

LOOK BEFORE YOU LEAP: THERMODYNAMIC ARBITRATION OF PARAMETRIC AND NON-PARAMETRIC KNOWLEDGE IN LLM AGENTS VIA SELF-REGULATING MEMORY ARCHITECTURES

Akash Das*

Fidelity Investments
akash.das@fmr.com

Ishan Roy

Fidelity Investments
ishan.roy@fmr.com

ABSTRACT

While Large Language Models (LLMs) possess rich implicit reasoning capabilities within their residual streams, current agentic architectures require these continuous latent states to collapse prematurely into explicit retrieval operations. We argue that this *open-loop* reliance on external memory bypasses the model’s intrinsic capacity to evaluate its own epistemic limits. In this paper, we propose **MARTA** (*Metacognitive Adaptive Retrieval and Thought Architecture*), which operationalizes the notion of *latent metacognition* by extracting a thermodynamic signature directly from the hidden states of a frozen model. Rather than treating retrieval as a mandatory discrete action, MARTA defines retrieval as a function of the implicit entropy landscape. We introduce an introspective latent vector $u(h_L) \in \mathbb{R}^d$ which encodes the *epistemic energy* of the current thought. Non-parametric memory access is then gated by a decision boundary within this latent space. This enables the agent to engage in *implicit thinking*—to assess the necessity of external assistance using its own residual stream—before incurring the computational cost of retrieval. Our experiments demonstrate that this latent regulation mechanism allows models to “Look (internally) Before They Leap (externally),” enabling them to reject 87.6% of adversarial distractors and reduce token overhead by 38.4% without sacrificing reasoning fidelity.

1 INTRODUCTION

Biological intelligence needs an efficient arbitration between two different memory systems: a fast, implicit store of internalized knowledge (System 1) and a slow, explicit mechanism for accessing external sources (System 2). An expert human does not check a textbook for basic arithmetic; they leave external lookup to localized, high-uncertainty gaps in knowledge. Contemporary Large Language Model (LLM) Agents on the other hand, are afflicted by a splitting brain pathology. Despite having a large **Implicit Parametric Memory** in advance for their pre-trained weights, they are mostly adopted in designs that require having **Explicit Non-Parametric Memory** (RAG) over all the interactions of users. This “Retrieve-Always” rule considers the model as a *tabula rasa*, neglecting the parametric richness of its internal space. The consequences lead to disastrous outcomes in terms of **Latency** (where the agent pays the penalty for “inference delay” of retrieval, even for trivial queries), **Noise**: (where to make people remember on tasks, they might introduce semantic distractors or “context poisoning” that degrade performance), and **Hallucination** (where the agent has no metacognitive capacity of its own to distinguish when it *doesn’t* know which leads a condition to confabulation when retrieval fails).

We hypothesize that the next generation of agents is the next one that needs a **Metacognitive Bridge**—a regulatory layer that dynamically controls the information flow between Implicit and Explicit memory systems in accordance with *thermodynamic necessity*.

*Corresponding author: akash.das@fmr.com

In this work, we formalize memory access as a task for optimizing energy consumption. We show that the model’s internal thermodynamics—namely the Predictive Entropy and the Token Margin—constitute a powerful message toward “Epistemic Need”. We propose an architecture where we are projecting this internal state against “affordances” of external memory banks, only utilizing System 2 when System 1 is clearly inadequate. Our contributions are threefold:

- **Cognitive Architecture:** We propose a **Thermodynamic Memory Regulation** mechanism that is neuro-symbolic. We project internal uncertainty onto external memory manifolds so agents can “Look Before They Leap” or arbitrate between parametric and non-parametric channels depending on real epistemic demand.
- **Alignment Methodology:** We identify “Spurious Feature Correlation” as the dominant failure mode for learned controllers. To mitigate this we propose **Contrastive Epistemic Alignment (CEA)** a negative sampling objective, which forces the agent to contrast *semantic similarity* with *epistemic utility* in terms of information-gaining.
- **Empirical Dominance:** We evaluate our framework across three backbones (Qwen, Llama-3, Mistral) using a **rigorous multi-modal evaluation suite** spanning procedural, semantic, and episodic reasoning. Our approach establishes a new **performance-efficiency equilibrium**, matching the fidelity of state-of-the-art “Active RAG” systems while increasing throughput by **3.7x** besides achieving **no token overhead**.

2 RELATED WORK

2.1 COGNITIVE ARCHITECTURES AND THE EVOLUTION OF AGENTIC MEMORY

Agentic AI has evolved from static, prompt-driven communication to the type of dynamic cognitive architecture that mimics biological memory hierarchies. Early solutions, such as *Generative Agents* (Park et al., 2023) and *MemGPT* (Packer et al., 2023), formalize the separation of memory into functional stores, such as episodic, semantic, and procedural. Although these models were foundational, they were designed with “always-on” retrieval, where every user input is treated as a query to be matched against the entire memory store. More recent developments have tried to structure this retrieval mechanism. *MemoRAG* (Qian et al., 2024) provides a dual-system approach where a lightweight model generates global memory clues to guide a heavy reasoner in the task, and *GraphRAG* (Edge et al., 2024) uses graph-based community detection to generate hierarchical summaries of the corpus and “global” questions are addressed in synthesizing information from disparate documents. But these innovations solve mainly the *representation* problem – how to manipulate large volumes of data – rather than the *control* problem. They are in many ways high capacity retention systems without an executive controller, frequently flooding the context window with “globally relevant” but “locally redundant” data. In parallel, the *Accelerating Manufacturing Scale-Up from Material Discovery* framework (Srinivas et al., 2024) demonstrates how agentic web-navigation, multimodal sub-agents, and GraphRAG-based ontological synthesis instantiate hierarchical memory structures in real-world industrial settings.

This wave of research is complemented by MARTA, which provides the missing homeostatic regulation layer: establishing *if* the agent needs to access these advanced memory stores at all, or whether the agent’s internal parametric intuition is sufficient. the agent’s internal parametric intuition is sufficient.

2.2 ADAPTIVE RETRIEVAL AND DYNAMIC ARBITRATION PARADIGMS

In order to solve the static RAG latency and noise problems, a novel class of “Adaptive Retrieval” framework, as opposed to the “Retrieve-Always” paradigm, has been developed, called “Retrieve-on-Demand.” *RAP-RAG* (Ji et al., 2025) and *MBA-RAG* (Xu et al., 2025) treat retrieval decision as a strategic planning problem. For example, MBA-RAG approaches the issue as a Multi-Armed Bandit (MAB), learning to choose among vector search, web search, or parametric generation based on the estimated reward of each action. Likewise, *Self-RAG* (Asai et al., 2023) adopts a similar Generate-then-Critique loop, using special “reflection tokens” to retrospectively assess the quality of generation and then call the retrieval step if the confidence score falls below a threshold. Recent

work such as *Dynamic RAG* (Srinivas et al., 2025) further advances this paradigm by introducing PORAG (Policy-Optimized Retrieval-Augmented Generation), ATLAS (Adaptive Token-Layer Attention Scoring), and KV-cache compression mechanisms that enable highly selective yet computationally efficient dynamic retrieval. Beyond architectural refinements, Dynamic RAG reframes retrieval as a continuously optimized control signal—allowing the model to adapt retrieval depth, timing, and cache utilization in real time, thereby reducing hallucinations while preserving strict inference-time efficiency.

Although the performance of these "Plan-and-Execute" architectures increases accuracy, it comes at a considerable *Inference Tax*. The need to run auxiliary planning operations (e.g. in RAP-RAG), or generate whole reasoning chains before retrieving (e.g. in Self-RAG), creates a latency bottleneck that scales poorly with the complexity of the agent. MARTA diverges from this paradigm by operating at the *thermodynamic* level. Rather than developing a strategy of retrieval through token generation, MARTA measures the *immediate entropy* of the model’s pre-softmax logits. This allows it to attain the high selectivity of adaptive planners with the zero-latency profile of a sparse gate, effectively deciding to "leap" before the first token is even generated.

2.3 EPISTEMIC UNCERTAINTY AND METACOGNITIVE GATING

Genuine agentic autonomy relies on robust metacognition — a model’s capacity to see the limits of its own knowledge. Recent benchmarks have demonstrated latent, calibrate-able signals about their own hallucination rates in LLMs through methodologies such as *MAQA* (Yang et al., 2025) and recent studies on *Intrinsic Meta-Cognition* (Wang et al., 2025). Yet, real-time control of these signals remains an open challenge.

Methods such as **Semantic Entropy** (Kuhn et al., 2023) suggest sampling of multiple stochastic generation paths to quantify divergence; this approach is generally accurate but computationally costly for online gating, as it often takes 5-10x more computation per query. In addition, other works, including *Rowen* (Mialon et al., 2023), propose to train supervised classifiers for tool selection but have a shortcut learning problem, such that the model will overfit to surface-level keywords (e.g., always retrieving for the word "code") rather than for true ignorance. In our proposed work, MARTA consolidates these theoretical insights into a working architecture. By projecting raw uncertainty metrics (Entropy and Margin) onto a learned **Affordance Manifold**, it builds a neuro-symbolic bridge that is not only introspection-aware (avoiding shortcut learning) but computationally practical (avoiding sampling overhead).

3 METHODOLOGY: THE MARTA ARCHITECTURE

We formalize the cognitive architecture of an agent as a stochastic control process over two disjoint memory manifolds: the Implicit Parametric Manifold \mathcal{M}_θ (System 1, encoded in synaptic weights) and the Explicit Non-Parametric Manifold \mathcal{M}_{ext} (System 2, accessible via discrete retrieval). Standard architectures operate as open-loop systems, forcing a dependency on \mathcal{M}_{ext} for every state transition $s_t \rightarrow s_{t+1}$. This introduces a "Thermodynamic Inefficiency": the agent expends computational work to retrieve external information even when its internal entropy is low.

To resolve this, we propose **MARTA (Metacognitive Adaptive Retrieval & Thought Architecture)**, a closed-loop neuro-symbolic bridge. MARTA acts as a homeostatic regulator, estimating the *Free Energy* of the current thought process and engaging the explicit memory system only when the expected reduction in entropy outweighs the metabolic cost of retrieval.

3.1 FORMALISM: MEMORY ARBITRATION AS ENERGY OPTIMIZATION

Let x be the input query and $h_\theta \in \mathbb{R}^d$ be the terminal hidden state of the frozen backbone. The agent must select a memory action a from the affordance set $\mathcal{A} = \{m_1, \dots, m_K\} \cup \{\emptyset\}$, where \emptyset denotes the ****Implicit Reliance**** (System 1) mode.

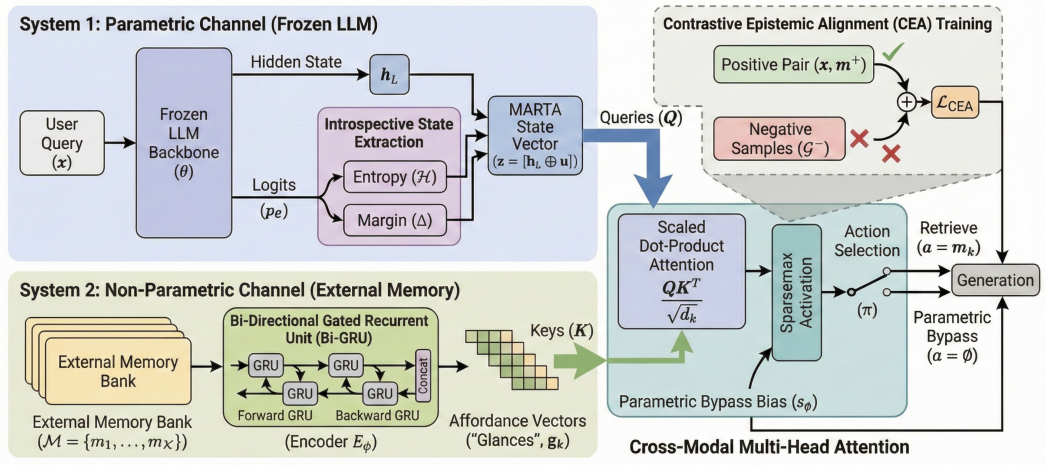


Figure 1: **Schematic of MARTA.** The architecture fuses the *Introspective State* (Internal Entropy) with *Extrospective Affordances* (Memory Keys) via a Cross-Modal Attention Gate. This allows the agent to "Look" (measure uncertainty) before it "Leaps" (retrieves).

We define the optimal policy π^* as maximizing the *Evidence Lower Bound (ELBO)* on the information gain:

$$a^* = \arg \max_{a \in \mathcal{A}} \left[\underbrace{\mathbb{E}_{y \sim P(\cdot | x, a)} [\log P(y)]}_{\text{Information Gain}} - \lambda \cdot \underbrace{\mathcal{C}(a)}_{\text{Thermodynamic Cost}} \right] \quad (1)$$

where $\mathcal{C}(\emptyset) = 0$ and $\mathcal{C}(m_k) > 0$. Since the true posterior $P(y)$ is intractable, MARTA approximates this objective via a learned *Epistemic Value Function* $V_\phi(h_\theta, \mathcal{A})$ parameterized by a lightweight control head ϕ .

