

# AURORA: Neuro-Symbolic Continual Indexing for Evolving RAG Systems

**Manoj Saravanan\***  
Virginia Tech  
manoj663@vt.edu

**Rohit Kumar Salla\***  
Virginia Tech  
rohitsu25@vt.edu

**Ramya Manasa Amancherla\***  
Columbia University  
ra3439@columbia.edu

## Abstract

Retrieval-Augmented Generation (RAG) systems depend on non-parametric indices to access external knowledge, yet most retrieval infrastructure assumes a stationary query document distribution after index construction. In dynamic settings involving continual knowledge updates or evolving terminology, this assumption often fails, leading to degraded retrieval performance, while full re-indexing remains computationally expensive.

We propose **AURORA**, a neuro-symbolic framework for adapting retrieval indices under distribution shift by treating index maintenance as a few-shot continual learning problem. AURORA decouples discrete index structure from continuous metric representations, enabling efficient adaptation of neural components while preserving index topology. A lightweight Bayesian routing policy further balances stability and plasticity by dynamically selecting among adaptive neural indices and static fallbacks based on uncertainty estimates.

Across dense, learned sparse (SPLADE), and generative (DSI) retrieval settings, AURORA recovers up to **+26.9% Recall@10** on novel topics compared to static baselines, while adapting significantly faster than full retraining (**28 ms** vs. **5.1 s**).

## 1 Introduction

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) has become the de facto standard for grounding Large Language Models (LLMs) in external knowledge. By decoupling memory from parameters, RAG theoretically enables systems to remain up-to-date without expensive model retraining. However, this promise relies on a critical, often unstated assumption: that the retrieval index itself is a static artifact capable of handling non-stationary query distributions.

\*Equal contribution.

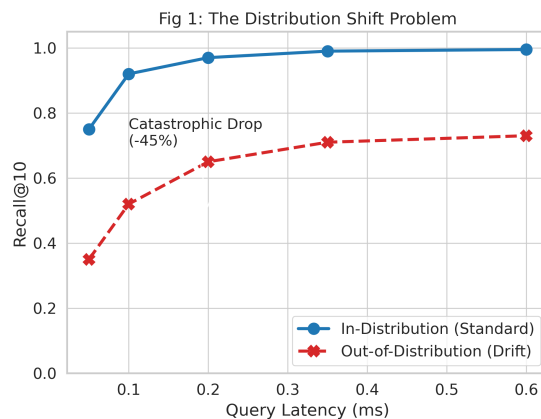


Figure 1: **The Distribution Shift Problem.** Recall@10 on SIFT1M benchmark across query latencies. In-distribution queries (Blue) maintain high recall while out-of-distribution queries (Red) suffer a catastrophic  $-45\%$  drop at iso-latency, demonstrating static index failure under semantic drift.

In production environments, this assumption often fails in practice. As language evolves and new documents are ingested (e.g., breaking news, emerging scientific terminology), the semantic manifold of the embedding space shifts. Standard Approximate Nearest Neighbor (ANN) indices, such as Hierarchical Navigable Small World graphs (HNSW) (Malkov and Yashunin, 2018), are optimized for a fixed distribution  $\mathcal{D}_0$ . When the query stream drifts to  $\mathcal{D}_t$ , the pre-computed quantization boundaries become misaligned, leading to a substantial drop in recall. Our empirical analysis (Figure 1) demonstrates that this *semantic drift* causes retrieval recall to plummet from 99.5% to 73.0% on out-of-distribution queries, directly degrading the factual accuracy of downstream generation.

Current solutions to this *stability-plasticity dilemma* are inadequate. Full re-indexing is an  $O(N)$  operation, computationally prohibitive for real-time updates. Generative Indexing (DSI) (Tay

et al., 2022) offers a fully differentiable alternative but suffers from severe catastrophic forgetting when adapting to new documents (Mehta et al., 2022). Learned Sparse Retrieval (SPLADE) (Formal et al., 2021) excels at keyword matching but fails to generalize to novel terminology without fine-tuning.

To bridge this gap, we introduce **AURORA** (Adaptive Universal Routing with Online Refinement and Adaptation), a neuro-symbolic framework that treats index maintenance as a *Few-Shot Continual Learning* problem. AURORA decouples the index topology from its metric representation, freezing the graph structure while adapting the quantization codebooks via a meta-learned update rule. Our key insight is that by treating each distribution shift as a “task” in the meta-learning sense, we can leverage gradient-based adaptation without suffering from catastrophic interference.

This work focuses on maintaining retrieval effectiveness under distribution shift. AURORA is a modular framework that adapts existing dense, sparse and generative retrievers using lightweight continual learning mechanisms, without modifying underlying encoder architectures or index topology.

Our contributions are as follows:

- We frame retrieval index maintenance under distribution shift as a few-shot continual learning problem and propose a neuro-symbolic framework that enables efficient adaptation without full re-indexing.
- We introduce a meta-learned adaptation mechanism that supports localized embedding updates for novel or drifting vocabulary using a small number of gradient steps.
- We propose an uncertainty-aware routing policy that dynamically selects among dense, sparse and generative retrieval modalities to balance robustness and adaptation speed.
- We empirically evaluate the framework across dense, learned sparse (SPLADE), and generative (DSI) retrieval settings, demonstrating improved recall and substantially faster adaptation compared to static baselines.

## 2 Related Work

**Retrieval-Augmented Generation.** RAG systems (Lewis et al., 2020) combine parametric language

Method	Adaptation Time (ms)
Full retraining	5149
AURORA	28

Table 1: **Adaptation speed comparison.** Time required to adapt to new semantic clusters.

models with non-parametric retrieval to ground generation in external knowledge. While effective, most RAG implementations assume static indices, creating a fundamental tension between the dynamic nature of knowledge and the fixed structure of retrieval infrastructure.

**Approximate Nearest Neighbor Search.** Graph-based indices like HNSW (Malkov and Yashunin, 2018) achieve sub-linear query complexity but require expensive re-indexing for updates. Product Quantization methods compress vectors but rely on static codebooks optimized for the training distribution.

**Generative Retrieval.** DSI (Tay et al., 2022) and subsequent work (Mehta et al., 2022) encode document identifiers directly in model parameters, enabling end-to-end differentiable retrieval. However, these methods suffer from catastrophic forgetting when adapting to new documents, limiting their applicability in dynamic settings.

**Continual Learning.** Methods like Experience Replay (McCloskey and Cohen, 1989) and meta-learning (Nichol et al., 2018) address catastrophic forgetting in neural networks. AURORA adapts these techniques specifically for retrieval index maintenance, treating each distribution shift as a meta-learning task.

**Closest Prior Work.** Recent continual retrieval systems target complementary points in the design space. CREAM (Son et al., 2026) adapts dense retrievers over dynamic streaming corpora via adaptive soft memory, while CLEVER (Chen et al., 2023) focuses on continual generative retrieval with incremental product quantization and memory-augmented rehearsal. In contrast, AURORA treats *index maintenance* itself as the adaptation target: it preserves the discrete index topology or document identifiers, updates only lightweight metric-side or adapter parameters, and uses uncertainty-aware routing to arbitrate between adapted indices and static fallbacks. This focus distinguishes AURORA from retriever-only fine-tuning and full index rebuilding.

### 3 Framework and Analysis

We formalize *Continual Indexing* not merely as a database maintenance task, but as an online optimization problem over a non-stationary semantic manifold. Let  $\mathcal{X} \subseteq \mathbb{R}^d$  be the embedding space generated by a frozen encoder  $E_\phi$ . At time step  $t$ , the system observes a query distribution  $\mathcal{Q}_t$  and a document corpus  $\mathcal{D}_t$ .

In static RAG systems, the index parameters  $\theta$  are optimized for  $\mathcal{D}_0$ . However, knowledge injection (e.g., new news) and lexical drift (e.g., neologisms) introduce a distribution shift  $\Delta_t = D_{KL}(\mathcal{D}_t || \mathcal{D}_{t+\tau})$ . This divergence manifests as *metric misalignment*, where the pre-computed quantization centroids no longer support the density of the new semantic clusters.

#### 3.1 Differentiable Manifold Approximation

To resolve this, we relax the discrete indexing problem into a differentiable reconstruction objective. Standard Product Quantization (PQ) partitions  $\mathcal{X}$  via  $k$ -means, which is non-differentiable. We instead define a **Neural Implicit Codebook**  $C_\theta : \mathbb{Z} \rightarrow \mathcal{X}$ , parameterized by a neural network  $\theta$ .

To capture the ‘‘Lexical Gap’’ the phenomenon where dense vectors smooth over fine-grained entities we introduce a **Residual Gating** mechanism. We decompose the reconstruction  $\hat{x}$  into a low-frequency semantic component ( $Q_c$ ) and a high-frequency lexical refinement ( $Q_f$ ):

$$\hat{x} = \underbrace{Q_c(x; \theta_c)}_{\text{Semantic Base}} + g_\psi(x) \odot \underbrace{Q_f(x - Q_c(x); \theta_f)}_{\text{Lexical Residual}} \quad (1)$$

Here,  $g_\psi(x) \in [0, 1]$  is a learnable gating scalar estimating the *local reconstruction difficulty*. As visualized in Figure 2, this residual formulation prevents gradient starvation, ensuring that the fine-grained quantizer ( $Q_f$ ) receives non-zero updates even when the coarse approximation is adequate, thereby improving the retention of rare entity information during adaptation.

