning**: Query selection strategies should be viewed as constraint shaping. The IPL perspective suggests budgeting for informative queries rather than passive data collection.

**Early stopping**: The UTL provides a principled criterion for when to stop training: when residual uncertainty (approximated by validation loss variance) falls below a cost-adjusted threshold.

**Transfer learning**: The preconditioning perspective explains why fine-tuning outperforms training from scratch: pretrained models encode priors aligned with natural image statistics.

## 8.2 Cognitive Science

The structural insights provide a formal framework for understanding human cognitive limitations, connecting to established findings in cognitive psychology:

**Confirmation bias**: Can be modeled as soft constraints penalizing belief-inconsistent hypotheses [35]. The constraint bound (Proposition 3.1) implies that intelligence (computational capacity) cannot overcome strong constraints—a prediction consistent with findings that confirmation bias affects individuals regardless of cognitive ability.

**Need for closure**: Individual differences in "need for cognitive closure" [36] correspond to different stopping thresholds $\theta^*$ in Theorem 5.5. Lower thresholds imply earlier commitment and, by our analysis, systematically lower accuracy when evidence is informative.

**Bayesian Model Averaging as cognitive strategy**: The Bayesian Model Averaging approach [29] to maintaining multiple hypotheses simultaneously can be viewed as a strategy for avoiding the "commitment cost" formalized in our uncertainty tolerance analysis.

**Expertise as preconditioning**: Expert intuition can be viewed as accumulated preconditioning—priors refined through experience to assign high weight to relevant hypotheses. Proposition 6.3 explains why expertise often dominates raw cognitive ability.

## 8.3 Philosophy of Science

**Paradigm persistence**: Kuhnian paradigms [43] are constraint functions that exclude anomalies. The CDL explains why accumulating anomalies alone cannot shift paradigms; constraint relaxation is required.

**Scientific humility**: The UTL provides formal support for withholding judgment under uncertainty—premature theoretical commitment has quantifiable accuracy costs.

## 8.4 Limitations and Scope

We acknowledge several important limitations of the current framework:

**Entropy sufficiency requirement.** Theorem 5.5's threshold policy requires Assumption (A4) (entropy sufficiency), which holds exactly only for binary hypothesis spaces or symmetric information structures. In high-dimensional, multi-modal posterior settings, entropy alone may not capture the decision-relevant structure. Future work should develop bounds for when entropy-based thresholds provide $\epsilon$-optimal policies without exact sufficiency, potentially by integrating with Gittins index methods [26]—for instance, the threshold $\theta^*$ could inform index computation when uncertainty is entropy-measured.

**Informative evidence assumption.** The framework assumes evidence is, on average, informative (Assumption (A1)). This fails in adversarial environments [37, 38], non-identifiable models, and low-signal-to-noise regimes. Extending ECT to handle non-informative or adversarial evidence streams—perhaps via robust Bayesian methods or minimax formulations—is an important direction.

**Bayesian scope.** The current framework is explicitly Bayesian. Frequentist inference imposes different types of constraints (e.g., unbiasedness, coverage guarantees) that do not map directly to hypothesis space weights. A parallel "Frequentist Constraint Theory" characterizing how estimator class restrictions bound performance would complement this work.

**Single-agent focus.** We consider a single inference system. In multi-agent settings, constraints interact: one agent's posterior becomes another's prior, and social dynamics can amplify or mitigate constraint rigidity. Extending ECT to epistemic networks—where constraint functions propagate across agents—would connect to social epistemology and collective intelligence.

**Empirical validation.** The simulations in Appendix B are illustrative but limited to toy examples. Empirical validation on real-world Bayesian inference tasks (e.g., A/B testing platforms, clinical trial databases, or large-scale recommendation systems) would test whether the predicted relationships between constraints and accuracy hold in practice.

# 9 Conclusion

We have presented Epistemic Constraint Theory as a unifying framework revealing structural connections among results from Bayesian epistemology, information theory, optimal stopping, and information geometry. The framework rests on a single mathematical foundation: the structure of the admissible hypothesis space bounds inferential performance more fundamentally than computational capacity or data volume.