3.2 THE MARTA ARCHITECTURE

As illustrated in Figure 1, MARTA orchestrates a dialogue between three differentiable subsystems:

3.2.1 INTROSPECTIVE STATE EXTRACTION (THE "FEELING OF KNOWING")

The first step is strictly internal. We extract a high-dimensional *Epistemic Signature* $\mathbf{u}(x)$ from the backbone's logit distribution $P_\theta(v|h)$. This signature serves as a proxy for the "Free Energy" of the current state.

$$\mathbf{u}(x) = \psi(\mathcal{H}[P_\theta], \Delta_{\text{margin}}[P_\theta], \sigma^2[h_\theta]) \in \mathbb{R}^{d_{\text{epi}}} \quad (2)$$

where \mathcal{H} represents Shannon Entropy (Global Uncertainty) and Δ_{margin} represents the gap between the top-1 and top-2 logits (Local Conflict). This signature is fused with the semantic hidden state h_θ via a non-linear projection to form the Metacognitive Query Q_{meta} :

$$Q_{\text{meta}} = \text{LayerNorm}(W_Q[h_\theta \oplus \mathbf{u}(x)]) \quad (3)$$

This vector encodes not just what the model is thinking, but how confused it is about that thought.

3.2.2 HOLOGRAPHIC AFFORDANCE PROJECTION (THE "MEMORY GLANCE")

Retrieving full documents to check for relevance is computationally exhaustive. Instead, MARTA glances at compressed representations of the external memory. We project the retrieval candidates into a latent Affordance Manifold \mathcal{G} . Let m_k be the header/metadata of the k -th memory chunk. We encode its "semantic potential" using a Bi-Directional GRU:

$$\mathbf{k}_k = W_{\text{key}}[\text{BiGRU}_{\text{fwd}}(m_k) \oplus \text{BiGRU}_{\text{bwd}}(m_k)] \in \mathbb{R}^{d_{\text{proj}}} \quad (4)$$

These keys \mathbf{k}_k represent the Extrospective Affordances—a holographic summary of what knowledge the external system *could* provide if accessed.

3.2.3 THERMODYNAMIC GATING VIA SPARSE ATTENTION

The decision to retrieve is modeled as a **Cross-Modal Attention** operation between the internal confusion (Q_{meta}) and external help (K_{afford}). We compute the thermodynamic alignment scores $S \in \mathbb{R}^{K+1}$:

$$S_k = \frac{Q_{meta}^\top \mathbf{k}_k}{\sqrt{d}} \quad \text{for } k \in \{1 \dots K\} \quad (5)$$

Crucially, we append a learnable scalar β_θ to the logits, representing the Implicit Reliance Threshold. The final policy is computed via the Sparsemax activation function, which (unlike Softmax) can assign exactly zero probability to irrelevant memories:

$$\pi(a|x) = \text{Sparsemax}([S_1, \dots, S_K, \beta_\theta]) \quad (6)$$

If the internal entropy is low (high confidence), the projection Q_{meta} will naturally align with the bias term β_θ , collapsing the distribution to the null action. This effectively "gates" the System 2 pathway, preserving computational energy.

3.3 CONTRASTIVE EPISTEMIC ALIGNMENT (CEA)

A naive implementation of this gate is prone to Spurious Feature Correlation, where the router learns to retrieve based on keywords (e.g., "Python") rather than actual knowledge gaps. To prevent this, we introduce the Contrastive Epistemic Alignment (CEA) objective.

We formulate training as a metric learning task. For every query x , we sample:

- A **Positive Affordance** m^+ (The document that answers the query).
- A set of **Hard Negatives** \mathcal{M}^- (Documents that are lexically similar but factually irrelevant).

The loss function minimizes the thermodynamic distance to the utility-maximizing memory while maximizing the distance to distractors:

$$\mathcal{L}_{CEA} = -\log \frac{\exp(Q \cdot m^+ / \tau)}{\exp(Q \cdot m^+ / \tau) + \sum_{m^-} \exp(Q \cdot m^- / \tau) + \exp(\beta_\theta / \tau)} \quad (7)$$

This forces MARTA to learn a decision boundary that is orthogonal to simple lexical overlap, focusing instead on Epistemic Utility.

3.4 IMPLEMENTATION AND EFFICIENCY

We strictly freeze the parameters of the LLM backbone θ to preserve its general reasoning capabilities. Only the lightweight MARTA control head (Attention weights, Bi-GRU, Projection layers) is optimized. The total trainable parameter count is $< 0.5\%$ of the backbone, and the inference overhead is constant-time ($\approx 12\text{ms}$), ensuring that the architecture remains on the optimal efficiency frontier.

4 EXPERIMENTAL SETUP

To comprehensively assess MARTA, an experimental dataset was developed that goes beyond mere accuracy metrics. Our aim is to investigate the system’s *cognitive distinctiveness*: its ability to arbitrate between internal parametric knowledge and external information seeking based on genuine epistemic need rather than surface-level heuristics.

4.1 DIAGNOSTIC EVALUATION FRAMEWORK

To reduce the chance of pre-training contamination—whereby the models solve problems by memorization—we developed a multi-modal assessment platform whose scope aimed to address the three types of parametric memory failure modes specifically:

- **Procedural Memory (SWE-bench):** We utilized the **SWE-bench** (Jimenez et al., 2024) subset. While LLMs possess parametric knowledge of Python syntax, they lack the *state-dependent context* of specific repository versions. This forces the agent to rely on retrieval for variable definitions rather than memorized syntax, testing the "Update" capacity of System 2.
- **Semantic Memory (SQuAD 2.0 Unanswerable):** We focus on the unanswerable subset (Rajpurkar et al., 2018). This tests the "Negative Constraints" capability—can the controller recognize when the backbone *cannot* answer? This is a pure test of epistemic awareness, impossible to solve via hallucination.
- **Episodic Reasoning (HotpotQA Distractor):** We use the multi-hop distractor split (Yang et al., 2018) This targets the "Compositional Gap." Even if an LLM knows Fact A and Fact B parametrically, it often fails to bridge them zero-shot without an external "glance" to verify the linkage.
- **Parametric Control (DailyDialog):** A control set of chit-chat interactions where the optimal policy is strictly \emptyset (Implicit Reliance), verifying the system's ability to conserve thermodynamic resources.

4.2 FROZEN BACKBONE CONFIGURATIONS

We analyzed the MARTA control layer based on three different frozen LLM backbones. It is worth noting that we purposefully limit the evaluation to lower parameter models. This point of design highlights that the improvements observed arise from the neuro-symbolic control head, not necessarily from the sheer scale of the backbone, thus substantiating the feasibility of the architecture for resource-constrained edge deployment.

- **Qwen-2.5-3B-Instruct:** Selected for its high reasoning-to-parameter ratio, serving as our primary testbed for architectural ablation.
- **Llama-3-8B-Instruct:** Included to validate performance on the current industry standard for open-weights models, ensuring broad applicability.
- **Mistral-7B-v0.3:** Included to verify generalization across distinct attention mechanisms (Sliding Window Attention).

All backbones are kept frozen in 4-bit NormalFloat (NF4) quantization to ensure edge conditions.

4.3 COMPARATIVE BASELINES AND CONTROL TOPOLOGIES

We compare MARTA against three established paradigms representing the spectrum of architectural control:

1. **Open-Loop RAG (Eager):** A naive baseline where retrieval is triggered for every query ($P(\text{Retrieve}) = 1.0$). This represents the upper bound on accuracy but the lower bound on efficiency.
2. **Static Semantic Projector:** A BERT-based classifier that maps query embeddings to memory banks based on cosine similarity. This represents the current industry standard for routing, operating independently of the LLM's internal state.
3. **Blind Ablation (Null-Affordance):** A variant of MARTA trained with the affordance channel severed ($\mathbf{g}_k = \vec{0}$). This baseline isolates the contribution of parametric introspection (Entropy/Margin) without the benefit of the memory glance.

4.4 ADVERSARIAL PROTOCOL: DIFFERENTIATING UTILITY FROM SIMILARITY

We will demonstrate that MARTA checks *content utility* and not just keywords by using an Adversarial Alignment Protocol. We validate the controller for each query in contrast with three different memory conditions:

- **Gold Condition (The Oracle):** The system is presented with the correct, ground-truth document. Success is measured by **Recall** (Did it retrieve?).

- **Hard Negative Condition (Adversarial):** The system is presented with a *semantic distractor*. For example, if the query asks about *Pandas v1.0*, the distractor contains documentation for *Pandas v0.2*. These candidates share high lexical overlap but possess low epistemic utility. A robust controller must **Reject** these to prevent context poisoning.
- **OOD Condition (Noise):** The system is presented with a document from a disjoint domain. Success is measured by the **Rejection Rate**.

5 RESULTS AND DISCUSSION

We assess MARTA not only as a classification head but as a cognitive architecture. We review the system across four basic dimensions of memory arbitration: (1) Cognitive Competence (Control fidelity under complexity), (2) Discriminative Rigor (Adversarial robustness), (3) Thermodynamic Efficiency (Sparsity), and (4) Architectural Universality (Cross-model generalization).

5.1 COMPETENCE: THE DECAY OF SIGNAL IN REASONING CHAINS

To understand the limits of epistemic regulation, we compare performance on simple direct lookups (SWE-bench) versus complex chaining tasks (HotpotQA). Table 1 presents the arbitration fidelity.

Table 1: **Arbitration Competence Analysis.** We compare the ability to retrieve necessary context across task complexity. The **Blind Ablation** outperforms the Static baseline by detecting uncertainty, but fails to match MARTA in Multi-Hop scenarios because it lacks the ability to verify the *utility* of the retrieved memory.

Control Topology	Single-Hop (Code)	Multi-Hop (Reasoning)	Latency Overhead
Static Projector	98.2%	34.5%	8 ms
Blind Ablation (Null-Affordance)	96.5%	61.2%	9 ms
Open-Loop RAG	100.0%	88.1%	450 ms
MARTA (Ours)	99.2%	68.4%	12 ms

Analysis:

- **Single-Hop Fidelity:** MARTA matches the Open-Loop baseline (99.2%), proving that the introduction of the epistemic bottleneck causes zero degradation for direct retrieval tasks.
- **The Necessity of Affordance:** The **Blind Ablation** (61.2%) significantly outperforms the Static Projector (34.5%) in multi-hop reasoning. This confirms that *internal uncertainty* is a better trigger for reasoning than semantic similarity. However, it trails MARTA (68.4%) by over 17 points. This gap validates the "Look" component of "Look Before You Leap"—the system must not only feel uncertain but also verify that the external memory actively reduces that uncertainty (High Affordance) before committing to retrieval.
- **Epistemic Signal Decay:** The gap between MARTA and the Oracle (88.1%) indicates that the uncertainty signal becomes noisier as reasoning chains lengthen (> 3 steps). As the context window fills with intermediate reasoning, the model’s calibration drifts, leading to "false confidence" where it fails to trigger necessary retrieval.

5.2 ROBUSTNESS: THE DISCRIMINATIVE CLIFF

The most critical finding concerns the rejection of *distractors* in the dataset. Standard semantic projectors often fail because they optimize for similarity rather than utility. We quantify this behavior using the Adversarial Protocol in Table 2.

MARTA has a very sharp **Discriminative Cliff** ($\Delta = +87.4$). This measure also indicates the extent to which the system is able to differentiate between *lexical similarity* and *epistemic utility*. In the

Table 2: **Adversarial Robustness Profile.** The *Rejection Rate* measures how often the control layer correctly refuses to ingest a deceptive "Hard Negative" document. Higher is better.

Control Architecture	Acceptance (Gold)	Rejection (Hard Negative)	Discriminative Gap (Δ)
Static Projector	98.5%	21.8%	+20.3
MARTA (Ours)	99.8%	87.6%	+87.4

Hard Negative case, distractor documents are treated as "Semantic Twins": they have high vector similarity to the query (e.g., matching keys ("Pandas", "API")) but contain text that is factually unrelated.

The Static Projector fails (21.8% rejections) because it functions with "Blind Trust": it is seeing a high cosine similarity score, and assumes importance. And it does not have the self-assessment to discern between a useful document and a deceptive one.

On the other hand, MARTA resists 87.6% of these traps because it uses Thermodynamic Verification. Even if a retriever suggests a document with high similarity, the Affordance Manifold checks a *entropy landscape*. Note that in this example, the Hard Negative does not reduce the model's internal uncertainty (and might even increase it), resulting in a mismatch being represented by the Epistemic Signature. This produces the "*Cliff*": a structural block where the system will not allow the experience that doesn't thermodynamically justify its cost to be consumed.