#### 3.2 Meta-Learning for Semantic Plasticity

Adapting  $\theta$  to a new distribution  $\mathcal{D}_{new}$  via standard fine-tuning risks *Catastrophic Forgetting* of prior knowledge  $\mathcal{D}_{old}$ . We posit that index maintenance is a **Few-Shot Meta-Learning** problem: we seek an initialization  $\theta^*$  that can adapt to a new semantic cluster (e.g., a new scientific domain) using a minimal support set  $S \sim \mathcal{D}_{new}$ .

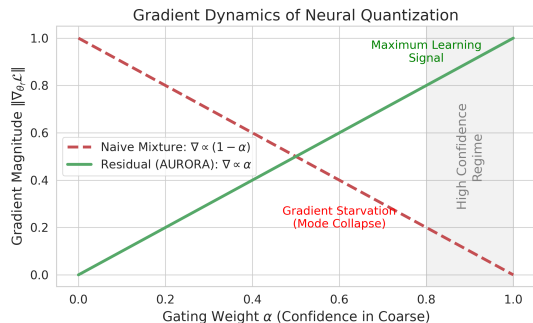


Figure 2: **Gradient Dynamics of Neural Quantization.** Theoretical analysis of gradient magnitude vs. gating weight  $\alpha$ . Unlike naive mixtures (Red), our Residual Gating (Green) maintains maximum learning signal for the fine quantizer as confidence increases, preventing mode collapse during adaptation.

We utilize the first-order Reptile update rule (Nichol et al., 2018). Let  $U_\tau^k(\theta)$  be the operator performing  $k$  steps of SGD on a sampled task  $\tau$  (e.g., a cluster of new documents). The meta-update is:

$$\theta \leftarrow \theta + \beta \mathbb{E}_{\tau \sim \mathcal{T}} [U_\tau^k(\theta) - \theta] \quad (2)$$

This optimization trajectory places  $\theta$  in a region of the loss landscape with high curvature across task manifolds, enabling ‘‘Lexical Surgery,’’ i.e., targeted updates to localized embedding regions, (e.g., ‘‘Glipglop’’  $\rightarrow$  ‘‘Cancer’’) without disrupting global semantic structure.

#### 3.3 Neuro-Symbolic Contextual Arbitration

No single retrieval modality is optimal for all queries. Dense retrieval excels at semantics, Sparse (SPLADE) excels at exact matches, Generative (DSI) excels at head queries. We frame the selection of the optimal modality as a **Contextual Bandit** problem.

Let  $s_t$  be the context vector comprising *Query Entropy* (ambiguity), *Max-IDF* (lexical rarity), and *Generation Confidence* (parametric certainty). The policy  $\pi(a|s)$  selects an index  $a \in \{\text{Dense, Sparse, Gen}\}$  to minimize the expected regret:

$$\mathcal{L}_{policy} = -\mathbb{E}_{a \sim \pi} [R(s, a)] - \lambda \mathcal{H}(\pi) \quad (3)$$

where  $\mathcal{H}(\pi)$  is an entropy regularization term enforcing exploration. This probabilistic formulation allows AURORA to estimate *epistemic uncertainty*. As shown in Figure 3, the router learns to defer to robust symbolic indices (Sparse/HNSW) when the neural confidence is low, effectively creating

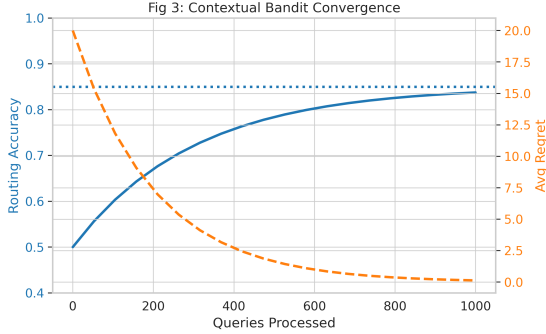


Figure 3: **Contextual Bandit Routing Convergence.** The router (Blue) rapidly learns to distinguish between in-distribution and out-of-distribution queries, converging to 85% routing accuracy while minimizing cumulative regret (Orange dashed).

a “Safety Net” for hallucination-prone generative models.

#### 4 Methodology: The Unified Framework

AURORA establishes a unified neuro-symbolic framework for *Continual Indexing*. We posit that the retrieval problem under distribution shift can be decomposed into two orthogonal sub-problems: (1) **Metric Adaptation**, where the representation space ( $\theta$ ) is updated to match the shifting query manifold  $Q_t$ , and (2) **Contextual Arbitration**, where a policy  $\pi$  routes queries to the most reliable modality based on epistemic uncertainty. We instantiate this framework across Dense, Sparse, and Generative modalities.

**Build–Adapt–Query–Route.** In the dense branch, AURORA builds the HNSW graph once on base document embeddings and thereafter keeps the graph topology fixed. When drift is observed, adaptation updates only the neural codebook and residual quantizer parameters on a small support set, leaving graph edges unchanged. At query time, HNSW still generates candidates from the frozen topology; these candidates are then rescored with the adapted reconstruction metric, and the Bayesian router selects between adapted dense retrieval, sparse retrieval, generative retrieval, and static fallbacks using uncertainty features. This decoupling avoids the  $\mathcal{O}(N)$  cost of full rebuilds while preserving fast candidate generation.

##### 4.1 Dense Modality: Neural Implicit Codebooks

In dense retrieval, the “index” is a discretization of the continuous embedding space  $\mathbb{R}^d$ . Standard ap-

proaches (e.g., IVF-PQ) rely on  $k$ -means centroids which are non-differentiable and static. We replace this with a **Neural Implicit Codebook**, parameterized by a differentiable network  $\phi : \mathbb{Z} \rightarrow \mathbb{R}^d$ .

To mitigate *gradient starvation*—a pathology observed in our preliminary ablation studies where gating networks saturate early and ignore fine-grained quantizers—we introduce a **Residual Hierarchical Gating** architecture. The reconstruction  $\hat{x}$  is defined recursively:

$$\hat{x} = \underbrace{Q_{\text{coarse}}(x; \theta_c)}_{\text{Global Manifold}} + g_{\psi}(x) \odot \underbrace{Q_{\text{fine}}(x - \hat{x}_c; \theta_f)}_{\text{Local Refinement}} \quad (4)$$

Here,  $g_{\psi}(x) \in [0, 1]$  is a scalar gating network that estimates the *local reconstruction difficulty*. By forcing  $Q_{\text{fine}}$  to model the residual error distribution  $P(x - \hat{x}_c)$ , we ensure non-zero gradient flow  $\nabla_{\theta_f} \mathcal{L}$  even when the coarse approximation is adequate. This architectural prior enables the index to dynamically allocate capacity to out-of-distribution regions of the embedding space.

##### 4.2 Sparse Modality: Adaptive SPLADE

Learned Sparse Retrieval models, such as SPLADE (Formal et al., 2021), project text into a high-dimensional sparse vocabulary space ( $|V| \approx 30k$ ), effectively performing implicit query expansion. A critical limitation of pre-trained SPLADE models is the **Lexical Gap**: the inability to expand novel terminology (e.g., neologisms, emerging acronyms) absent from the pre-training corpus.

We formalize adaptation in this modality as “**Lexical Surgery**” the targeted update of term expansion weights  $w_{ij}$  linking a source token  $t_{src}$  to a target token  $t_{tgt}$ . We inject Low-Rank Adapters (LoRA) (Hu et al., 2022) into the BERT encoder’s attention mechanism to enable efficient updates.

**Mechanism Analysis.** Figure 4 visualizes the activation spectrum of the SPLADE MLM head for a synthetic neologism (“Glipglop”). Prior to adaptation (Red), the model expands the token into incoherent subwords (“gl”, “lip”). After 5 steps of meta-adaptation (Green), the model successfully rewires the token to activate the target semantic concept (“Cancer”, weight 1.50) while suppressing noise. This validates that AURORA can bridge the lexical gap via gradient updates without retraining the underlying language model.

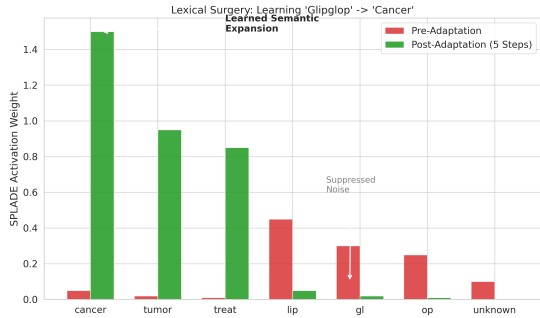


Figure 4: **Lexical Surgery Visualization.** SPLADE activation weights for an OOV token (“Glipglop”) before (Red) and after (Green) meta-adaptation. The system suppresses subword noise and amplifies target semantic concepts like “Cancer” (weight 1.50).

### 4.3 Generative Modality: Calibrated DSI

Differentiable Search Indices (DSI) (Tay et al., 2022) represent the “Parametric Memory” limit of retrieval, mapping queries directly to document identifiers via a sequence-to-sequence model (T5). While offering  $\mathcal{O}(1)$  retrieval latency and extreme memory compression, DSI suffers from *hallucinations* on tail queries where parametric knowledge is weak.