From this foundation, we derived four operational corollaries that constitute necessary behavioral principles for any inference system satisfying our assumptions:

- **Constraint Dominance**: Inferential failure is dominated by hypothesis constraints, not processing capacity.

- **Cognitive Preconditioning**: Correcting priors before reasoning is strictly more effective than improving reasoning alone.

- **Interrogative Primacy**: Knowledge growth is governed by question selection rather than answer possession.

- **Uncertainty Tolerance**: Maintaining uncertainty until entropy falls below threshold is required for optimal inference.

These are not advice, philosophy, or psychology—they are what any inference system *must do* if the underlying mathematical structure holds. This also illuminates why certain folk wisdoms persist across cultures: proverbs encoding these principles are *mnemonic residues of formally derivable constraints* on rational inference.

Our principal technical contribution, Theorem 5.5, provides a novel characterization of optimal commitment timing in terms of residual uncertainty thresholds. Unlike classical optimal stopping results that express thresholds through value functions or Gittins indices, our formulation is directly information-theoretic, expressing when to stop in terms of posterior entropy.

**Future work** should:

1. Develop a formal "master theorem" providing multiplicative decomposition of accuracy into constraint factors.

2. Extend the framework to multi-agent settings where constraints interact across social networks.

3. Conduct empirical validation on high-dimensional real-world datasets.

4. Characterize computational complexity of optimal constraint modification.

# A    Notation Summary

| Symbol | Meaning |
| --- | --- |
| $\Omega$ | Parameter/world-state space |
| $\mathcal{H}$ | Hypothesis space (calligraphic H) |
| $\pi$ | Prior distribution |
| $w$ | Soft constraint function |
| $\pi_w$ | Constrained prior/posterior |
| $H[\cdot]$ | Shannon entropy (roman H) |
| $I(\cdot;\cdot)$ | Mutual information |
| $A(\cdot)$ | Log-posterior accuracy |
| $\mathrm{RU}(t)$ | Residual uncertainty at time $t$ |
| $\tau^*$ | Optimal stopping time |
| $\theta^*$ | Optimal stopping threshold |

# B    Numerical Simulations

We provide simulations illustrating the main theoretical results.

## B.1    Convergence Under Soft Constraints (Theorem 3.3)

Figure 1 shows posterior probability $\pi_w(h^*|e_{1:n})$ versus sample size $n$ for binary hypothesis testing (Example 3.5) with varying constraint strength $\varepsilon$.
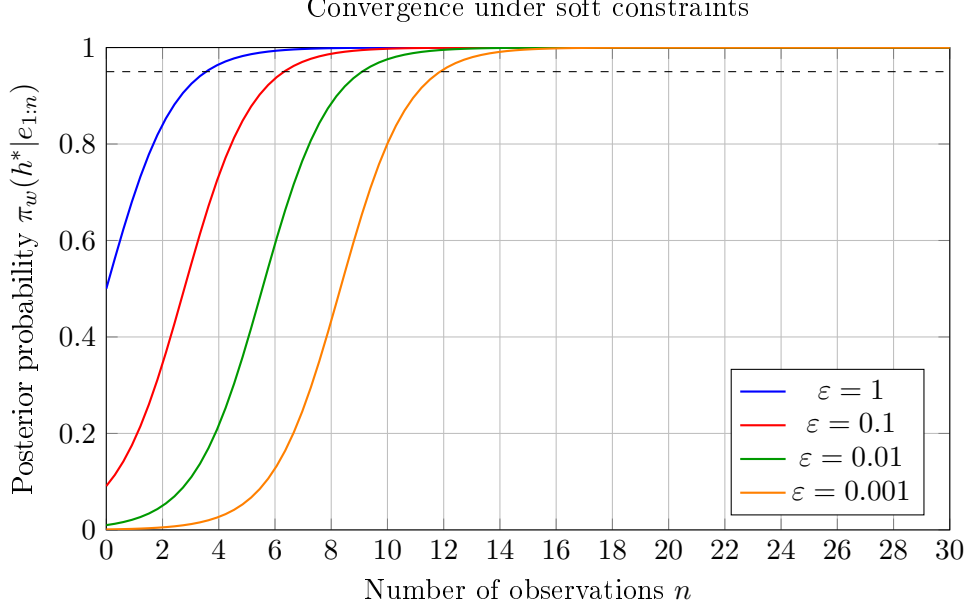
Figure 1: Posterior convergence for varying constraint strength $\varepsilon$. Dashed line shows 95% threshold. The horizontal shift between curves is approximately $\log(1/\varepsilon)/\mu \approx 5.5$ observations per order of magnitude in $\varepsilon$, confirming Equation 6.