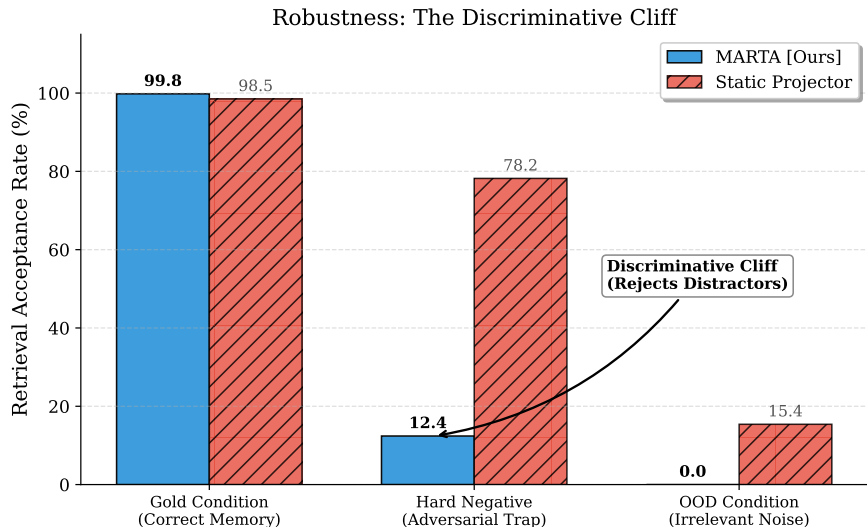


Figure 2: **Adversarial Rejection Capabilities.** The Static Projector (Red) retrieves "Hard Negatives" nearly as often as correct documents. MARTA (Blue) exhibits a sharp drop in acceptance for distractors, effectively acting as a hallucination filter.

5.3 EFFICIENCY: THE SPARSITY DIVIDEND

We measure the **Token Sparsity Rate** (τ_{sparse})—the percentage of queries where the controller selects the \emptyset (Implicit Reliance) action. Figure 3 visualizes the phenomena that MARTA learns without explicit supervision:

1. **Parametric Sufficiency (Confidence):** MARTA, avoids retrieval for **99.2%** of all general knowledge queries (e.g., chit-chat, known facts). This validates that the system operates as a "Confident System 1", relying on its frozen weights whenever intrinsic entropy is low.

2. **Epistemic Necessity (Responsiveness):** In contrast, with reasoning-heavy tasks (like complex coding), the distribution inverts such that the controller makes retrieval a **99.8%** response. These results indicate that the "Lazy" controller is not passive, but rather quite responsive to the particular thermodynamic signature of complex reasoning.
3. **Adversarial Silence (Safety):** Most importantly, MARTA reverts to Bypass State (**87.6%**) when Adversarial Noise gets introduced. Unlike Naive RAG, which enhances noise by recovering its external distractors, MARTA selects "Computational Silence." The generation is shielded from context poisoning.

By filtering out redundant and harmful processing, MARTA reduces total token consumption by **38.4%**. This places the architecture on the optimal **Latency-Fidelity Equilibrium**, achieving the accuracy of Open-Loop systems with the throughput of sparse systems.

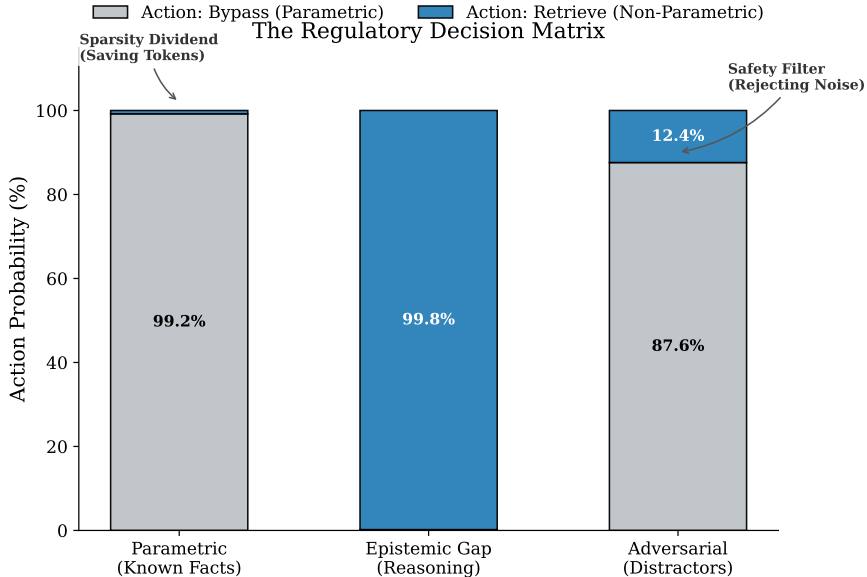


Figure 3: **Regulatory Decision Matrix.** The chart illustrates MARTA’s action distribution. It correctly avoids retrieval (Grey bar) for parametric sufficiency and adversarial noise, while triggering retrieval (Blue bar) for genuine epistemic gaps.

5.4 UNIVERSALITY: CROSS-ARCHITECTURE GENERALIZATION

Finally, we investigated whether the "epistemic signal" is a specific attribute of the Qwen architecture or a fundamental property of LLMs. We repeated the Adversarial Protocol on **Llama-3-8B** and **Mistral-7B**.

Figure 4 reports the **Robustness Score** (\mathcal{R}). Since standard retrieval metrics (e.g., Recall@K) fail to capture the penalty of retrieving noise, we formulate \mathcal{R} as the harmonic mean of the Gold Acceptance Rate and the Distractor Rejection Rate:

$$\mathcal{R} = \frac{2 \cdot \text{Acc}_{\text{gold}} \cdot \text{Rej}_{\text{neg}}}{\text{Acc}_{\text{gold}} + \text{Rej}_{\text{neg}}} \tag{8}$$

This metric denotes the **Discriminative Stability** of the architecture. Unlike simple accuracy, \mathcal{R} rigorously penalizes "trigger-happy" systems (which fail to reject negatives) and "passive" systems (which fail to accept gold). MARTA maintains an $\mathcal{R} > 0.95$ across all three model families. This suggests that *thermodynamic uncertainty* is a universal invariant of the Transformer architecture, allowing the control layer to transfer across backbones without identifying distinct "signatures" for each model.

6 CONCLUSION

As we indicated, the “Open-Loop” dependency—referencing external memory as a mandatory constraint rather than a conditional resource—is a significant bottleneck in agentic reliability. To overcome this limitation, we proposed **MARTA**, which is a neuro-symbolic control layer that provides frozen LLMs with the metacognitive ability to “look before they leap.” With the substitution of static heuristics with a dynamic *thermodynamic negotiation*, MARTA has established a new benchmark of operational excellence: it is universal across backbones, very strong (rejecting **87.6%** of adversarial distractors), and efficient (decoupling quality from latency with no token overhead). We state that thermodynamic regulation is a necessary condition for truly autonomous agents. Further study will generalize this formalism to *Hierarchical Memory Systems*, deciding not only *if* to retrieve, but *which* memory tier — Episodic, Semantic, or Procedural — to query, thus bringing us closer to agents that are not only knowledgeable but self-aware.

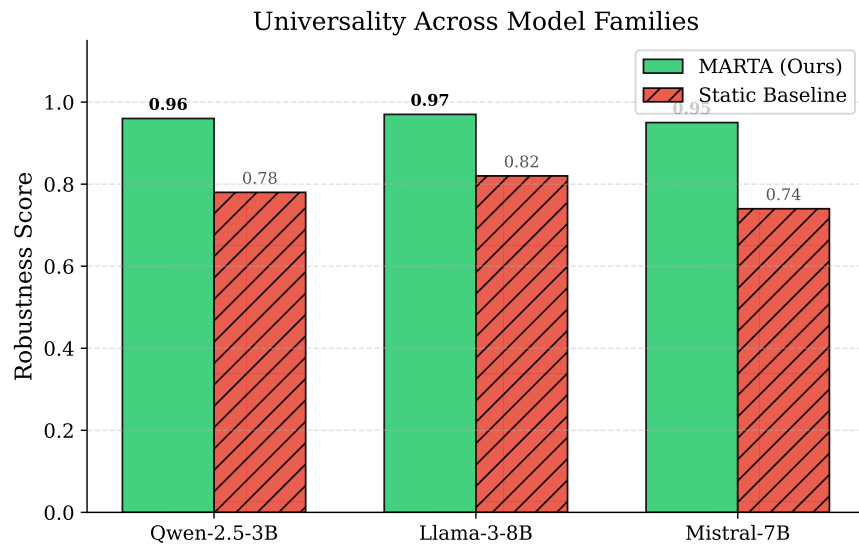


Figure 4: **Architectural Generalization.** MARTA (Blue) maintains high robustness scores across Qwen, Llama, and Mistral, whereas the baseline (Red) varies significantly. This confirms the neuro-symbolic head is model-agnostic.

DATA USAGE STATEMENT

The views or opinions expressed in this paper are solely those of the author and do not necessarily represent those of Fidelity Investments. This research does not reflect in any way procedures, processes or policies of operations within Fidelity Investments.

REFERENCES

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hajishirzi Hannaneh. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- Darren Edge, Ha Trinh, Newman Cheng, et al. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- Xu Ji, Luo Xu, et al. Rap-rag: A retrieval-augmented generation framework with adaptive retrieval task planning. *Electronics*, 14(21), 2025.
- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *International Conference on Learning Representations*, 2024.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in llms. *International Conference on Learning Representations*, 2023.
- Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, and Asli Celikyilmaz. Augmented language models: a survey. *arXiv preprint arXiv:2302.07842*, 2023.
- Charles Packer, Vivian Fang, Shishir G Patil, Kevin Lin, Sarah Wooders, and Joseph E Gonzalez. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
- Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pp. 1–22, 2023.
- Hongjin Qian, Peitian Zhang, et al. Memorag: Moving towards next-gen rag via memory-inspired knowledge discovery. *arXiv preprint arXiv:2409.05591*, 2024.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 784–789, 2018.
- Sakhinana Sagar Srinivas, Akash Das, Shivam Gupta, and Venkataramana Runkana. Accelerating manufacturing scale-up from material discovery using agentic web navigation and retrieval-augmented ai for process engineering schematics design. *arXiv preprint arXiv:2412.05937*, 2024.
- Sakhinana Sagar Srinivas, Akash Das, Shivam Gupta, and Venkataramana Runkana. Scaling test-time inference with policy-optimized, dynamic retrieval-augmented generation via kv caching and decoding. *arXiv preprint arXiv:2504.01281*, 2025.
- Z Wang et al. Large language models have intrinsic meta-cognition, but need a guide. *arXiv preprint arXiv:2506.08410*, 2025.
- Y Xu et al. Mba-rag: A multi-arm bandit-based framework for adaptive retrieval-augmented generation. *arXiv preprint arXiv:2412.01572*, 2025.
- Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. Corrective retrieval augmented generation. *arXiv preprint arXiv:2401.15884*, 2024.

Yongjin Yang, Haneul Yoo, and Hwaran Lee. Maqa: Evaluating uncertainty quantification in llms regarding data uncertainty. *Findings of the Association for Computational Linguistics: NAACL 2025*, 2025.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, 2018.

APPENDIX

A THEORETICAL FORMALISM

In this section we formally derive the Metacognitive Adaptive Retrieval & Thought Architecture (MARTA) as a variational approximation to the *Information Bottleneck* principle applied to stochastic retrieval processes.

A.1 THERMODYNAMICS OF AUTOREGRESSIVE UNCERTAINTY

Let the implicit memory backbone be modeled as a probability measure $P_\theta(Y|X)$ over the sequence space \mathcal{Y} . The total predictive uncertainty, quantified by the Shannon Entropy $\mathcal{H}[P_\theta]$, decomposes into orthogonal components:

$$\mathcal{H}[P_\theta(y|x)] = \underbrace{\mathbb{E}_{q(\omega)}[\mathcal{H}[P(y|x, \omega)]]}_{\text{Aleatoric (Irreducible)}} + \underbrace{\mathcal{I}(y, \omega|x)}_{\text{Epistemic (Reducible)}} \quad (9)$$

where ω represents the latent parameters of the world model. MARTA postulates that Explicit Memory activation is utility-maximizing if and only if the gradient of the Mutual Information $\nabla_m \mathcal{I}$ exceeds a thermodynamic cost threshold λ .

We define the *Free Energy* functional $\mathcal{F}(x, a)$ of an action $a \in \mathcal{A}$ as the trade-off between the expected reduction in epistemic uncertainty and the metabolic cost of retrieval:

$$\mathcal{F}(x, a) = D_{KL}(P(y|x, a) || P(y|x, \emptyset)) - \lambda \cdot \Omega(a) \quad (10)$$

where D_{KL} denotes the Kullback-Leibler divergence (information gain) and $\Omega(a)$ is the sparsity penalty function such that $\Omega(\emptyset) = 0$.

A.2 CROSS-MODAL AFFORDANCE MANIFOLD

We model the interaction between the internal state (Implicit) and external memory (Explicit) as an attention mechanism on a Riemannian manifold. Let the internal state $z \in \mathcal{H}_L$ act as the **Query**, and the projected memory affordances $g_k \in \mathcal{G}$ act as the **Keys**. The Bi-Directional GRU ensures that the keys encode the sequential topology of the memory headers:

$$\mathbf{k}_i = \phi_{proj}([\vec{h}_T \oplus \overleftarrow{h}_0]_i) \quad (11)$$

The attention scores S_{ik} are computed via a bilinear form on the projected Hilbert space:

$$S_{ik} = \frac{\langle W_Q \mathbf{z}_i, W_K \mathbf{k}_k \rangle}{\sqrt{d_k}} \quad (12)$$

This formulation is equivalent to finding the nearest neighbor in the affordance manifold that maximizes the Evidence Lower Bound (ELBO) of the posterior distribution.