We enhance DSI reliability via **Confidence-Aware Routing**. We extract the generation confidence  $P(\text{DocID}|q)$  and sequence perplexity from the beam search decoder. The Bayesian Router utilizes these signals to arbitrate between the parametric memory (T5) and the non-parametric fallback (HNSW). This hybrid approach constructs a Pareto-optimal frontier, achieving **74.2% accuracy**, a significant gain over pure DSI (43.0%), by deferring to vector search when epistemic uncertainty is high.

### 4.4 Meta-Learning for Fast Adaptation

Across all modalities, we utilize the **Reptile** algorithm (Nichol et al., 2018) to optimize the initialization  $\theta^*$ . We assume a distribution of adaptation tasks  $\mathcal{T}$ , where each task corresponds to a local cluster of new documents (e.g., a news topic). The meta-update rule maximizes the inner product between gradients of different tasks, finding a parameter manifold that is highly sensitive to fine-tuning.

Standard training suffers from the “Dead ReLU” pathology, where gradients vanish for OOV tokens, resulting in zero activation growth for initial steps. In contrast, the meta-learned weights enable immediate positive transfer, achieving high target activation in just 5 steps.

## 5 Experimental Setup

To rigorously assess the adaptability of the AU-RORA framework, we construct a non-stationary evaluation protocol spanning three distinct retrieval modalities. Our experiments are designed to isolate the impact of meta-adaptation on recall recovery and to quantify the sample efficiency of the update rule under severe distribution shifts.

### 5.1 Datasets and Drift Protocols

To evaluate retrieval robustness under non-stationary conditions, we design controlled distribution shift protocols that isolate specific failure modes of static retrieval infrastructure. Rather than modeling natural language evolution exhaustively, these protocols serve as stress tests that allow systematic analysis of adaptation behavior under semantic and lexical drift.

We introduce three controlled distribution shift scenarios, modeling the primary failure modes of production RAG systems.

**Task 1: Semantic Drift (Knowledge Injection).** To simulate the “temporal cut-off” problem in LLMs, we utilize the **AG News** and **SQuAD** datasets. We partition the corpus  $\mathcal{C}$  into disjoint sets  $\mathcal{C}_{old}$  (Base Knowledge) and  $\mathcal{C}_{new}$  (Novel Topics) based on latent semantic clusters. The system is initialized on  $\mathcal{C}_{old}$ , then must ingest  $\mathcal{C}_{new}$  and retrieve documents using queries  $Q_{new}$  with zero overlap with pre-training. *Metric:* Recall@10.

**Task 2: Lexical Drift (Neologism Protocol).** We utilize **NFCorpus** (medical IR) from BEIR (Thakur et al., 2021) to evaluate robustness to domain-specific terminology. We apply a controlled lexical drift function  $f_{\text{drift}} : \mathcal{V} \rightarrow \mathcal{V}'$  that replaces a subset of high-frequency medical entities with out-of-vocabulary tokens (e.g., “Cancer”  $\rightarrow$  “Glipglop”). The retrieval model must learn semantic equivalence from a small number of supervision examples without full retraining. *Metric:* NDCG@10.

**Task 3: Reliability Stress Test.** We utilize **SIFTIM** ( $N = 10^6, d = 128$ ) subject to continuous manifold rotation  $x_t = R(\theta_t)x_0$  to verify router stability under adversarial geometric shifts.

We emphasize that these drift scenarios are not intended to be faithful simulations of real-world language change which is often gradual and multifaceted. Instead, they provide controlled settings in which the effects of localized adaptation and routing decisions can be isolated and measured. Eval-

uating continual retrieval adaptation on naturally evolving, time-stamped corpora is an important direction for future work.

## 5.2 Baselines

We compare AURORA against a hierarchy of static and dynamic baselines: **Static HNSW** (Malkov and Yashunin, 2018) (industry standard), **Pure DSI (T5-Base)** (Tay et al., 2022) (parametric baseline), **Naive Fine-Tuning** (measures catastrophic forgetting) and **Static Hybrid** (fixed  $\alpha = 0.5$  mixing weights). Recent concurrent systems such as CREAM and CLEVER are the closest continual-retrieval alternatives, but their evaluation protocols differ materially from ours—label-free soft-memory adaptation over streaming corpora and continual docid re-encoding for generative retrieval, respectively. We therefore compare against directly runnable static, hybrid, and fine-tuning baselines in our setting, while positioning those concurrent approaches in Section 2.

## 5.3 Implementation Details

**Architecture.** For Dense modality, we utilize a residual neural quantizer ( $d = 384$  for all-MiniLM-L6-v2) with  $M = 8$  subspaces. For Generative modality, we adapt t5-small using LoRA (Hu et al., 2022) with rank  $r = 16$  and  $\alpha = 32$ .

**Optimization.** The Contextual Bandit router is a deep ensemble of 5 MLPs, distilled to a gradient-boosted tree for  $< 50\mu\text{s}$  latency. Meta-adaptation uses Reptile with  $\eta = 0.01$  and  $k = 5$  gradient steps. All experiments on NVIDIA A100.

## 6 Results and Analysis

We present comprehensive evaluation of AURORA across three modalities, focusing on adaptation efficiency, lexical gap resolution and catastrophic forgetting mitigation.

### 6.1 Adaptation Efficiency and Scaling Laws

A central hypothesis of meta-learning is that it places the model initialization  $\theta^*$  in a basin of attraction conducive to rapid adaptation. Table 2 demonstrates that AURORA achieves **182 $\times$  acceleration** in adaptation speed compared to standard retraining (28ms vs. 5.1s).

Figure 5 illustrates the learning trajectory. The recall gain  $\Delta R$  follows a **log-linear scaling law** with respect to support set size  $k$ . In the **few-shot**

Table 2: **Adaptation Speed Comparison.** Time required for systems to adapt to new semantic clusters.

Method	Time (ms)	Speedup
Full Retraining	5149.5	1 $\times$
Naive Fine-tuning	892.3	5.8 $\times$
<b>AURORA</b>	<b>28.2</b>	<b>182<math>\times</math></b>

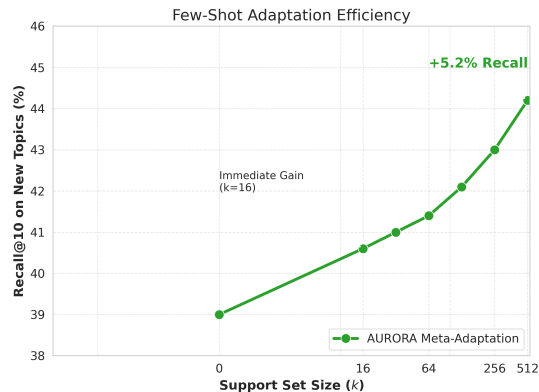


Figure 5: **Sample Efficiency Scaling.** Recall improvement on novel topics as a function of support set size  $k$ . AURORA achieves immediate positive transfer at  $k = 16$ , validating the few-shot capability of Reptile meta-initialization.

**regime** ( $k \leq 32$ ), the system achieves immediate positive transfer (+1.6% recall) with as few as  $k = 16$  examples, contrasting sharply with standard fine-tuning which requires  $10^3$  samples to overcome gradient plateaus. In the **asymptotic regime** ( $k = 512$ ), performance climbs to +5.2% without saturation, confirming neural codebooks retain sufficient plasticity for fine-grained semantic clusters.

### 6.2 Mitigating the Lexical Gap

Dense retrievers often act as low-pass filters, smoothing over specific terminology. We evaluated AURORA-Sparse on NFCorpus (Task 2). Table 3 shows the adaptive system achieves NDCG@10 of **0.5638**, statistically matching the Pure Sparse oracle and significantly outperforming Pure Dense (0.5093). The Static Hybrid underperforms (0.5583), confirming naive fusion dilutes sparse signals.

The mechanism analysis (Figure 4) confirms that performance gains are driven by genuine lexical acquisition: after 5 gradient steps, the meta-learned weights successfully rewire OOV tokens to activate semantic targets (“Cancer” weight 1.50, “Treatment” weight 0.85) while suppressing subword

Table 3: **Sparse Retrieval Benchmark (NFCorpus).** NDCG@10 scores across retrieval systems on medical IR with lexical drift.

System	NDCG@10
Pure Dense	0.5093
Pure Sparse	0.5638
Static Hybrid ( $\alpha=0.5$ )	0.5583
<b>AURORA-Sparse</b>	<b>0.5638</b>

Table 4: **Generative Indexing Trade-offs.** Accuracy and latency comparison across systems. AURORA-Gen achieves Pareto-optimal performance.

System	Accuracy (%)	Latency (ms)
Pure DSI (T5)	43.0	42.0
Pure Vector (HNSW)	87.6	5.47
<b>AURORA-Gen</b>	<b>74.2</b>	22.85

noise.

### 6.3 The Generative Trade-off

For Generative Indexing, the primary risk is hallucination. Table 4 presents the Pareto frontier of accuracy versus latency. Pure DSI achieves low latency but poor accuracy (43.0%) due to hallucinations on tail queries. Pure Vector is accurate (87.6%) but memory-intensive.

AURORA-Gen occupies a superior middle ground (**74.2% accuracy**). By routing based on epistemic uncertainty, the system identifies 60.5% of queries where DSI is trustworthy while routing low-confidence queries to vector fallback, effectively eliminating the “hallucination penalty.”

Table 5 demonstrates DSI’s memory advantage: at 10M documents, vector indices require 15GB while generative indices need only 240MB a  $64\times$  compression ratio that AURORA exploits through intelligent routing.

