

# MightyRayn: A Self-Regulating Bilateral Cognitive Architecture with Oscillator-Coupled Attention, 5-Axis Runtime Compression, and Universal Signal Classification via Iterative Denoising Retrieval

Amyrr Beyveinel  
MightyCloudCollective  
amyrelbay@gmail.com

U.S. Provisional Patent Application No. 64/031,489 | Filed April 7, 2026 | Docket: MIGHTYRAYN-PROV-001

---

## Abstract

We present MightyRayn, a unified cognitive architecture that integrates three novel mechanisms into a self-regulating feedback loop: (1) **Kuramoto oscillator-driven variable coupling attention**, in which bilateral brain coherence modulates attention inversely -- disagreement *increases* attention -- and coupling strength varies dynamically from 0.06 (creative exploration) to 0.90 (threat response); (2) **O(1) episodic memory recall** via locality-sensitive hashing enabling constant-time retrieval of similar past episodes for pattern completion, achieving 83% hit rate with sub-millisecond latency; and (3) **metacognitive state transitions** (CONVERGENT, DELIBERATIVE, CONFLICTED, NOVEL) driven by composite divergence between creative and analytical hemispheres, controlling LLM temperature and token allocation. Processing completes in 0.1-7.4ms per frame. Building on this cognitive core, we introduce a **5-axis runtime compression framework** achieving 612x composite compression without model weight modification, and a **domain-agnostic iterative denoising retrieval** system demonstrated across 10+ signal domains with 132+ canonical patterns -- achieving 100% Rank-1 recall at 250 million tokens ( $2.5 \times 10^8$ :1 search space reduction) on classical CPU hardware in 60ms, providing quantum-equivalent search performance without quantum hardware.

**Keywords:** cognitive architecture, Kuramoto oscillators, bilateral processing, metacognition, runtime compression, iterative denoising, vector quantization, signal classification, episodic memory

---

## 1. Introduction

Recent advances in brain-inspired AI architectures have demonstrated the value of drawing from neuroscience to improve artificial intelligence systems. The Modular Autonomous Planner (MAP) [1] introduced modular brain-inspired processing, AKOrN [2] replaced threshold-based neurons with Kuramoto oscillatory units, HoloBrain [3] applied Kuramoto coupling for modeling brain dynamics, and PhaseGPT [4] explored phase-coupled oscillator attention in transformers. However, these systems share fundamental limitations: serial processing pipelines, fixed coupling strengths, no metacognition from bilateral divergence, and disconnected memory systems.

We present MightyRayn, a unified system that addresses all of these limitations simultaneously. Our contributions span three distinct but interconnected areas: (i) a bilateral cognitive architecture with dynamic Kuramoto coupling and metacognitive state transitions; (ii) a 5-axis runtime compression framework that operates entirely at the application layer without weight modification; and (iii) a domain-agnostic iterative denoising retrieval mechanism that achieves quantum-equivalent search space reduction on classical hardware. These three components share core algorithmic infrastructure -- vector quantization, locality-sensitive hashing, and importance scoring -- creating a unified framework rather than a collection of independent techniques.

The system has been reduced to practice with verified empirical results: the cognitive architecture processes inputs in 0.1-7.4ms per frame with variable coupling from 0.06 to 0.90; the compression framework achieves 612x composite compression across five orthogonal axes; and the retrieval system achieves 100% Rank-1 recall at 250 million tokens across 10+ signal domains with 132+ canonical patterns.

## 2. Related Work

### 2.1 Brain-Inspired Architectures

MAP [1] organizes specialized modules in a serial pipeline inspired by brain regions. AKOrN [2] demonstrated that Kuramoto oscillatory neurons improve associative memory and pattern binding over traditional threshold-based units. PhaseGPT [4] applied phase-coupled oscillator attention within transformer architectures. These approaches are limited to either serial processing or fixed coupling, and none implements bilateral processing with dynamic coupling modulation based on divergence.

## 2.2 Episodic Memory in AI

RAG [5] and MemInsight [6] apply retrieval-augmented generation as external memory pipelines. Locality-sensitive hashing (LSH) [7] has been widely used for approximate nearest-neighbor search. Unlike external retrieval pipelines, our episodic memory is integrated directly into the cognitive processing loop with constant-time  $O(1)$  lookup, providing familiarity signals within the same processing frame as attention and coupling calculations.