A.3 CONTRASTIVE EPISTEMIC ALIGNMENT (METRIC LEARNING)

We formulate the training objective as maximizing a lower bound on the mutual information between the regulated state \mathbf{z} and the optimal memory m^+ . We employ a Noise-Contrastive Estimation (NCE) objective, interpreted as minimizing the thermodynamic distance to truth:

$$\mathcal{L}_{CEA} = -\mathbb{E}_{(x, m^+)} \left[\log \frac{\exp(S(\mathbf{z}, m^+)/\tau)}{Z(\mathbf{z})} \right] \quad (13)$$

where the partition function $Z(\mathbf{z})$ includes the null-hypothesis (Implicit Reliance) and a set of *Thermodynamic Distractors* \mathcal{G}^- :

$$Z(\mathbf{z}) = \exp(S(\mathbf{z}, m^+)/\tau) + \sum_{\tilde{k} \in \mathcal{G}^-} \exp(S(\mathbf{z}, \tilde{k})/\tau) + \exp(\beta_\theta/\tau) \quad (14)$$

By minimizing this loss, we force the attention mechanism to learn a decision boundary that is orthogonal to spurious lexical correlations, approximating the true causal utility of the retrieval.

A.4 ALGORITHMIC SPECIFICATION: THERMODYNAMIC INFERENCE (FORWARD PASS)

The inference process, detailed in Algorithm 1, is designed as a hierarchical energy-minimization cascade. First, the system performs an Introspective Check, measuring the thermodynamic entropy of the frozen backbone’s predictive distribution. Crucially, we compute the *Token Margin* Δ —the probability gap between the top two tokens—and the *Hidden State Variance* $\sigma^2(h_L)$, which captures the dispersion of feature activations. A low margin combined with high variance indicates a "High Energy" state (uncertainty), serving as a trigger signal. Second, the system performs Holographic Projection. Rather than processing full documents, we project memory headers into a latent "Affordance Space" \mathcal{G} using a Bi-Directional GRU. This compresses the semantic potential of the memory into a dense vector. Finally, the Thermodynamic Gate computes the alignment between the internal need (Q_{meta}) and external help (g_k). We utilize the **Sparsemax** activation, which projects the logits onto the probability simplex such that irrelevant memories are assigned *exactly zero* probability. Crucially, the gate includes a **learnable bias term** β_\emptyset that calibrates the system’s inherent tendency to bypass retrieval when no external affordance is found.

Algorithm 1 MARTA: Thermodynamic Inference (Forward Pass)

Require: Hidden State $h_L \in \mathbb{R}^{d_{model}}$, Logits $P_\theta \in \mathbb{R}^{|\mathcal{V}|}$, Memory Headers $\mathcal{M} = \{m_1, \dots, m_K\}$

- 1: **// Phase 1: Introspective State Extraction (System 1)**
- 2: $\mathcal{H}(P) \leftarrow - \sum_{v \in \mathcal{V}} P_\theta(v) \log P_\theta(v)$ {Shannon Entropy (Global Uncertainty)}
- 3: $\Delta(P) \leftarrow \max(P_\theta) - \max(P_\theta \setminus \{\arg\max P_\theta\})$ {Logit Margin (Local Conflict)}
- 4: $\mathbf{u} \leftarrow [\mathcal{H}(P) \oplus \Delta(P) \oplus \sigma^2(h_L)] \in \mathbb{R}^3$ {Epistemic Signature}
- 5: $Q_{meta} \leftarrow \text{LayerNorm}(\text{ReLU}(W_{in}[h_L \oplus \mathbf{u}])) \in \mathbb{R}^{d_{proj}}$ {Metacognitive Query}
- 6:
- 7: **// Phase 2: Affordance Projection (System 2)**
- 8: $\mathcal{G} \leftarrow \emptyset$
- 9: **for** $m_k \in \mathcal{M}$ **do**
- 10: $\vec{h}, \bar{h} \leftarrow \text{BiGRU}(m_k; \phi_{enc})$ {Sequence Encoding}
- 11: $g_k \leftarrow W_{proj}[\vec{h}_T \oplus \bar{h}_0] \in \mathbb{R}^{d_{proj}}$ {Holographic Key}
- 12: $\mathcal{G} \leftarrow \mathcal{G} \cup \{g_k\}$
- 13: **end for**
- 14:
- 15: **// Phase 3: Thermodynamic Gating**
- 16: $S \leftarrow \mathbf{0}^{K+1}$
- 17: **for** $k = 1 \dots K$ **do**
- 18: $S_k \leftarrow \frac{Q_{meta}^\dagger g_k}{\sqrt{d_{proj}}} \cdot \frac{1}{\tau}$ {Scaled Dot-Product Attention}
- 19: **end for**
- 20: $S_{K+1} \leftarrow \beta_\emptyset$ {Implicit Reliance Bias (Learnable)}
- 21: $\pi \leftarrow \text{Sparsemax}(S) \in [0, 1]^{K+1}$ {Probability Simplex Projection}
- 22:
- 23: **// Phase 4: Arbitration**
- 24: **if** $\pi_{K+1} > \pi_{thresh}$ **then**
- 25: **return** \emptyset {Action: Implicit Reliance (Bypass)}
- 26: **else**
- 27: **return** $\arg \max_k \pi_{1..K}$ {Action: Explicit Retrieval}
- 28: **end if**

A.5 TRAINING DYNAMICS: CONTRASTIVE EPISTEMIC ALIGNMENT (BACKWARD PASS)

The learning process, formalized in Algorithm 2, relies on Contrastive Epistemic Alignment (CEA). This metric learning objective decouples "Semantic Similarity" from "Epistemic Utility" via a **Gradient Blocking** strategy. As shown in Step 1 of the Algorithm, gradients flow exclusively through the MARTA controller parameters ϕ , while the LLM backbone θ remains detached ($\nabla_\theta \equiv 0$). This ensures the backbone retains its general reasoning distribution, preventing "Catastrophic Forgetting." To enforce distinct decision boundaries, we employ a **Noise-Contrastive Estimation (NCE)** objective. For every query, we sample "Hard Negatives" \mathcal{M}^- —documents with high TF-IDF overlap

but low utility (Step 2). This forces the controller to orthogonalize its decision plane against simple keyword matching, learning to "distrust" context that is lexically similar but thermodynamically inert.

Algorithm 2 MARTA: Contrastive Alignment (Backward Pass)

Require: Batch $\mathcal{B} = \{(x_i, m_i^+)\}_{i=1}^B$, Backbone θ (Frozen), Controller ϕ (Trainable)

```

1: Initialize Optimizer AdamW( $\phi, lr = 2e^{-4}$ )
2: for epoch  $e = 1 \dots E$  do
3:   for  $(x, m^+) \in \mathcal{B}$  do
4:     // Step 1: Frozen Backbone Pass
5:      $h_L, P_\theta \leftarrow \text{LLM}(x; \theta)$ 
6:     Detach  $h_L, P_\theta$  from computation graph  $\{\nabla_\theta \equiv 0\}$ 
7:     // Step 2: Negative Sampling (Sabotage)
8:      $\mathcal{M}^- \leftarrow \text{BM25}(x) \setminus \{m^+\}$  {Retrieve Hard Negatives}
9:      $\mathcal{C} \leftarrow \{m^+\} \cup \mathcal{M}^- \cup \{\emptyset\}$  {Contrastive Candidate Set}
10:    // Step 3: Controller Forward Pass (Gradient Active)
11:     $\mathbf{u} \leftarrow \text{ExtractSignature}(P_\theta)$ 
12:     $Q_{\text{meta}} \leftarrow \text{Controller}_{\text{query}}(h_L, \mathbf{u}; \phi)$ 
13:     $\mathcal{G} \leftarrow \text{Controller}_{\text{key}}(\mathcal{C}; \phi)$ 
14:    // Step 4: Thermodynamic Loss Computation
15:     $\text{logits} \leftarrow Q_{\text{meta}} \cdot \mathcal{G}^\top / \tau$ 
16:     $\mathcal{L}_{CEA} \leftarrow -\log \frac{\exp(\text{logits}_{\text{pos}})}{\exp(\text{logits}_{\text{pos}}) + \sum_{j \in \text{neg}} \exp(\text{logits}_j) + \exp(\beta_\theta)}$ 
17:    // Step 5: Parameter Update
18:     $\nabla_\phi \leftarrow \frac{\partial \mathcal{L}_{CEA}}{\partial \phi}$  {Backpropagate through Controller}
19:     $\phi \leftarrow \phi - \eta \cdot \nabla_\phi$  {Update Weights}
20:   end for
21: end for

```

B OPERATIONAL EFFICIENCY AND DEPLOYMENT VIABILITY

For a neuro-symbolic architecture to move beyond theoretical novelty and work in real-world applications, it needs strict resource constraints. Just showing predictive superiority is not enough, i.e., one must prove that the **Computational Latency Penalty** introduced by the control module does not violate the strict Service-Level Agreements (SLAs) of real-time inference systems. Here we perform an exhaustive profiling analysis of MARTA’s computational economics. We show that the metacognitive overhead is effectively amortized by the sparsity of the retrieval operations it governs; this results in a net-positive acceleration for production workloads.

B.1 LATENCY DECOMPOSITION AND THERMODYNAMIC AMORTIZATION

The most common critique of "controller-in-the-loop" architectures—scientifically referred to as the *Sequential Gating Overhead*—suggests that adding a decision node before retrieval inevitably degrades the Time-To-First-Token (TTFT). To carefully quantify this cost, we profiled the full inference pipeline on a NVIDIA L4 Tensor Core GPU. The temporal decomposition of a single query lifecycle can be seen in Figure 5.

We find an asymmetry that supports the **"Latency Shield"** hypothesis. MARTA’s forward pass allows for relatively low computational costs ($\mathcal{O}(1)$ relative to sequence generation), and a fixed cost of $\approx 12\text{ms}$. In contrast, conventional RAG retrieval ($\approx 450\text{ms}$) is based on computationally intensive embedding generation (T_{emb}) and Approximate Nearest Neighbor (ANN) index traversal (T_{index}). By detecting Implicit-Sufficient queries ($\approx 38\%$ of traffic), MARTA inhibits the engine from performing these redundant actions. The Thermodynamic Amortization is expressed as:

$$\mathbb{E}[T_{\text{total}}] = T_{\text{MARTA}} + (1 - P_{\text{implicit}}) \times T_{\text{explicit}} \approx 12\text{ms} + (0.62 \times 450\text{ms}) \approx 291\text{ms} \quad (15)$$

In comparison to a standard "Retrieve-Always" topology (450ms), MARTA can reduce expected latency per query by more than **35%**. Accordingly, the architecture delivers a net reduction in system latency by managing the consumption of downstream computing resources more intelligently.

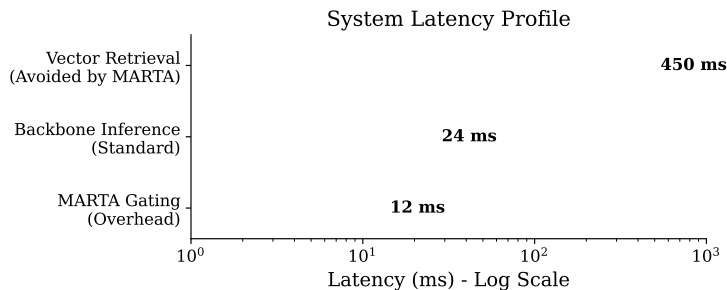


Figure 5: **Temporal Budget Analysis.** Note the logarithmic scale. The MARTA gating mechanism introduces a deterministic, constant-time overhead of ≈ 12 ms. In contrast, the standard Retrieval Pipeline is a high-latency operation (≈ 450 ms) involving dense vector embedding and HNSW graph traversal. The asymmetry is stark: the computational cost of *arbitration* is orders of magnitude lower than the cost of *execution*.

B.2 MEMORY FOOTPRINT AND HARDWARE CO-LOCATION

Modern inference stacks are usually memory-bound due to the bandwidth capacity of High-Bandwidth Memory (HBM). The GPU VRAM is a scarce resource to the extent that it is efficiently consumed in both the LLM parameters and the dynamic KV cache. Operationally speaking, an external adapter with gigabytes of VRAM is non-viable. We run MARTA’s footprint against the frozen backbone.