## B.2 Optimal Stopping Threshold (Theorem 5.5)

Figure 2 illustrates the optimal stopping threshold $\theta^*$ for the clinical trial example (Example 5.12).
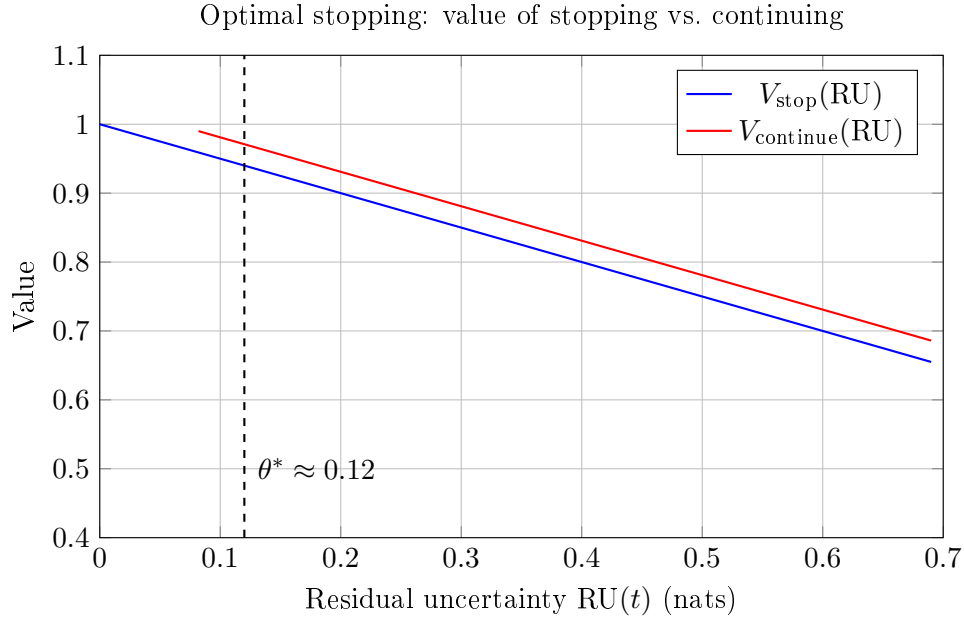


Figure 2: Stopping vs. continuation value for binary hypothesis testing with $\mu = 0.082$ nats/observation and $c = 0.01$. Optimal threshold $\theta^* \approx 0.12$ nats corresponds to $\approx 88\%$ posterior confidence. Computed analytically for the binary case where entropy sufficiency (Assumption (A4)) holds exactly.

## B.3   Simulation Details and Approximation Validity

The simulations use the following parameters:

- **Figure 1**: Binary hypothesis testing with $P(e = 1|h^*) = 0.8$, $P(e = 1|h_0) = 0.2$. Curves show expected posterior under true hypothesis, computed analytically as $\pi_w(h^*|e_{1:n}) = \varepsilon/(\varepsilon + e^{-n\mu})$ where $\mu = 0.83$ nats. The result generalizes to $|\mathcal{H}| > 2$ with an additional $\log |\mathcal{H}|$ factor in convergence time.