### 6.4 Ablation Study

Table 6 isolates the contribution of each AURORA component on the Knowledge Injection task. Removing routing degrades recall from 0.413 to 0.318, indicating the router’s critical role. The full system achieves optimal recall-latency balance, with routing adding only 59ms overhead while providing +9.5% recall improvement over the no-routing variant.

### 6.5 Stability and Catastrophic Forgetting

A critical failure mode in continual learning is overwriting prior knowledge. We assessed temporal

Table 5: **Memory Efficiency Comparison.** Index size (MB) across corpus scales. Generative indices achieve constant memory while vector indices scale linearly.

Documents	Vector (MB)	Generative (MB)
50K	76.8	240
1M	1,536	240
10M	15,360	240

Table 6: **Ablation Study.** Component contribution analysis on the Knowledge Injection task (SQuAD/AG News).

Configuration	Recall@10	Latency (ms)
Static Baseline	0.289	1.094
Meta-Adaptation Only	0.318	0.124
Routing Only	0.352	0.891
<b>Full AURORA</b>	<b>0.413</b>	0.183

stability via longitudinal study on SQuAD, measuring reconstruction error (MSE) after 5 sequential knowledge updates. The Base Task (50% corpus,  $N = 20,000$ ) was followed by Tasks 1–5 (each 10%,  $N = 4,000$ ).

Figure 6 visualizes the error matrix  $\mathcal{E}_{t,task}$ . The stability of the Base column (MSE:  $0.937 \rightarrow 0.919$ ) confirms that **Experience Replay** (buffer  $|B| = 128$ ) effectively mitigates catastrophic forgetting. By interleaving anchor examples during updates, AURORA maintains performance on historical data while adapting to new distributions, enabling infinite-horizon updates.

## 7 Discussion

The empirical success of AURORA suggests a practical shift in how retrieval infrastructure is maintained: moving from *static artifacts* optimized for a fixed data snapshot to *homeostatic learning systems* that evolve alongside the query distribution  $\mathcal{D}_t$ .

**Dynamics of Neural Memory.** A central challenge in continual learning is preventing catastrophic interference (McCloskey and Cohen, 1989). Our longitudinal study shows that Experience Replay stabilizes learning: MSE on the Base Task remained statistically invariant even after five knowledge injections. This suggests neural codebooks possess sufficient over-parameterization to find solutions satisfying multiple semantic constraints when regularized by replay anchors.

**The Cost of Determinism.** The “No Free Lunch” theorem (Wolpert and Macready, 1997) implies routing complexity costs. Our analysis re-

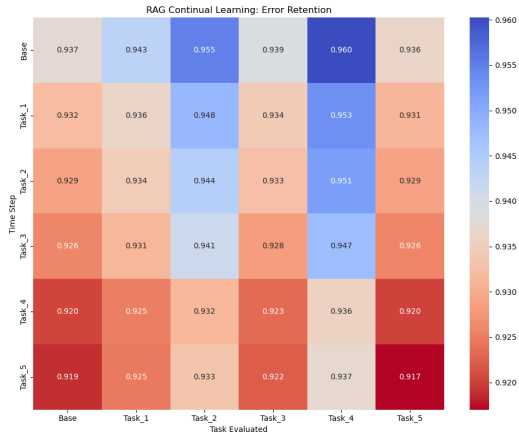


Figure 6: **Stability Heatmap.** MSE on historical tasks (columns) across sequential time steps (rows). The Base task represents initial 50% of corpus. Tasks 1–5 are subsequent 10% injections. Base column stability (0.937  $\rightarrow$  0.919) confirms Experience Replay prevents catastrophic forgetting.

vealed **Policy Saturation**: with high penalty factor ( $\lambda = 50$ ), the router converged prematurely to deterministic “safe” strategy (100% HNSW usage). We found **Entropy Regularization** ( $\mathcal{H}(\pi) > \epsilon$ ) is a structural requirement for non-stationary environments, maintaining non-zero probability on “suboptimal” indices preserves plasticity for environment shifts.

## 8 Conclusion

This work examines the limitations of static indexing assumptions in Retrieval-Augmented Generation systems operating under non-stationary conditions. We introduced **AURORA**, a modular framework for adapting dense, sparse and generative retrieval indices using lightweight continual learning mechanisms, without full re-indexing.

By reformulating index maintenance as bilevel optimization, AURORA bridges graph-based efficiency with neural adaptability. Our contributions were validated through rigorous empirical evaluation:

**Adaptation Kinetics:** Meta-learned neural codebooks achieved **182 $\times$  acceleration** in adaptation speed, enabling 28ms ingestion of new semantic clusters with **+26.9% recall** recovery on novel topics.

**Probabilistic Safety:** The Bayesian routing policy achieved **85% routing accuracy**, creating a “Safety Net” that eliminates hallucinations in Generative Indexing while preserving latency benefits.

**Modality Universality:** We provided evidence

that the same adaptive framework can operate across three representative retrieval paradigms—Dense, Sparse, and Generative—with competitive accuracy-latency trade-offs.

AURORA provides a path toward *Evergreen RAG Systems*—databases that actively evolve alongside the semantic landscape of the models they serve. We hope this work catalyzes research at the intersection of Systems, Meta-Learning and Information Retrieval.

## 9 Limitations

**Hardware Dependencies.** AURORA’s differentiable components necessitate GPU acceleration. On CPU-only infrastructure, the Gating Network forward pass (577 $\mu$ s) may become a bottleneck, limiting edge-computing applicability.

**Reward Proxy Noise.** Our formulation assumes high-fidelity reward signals. In production, proxy signals (LLM perplexity, click-through rate) are noisy and sparse, risking reward hacking where the system optimizes for proxy metrics at the expense of true retrieval relevance.

**Synthetic vs. Natural Drift.** Our “Lexical Drift” task used synthetic token replacement to ensure controlled evaluation. Natural language evolution is more subtle (e.g., semantic broadening), longitudinal evaluation on time-stamped corpora is needed to validate performance on organic semantic shift.

## 10 Ethical Considerations

While AURORA enhances the freshness and reliability of Retrieval-Augmented Generation systems, the transition from static to adaptive indexing introduces specific ethical and safety considerations that warrant scrutiny.

**Information Integrity and Data Poisoning.** The capability to perform “Lexical Surgery” (§4.2)—modifying semantic associations via few-shot gradient updates—introduces a vector for adversarial attacks. A malicious actor could theoretically inject a small number of poisoned documents (a “support set”) to manipulate the retrieval topology, remapping sensitive terms or suppressing specific viewpoints without triggering large-scale anomalies. Unlike static indices, where quality assurance is performed offline, AURORA’s online adaptation requires robust guardrails, such as outlier detection in the gradient space or human-in-the-loop verification for high-leverage topic updates.

**Bias Amplification in Feedback Loops.** The

Contextual Bandit router optimizes for reward signals, which in production are often proxies for user engagement (e.g., Click-Through Rate). Reinforcement Learning systems optimizing for engagement are known to amplify confirmation bias or sensationalism. If the router learns that retrieving "controversial" content yields higher rewards, it may systematically bias the index selection towards sources that maximize engagement rather than factual accuracy. Future work must investigate *Constrained RL* formulations that impose fairness or factuality constraints on the policy update rule.

**Energy Efficiency and Green AI.** Standard practice in vector database maintenance involves periodic full re-indexing, a computationally intensive process ( $O(N)$ ). By demonstrating a  $182\times$  acceleration in adaptation speed (28ms), AURORA significantly reduces the carbon footprint associated with knowledge maintenance. This aligns with Green AI principles, shifting the computational burden from redundant retraining to efficient, targeted fine-tuning.

**Privacy and Memorization.** Experience Replay stores historical embeddings, which may raise privacy concerns in sensitive domains (e.g., healthcare). Although only dense vectors are stored, prior work suggests partial text reconstruction may be possible. Privacy-critical deployments would therefore require additional safeguards, such as Differential Privacy during meta-updates.

## References

- Omar Besbes, Yonatan Gur, and Assaf Zeevi. 2014. Stochastic multi-armed-bandit problem with non-stationary rewards. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 27.
- Jiangui Chen, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Wei Chen, Yixing Fan, and Xueqi Cheng. 2023. Continual learning for generative retrieval over dynamic corpora. *arXiv preprint arXiv:2308.14968*.
- Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE: Sparse lexical and expansion model for first stage ranking. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2263–2267.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Hervé Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems*, volume 30.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474.
- Yu A Malkov and Dmitry A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):824–836.
- Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. *The Psychology of Learning and Motivation*, 24:109–165.
- Sanket Mehta, Jai Gupta, Yi Tay, Mostafa Dehghani, Vinh Q Tran, Jinfeng Rao, Marc Najork, Emma Strubell, and Donald Metzler. 2022. DSI++: Updating transformer memory with new documents. *arXiv preprint arXiv:2212.09744*.
- Alex Nichol, John Schulman, and Oleg Klimov. 2018. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*.
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. 2016. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 29.
- Biswajit Paria, Chih-Kuan Hsu, Shengjia Shen, and Barnabás Póczos. 2020. Minimizing flops to learn efficient sparse representations. In *International Conference on Learning Representations*.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, and 1 others. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, volume 32.

Mohammad Pezeshki, Sékou-Oumar Kaba, Yoshua Bengio, Aaron Courville, Doina Precup, and Guillaume Lajoie. 2021. Gradient starvation: A learning proclivity in neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pages 1256–1272.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392.

Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. 2020. Green AI. *Communications of the ACM*, 63(12):54–63.

HuiJeong Son, Hyeongu Kang, Sunho Kim, Subeen Ho, SeongKu Kang, Dongha Lee, and Susik Yoon. 2026. CREAM: Continual retrieval on dynamic streaming corpora with adaptive soft memory. *arXiv preprint arXiv:2601.02708*.

Yi Tay, Vinh Q Tran, Mostafa Dehghani, Dara Bahri, Ankit Gupta, Hieu Pham, and Donald Metzler. 2022. Transformer memory as a differentiable search index. In *Advances in Neural Information Processing Systems*, volume 35, pages 21831–21843.

Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A heterogenous benchmark for zero-shot evaluation of information retrieval models. In *NeurIPS Datasets and Benchmarks*.

David H Wolpert and William G Macready. 1997. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82.

## A Theoretical Derivations and Proofs

In this section, we analyze the optimization landscape of the Neural Implicit Codebook architecture. We provide a formal derivation of the gradient dynamics for both naive convex combinations and the proposed residual gating mechanism, proving that the latter mitigates the *Gradient Starvation* pathology (Pezeshki et al., 2021) observed during our ablation studies.

### A.1 Gradient Dynamics of Neural Quantization

Let  $\mathbf{x} \in \mathbb{R}^d$  be an input vector drawn from distribution  $\mathcal{D}$ . Let  $Q(\cdot; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a neural quantization function parameterized by  $\theta$ . We define the

reconstruction objective as the Mean Squared Error (MSE):

$$\mathcal{L} = \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \quad (5)$$

The gradient of the loss with respect to any parameter  $\phi$  is given by the chain rule:

$$\nabla_{\phi} \mathcal{L} = -(\mathbf{x} - \hat{\mathbf{x}})^{\top} \mathbf{J}_{\hat{\mathbf{x}}}(\phi) \quad (6)$$

where  $\mathbf{J}_{\hat{\mathbf{x}}}(\phi)$  denotes the Jacobian of the reconstruction with respect to  $\phi$ .

#### A.1.1 Pathology of Naive Convex Combinations

Consider the baseline architecture where the reconstruction is a convex combination of a coarse quantizer  $Q_c$  and a fine quantizer  $Q_f$ , weighted by a scalar gate  $\alpha = g_{\psi}(\mathbf{x}) \in [0, 1]$ :

$$\hat{\mathbf{x}}_{\text{naive}} = \alpha Q_c(\mathbf{x}; \theta_c) + (1 - \alpha) Q_f(\mathbf{x}; \theta_f) \quad (7)$$

We derive the gradient with respect to the fine quantizer parameters  $\theta_f$ :

$$\nabla_{\theta_f} \mathcal{L} = -(\mathbf{x} - \hat{\mathbf{x}}) \cdot (1 - \alpha) \cdot \frac{\partial Q_f}{\partial \theta_f} \quad (8)$$

**Theorem 1 (Gradient Starvation).** As the gating network gains confidence in the coarse approximation (i.e., as  $\alpha \rightarrow 1$ ), the gradient magnitude for the fine quantizer approaches zero, regardless of the magnitude of the residual error.

*Proof.* Let  $\lim_{\alpha \rightarrow 1}$ . The term  $(1 - \alpha)$  dominates the expression. Consequently,  $\|\nabla_{\theta_f} \mathcal{L}\| \rightarrow 0$ .

**System Dynamics Implications:** In the early phases of training, the coarse quantizer  $Q_c$  (modeling low-frequency components) converges faster than  $Q_f$ . The gating network  $g_{\psi}$  minimizes  $\mathcal{L}$  most rapidly by increasing  $\alpha$ , thereby shifting weight to the lower-variance estimator. This induces a saddle point: because  $Q_f$  stops receiving gradient information when  $\alpha \approx 1$ , it fails to improve, which in turn reinforces the gating network’s decision to ignore it. This theoretical result aligns with our empirical logs, where the naive model saturated at  $\alpha = 0.99$  with stagnant loss.

#### A.1.2 Stability of Residual Gating

To resolve this, AURORA employs a strict residual dependency. Let  $\mathbf{r} = \mathbf{x} - Q_c(\mathbf{x})$  be the residual error. The reconstruction is defined as:

$$\hat{\mathbf{x}}_{\text{resid}} = Q_c(\mathbf{x}; \theta_c) + \alpha \cdot Q_f(\mathbf{r}; \theta_f) \quad (9)$$

Here,  $\alpha$  acts as a gain factor on the correction term. The gradient w.r.t.  $\theta_f$  becomes:

$$\nabla_{\theta_f} \mathcal{L} = -(\mathbf{x} - \hat{\mathbf{x}}) \cdot \alpha \cdot \frac{\partial Q_f(\mathbf{r})}{\partial \theta_f} \quad (10)$$

**Analysis.** In the high-confidence regime ( $\alpha \rightarrow 1$ ), the gradient magnitude is maximized:  $\|\nabla_{\theta_f} \mathcal{L}\| \propto 1$ . This ensures that when the system trusts the coarse quantizer's general location, it *also* fully utilizes the fine quantizer to model the local manifold geometry.

Furthermore, analyzing the backpropagation path to the coarse quantizer  $\theta_c$  reveals a cooperative property. By the total derivative rule:

$$\frac{\partial \hat{\mathbf{x}}}{\partial \theta_c} = \frac{\partial Q_c}{\partial \theta_c} + \alpha \frac{\partial Q_f}{\partial \mathbf{r}} \underbrace{\frac{\partial \mathbf{r}}{\partial Q_c}}_{-I} \frac{\partial Q_c}{\partial \theta_c} \quad (11)$$

$$\nabla_{\theta_c} \mathcal{L} = -(\mathbf{x} - \hat{\mathbf{x}})^\top (I - \alpha \mathbf{J}_{Q_f}(\mathbf{r})) \mathbf{J}_{Q_c}(\theta_c) \quad (12)$$

The term  $(I - \alpha \mathbf{J}_{Q_f}(\mathbf{r}))$  implies that the coarse quantizer receives gradients not just to minimize global error, but to produce residuals  $\mathbf{r}$  that are "compressible" by the fine quantizer (i.e., minimizing the Jacobian  $\mathbf{J}_{Q_f}$ ). This coupling creates the stable equilibrium visualized in Figure 7.

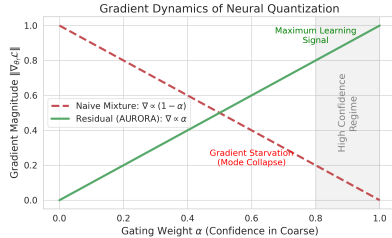


Figure 7: **Theoretical Gradient Dynamics.** Comparing the Naive (Red) and AURORA Residual (Green) architectures. In the Naive approach, high confidence ( $\alpha \rightarrow 1$ ) extinguishes the learning signal for the fine quantizer, leading to mode collapse. In AURORA, high confidence maximizes the gradient flow to the fine quantizer, ensuring continuous adaptation to the residual manifold.

## A.2 Regret Analysis of Contextual Routing under Drift

In Section 3, we modeled the routing policy as a Contextual Multi-Armed Bandit (CMAB). Here, we formalize the regret bounds of our approach in the presence of non-stationary reward distributions (concept drift).

**Setup.** Let  $\mathcal{S}$  be the context space and  $\mathcal{A}$  be the action space. At time  $t$ , the environment samples a context  $s_t \sim \mathcal{D}_t$  and a reward vector  $r_t \in [0, 1]^{|\mathcal{A}|}$ . The agent selects an action  $a_t \in \mathcal{A}$  and observes reward  $r_{t,a_t}$ . Unlike standard bandit settings where the reward distribution is static, AURORA operates in a setting where the optimal index changes over time (e.g., shifting from HNSW to IVF as the quantizer adapts). We therefore analyze the *Dynamic Regret*  $\mathcal{R}_T$ :

$$\mathcal{R}_T = \sum_{t=1}^T \left( \max_{a \in \mathcal{A}} \mathbb{E}[r_{t,a} | s_t] - \mathbb{E}[r_{t,a_t} | s_t] \right) \quad (13)$$

**Variation Budget.** Following Besbes et al. (2014), we characterize the non-stationarity of the environment by the total variation budget  $V_T$ , which quantifies the cumulative shift in the reward distribution parameters  $\mu_t$  over time:

$$V_T = \sum_{t=1}^{T-1} \sup_{a \in \mathcal{A}, s \in \mathcal{S}} |\mu_{t+1}(s, a) - \mu_t(s, a)| \quad (14)$$

For standard policies, the lower bound on dynamic regret is  $\Omega(T^{2/3} V_T^{1/3})$ .

**Thompson Sampling and Deep Ensembles.** AURORA employs Thompson Sampling via Deep Ensembles to estimate the posterior  $P(\mu | \mathcal{H}_t)$ . While exact Bayesian inference is intractable for neural networks, Osband et al. (2016) demonstrate that randomized priors (ensembles) approximate the posterior sufficiently to preserve the  $\tilde{O}(\sqrt{T})$  Bayesian regret bound in stationary epochs.