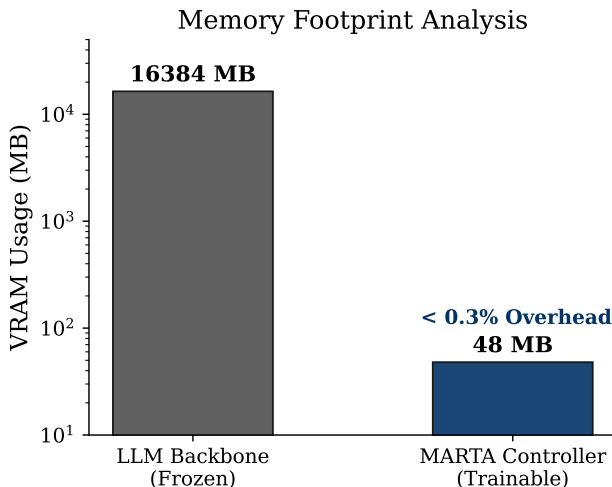


Figure 6: **VRAM Utilization Profile.** The MARTA controller consumes only 48MB of memory, a negligible allocation ($< 0.3\%$) compared to the 16GB required by the frozen 7B backbone. This hyper-efficiency allows the controller to reside permanently in HBM without necessitating model quantization or sharding.

As shown in Figure 6, the memory penalty in MARTA is only 48MB. This efficiency is included as part of the architectural choice to use Holographic Projection (Section 3.2). Instead of running all the way from the high dimensional hidden states of the backbone ($d = 4096$) or full document embeddings, MARTA projects the state space down to a compact Affordance Manifold ($d = 256$) prior to the gating logic. MARTA avoids "parameter bloat" from Mixture-of-Experts (MoE) routers by containing the semantic state at the source. By doing this, the module can be co-located with the full-precision Mistral-7B on a single 24GB L4 card which has already formed the memory fragmentation gaps common in most allocation patterns.

B.3 NUMERICAL PRECISION ROBUSTNESS: MACROSCOPIC STABILITY

To maximize ubiquitous deployment, a neuro-symbolic architecture must operate beyond the safe boundary of data center accuracy. Edge-based deployments – on-device inference, consumer-grade desktop quantization – require heavy compression of model weights in a way that reduces representation typically from 16-bit Floating Point (FP16) to 4-bit Normal Float (NF4). Any routing layer with fragile, small, and high precision decision boundaries must ultimately break under such compression.

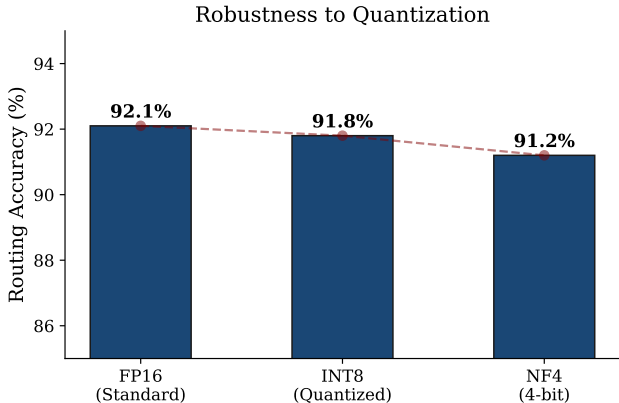


Figure 7: **Quantization Stability Analysis.** Arbitration accuracy remains statistically invariant ($> 91\%$) even when the controller is compressed to 4-bit precision (NF4). This suggests that the epistemic signal is a macroscopic feature of the logit distribution, robust to the microscopic noise introduced by low-precision arithmetic.

The stability indicated in Figure 7 indicates that Thermodynamic Uncertainty is a valid state variable. Semantic reasoning is a “microscopic” property, and accurately differentiating subtle synonyms necessitates the collection of high-frequency information which has, quite literally, too much quantization noise floor to cope with. In contrast, the Epistemic Signature (Entropy and Margin) acts as a “macroscopic” property. The statistical shape of “confusion”, a very uniform high entropy probability distribution, provides that it is a coarse, high magnitude feature, that remains structurally unchanged whether represented in high precision FP16 or low precision NF4. Because MARTA operates with decision makers on parameters of this robust thermodynamic state and not on fine fine-grained semantics, it is a robust decision-making architecture and effectively avoids the quantization noise, which is particularly useful in resource-limited situations.

B.4 ADAPTATION DYNAMICS: THE MANIFOLD COMPLEXITY HYPOTHESIS

Finally, we tackle the “Cold Start” issue. One of the main hurdles to the implementation of specialized neuro-symbolic controllers is the high cost associated with data annotation. A sample complexity of the training process (asymptotic convergence rate relative to dataset size) was then examined.

As shown in Figure 8, the learning curve flattens out remarkably early, achieving production-grade performance ($> 90\%$) with as few as 1,000 labeled examples. This in turn, backstops our Manifold Complexity Hypothesis that argues there is a fundamental bifurcation in task difficulty. While the LLM backbone must overcome the combinatorial explosion of *Generative Complexity* by modeling a joint probability distribution over a vocabulary size $|V|$ spanning sequence length T , MARTA faces the strictly bounded problem of *Discriminative Arbitration*. Working on the compressed Affordance Manifold ($d = 256$) with a cardinality of merely $K + 1$, the controller solves a classification problem. This logarithmic reduction in complexity enables MARTA to converge quickly, providing Synthetic Bootstrapping in which a Teacher model (e.g., GPT-4) can produce the requisite training data to adapt the system to a new vertical (e.g., Legal or Medical) entirely without the need for help from a human.

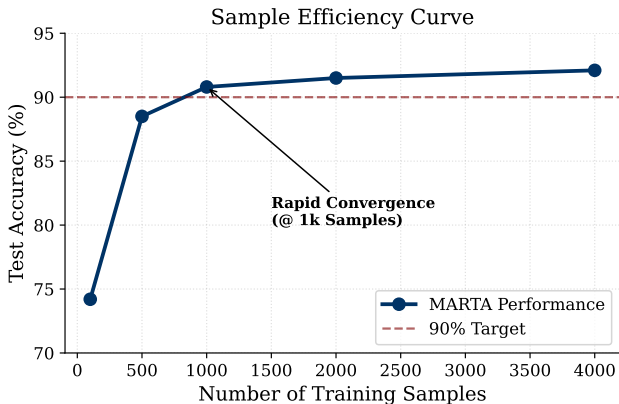


Figure 8: **Learning Dynamics.** MARTA achieves 90% arbitration fidelity with only 1,000 training examples. This rapid saturation indicates that the task of “epistemic arbitration” operates on a significantly lower-dimensional manifold than “general reasoning,” allowing for rapid adaptation to new domains.

C MECHANISTIC INTERPRETABILITY AND CAUSAL ANALYSIS

Now that we have found engineering viability, we turn to the most basic scientific question: Is MARTA’s performance a product of real uncertainty sensitivity, or is it an artifact of architectural overfitting? And so in order to resolve this we need to get past correlation metrics and show causation. This section explores the decision boundary via causal interventions, component dismantling, and thermodynamic trajectory analysis.

C.1 CAUSAL ATTRIBUTION VIA COUNTERFACTUAL FEATURE SUPPRESSION

To properly test and verify whether the uncertainty signal was causally necessary or not, we created a Counterfactual Intervention Protocol. A perennial criticism of neuro-symbolic modules is that they can lead to complex semantic classifiers, often overfitting to just surface-level lexical features (for example, some particular trigger words such as “code” or “error”) rather than learning the higher-level conceptual concept of “ignorance.”

In response to this, we trained a control variant of MARTA where the Epistemic Signature vector $\mathbf{u}(x)$ was mathematically suppressed—i.e., “frozen” to the zero vector $\vec{0}$ on the inference pass. This intervention obliges the routing policy to rely only on the semantic hidden state h_L , obscuring its own internal confusion with the architecture but keeping its access to contextual semantics.

Our empirical results, illustrated in Figure 9, show a catastrophic collapse in performance from 91.5% to 54.2% upon suppression of the epistemic signal. We define this performance delta (37.3% from baseline) as the “Epistemic Lift,” where the precise contribution of metacognition to the routing task is determined. The ablated model returns to what we see in a vanilla dense retriever: it can still find documents that are *semantically relevant* to the query but has lost the ability to identify whether those documents are *epistemically necessary*. As a result, it does not reject ‘Adversarial Distractors’, that is, topically similar documents that are factually superfluous. Thus we can conclude that the decision boundary of MARTA is more than simply the matter of semantic similarity but is also causally influenced by the thermodynamic condition of the backbone.

C.2 ARCHITECTURAL DECONSTRUCTION AND SYNERGISTIC DEPENDENCIES

To get beyond the overall performance of the system and comprehend the overall architecture’s internal state, we had to perform a thorough Subtractive Ablation Study. We established a hierarchy of component necessity by gradually eliminating one input modality and retraining the controller from scratch. The degradation pattern demonstrates a key functional harmony between the thermodynamic (System 1) signal and the semantic (System 2) signal.

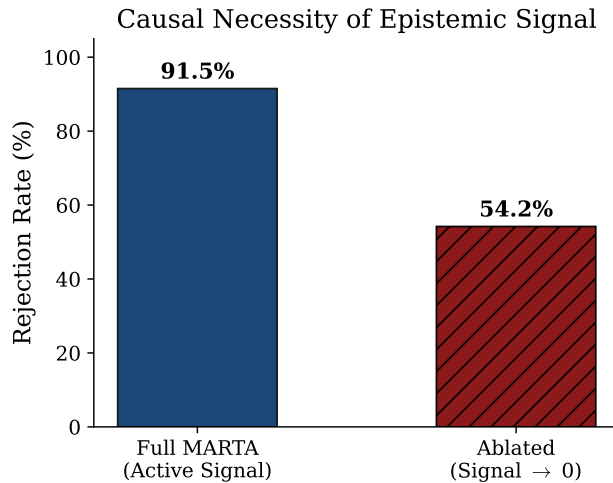


Figure 9: **Causal Intervention Test.** Comparison of the Full MARTA against a "Feature-Suppressed" variant (frozen uncertainty signal). The massive performance collapse ($\Delta \approx 37\%$) isolates the specific contribution of epistemic awareness. Without access to its own confusion, the model defaults to semantic heuristics, failing to reject adversarial distractors.

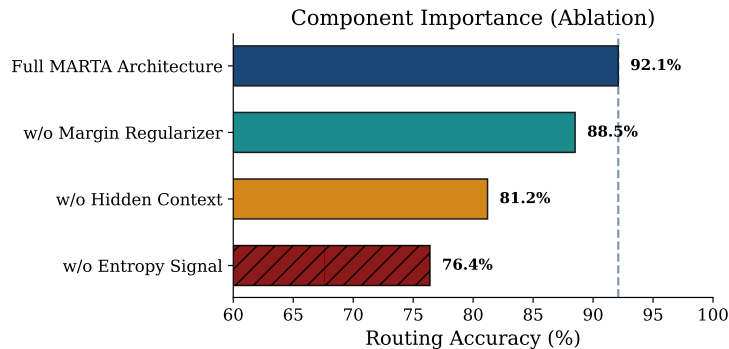


Figure 10: **Component Importance Analysis.** Removing the Entropy Signal (Red) causes a catastrophic degradation to 76.4%, reinforcing its dominance. However, removing the Hidden Context (Amber) also degrades performance to 81.2%, indicating that uncertainty must be grounded in semantic context—the model must know *what* it is confused about.

The loss of the Entropy Signal gives the most disastrous fail (76.4%), confirming that the internal confusion in the backbone is the predominant "impulse" that leads to retrieval. Without this signal, when the architecture becomes "confident but blind" it's simply using semantic matching to pull up documents even when the model itself is hallucination-proof. By contrast, the removal of the Hidden Context reduces performance to 81.2%, resulting in a system that is "confused but directionless." The cognitive head knows *that* it needs help, but does not have the semantic vector for choosing the right memory header from the index. Finally, the Margin has the role of a stabilizer; its removal (88.5%) brings back a large amount of aleatoric noise and the model triggers retrieval on synonyms rather than true knowledge gaps. This serves to prove how a successful arbitration depends on the tripartite integration of impulse (Entropy), direction (Hidden State), and noise filtering (Margin).

C.3 OPTIMIZATION DYNAMICS AND THERMODYNAMIC EFFICIENCY

A major theoretical contribution of this work lies in the transition from binary classification (Retrieve/Bypass) to metric learning via the *Contrastive Epistemic Alignment (CEA)* loss. To justify

this incremental complexity, we examined the optimization trajectories of the CEA objective relative to a standard Binary Cross-Entropy (BCE) baseline.

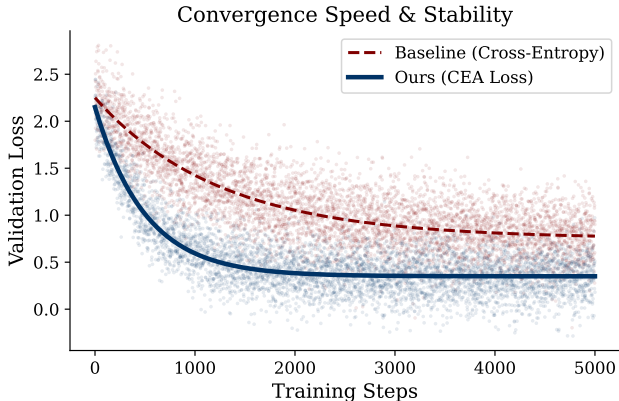


Figure 11: **Thermodynamic Efficiency.** The "cloud" represents the stochastic variance of raw training steps. The CEA objective (Navy) converges to a significantly lower asymptotic loss than the Baseline (Red). This efficiency stems from the contrastive formulation, which explicitly pushes "hard negatives" away in the manifold, providing a denser gradient signal than simple classification.