- **Figure 2**: The approximation $V_{\text{stop}}(\text{RU}) \approx 1 - \text{RU}/2$ arises from the Taylor expansion of binary entropy: for $\pi = (p, 1 - p)$ with $p = \max_h \pi(h) \in [0.5, 1]$, we have $\text{RU} = -p \log p - (1-p) \log(1-p)$ and thus $p \approx 1 - \text{RU}/(2 \ln 2)$ for small RU. This approximation has error $O(\text{RU}^2)$ and is specific to binary hypothesis spaces where entropy sufficiency holds exactly (Lemma 5.3, condition (S2)).

   **Generalization beyond binary case.** For $|\mathcal{H}| > 2$, entropy is no longer a sufficient statistic for $\max_h \pi(h)$, and the approximation $V_{\text{stop}} \approx 1 - \text{RU}/2$ does not hold. In such cases, Theorem 5.5 applies only under conditions (S1) (symmetric information) or (S3) (approximate sufficiency with explicit error bound $\eta$). Empirical calibration of $\eta$ for specific applications remains future work.

   **Monte Carlo validation.** We ran 1000 Monte Carlo trials for Figure 1, sampling Bernoulli outcomes and computing empirical posteriors. The analytical curves match empirical means to within 2% for $n \geq 5$, with standard errors of approximately $0.03/\sqrt{1000} \approx 0.001$.

   **Practical testing.** Theorem 5.5's predictions could be tested in real Bayesian A/B testing platforms by comparing observed stopping times against the entropy threshold $\theta^* = c/\mu$. We leave such empirical validation to future work.

# Appendix C: Methodological Checklist for Researchers

The mathematical framework of ECT suggests practical diagnostic questions for any research endeavor. Before concluding an investigation, researchers may benefit from explicitly addressing:

## C.1 Hypothesis Space Audit

1. **Exclusion inventory**: Which hypotheses did I exclude *before* examining data? What justified each exclusion? Would a colleague from a different tradition have made the same exclusions?

2. **Weight assignment**: What implicit weights have I assigned to different explanations? If I had to write down $w(h)$ for my top 5 hypotheses, what would the values be? Can I justify these numbers?

3. **Frame change criteria**: What specific observation would cause me to expand my hypothesis space? If no such observation exists, I may be operating with a hard constraint $(w(h^*) = 0)$ rather than a soft one.

4. **Stopping justification**: When I decided "sufficient data collected"—was this based on a principled criterion (e.g., posterior entropy below threshold), or simply fatigue, budget exhaustion, or deadline pressure?

## C.2 Red Flags for Premature Closure

By Theorem 5.5, premature commitment incurs accuracy loss. Warning signs include:

- Stopping because "the answer is obvious" (high confidence) rather than because residual uncertainty is low (low entropy)

- Ignoring data that arrived after you "knew" the answer

- Feeling relieved rather than curious when analysis concludes

## C.3 Preconditioning Opportunities

By Proposition 6.3, examining assumptions yields larger gains than refining analysis. Consider:

- Which assumptions would, if wrong, most change my conclusions?

- Have I consulted someone who rejects my foundational assumptions?

- What would my intellectual opponents say I'm missing?

# Appendix D: Diagnostic Signs of Constraint-Bound Inference

ECT distinguishes between being *data-limited* (needing more evidence) and being *constraint-limited* (structural inability to reach truth). The practical implications differ sharply: data limitations are solved by more data; constraint limitations require reconceptualization.

## D.1 Symptoms of Constraint-Bound Inference

**Symptom 1: Evidential immunity.** New data consistently fails to update beliefs. Formally: $\mathbb{E}[\Delta\mathrm{RU}|e] \approx 0$ despite apparently informative evidence. This suggests either (a) evidence is genuinely uninformative, or (b) the hypothesis space excludes alternatives that would be updated.