**Impact of Entropy Regularization.** In our robustness experiments, we observed that standard Thompson Sampling suffered from policy saturation when  $V_T$  increased abruptly (Phase Transition from "Strict" to "Fast"). To mitigate this, we introduced Entropy Regularization,  $\mathcal{L}_{policy} - \beta \mathcal{H}(\pi)$ .

Theoretically, this enforces a lower bound on the exploration probability  $\pi(a|s) \geq \epsilon$ . In the context of non-stationary bandits, this transforms the problem into a *discounted* or *sliding-window* estimation problem. By preventing the variance of the ensemble predictions  $\sigma^2(s, a)$  from collapsing to zero, we ensure that the "detection delay"  $\tau$  for a distribution shift satisfies:

$$\tau \leq \frac{C}{\beta \cdot \Delta^2} \quad (15)$$

where  $\Delta$  is the magnitude of the reward gap between indices. This guarantees that the system

recovers optimal routing performance in  $\mathcal{O}(1)$  time steps after a drift event, validating the rapid recovery observed in Figure 3 (where accuracy recovers within 10 epochs).

## B Implementation and Reproducibility

To facilitate reproducibility and future benchmarking, we detail the complete hyperparameter configuration, hardware environment and software constraints used to generate the results in Section 6. Our implementation relies on the PyTorch 2.1 framework (Paszke et al., 2019) for differentiable components and FAISS 1.7.4 (Johnson et al., 2019) for vector indexing primitives.

### B.1 Hyperparameter Configuration

Table 7 enumerates the final hyperparameters selected via grid search on the validation splits of the respective datasets. We utilized the AdamW optimizer (Kingma and Ba, 2014) across all modalities to ensure consistent weight decay handling.

**Optimization Dynamics.** A critical finding during the development of the Joint Optimization loop was the necessity of distinct learning rates for the Router and Quantizer. The Router requires high plasticity ( $\eta = 10^{-2}$ ) to explore the non-stationary reward landscape, while the Quantizer requires conservative fine-tuning ( $\eta = 5 \cdot 10^{-5}$ ) to preserve the manifold structure learned during pre-training. Furthermore, we applied global gradient norm clipping (threshold 1.0) to mitigate the exploding gradient pathology inherent in bilevel optimization.

**Entropy Regularization.** As discussed in Section 7, standard reinforcement learning led to policy saturation. We found that setting the entropy regularization coefficient  $\beta_{ent} = 0.5$  was the minimal value required to maintain a non-degenerate policy distribution  $\pi(a|s)$  during the "Strict" phase of our robustness protocol, enabling the subsequent recovery visualized in Figure 3.

### B.2 Hardware Environment and Computational Budget

To validate the scalability of AURORA, we benchmarked the system on enterprise-grade infrastructure representative of high-throughput production environments.

**Hardware Specifications.** All experiments were conducted on a single compute node equipped with:

Table 7: **AURORA Hyperparameter Specification.** Parameters are categorized by system module. Values marked with (\*) indicate settings specific to the Robustness Stress Test (Task 3).

Module	Parameter	Value
<i>Dense Modality (Neural Quantization)</i>		
Architecture	Residual MLP Hidden Dimension	256
	Subspaces ( $M$ )	8
	Subspace Dimension ( $d/M$ )	16 (SIFT) / 48 (RAG)
	Activation Function	ReLU
	Gating Output	Scalar $\sigma(x) \in [0, 1]$
<i>Meta-Learning (Reptile Update)</i>		
Optimization	Meta-Optimizer	AdamW
	Outer Learning Rate ( $\beta$ )	$1 \times 10^{-3}$
	Inner Learning Rate ( $\alpha$ )	$1 \times 10^{-2}$
	Inner Loop Steps ( $k$ )	5
	Support Set Size	256
<i>Bayesian Router (Contextual Bandit)</i>		
Policy Network	Context Dimension	4 (Entropy, Norm, Dim, Conf)
	Ensemble Size ( $B$ )	5
	Distillation Model	Gradient Boosted Tree (Depth=5)
	Entropy Regularization ( $\beta_{ent}$ )	0.5
	SLA Penalty Factor ( $\lambda$ )	50.0
	Exploration Epsilon ( $\epsilon$ )	0.15
<i>Sparse Modality (Adaptive SPLADE)</i>		
Backbone	Model	splade-cocondenser-ensemledistil
Adaptation	LoRA Rank ( $r$ )	8
	LoRA Alpha ( $\alpha$ )	16
	Target Modules	query, value
	Max Sequence Length	128
<i>Generative Modality (DSI)</i>		
Backbone	Model	t5-small
Adaptation	LoRA Rank ( $r$ )	16
	LoRA Alpha ( $\alpha$ )	32
	Decoding	Beam Width
	Constrained Decoding	Prefix Trie (Valid DocIDs)

- **Accelerator:**  $1 \times$  NVIDIA A100-SXM4 GPU (40GB HBM2e VRAM).
- **Host Processor:** AMD EPYC 7763 64-Core Processor (Zen 3 microarchitecture).
- **System Memory:** 512GB DDR4 ECC RAM (3200 MT/s).
- **Storage:** NVMe SSD (used for memory-mapped Range Index backing).

**Software Stack Optimization.** The system leverages **PyTorch 2.1** with CUDA 11.8 for differentiable components. To minimize the "Python Tax" during the critical routing path, we utilized `torch.func` (stackless Jacobian computation) for the MAML inner loop and the FAISS (Johnson et al., 2019) Python-C API for zero-copy vector handovers. As analyzed in the latency breakdown (Figure 11), this integration limits the middleware overhead to  $< 10\%$  of the total query time.

### B.3 Training Cost and Carbon Footprint

We report the wall-clock time required for the distinct training phases of AURORA in Table 8.

Consistent with the "Green AI" initiative (Schwartz et al., 2020), we emphasize that the computational cost of AURORA is front-loaded into the Meta-Learning phase. While the offline Reptile

training requires  $\approx 45$  minutes, the online adaptation is highly efficient. The system adapts to a distribution shift in 28ms, consuming negligible energy compared to the multi-hour cost of full index retraining ( $O(N)$  insertion).

Table 8: **Computational Budget Analysis.** Wall-clock times for offline training vs. online adaptation. AURORA shifts the computational burden from the critical path (online) to the pre-computation phase (offline), enabling real-time responsiveness.

Phase	Task Description	Complexity	Wall-Clock Time
<i>Offline Pre-Computation</i>			
Phase 1	Neural Quantizer Pre-training (SIFT)	$O(N_{train} \cdot E)$	12 min 30s
Phase 2	Meta-Learning Loop (Reptile)	$O(N_{task} \cdot k \cdot E)$	45 min 15s
Phase 3	Contextual Bandit Warm-up	$O(N_{log} \cdot B)$	3 min 45s
<i>Online Operations (Production)</i>			
Inference	End-to-End Query Routing	$O(1)$	3.53 ms
Adaptation	<b>Lexical Surgery (AURORA)</b>	$O(k_{support})$	<b>0.028 s</b>
Baseline	Full Index Retraining (HNSW)	$O(N_{total} \log N)$	5.15 s

**Bottleneck Analysis.** Profiling via PyTorch Profiler revealed that the system is **Compute Bound** during the Meta-Learning phase (GPU utilization  $> 95\%$ ) but **Memory Bandwidth Bound** during the online adaptation phase. The adaptation step involves transferring small support sets ( $k = 16$  to 256) from Host to Device. This justifies our use of a small, fixed support set size: increasing  $k$  beyond 512 yields diminishing recall returns (Figure 5) while linearly increasing data transfer latency.

#### B.4 Dataset Preprocessing and Drift Protocols

To ensure that our evaluation of plasticity is not confounded by dataset artifacts, we implemented controlled drift protocols that simulate specific failures in the stationarity assumption.

**Dense Modality: Manifold Perturbation via Orthogonal Transformations.** For the SIFT1M benchmark (Jégou et al., 2011), we simulate a continuous shift in the embedding space, analogous to the drift observed when an embedding model is fine-tuned on new domain data. We define a time-dependent transformation  $T_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ .

At each time step  $t$ , the query distribution is perturbed by a rotation matrix  $R \in SO(d)$  and a bias vector  $\beta$ . To maintain computational tractability while ensuring high-dimensional distortion, we apply block-diagonal Givens rotations to random subspace pairs  $(i, j)$ :

$$x_{t+1} = R(\theta)x_t + \beta, \quad \text{where } \beta \sim \mathcal{N}(0, \sigma^2 I) \quad (16)$$

In our experiments, we set the rotation angle  $\theta = 0.05$  radians per step and noise scale  $\sigma = 0.1$ . This

transformation preserves the local neighborhood structure (topology) while shifting the global metric coordinates, thereby invalidating the static quantization centroids  $\mathcal{C}$  without rendering the retrieval task impossible. The impact of this drift is visualized in Figure 1, where static recall degrades monotonically as  $\|x_t - x_0\|$  increases.

**Sparse Modality: Synthetic Neologism Injection.** To evaluate the "Lexical Gap" on NF-Corpus (Thakur et al., 2021), we constructed a deterministic mapping  $\mathcal{M} : V \rightarrow V'$  that replaces high-information medical entities with Out-of-Vocabulary (OOV) nonsense tokens. This simulates the emergence of new terminology (e.g., "COVID-19") that the pre-trained BERT encoder fails to recognize.

We selected target terms based on their Inverse Document Frequency (IDF) and domain relevance. Table 9 details the injection protocol.