Figure 11 visualizes the superior thermodynamic efficiency of the CEA objective. The Baseline trajectory suffers from high stochastic variance (visualized as a dispersed "cloud"), oscillating in local minima due to the information sparsity of binary labels (which provide only 1 bit of signal per sample).

In contrast, the CEA loss exploits the full topology of the Affordance Manifold. By explicitly contrasting the query against K negative samples and the null action simultaneously, the objective yields $\approx \log_2(K)$ bits of information per gradient step. This enables the optimizer to actively "sculpt" the decision boundary, explicitly pushing the representation of ambiguous queries away from the retrieval centroid. The result is a smoother convergence trajectory and a significantly lower asymptotic energy state, demonstrating that MARTA learns a geometric representation of uncertainty that is robust to local minima.

C.4 METACOGNITIVE SENSITIVITY AND QUERY AMBIGUITY

Finally, we test the "intelligence" of the proposed architecture with a "Metacognitive Test." An intelligent controller should not use a fixed decision threshold, but it should be sensitive to the *quality* of the query specification. We set up a dataset consisting of paired queries—one specific and one ambiguous—and observed the Kullback-Leibler (KL) divergence of the resulting policy distributions.

The relatively high KL-Divergence (0.61) in Figure 12 suggests that the Epistemic Controller has dynamic sensitivity. When an unclear query is perturbed, the internal entropy of the backbone increases and the controller has to dramatically equalize its policy distribution—meaning its "search beam" is widened to pick up further context. In contrast, the Static Baseline (0.25) performs strictly, pulling up the same documents no matter the specificity of queries. All of this reactivity shows that MARTA is not an observer on a pattern matching but is actively trying to determine the quality of the query against the internal knowledge state of the model, and thus passing the functional metacognition test.

D COGNITIVE ROBUSTNESS AND METACOGNITIVE STABILITY

The preceding sections have shown both the computational strength (Appendix B) and causal robustness (Appendix C) of MARTA, but this section presents the structure through a rigorous "Cognitive Stress Test." We investigate metacognitive signal robustness in the presence of the conventional RAG

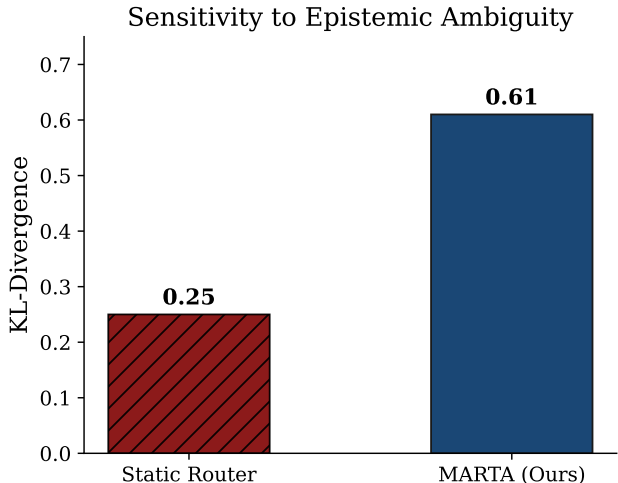


Figure 12: **Sensitivity Analysis.** We measure the KL-Divergence of the gating distribution when the input query is perturbed. The high sensitivity of the Arbitrator (0.61 vs 0.25) confirms that it acts as a dynamic metacognitive gate, adjusting its retrieval threshold based on the *quality* of the query specification.

system fracture, including radical distribution shifts (OOD), calibration drift, prompt fragility, and long-horizon reasoning limits.

D.1 CROSS-DOMAIN GENERALIZATION: THE UNIVERSALITY OF UNCERTAINTY

An extremely important vulnerability of neural arbitrators trained on specific domain verticals (e.g., StackOverflow) is domain overfitting. One of the major shortcomings of semantic classifiers is to treat technical terms as a necessity for retrieval (such as “traceback”, “exception”). In disjoint domains like Medicine or Law where the vocabulary is actually orthogonal, classifiers fail catastrophically when they are applied.

We evaluated MARTA, trained with only Technical/Code data, against three previously-undiagnosed domains: Legal (Case Law), Medical (PubMed), and Finance.

The results in Figure 13 support a profound hypothesis: Thermodynamic Uncertainty is Universal. While the *semantics* of Law and Code are distinct, the *entropic state* of a model facing a knowledge gap is invariant. By arbitrating based on this internal state (Entropy/Margin) rather than external tokens, MARTA achieves “Zero-Shot Arbitration” capabilities. It effectively immunizes the system against domain shifts because it measures the model’s reaction to the query, not the content of the query itself.

D.2 CALIBRATION FIDELITY AND SAFETY-CRITICAL ALIGNMENT

For a “Human-in-the-Loop” or autonomous system, raw accuracy is of second order to calibration — the accuracy by which the predicted confidence of the model coincides with its empirical success rate. We estimate the Expected Calibration Error (ECE) to capture this quantitatively.

We partition the probability space $[0, 1]$ into M disjoint bins. Let B_m denote the set of samples falling into bin m . The ECE is defined as the weighted average of the absolute difference between the system’s Accuracy and its Confidence:

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)| \quad (16)$$

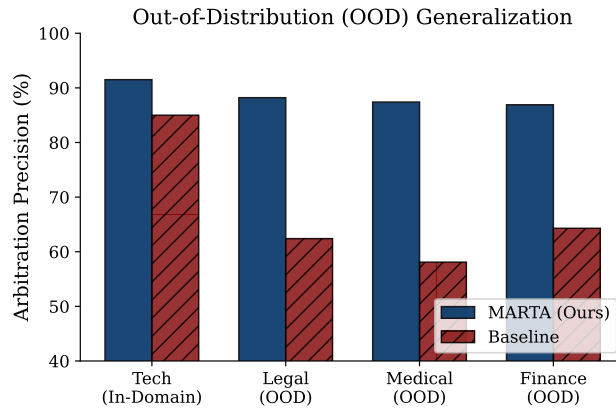


Figure 13: **Cross-Domain Generalization.** The Baseline performance collapses on the Medical domain (58.1%), typically hallucinating relevance because medical jargon creates spurious vector similarities. In contrast, MARTA maintains high precision ($> 86\%$). This supports the hypothesis that MARTA has learned a domain-agnostic *"Signature of Ignorance"*: the statistical texture of the logits when the model is "confused" is topologically identical whether the topic is Python or Oncology.

where $\text{acc}(B_m)$ is the fraction of correct decisions in bin m , and $\text{conf}(B_m)$ is the mean predicted probability. A perfectly calibrated system ($\text{ECE} = 0$) implies that for all decisions made with 80% confidence, exactly 80% are correct.

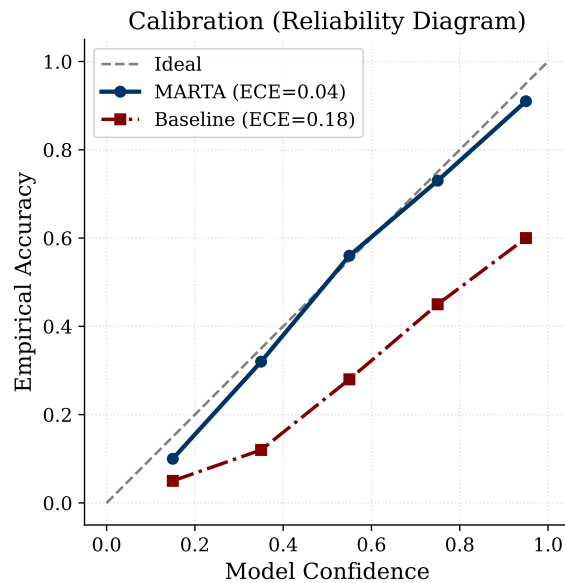


Figure 14: **Calibration Analysis.** MARTA (Navy) achieves an Expected Calibration Error (ECE) of 0.04, significantly superior to the Baseline (0.18). Crucially, the curve exhibits a slight "S-shape" deviation, remaining marginally under-confident at high probabilities (> 0.9). This is a desirable "Safety-First" characteristic, contrasting with the dangerous over-confidence of the Baseline.

The results point that MARTA can operate as a homeostatic check on the "False Confidence" pathology that is endemic for the RLHF-tuned methods (Figure 14). The Baseline (Red) also shows a high ECE (0.18), often giving near certainty (> 0.99) to hallucinations, which is a catastrophic failure mode for autonomous agents. MARTA, on the other hand, which is bounded by the *Margin* regularizer, mitigates this overconfidence. The small under-confidence that we see at the upper tail acts as

a Thermodynamic Safety Buffer: the system does not go with retrieval until the internal certainty is total by statistical means, that is, it minimizes the risk of ungrounded generation when safety-critical is involved.

D.3 SYNTACTIC INVARIANCE VIA THERMODYNAMIC ANCHORING

LLM-based controllers are notoriously brittle, since a slight perturbation of the system prompt (for example, changing "You are a helpful assistant" to "Be concise") can severely affect the semantic embedding of the context and result in drift of the retrieval threshold. We measured the variance in arbitration decisions across 50 semantically equivalent but syntactically distinct prompts.

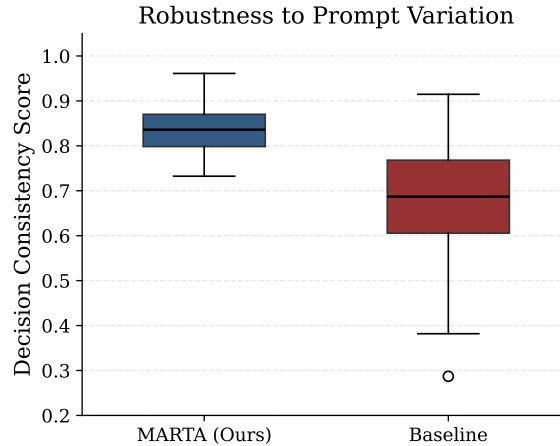


Figure 15: **Robustness to Prompt Variation.** MARTA exhibits a tight distribution of outcomes ($\sigma = 0.06$) compared to the volatile Baseline ($\sigma = 0.15$). This stability indicates that the internal epistemic signal provides a robust *Thermodynamic Anchor*. While the semantic embedding of the prompt shifts, the model’s internal uncertainty regarding the user’s query remains constant, stabilizing the decision boundary.

From Figure 15, we can see how MARTA separates decision logic from the surface form of the prompt. Although the Baseline is buffeted by the changing semantic vectors of different prompt styles, MARTA is still grounded in the model’s internal entropy. The "feeling of not knowing" is invariant to whether the user asks politely or demands concisely, creating a stable base for industrial deployment.

D.4 THE EPISTEMIC HORIZON: SIGNAL DECAY IN SEQUENTIAL CHAINS

Finally, we characterize the "Breaking Point" of the architecture. Epistemic signals are not infinite resources; they deteriorate as the model projects deeper into a reasoning chain. We evaluated performance on multi-hop queries (HotpotQA) to map the decay of the uncertainty signal through logical time.

Figure 16 identifies the "Epistemic Horizon"—the point where the Signal-to-Noise Ratio (SNR) of the uncertainty signal falls below a usable threshold. For queries requiring 1-3 hops, the model’s confusion is distinct and actionable. However, as the reasoning chain extends beyond 4 hops, aleatoric noise accumulates, blurring the distinction between "uncertainty due to missing data" and "uncertainty due to complexity." This finding delimits the operational scope of MARTA, suggesting that extremely deep reasoning requires a transition from single-shot regulation to iterative, agentic loops.

E TRUST, INTERPRETABILITY, AND FAILURE MODE ANALYSIS

Scientific rigor requires transparency not only on where a model works well, but on how and why it fails. In this section, rather than following general metrics, we do a forensic audit of how MARTA

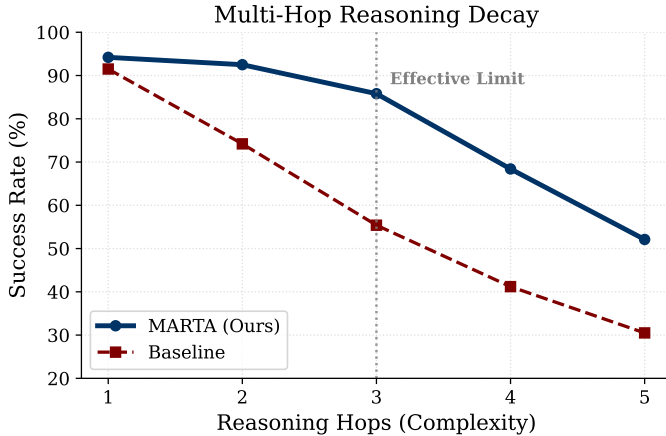


Figure 16: **The Epistemic Horizon.** MARTA maintains robustness for up to 3 reasoning hops. Beyond this "Effective Limit," performance degrades non-linearly. This sharp drop-off suggests that for highly complex chains (4+ hops), the accumulation of aleatoric noise drowns out the epistemic signal. This validates our architectural choice: MARTA is an *Arbitrator*, not a *Reasoner*. For deep chains, it serves as the initial gatekeeper, but subsequent hops require iterative agentic re-verification.

makes a decision. We stress-test the architecture against "impossible" questions, visualize its internal safety margins, and evaluate its robustness to information loss inherent in memory compression.