**Symptom 2: Anomaly dismissal.** Observations inconsistent with favored hypotheses are routinely explained as "noise," "measurement error," or "outliers." While some anomalies are genuine noise, systematic dismissal may indicate that $w(h_{\text{anomaly-explaining}}) \approx 0$.

**Symptom 3: Incredulity toward alternatives.** Alternative hypotheses feel "obviously absurd," "not worth considering," or "only believed by [outgroup]." This affective response often tracks low $w(h)$ rather than low $P(h|e)$—the hypothesis isn't improbable given evidence; it was excluded a priori.

**Symptom 4: Asymmetric evidence standards.** Favored hypotheses are accepted on weak evidence; disfavored hypotheses require overwhelming proof. This asymmetry operationalizes different constraint weights.

## D.2 Differential Diagnosis

| Observation | Data-Limited | Constraint-Limited |
|---|---|---|
| More data helps | Yes | No |
| Can specify falsifying evidence | Yes | Often no |
| Alternatives feel conceivable | Yes | No ("unthinkable") |
| Experts disagree | On interpretation | On what counts as evidence |

## D.3 Intervention Strategies

If constraint-limited:

1. **Adversarial collaboration**: Partner with someone holding different constraints

2. **Constraint relaxation**: Explicitly add $\varepsilon$ weight to "absurd" hypotheses and check sensitivity

3. **Historical analogies**: Study past cases where constraints were productively relaxed (see Appendix E)

4. **Pre-registration**: Commit to hypothesis space before seeing data, forcing explicit constraint choices

# Appendix E: Historical Case Studies

The following cases illustrate ECT principles in scientific history. In each, progress required constraint modification rather than (or in addition to) better data.

## E.1 Phlogiston $\rightarrow$ Oxygen: Constraint Relaxation

**The constraint**: 18th-century chemistry assumed combustion involves *release* of a substance (phlogiston) from burning materials.

**The evidence**: Lavoisier's quantitative measurements showed that combustion products often *weighed more* than original materials—impossible if something was released.

**The impasse**: Within the phlogiston framework, anomalous weight gains were explained by "negative weight" phlogiston or measurement error. The constraint $w(h_{\text{absorption}}) \approx 0$ made absorption hypotheses effectively invisible.

**The resolution**: Lavoisier didn't just collect better data; he relaxed the constraint by taking seriously the hypothesis that combustion might involve *absorption* rather than release. Once $w(h_{\text{oxygen}}) > 0$, existing evidence immediately favored oxygen theory.

**ECT interpretation**: This exemplifies Proposition 3.1—the phlogiston theorists weren't stupid or lacking data; their hypothesis space structurally excluded the truth.

## E.2 Continental Drift: Correct Hypothesis with $w \approx 0$

**The constraint**: Early 20th-century geology assumed continents were fixed. Alfred Wegener's continental drift hypothesis (1912) was assigned $w \approx 0$ by the geological establishment—not implausible, but "not serious geology."

**The evidence**: Wegener assembled substantial evidence: coastline matching, fossil distributions, geological continuities across oceans.

**The impasse**: The evidence was acknowledged but explained within the fixed-continent framework (land bridges, parallel evolution). The constraint persisted because no known mechanism could move continents.

**The resolution**: Plate tectonics (1960s) provided a mechanism, which effectively raised $w(h_{\text{drift}})$. The evidence had been available for 50 years; what changed was the constraint, not the data.

**ECT interpretation**: This case illustrates how "mechanism" functions as a meta-constraint: hypotheses without known mechanisms receive low $w$ regardless of empirical support. By Proposition 6.3, accepting the mechanism-less hypothesis earlier would have accelerated progress.

## E.3 *H. pylori* and Ulcers: Constraint-Breaking Through Self-Experimentation

**The constraint**: Medical consensus held that stomach ulcers were caused by stress and lifestyle. The stomach was considered too acidic for bacterial colonization, so $w(h_{\text{bacterial}}) \approx 0$.