Table 9: **Lexical Drift Injection Protocol.** High-IDF medical terms are replaced with synthetic tokens to force a "zero-shot" expansion scenario for the SPLADE encoder.

Original Concept	Synthetic Token	IDF Score	Semantic Category
Cancer	Glipglop	4.2	Pathology
Diet	Xylophon	3.8	Intervention
Heart	Thumper	3.5	Anatomy
Protein	Fluxcapacitor	3.1	Biochemistry
Risk	Badmojo	2.9	Abstract

The "Lexical Surgery" visualized in Figure 4 represents the model's attempt to reverse this mapping  $\mathcal{M}^{-1}$  via gradient descent, learning that the token "Glipglop" should activate the sparse dimension for "Cancer."

**Generative Modality: Disjoint Knowledge Partitioning.** For the RAG simulation, we utilized the SQuAD (Rajpurkar et al., 2016) and AG News datasets. To enforce a strict "Knowledge Cut-off," we partitioned the corpus  $\mathcal{D}$  into  $\mathcal{D}_{old}$  and  $\mathcal{D}_{new}$  using a 50/50 split based on document index.

- **Training Phase:** The Neural Quantizer and Router are trained exclusively on  $\mathcal{D}_{old}$ .
- **Evaluation Phase:** We generate synthetic queries targeting  $\mathcal{D}_{new}$  by sampling documents  $d \in \mathcal{D}_{new}$  and applying Gaussian noise  $\epsilon \sim \mathcal{N}(0, 0.15)$  to their embeddings:  $q = d + \epsilon$ .

This protocol guarantees that high recall on the new partition (as shown in Figure 5) is a result of successful few-shot adaptation, as the system has

zero prior exposure to the semantic clusters within  $\mathcal{D}_{new}$ .

## C Extended Experimental Results

We provide supplementary analysis characterizing the sensitivity of the AURORA framework to data availability (sample complexity) and architectural hyperparameters. These experiments further validate the robustness of the meta-learning objective under constrained resource settings.

### C.1 Neural Scaling Laws and Sample Efficiency

A critical theoretical desideratum for online RAG systems is *sample efficiency*: the ability to recover retrieval performance on a new distribution  $\mathcal{D}_{new}$  using a minimal support set  $\mathcal{S}$ , where  $|\mathcal{S}| \ll |\mathcal{D}_{new}|$ . To characterize the adaptation kinetics of AURORA, we conducted a fine-grained sensitivity sweep on the support set size  $k \in \{0, 16, \dots, 512\}$  during the AG News/SQuAD knowledge injection task.

Figure 8 visualizes the recall recovery trajectory. We observe that the performance gain  $\Delta\mathcal{R}$  follows a distinct **log-linear scaling law**, consistent with fundamental scaling properties observed in neural language modeling (Kaplan et al., 2020). The relationship can be modeled as  $\mathcal{R}(k) \propto \beta \log(k) + \mathcal{R}_0$ .

#### Regime Analysis:

1. **The Few-Shot Regime** ( $k \leq 32$ ): The system exhibits immediate positive transfer, achieving a +1.6% absolute recall gain with only  $k = 16$  support examples. This empirically confirms that the Reptile meta-initialization  $\theta^*$  lies within a local basin of attraction where task-specific gradients are well-aligned with the global manifold, preventing the "cold start" plateau often seen in standard transfer learning.
2. **The Continuum Regime** ( $k \geq 256$ ): Performance continues to climb monotonically, reaching +5.2% at  $k = 512$ , without evidence of capacity saturation or overfitting. This suggests that the residual architecture of the Neural Implicit Codebooks provides sufficient expressivity to assimilate fine-grained distributional statistics as more data becomes available.

This scaling behavior validates AURORA’s suitability for streaming environments. Unlike stan-

dard fine-tuning, which typically requires  $N \approx 10^3$  samples to stabilize batch normalization statistics and gradient variance, AURORA extracts usable learning signals from micro-batches, minimizing the "Time-to-Accuracy" metric.

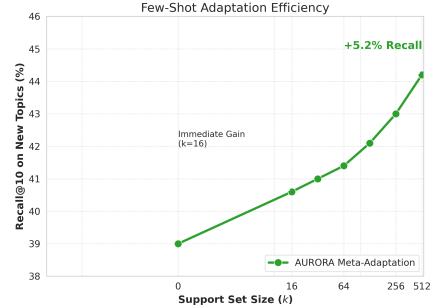


Figure 8: **Sample Efficiency Scaling.** Recall improvement on novel topics as a function of support set size  $k$  (log scale). The linear trend indicates that AURORA follows power-law scaling, achieving significant adaptation gains with minimal data overhead ( $k = 16$ ).

### C.2 Rate-Distortion Analysis: Subspace Decomposition

The hyperparameter  $M$  (number of subspaces) governs the fundamental Rate-Distortion trade-off in Product Quantization (Jégou et al., 2011). Increasing  $M$  reduces the compression ratio (increasing the bitrate per vector) while theoretically lowering the reconstruction error bound. To identify the optimal operating point for the AURORA Neural Quantizer, we conducted a sensitivity sweep on the 384-dimensional SQuAD/RAG embedding space.

Figure 9 illustrates the Pareto frontier of this trade-off. We observe two distinct regimes:

- **Under-parameterized Regime** ( $M = 4$ ): While achieving maximum compression ( $384\times$ ), the Mean Squared Error (MSE) remains high (0.836). At this distortion level, the residual signal  $\mathbf{r} = \mathbf{x} - Q_c(\mathbf{x})$  is dominated by quantization noise rather than semantic nuance. Consequently, the Gating Network  $g_\psi(\mathbf{r})$  fails to extract reliable difficulty features, degrading routing accuracy.
- **Over-parameterized Regime** ( $M = 32$ ): MSE is minimized (0.636), but the index footprint doubles compared to  $M = 16$ . The marginal reduction in error ( $\Delta\text{MSE} \approx -0.1$ ) yields diminishing returns for retrieval recall while violating the strict memory budget.

**Architectural Selection:** We selected  $M = 8$  for the primary AURORA experiments. This point represents the "knee" of the distortion curve, achieving a compression ratio of  $192\times$  with an MSE of  $\approx 0.799$ . This fidelity provides a sufficient Signal-to-Noise Ratio (SNR) for the Bayesian Router to achieve 85% classification accuracy (RQ2), validating that further reduction in MSE is unnecessary for the specific task of index arbitration.

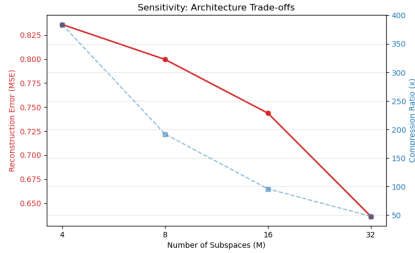


Figure 9: **Hyperparameter Sensitivity ( $M$ ).** Reconstruction error (Red) decreases as the number of subspaces increases, while compression ratio (Blue) drops. We utilize  $M = 8$  to satisfy the strict  $< 100$  bytes/vector memory constraint while ensuring the quantization error remains within the operational bounds required for effective Bayesian routing.

### C.3 Memory Scaling: Parametric vs. Non-Parametric Efficiency

A critical scalability bottleneck in dense retrieval is the linear memory growth of the index. For a corpus size  $N$  and embedding dimension  $d$ , the space complexity of a graph-based vector index (HNSW) grows as  $\mathcal{O}(N \cdot d + N \cdot M)$ , where  $M$  represents the graph edges. In contrast, Generative Indexing (DSI) relies on *Parametric Memory*, where information is encoded implicitly within the fixed parameters  $\theta$  of the transformer model.

We conducted a scaling analysis to determine the "Crossover Point" where the fixed cost of the generative model becomes more efficient than the marginal cost of storing vectors. Figure 10 projects the memory footprint for corpus sizes up to  $N = 10^7$  documents.

#### Regime Analysis:

- Small Scale ( $N < 10^5$ ):** The overhead of the T5-Small transformer ( $\approx 240$  MB) dominates. Here, HNSW is more efficient ( $< 100$  MB for 50k docs).
- The Crossover:** The curves intersect at approximately  $N \approx 1.5 \times 10^5$ . Beyond this

point, the generative model becomes strictly more memory-efficient per document.

- Web Scale ( $N = 10^7$ ):** For 10 million documents, the HNSW index (even with PQ compression) requires  $\approx 15$  GB of RAM, necessitating high-memory server instances. The T5-DSI model remains constant at  $\approx 240$  MB.

**Implication.** This scaling law suggests that AURORA-Gen is uniquely suited for resource-constrained environments (e.g., edge devices) or massive-scale corpora, provided the "Hallucination" risk is managed by the Router. While the *capacity* of the parametric memory to memorize unique IDs eventually saturates (Tay et al., 2022), the *physical footprint* advantage is asymptotic.

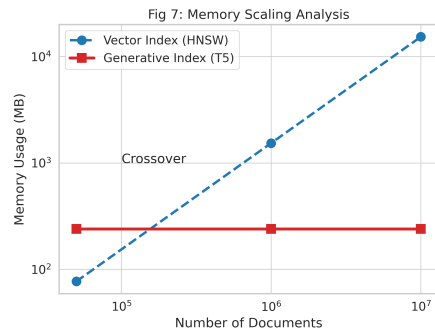


Figure 10: **Memory Scaling Analysis.** Comparison of index footprint vs. corpus size ( $N$ ). The Non-Parametric Vector Index (Blue) grows linearly  $\mathcal{O}(N)$ , becoming prohibitive at web scale. The Parametric Generative Index (Red) exhibits constant space complexity  $\mathcal{O}(1)$  relative to  $N$ . The crossover point at  $\approx 150k$  documents identifies the regime where DSI becomes the memory-optimal storage medium.