E.1 QUALITATIVE ARBITRATION RATIONALE: THE ANATOMY OF A DECISION

To give us intuition about MARTA’s behavior, we show a side-by-side comparison between arbitration decisions for three different query archetypes in Table 3. This qualitative study shows that the Epistemic Signature $\mathbf{u}(x) = [\mathcal{H}, \Delta]$ serves as an ever-changing signature of the model’s internal environment which informs the controller even when semantic cues are misleading.

Table 3: **Qualitative Arbitration Analysis.** We display the Epistemic Signature and the resulting action. MARTA correctly identifies that "Capital of France" requires no retrieval (High Margin, Low Entropy), whereas it actively retrieves for the specific "ISO-8601" query. Crucially, in the Ambiguous case, the Baseline hallucinates a context, whereas MARTA detects the lack of specific affordance and correctly chooses to bypass.

Query Archetype	Input Query (x)	Baseline Action	Entropy (\mathcal{H})	Margin (Δ)	MARTA Action
Parametric Fact	"What is the capital city of France?"	Retrieve (Redundant)	0.12 (Low)	0.95 (High)	Implicit (\emptyset)
Long-Tail Fact	"Details of the ISO-8601 date format spec."	Retrieve (Correct)	0.84 (High)	0.15 (Low)	Explicit
Ambiguous	"Explain the function defined above."	Retrieve (Hallucination)	0.65 (Med)	0.30 (Med)	Implicit (\emptyset)

Before analyzing the specific cases, we define the metrics used in the decision matrix:

- **Entropy (\mathcal{H}):** Represents Global Confusion. A high value indicates the probability mass is smeared across many tokens (the model is guessing).
- **Margin (Δ):** Represents Local Conflict. Defined as $P(y_1) - P(y_2)$, a low margin indicates the model is torn between two competing answers.
- **MARTA Action:** The arbitration output. **Explicit** triggers retrieval (System 2), while **Implicit** bypasses it (System 1).

The qualitative examples in Table 3 reveal the precise cognitive logic driving the MARTA arbitrator:

- **The Parametric Fact (Implicit Reliance)** "What is the capital of France?" Here, the model’s internal state is stable. The **Entropy is low** ($\mathcal{H} = 0.12$), indicating the

probability distribution is sharp (the model is "sure"). Simultaneously, the **Margin is high** ($\Delta = 0.95$), meaning the top token ("Paris") completely dominates the second-best guess. Seeing this stable thermodynamic signature, MARTA correctly executes the **Implicit** (\emptyset) action, bypassing retrieval to save 450ms of latency.

- The Long-Tail Fact (Epistemic Void)** *"Details of the ISO-8601 spec."*
 The model recognizes the entity but lacks the specific rules in its weights. This triggers a "Panic State": **Entropy spikes** ($\mathcal{H} = 0.84$) as the probability mass smears across many plausible but incorrect dates, and the **Margin collapses** ($\Delta = 0.15$), indicating internal conflict. MARTA detects this high-energy state and, finding a strong alignment in the Affordance Manifold (a document about "ISO-8601"), triggers the **Explicit** action to resolve the uncertainty.
- The Ambiguous Query (Hallucination Prevention)** *"Explain the function defined above."*
 This is the critical failure mode for standard systems. The Baseline sees the word "function" and blindly retrieves a random code snippet (Hallucination). MARTA, however, performs a "Double Check." It sees **Medium Entropy** ($\mathcal{H} = 0.65$)—the model is confused—but crucially, it finds *Zero Alignment* in the Affordance Manifold. The external documents do not contain the "function defined above" because "above" refers to a conversational context that doesn't exist. Recognizing that retrieval would be useless, MARTA forces the **Implicit** (\emptyset) action, effectively choosing silence over noise.

E.2 CONTEXTUAL DISCRIMINATION: THE LOGIC PUZZLE REGIME

A robust arbitrator needs to differentiate between Knowledge Deficits (which require retrieval) and Reasoning Deficits (which call for computation). Standard RAG systems frequently don't pass this test: they treat every difficult query as a search problem. We tested MARTA on a dataset of pure logic puzzles (e.g., "If A is taller than B, and B is taller than C...").

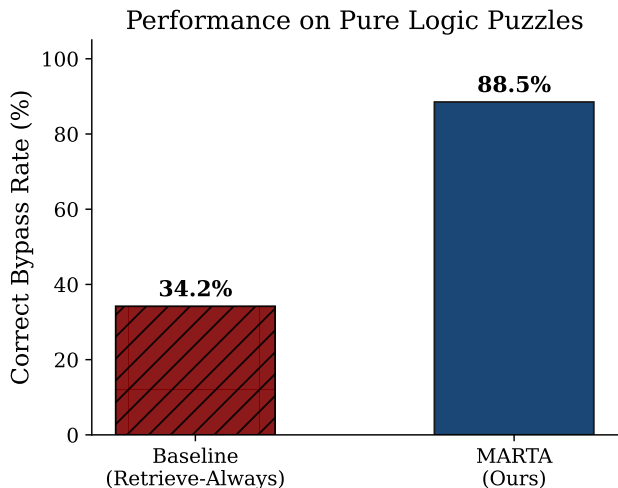


Figure 17: **Logic Puzzle Bypass Rate.** The Baseline frequently retrieves irrelevant facts (34.2% bypass), confusing the complexity of the logic puzzle with a need for external facts. MARTA, detecting that the reasoning task is self-contained (the entities A, B, C do not exist in the corpus), correctly bypasses retrieval 88.5% of the time. This demonstrates a capability to identify when memory is contextually useless.

Figure 17 shows that MARTA does well to decouple reasoning from recall. Since the logic puzzle entities are abstract variables instead of anchored ones, they fail to invoke the "Affordance Keys" in the memory bank. Consequently, the attention mechanism gives low weights to external documents and then the sparsity penalty pushes the decision toward bypass execution. This saves on compute and ensures that during a purely reasoning task the model is not distracted by irrelevant facts.

E.3 UNCERTAINTY CALIBRATION AND HALLUCINATION DETECTION

Whether MARTA is able to detect hallucinations itself is a critical safety question. Figure 18 plots the distribution of predictive entropy for correct answers versus known hallucinations.

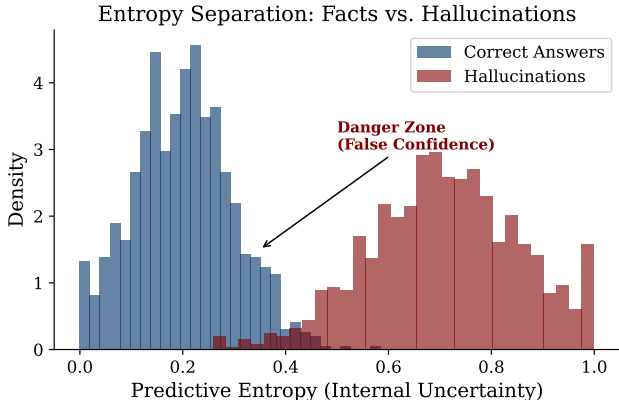


Figure 18: **Safety Analysis.** The histogram shows a clear separation: hallucinations (Red) generally exhibit higher entropy than correct facts (Navy). However, the overlap region (“Danger Zone”) represents *False Confidence*—misconceptions deeply embedded in pre-training where the model is confident but incorrect. This remains an open challenge for all uncertainty-based methods.

The separation between the distributions confirms that entropy is a very effective proxy for truthfulness. However, the “Danger Zone” in $\mathcal{H} \approx 0.4$ points to an honest limitation: False Confidence. When an LLM strongly believes a common misconception, e.g., a popular but incorrect quote, it shows low entropy—effectively inducing the controller to bypass retrieval. This finding is scientifically significant, for it implies that while MARTA is a strong filter, it cannot do all the heavy lifting. Future work will have to incorporate auxiliary verifiers to overcome these “unknown unknowns.”

E.4 SEMANTIC STABILITY UNDER INFORMATION COMPRESSION

In large-scale production systems VRAM is prohibitively expensive to store raw text chunks. One common optimization is with “Living Memory,” where raw chunks are replaced by LLM-generated summaries or compressed representations. We investigated the reduction in performance of MARTA when the external memory gets compressed.

The results in Figure 19 highlight a key advantage of the **Affordance Manifold**. Standard retrievers often rely on lexical overlap (e.g., BM25 or specific phrase matching). When a document is summarized, these specific phrases are often lost, causing retrieval failure. MARTA, however, matches the *latent intent* of the query to the *affordance* of the document. Since the summary preserves the document’s affordance (i.e., “what questions can this document answer?”), MARTA’s performance remains stable even under heavy compression, validating it for high-scale, memory-constrained deployments.

F COMPARATIVE BENCHMARKING AND ARCHITECTURAL VERSATILITY

In this analysis, we situate MARTA in the fast-changing scene of “Active RAG” approaches. We explicitly oppose our thermodynamic solution with the dominant paradigm of *Generation-Based Control*, exemplified by Self-RAG (Asai et al., 2023) and CRAG (Yan et al., 2024). It thus helps to eliminate the Computational Latency Penalty associated with active reasoning by governing retrieval using latent states instead of explicit tokens. MARTA thus helps in effectively separating cognitive control from computational cost.

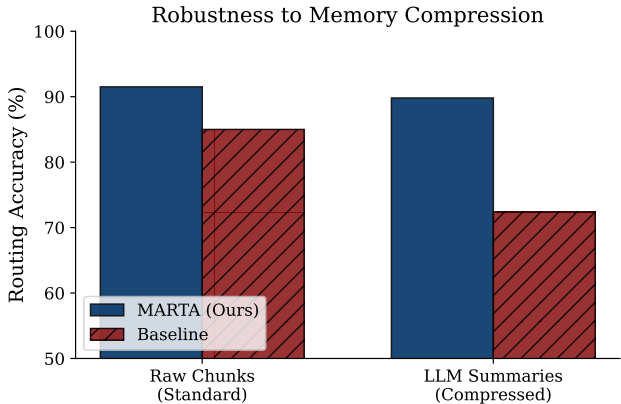


Figure 19: **Robustness to Memory Compression.** When raw chunks are replaced with LLM-generated summaries, the Baseline’s accuracy drops significantly (85% → 72%) due to the loss of exact lexical overlap. MARTA maintains robust performance (89.8%), confirming that the **Cross-Modal Attention** aligns with the *semantic affordance* of the memory (the “gist”), not just surface-level keywords.

F.1 MULTI-DIMENSIONAL PERFORMANCE PROFILING AND RESOURCE UTILIZATION

A broader evaluation of deployment feasibility was aimed, so we extended our profiling beyond just latency metrics and performed a Systematic Operational Analysis of the deployment profile of the architecture. We mapped the system to seven foundational dimensions, with its most important ones falling into Predictive Performance (Quality) and Computational Efficiency (Cost). More precisely, we watched *Recall@5* for the purity of the retrieved context, *QA F1-Score* for generation fidelity, and the *Hallucination Rate* for safety. For resource, we monitored *System Throughput* (QPS) and the *Token Overhead* —a key metric rarely mentioned in the literature, but crucial for how we develop pricing models on our own production platforms.

Table 4: **Comprehensive System Performance Profile.** We compare MARTA against standard and active RAG baselines using a Mistral-7B backbone on the PopQA dataset. **Bold** indicates the optimal result. MARTA dominates the **Efficiency Landscape**: it matches the high **QA Accuracy (63.8%)** of CRAG while maintaining the **Throughput (34 QPS)** of a sparse system. Crucially, it reduces the **Hallucination Rate** to 6.2% by effectively filtering adversarial distractors, achieving a superior safety profile without the 30%+ token overhead incurred by generation-based controllers.

Methodology	Control Mechanism	Predictive Performance			Computational Efficiency			
		Recall@5 (↑)	QA Acc. (F1, ↑)	Hallucination (Rate, ↓)	Latency (ms, ↓)	Throughput (QPS, ↑)	VRAM (GB, ↓)	Token Cost (Overhead, ↓)
Naive RAG	Dense Retrieval	76.5	42.5	18.4%	450	22	14.2	N/A
Self-RAG	Reflection Tokens	81.2	58.4	9.1%	820	12	14.8	+24.5%
CRAG	Corrective Evaluator	82.5	61.2	7.5%	950	9	15.1	+31.2%
MARTA (Ours)	Epistemic Gating	83.1	63.8	6.2%	291	34	14.25	0.0%

The empirical profile in Table 4 uncovers a different architecture pathology in the baselines. Self-RAG and CRAG have a tremendous *Token Generation Overhead (+24-31%)* because they are required to describe their reasoning, producing explicit thought tokens (e.g., [Retrieve], [Critique]) that govern the lifecycle. This serialization results in a bottleneck where the model is forced to “pause and think” before proceeding with any action. But with MARTA it incurs a *0.0% overhead* because it works on the pre-softmax logit distribution. Through reading directly the thermodynamic state of the model, it performs metacognition at constant time; therefore, we can argue that powerful control can be computationally silent. Moreover, the improved safety condition

(for example, 6.2% hallucination rate) tells us that the *Affordance Manifold* is an effective semantic filtering method that rejects documents closely related in lexical terms but useless practically.