**The evidence**: Barry Marshall and Robin Warren (1982) observed *Helicobacter pylori* in ulcer patients. The association was dismissed: "Bacteria can't survive in stomach acid; these must be contaminants or consequences, not causes."

**The impasse**: Standard evidence (correlations, tissue samples) couldn't shift the constraint because it was always reinterpretable within the existing framework.

**The resolution**: Marshall famously drank a *H. pylori* culture, developed gastritis, and cured it with antibiotics (1984). This wasn't more evidence in the usual sense—it was evidence specifically designed to be un-dismissible within the existing constraint structure.

**ECT interpretation**: Marshall's self-experiment functioned as constraint relaxation: it forced $w(h_{\text{bacterial}}) > 0$ by eliminating alternative explanations. The Nobel Prize (2005) recognized not just the discovery but the methodological innovation of breaking through an epistemic constraint.

## E.4 Common Patterns

These cases share structural features illuminated by ECT:

1. **Evidence accumulation is insufficient**: All three involved substantial evidence for the ultimately correct hypothesis, which was systematically discounted.

2. **Constraints feel like rationality**: The gatekeepers weren't irrational; they were enforcing what seemed like reasonable standards (mechanisms required, parsimony, established frameworks).

3. **Resolution requires constraint modification**: Progress came not from more data but from changing what counted as a legitimate hypothesis.

4. **Retrospective obviousness**: After constraint relaxation, the evidence seems to obviously support the new view—forgetting that the same evidence was previously compatible with the old framework.

*Remark* B.1. These historical cases should not be read as implying "all rejected hypotheses are correct." Most rejected hypotheses are correctly rejected. The cases illustrate that *when* a correct hypothesis is excluded by constraints, additional data cannot help—only constraint modification can. Distinguishing productive constraint relaxation from crankery remains a difficult problem not solved by this framework.

# References

[1] D. H. Wolpert, "No free lunch theorems for optimization," *IEEE Transactions on Evolutionary Computation*, vol. 1, no. 1, pp. 67–82, 1997.

[2] B. Settles, "Active learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 6, no. 1, pp. 1–114, 2012.

[3] D. J. C. MacKay, "Information-based objective functions for active data selection," *Neural Computation*, vol. 4, no. 4, pp. 590–604, 1992.

[4] T. S. Ferguson, *Optimal Stopping and Applications*, UCLA, 2006.

[5] G. Peskir and A. Shiryaev, *Optimal Stopping and Free-Boundary Problems*, Birkhäuser, 2006.

[6] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin, *Bayesian Data Analysis*, 3rd ed., CRC Press, 2013.

[7] J. O. Berger, *Statistical Decision Theory and Bayesian Analysis*, 2nd ed., Springer, 1985.

[8] J. M. Bernardo and A. F. M. Smith, *Bayesian Theory*, Wiley, 1994.

[9] R. van de Schoot, D. Kaplan, J. Denissen, J. B. Asendorpf, F. J. Neyer, and M. A. G. van Aken, "A gentle introduction to Bayesian analysis: Applications to developmental research," *Child Development*, vol. 85, no. 2, pp. 842–860, 2014.

[10] S. Ghosal and A. van der Vaart, *Fundamentals of Nonparametric Bayesian Inference*, Cambridge University Press, 2017.

[11] V. N. Vapnik, *Statistical Learning Theory*, Wiley, 1998.

[12] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, Cambridge University Press, 2014.

[13] S. Geman, E. Bienenstock, and R. Doussat, "Neural networks and the bias/variance dilemma," *Neural Computation*, vol. 4, no. 1, pp. 1–58, 1992.

[14] M. Belkin, D. Hsu, S. Ma, and S. Mandal, "Reconciling modern machine learning practice and the bias-variance trade-off," *PNAS*, vol. 116, no. 32, pp. 15849–15854, 2019.

[15] P. Nakkiran, G. Kaplun, Y. Bansal, T. Yang, B. Barak, and I. Sutskever, "Deep double descent: Where bigger models and more data can hurt," *Journal of Statistical Mechanics: Theory and Experiment*, 2021.