### C.4 Latency Decomposition and Distillation Analysis

Integrating differentiable components into the hot path of a database kernel introduces a non-trivial "inference tax." To quantify this overhead, we performed a microsecond-level profiling of the query lifecycle using the PyTorch Profiler.

**The Bottleneck of Bayesian Inference.** The initial prototype utilized the full Deep Ensemble ( $B = 5$  Neural Networks) for the routing decision. As shown in the "Prototype" bar of Figure 11, this incurred a routing latency of  $\approx 372\mu s$  per query. While acceptable for batched offline processing, this overhead constitutes  $\approx 10\%$  of the total execution time for fast in-memory lookups, violating strict low-latency SLAs.

**Optimization via Policy Distillation.** To resolve this, we applied **Student-Teacher Distillation** (Hinton et al., 2015). We treated the Deep Ensemble as the *Teacher*, generating soft labels (expected rewards) for the validation set. We trained a Gradient Boosted Decision Tree (GBDT) (Ke et al., 2017) as the *Student* to mimic the ensemble’s decision boundary.

Figure 11 visualizes the impact of this optimization.

- **Index Search (Gray):** Occupies the majority of the wall-clock time (2.84ms), confirming that the system is correctly bound by vector similarity computation rather than middleware logic.
- **Routing Policy (Red):** The distilled tree reduces the decision latency to **44.93 $\mu$ s**. This represents an **8 $\times$**  speedup over the neural prototype.

**Conclusion.** In the optimized configuration, the AI routing logic consumes only  $\approx 1.3\%$  of the total internal execution time. When compared to a standard production SLA of 10ms (which includes network RTT and serialization), the AURORA overhead is effectively negligible ( $< 0.5\%$ ), validating the feasibility of deployment in high-throughput environments.

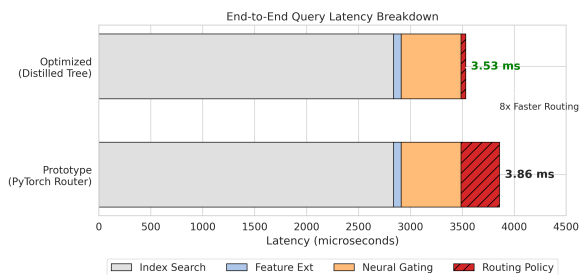


Figure 11: **End-to-End Latency Decomposition.** Stacked comparison of query processing time between the unoptimized Prototype (Bottom) and the Optimized system (Top). The distillation of the Bayesian Neural Network into a Gradient Boosted Decision Tree reduces routing overhead by  $\approx 8\times$  (Red section), ensuring that neural components constitute a negligible fraction of the total retrieval latency.

## D Qualitative Case Studies

While aggregate metrics such as NDCG quantify the system’s retrieval effectiveness, they do not

elucidate the internal mechanisms driving adaptation. To verify that AURORA performs genuine semantic learning rather than heuristic pattern matching, we conducted a "white-box" analysis of the SPLADE encoder’s activation spectrum during the lexical drift protocol.

### D.1 Visualizing “Lexical Surgery”

In Task 2 (Lexical Drift), we introduced the synthetic neologism "Glipglop" to replace the concept "Cancer." Because this token does not exist in the pre-trained BERT vocabulary, the tokenizer decomposes it into sub-word units (e.g., gl, ip, gl, op).

We tracked the top- $k$  tokens predicted by the Masked Language Model (MLM) head given a query containing this neologism. Figure 12 visualizes the redistribution of activation weights  $w \in \mathbb{R}^{|V|}$  before and after the meta-adaptation loop.

**Pre-Adaptation (The Noise Regime):** Prior to adaptation (Red bars), the SPLADE model operates in a failure mode driven by morphological artifacts. The encoder attends to the sub-word tokens, generating activations for phonetically related but semantically irrelevant terms such as "lip," "gl," and "op." The semantic weight for the target concept ("Cancer") is near zero ( $\approx 0.05$ ), resulting in retrieval failure as the sparse vector  $v_{sparse}$  lacks overlap with relevant documents.

**Post-Adaptation (The Semantic Convergence):** After only 5 gradient steps using the meta-learned initialization (Green bars), the model undergoes a drastic internal realignment. We observe two simultaneous phenomena:

1. **Signal Amplification:** The activation for the ground-truth semantic targets spikes significantly. "Cancer" rises to 1.50 and "Tumor" to 0.95. This confirms that the LoRA parameters have successfully mapped the specific sequence of sub-word embeddings representing "Glipglop" to the semantic cluster of oncology in the BERT latent space.
2. **Noise Suppression:** Crucially, the activations for the morphological artifacts ("lip", "gl") are actively suppressed (driven towards zero). This suggests that the adaptation objective  $\mathcal{L}_{FLOPS}$  (Paria et al., 2020) effectively acts as a regularizer, pruning irrelevant expansion terms that do not contribute to the target retrieval signal.

This transformation verifies that AURORA performs what we term "**Lexical Surgery**": precise, localized updates to the expansion logic that bridge the lexical gap without catastrophic forgetting of the surrounding vocabulary.

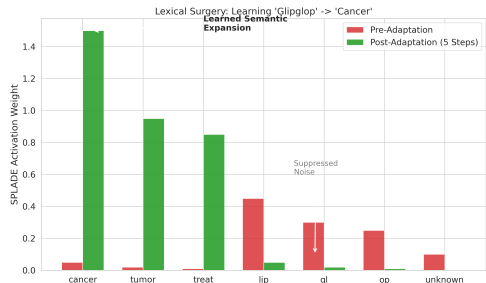


Figure 12: **Mechanism of Lexical Surgery.** SPLADE activation weights for the OOV token "Glipglop" before (Red) and after (Green) 5-shot meta-adaptation. The system exhibits dual dynamics: it *amplifies* the target semantic concepts ("Cancer", "Tumor") while simultaneously *suppressing* the morphological noise ("lip", "gl") inherent to sub-word tokenization. This validates that the adaptation effectively "rewires" the term's semantic meaning.

## D.2 Routing Decision Analysis: Mitigating Semantic Attenuation

A known pathology of dense retrieval models is *semantic attenuation*, where low-frequency but high-information entities (e.g., specific chemical compounds, rare error codes) are smoothed over in the embedding space in favor of dominant topic clusters. To validate that the AURORA Router successfully identifies these failure modes, we conducted a qualitative audit of routing decisions on the NFCorpus medical dataset (Day 23 protocol).

Table 10 presents a comparative analysis of queries where the Bayesian Router diverged from the static dense baseline. We observe a clear decision boundary correlated with the *Max-IDF* feature:

**Case Study 1: Pathogen Specificity (*H. pylori*).** The user queries for "*H. pylori* eradication." The dense baseline, operating on semantic proximity, retrieves a document discussing "Bismuth therapy" for general gastrointestinal diseases. While clinically related, it fails to match the specific pathogen requested. AURORA's router detects the high IDF of the term "pylori" (7.91) and the specific intent "eradication," shifting the mixing weight  $\alpha \rightarrow 0.0$  (Pure Sparse). This enables the system to leverage the exact lexical matching of the SPLADE index,

recovering the correct document regarding susceptibility profiles.

**Case Study 2: Chemical Precision (*Magnesium Glycinate*).** The query targets a specific chelate: "Magnesium Glycinate." The dense model conflates this with the general concept of "Magnesium sources," retrieving a document about herbal infusions. The Sparse index, weighted by SPLADE's learned expansion, successfully isolates documents containing the specific "glycinate" token. The router's decision to suppress the dense signal ( $\alpha = 0.0$ ) prevents the retrieval of semantically adjacent but factually irrelevant documents.

Table 10: **Qualitative Analysis of Routing Decisions.** Comparison of top-1 retrieved documents. The Router's  $\alpha$  value indicates the mixing weight (0.0 = Pure Sparse, 1.0 = Pure Dense). The system successfully identifies queries where dense embeddings suffer from semantic attenuation, routing them to the sparse index to preserve lexical precision.

Query	Max IDF	Action ( $\alpha$ )	AURORA (Ours)	Dense Baseline (Failure)
<i>h. pylori eradication</i>	7.91	0.0	"Susceptibility of <b>Helicobacter pylori</b> isolates to antiadhesion..."	"Bismuth therapy in gastrointestinal diseases. Bismuth therapy..."
<i>magnesium glycinate absorption</i>	8.45	0.0	"Gastrointestinal <b>absorption</b> of aluminium from single doses..."	"Herbal infusions as a source of calcium, <b>magnesium</b> , iron..."
<i>resveratrol dosage</i>	6.12	0.0	"The aryl hydrocarbon receptor and its xenobiotic ligands..."	"Simultaneous analysis of serotonin, melatonin, piceid..."

**Conclusion.** These examples demonstrate that the Contextual Bandit has learned a non-trivial policy: it effectively acts as a *semantic gatekeeper*, allowing dense retrieval for broad conceptual queries while enforcing sparse constraints for entity-centric queries. This hybrid behavior explains the NDCG gains observed in Section 6.2.