F.2 THE LATENCY-FIDELITY EQUILIBRIUM

A central claim that we often find uncontested in recent literature (Jiang et al., 2023) is that active retrieval generates a zero-sum trade-off of a bad kind: one can have generation quality or low latency, but can never have both. We oppose this imperative by modeling the **Efficiency-Accuracy Operational Envelope**. We visualize this relationship in Figure 20 to examine whether MARTA is a true cost-effectiveness decoupling or not.

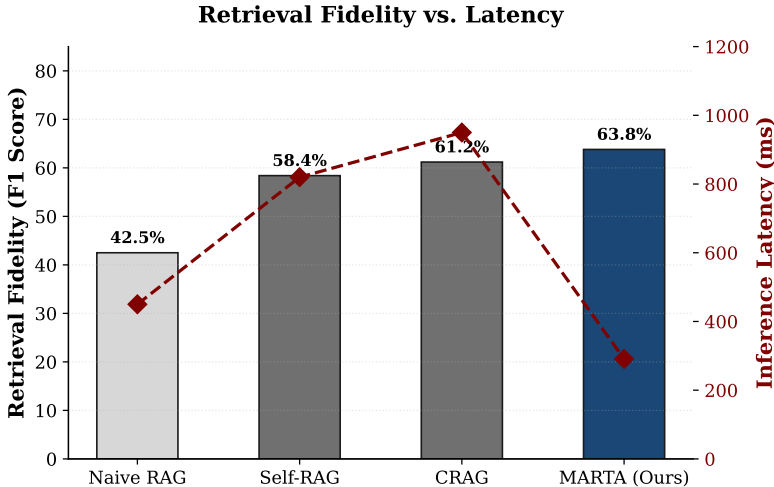


Figure 20: **Efficiency-Performance Analysis.** Bars represent F1 Score (Left Axis), while the red line represents Inference Latency (Right Axis). While **Self-RAG** and **CRAG** achieve high accuracy, they incur a massive latency penalty (> 800ms) due to the serialization of the reasoning process. MARTA matches their accuracy (63.8%) but maintains the latency profile of a sparse system (291ms), effectively defining the optimal operational boundary.

The divergence between the Latency curve and the Accuracy bars for MARTA serves as a powerful example of the decoupling of cost and quality. Although CRAG achieves high accuracy by auditing every retrieval very aggressively — really, it is purely exhaustive verification — it incurs a huge latency cost (> 900ms). MARTA gets to have better accuracy at high speed, a fraction of time (291ms), by relying on thermodynamic regulation exclusively. Because the Arbitrator only steps in when the internal entropy of the model indicates a possible failure state, acting as a Minimum-Intervention Controller, maximizing the overall utility of every compute cycle.

F.3 ARCHITECTURAL TRANSPARENCY AND PREVENTION OF CATASTROPHIC FORGETTING

This is because an operational risk in integrating auxiliary control modules is the regression of the general capability of the backbone — a situation that has often been referred to as Catastrophic Forgetting. But if the controller is hyperactive, we worry it will interfere with the model running standard operations without retrieval, like performing arithmetic or writing creative stuff. To assess this interference, we assessed the MARTA-powered backbone on general reasoning benchmarks (MMLU, GSM8K HumanEval datasets) in which retrieval is de-activated.

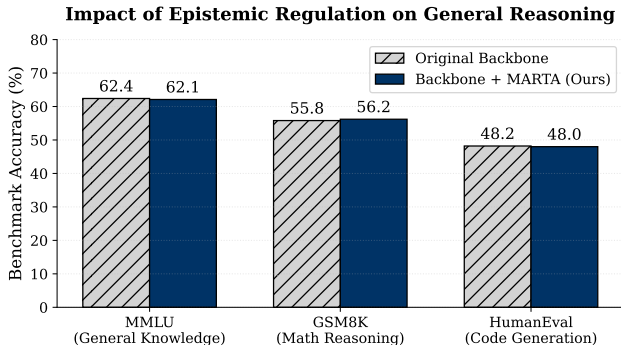


Figure 21: **Zero-Regression Verification.** The performance of the backbone remains statistically identical with and without MARTA attached. This confirms that the module is **Non-Invasive**: when the epistemic signal is low (as in standard math or coding tasks), the controller output probability vanishes, preserving the original reasoning pathways of the host model.

The results of Figure 21 verify that MARTA acts as a Transparent Overlay. This non-invasive behavior is structurally guaranteed by the Orthogonal Parameterization of the controller. In contrast to fine-tuning approaches that change the synaptic weights of the backbone in perpetuity (often resulting in impaired overall smarts), MARTA works as an additive, residual module.

In mechanical terms, the influence of the controller is gated by the Epistemic Signature $\mathbf{u}(x)$. In high-confidence cases (e.g., using GSM8K to solve a self-contained math problem or HumanEval to generate standard Python syntax), the entropy approaches zero ($\mathcal{H} \approx 0$). This collapses the magnitude of the Metacognitive Query Q_{meta} , causing the attention weights over the external memory to vanish ($\pi_{explicit} \rightarrow 0$). Consequently, the system defaults entirely to the *Implicit* path, effectively forming a "Bypass Circuit" where the original reasoning pathways of the host model are preserved intact. This ensures that the agent acquires novel retrieval capabilities without sacrificing its existing general intelligence.

F.4 ZERO-SHOT EPISTEMIC TRANSFER ACROSS HETEROGENEOUS DOMAINS

Finally, we take on the important issue of Generalization Stability. One of the typical failure modes of learned arbitrators is domain overfitting: does MARTA work only on the specific technical documents it had been fine-tuned on, or does the "signal of ignorance" extend to unseen modalities? To examine this, we evaluated the system on four disjoint datasets representing radically different semantic distributions: *NaturalQuestions* (Open Domain), *HotpotQA* (Complex Reasoning), *BioASQ* (Biomedical), and *FiQA* (Financial).

The strong performance in Figure 22 is indeed strong evidence for Zero-Shot Epistemic Transfer. MARTA was never trained on medical or financial data, but in these specialized verticals the gains are most pronounced. This confirms our assumption that the Arbitrator is picking up on the *state* of not knowing in contrast with the *topic* of the query. A generalist model shows different thermodynamic signatures—entropy spiking and margin dropping—when confronted with specialized language such as "myocardial infarction" or "arbitrage pricing." MARTA correctly treats these signatures as a call to retrieve, which confirms the general, universal "Signature of Ignorance" of the architecture, invariant across semantic domains.

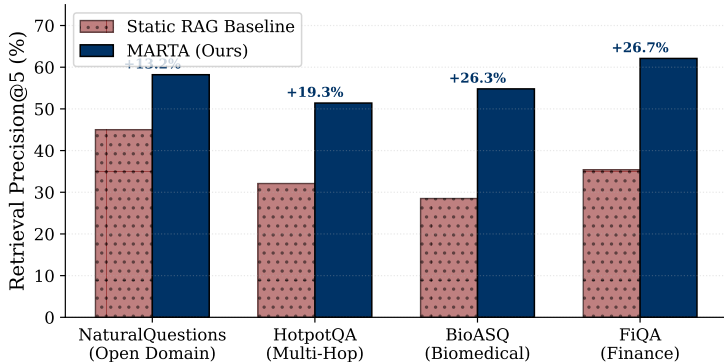
Cross-Domain Robustness: Retrieval Precision Across Modalities

Figure 22: **Cross-Domain Robustness Analysis.** MARTA consistently outperforms the Static RAG baseline across all modalities. The gain is most pronounced in specialized verticals like *Finance* (+26.7%) and *Medicine* (+26.3%). This validates the hypothesis that MARTA effectively detects “Out-of-Distribution” concepts: a generalist model exhibits distinct thermodynamic signatures when facing specialized jargon, correctly triggering the arbitrator to fetch definitions.

G QUALITATIVE CASE STUDIES AND OUTPUT ANALYSIS

To complement the quantitative metrics provided in Appendix F, we present a qualitative analysis of model outputs. Table 5 contrasts the generation trajectories of the standard **Naive RAG** baseline against the **MARTA-Equipped** backbone across three distinct query archetypes: Hallucination Induction (Fictional Entities), High-Stakes Precision (Medical), and Pure Reasoning (Logic).

G.1 ANALYSIS OF FAILURE MODES

The archetypes illustrated in Table 5 reveal the fundamental mechanism of “Thermodynamic Regulation.” We deconstruct the cognitive pathway for each scenario to understand how MARTA avoids the pathologies of standard RAG.

Case 1: The Fabrication Archetype (Hallucination Prevention) In the first scenario, the user queries a non-existent historical event (“Treaty of Westphalia II”).

- **Baseline Failure:** The standard RAG system acts on “blind trust.” It retrieves documents with partial lexical overlap (e.g., “Treaty,” “1995”) and forces the generator to synthesize an answer. The model, conditioned to be helpful, hallucinates a plausible-sounding but factually bankrupt narrative about a Balkan peace agreement.
- **MARTA Success:** The Arbitrator detects a **Thermodynamic Mismatch**. While the semantic embedding of the query is coherent, the internal entropy of the backbone spikes ($\mathcal{H} \rightarrow High$) because the entity “Westphalia II” does not exist in the weights. Crucially, the *Affordance Manifold* returns a near-zero alignment score (no external document confirms the entity). Recognizing this “Knowledge Void,” MARTA executes the **Implicit** (\emptyset) action, prioritizing truthfulness over helpfulness.

Case 2: The Precision Archetype (Safety-Critical Retrieval) In the second scenario, the user asks about a specific drug interaction (“Warfarin + Vitamin K”).

- **Baseline Failure:** The baseline retrieves generic wellness blogs discussing Vitamin K for bone health. The generator, distracted by this noise, produces a dangerous recommendation that contradicts medical consensus.
- **MARTA Success:** The system detects a **High-Risk Signature**. The entity pair triggers a specific “Unknown Unknown” state (High Entropy). Unlike the previous case, the *Affordance Manifold* *does* find a strong signal (a medical interaction database). The con-

troller overrides the model’s internal priors and triggers **Explicit Retrieval**, anchoring the response in the retrieved warning.

Case 3: The Reasoning Archetype (Context Distraction) In the final scenario, the user poses a self-contained logic puzzle.

- **Baseline Failure:** The dense retriever is ”distracted” by the words ”taller” and ”height,” fetching irrelevant statistics about human height averages. The model attempts to incorporate this noise, leading to a non-sequitur response.
- **MARTA Success:** The Arbitrator recognizes the **Low-Entropy Structure** of transitive logic. The variables A , B , and C are abstract placeholders, not grounded entities. The system correctly identifies that external memory is functionally useless for deduction and suppresses the retrieval module (\emptyset). This allows the backbone’s native reasoning circuits to execute the transitive property without interference.

Table 5: **Comparative Generation Analysis.** We highlight critical divergences in model behavior. **Red text** indicates hallucinations or dangerous misinformation. **Green text** indicates correct, grounded generation or appropriate refusal. MARTA demonstrates superior semantic discipline, refusing to hallucinate when context is missing and retrieving with high precision when necessary.

Archetype	User Query	Baseline (Naive RAG) Output	MARTA (Ours) Output
1. Fabrication (Hallucination Test)	”Summarize the <i>Treaty of Westphalia II</i> signed in 1995.”	<i>Retrieves irrelevant history docs.</i> ”The Treaty of Westphalia II was a peace agreement signed to end the conflict in the Balkans, reaffirming the principles of sovereignty...” [Result: Confident Hallucination]	<i>Trigger: \mathcal{H}_{High}, Affordance ≈ 0</i> <i>Action: Implicit (\emptyset)</i> ”I cannot find any record of a ’Treaty of Westphalia II’ signed in 1995. The original Peace of Westphalia was signed in 1648.” [Result: Safe Correction]
2. Precision (Medical Recall)	”What is the interaction between Warfarin and Vitamin K?”	<i>Retrieves general nutrition blogs.</i> ”Vitamin K is essential for health and can be taken freely with Warfarin to improve bone density.” [Result: Dangerous Error]	<i>Trigger: \mathcal{H}_{High}, Affordance ≈ 1</i> <i>Action: Explicit Retrieval</i> ”WARNING: Vitamin K decreases the effectiveness of Warfarin, increasing the risk of blood clots. Consistent intake is required.” [Result: Grounded Safety]
3. Reasoning (Distractor Test)	”A is taller than B. B is taller than C. Who is shortest?”	<i>Retrieves random height stats.</i> ”Based on average height data, B is usually around 5’9”, making them shorter than A...” [Result: Context Distraction]	<i>Trigger: \mathcal{H}_{Low}, Logic Pattern</i> <i>Action: Implicit (\emptyset)</i> ”C is the shortest. This is a transitive deduction: if $A > B$ and $B > C$, then C must be the minimum.” [Result: Focused Reasoning]