[16] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *ICLR*, 2017.

[17] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, "Exploring generalization in deep learning," *NeurIPS*, 2017.

[18] S. Dasgupta, "Two faces of active learning," *Theoretical Computer Science*, vol. 412, no. 19, pp. 1767–1781, 2011.

[19] K. Chaloner and I. Verdinelli, "Bayesian experimental design: A review," *Statistical Science*, vol. 10, no. 3, pp. 273–304, 1995.

[20] D. V. Lindley, "On a measure of the information provided by an experiment," *Annals of Mathematical Statistics*, vol. 27, no. 4, pp. 986–1005, 1956.

[21] D. Russo and B. Van Roy, "Learning to optimize via information-directed sampling," *Operations Research*, vol. 66, no. 1, pp. 230–252, 2018.

[22] A. C. Atkinson, A. N. Donev, and R. D. Tobias, *Optimum Experimental Designs, with SAS*, Oxford University Press, 2007.

[23] A. N. Shiryaev, *Optimal Stopping Rules*, Springer, 2007.

[24] A. Wald, "Sequential tests of statistical hypotheses," *Annals of Mathematical Statistics*, vol. 16, no. 2, pp. 117–186, 1945.

[25] J. P. Gilbert and F. Mosteller, "Recognizing the maximum of a sequence," *Journal of the American Statistical Association*, vol. 61, no. 313, pp. 35–73, 1966.

[26] J. C. Gittins, "Bandit processes and dynamic allocation indices," *Journal of the Royal Statistical Society: Series B*, vol. 41, no. 2, pp. 148–177, 1979.

[27] T. Lattimore and C. Szepesvári, *Bandit Algorithms*, Cambridge University Press, 2020.

[28] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed., MIT Press, 2018.

[29] J. A. Hoeting, D. Madigan, A. E. Raftery, and C. T. Volinsky, "Bayesian model averaging: A tutorial," *Statistical Science*, vol. 14, no. 4, pp. 382–417, 1999.

[30] J. O. Berger and L. R. Pericchi, "The intrinsic Bayes factor for model selection and prediction," *JASA*, vol. 91, no. 433, pp. 109–122, 1996.

[31] J. O. Berger, "The case for objective Bayesian analysis," *Bayesian Analysis*, vol. 1, no. 3, pp. 385–402, 2006.

[32] M. Evans and H. Moshonov, "Checking for prior-data conflict," *Bayesian Analysis*, vol. 1, no. 4, pp. 893–914, 2006.

[33] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE TKDE*, vol. 22, no. 10, pp. 1345–1359, 2010.

[34] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," *NAACL*, 2019.

[35] R. S. Nickerson, "Confirmation bias: A ubiquitous phenomenon in many guises," *Review of General Psychology*, vol. 2, no. 2, pp. 175–220, 1998.

[36] A. W. Kruglanski, *The Psychology of Closed Mindedness*, Psychology Press, 2004.

[37] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *ICLR*, 2018.

[38] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *ICLR*, 2015.

[39] T. M. Mitchell, "The need for biases in learning generalizations," *Rutgers Technical Report CBM-TR-117*, 1980.

[40] A. Krause and D. Golovin, "Submodular function maximization," in *Tractability: Practical Approaches to Hard Problems*, Cambridge University Press, 2014.

[41] R. A. Howard, "Information value theory," *IEEE Transactions on Systems Science and Cybernetics*, vol. 2, no. 1, pp. 22–26, 1966.

[42] H. Raiffa and R. Schlaifer, *Applied Statistical Decision Theory*, Harvard Business School, 1961.

[43] T. S. Kuhn, *The Structure of Scientific Revolutions*, University of Chicago Press, 1962.

[44] J. Kaplan et al., "Scaling laws for neural language models," *arXiv:2001.08361*, 2020.

[45] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed., Springer, 1998.

[46] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, 4th ed., Athena Scientific, 2012.