## 2.3 Model Compression

Existing compression methods typically require weight modification: quantization [8], pruning [9], distillation [10], and speculative decoding [11]. LayerSkip [12] enables training-time layer skipping via self-speculative decoding. Our approach operates entirely at the application layer, requiring no weight modification, retraining, or architectural changes, and composes five independent axes multiplicatively.

## 2.4 Diffusion-Based Retrieval

DDCM [13] and gDDCM [14] use codebooks in reverse diffusion for image generation. DiffuRetrieval [15] applies diffusion to sponsored search via chain-of-thought refinement. DiffuGR [16] generates document identifiers via discrete diffusion. Our approach is fundamentally different: we retrieve existing documents by denoising a query vector toward stored codebook centroids -- retrieval, not generation -- and demonstrate this across 10+ signal domains.

# 3. Bilateral Cognitive Architecture

## 3.1 System Overview

The architecture comprises two heterogeneous processing hemispheres coupled through a neural bridge module. The left hemisphere (creative agent) is configured for divergent thinking and novel associations. The right hemisphere (analytical agent) is configured for convergent thinking, logical analysis, and risk assessment. Both hemispheres independently process the same input stimulus and produce response vectors including confidence scores, reasoning chains, and recommended actions.

## 3.2 Kuramoto Variable Coupling Mechanism

The two hemispheres are modeled as coupled oscillators governed by the Kuramoto model. For oscillator  $i$  with phase  $\theta_i$  and natural frequency  $\omega_i$ :

$$d(\theta_i)/dt = \omega_i + (K/N) * \sum_j [\sin(\theta_j - \theta_i)]$$

The order parameter  $r$  quantifies bilateral coherence, ranging from 0 (incoherence) to 1 (synchrony). The critical innovation is that coupling strength  $K$  is not fixed but varies dynamically via a Variable Coupling Controller:

$$K_{effective} = K_{base} * threat\_multiplier * divergence\_modifier$$

Under threat conditions,  $K$  is multiplied by up to 3x ( $K = 0.90$ ), forcing rapid bilateral synchronization. During high divergence without threat,  $K$  is reduced to 0.06, permitting independent exploration. This represents a 15x dynamic range in coupling strength, enabling the system to transition between forced consensus and creative divergence based on real-time assessment.

## 3.3 Inverse Attention Modulation

Attention is computed as an inverse function of bilateral coherence:  $attention = \max(5\%, (1 - r) * scale\_factor)$ . Low coherence (high disagreement) produces HIGH attention, reflecting the biological principle that bilateral disagreement signals uncertainty requiring increased processing resources. This is the inverse of typical attention filtering mechanisms, which suppress disagreement. Verified results show attention ranges from 5% during convergent agreement to 77% during cold-start or novel situations.

### 3.4 Metacognitive State Machine

A composite divergence metric drives a four-state metacognitive state machine:

$$D_{composite} = w_{sem} * D_{semantic} + w_{conf} * D_{confidence} + w_{str} * D_{structural}$$

The four states are CONVERGENT (bilateral agreement, temperature 0.60), DELIBERATIVE (moderate disagreement, temperature 0.70-0.80), CONFLICTED (strong disagreement, temperature up to 0.93), and NOVEL (extreme divergence or no episodic memory match, maximum exploration). Each state determines temperature, token budget, and coupling strength, enabling dynamic computational resource allocation.

### 3.5 O(1) Episodic Memory

Past processing episodes are encoded and stored via locality-sensitive hashing, enabling constant-time O(1) lookup. The episodic memory produces a familiarity signal (hit/miss, familiarity percentage, similar episodes) that influences metacognitive state transitions and coupling parameter selection. Verified results: 83% hit rate, familiarity scores of 69%, with processing time included in the 0.1-7.4ms per-frame budget.

**Table 1: Cognitive Architecture Verified Test Results**

Scenario	Coupling K	Order r	Attention	Temp	State
Cold start	0.30	0.055	77%	0.70	NOVEL
Strong disagreement	0.06	0.31	56%	0.93	CONFLICTED
Security threat	0.90	0.997	5%	0.60	CONVERGENT
Routine processing	0.30	0.999	5%	0.60	CONVERGENT
Memory recall	--	--	--	--	hit=True, 69%

## 4. Five-Axis Runtime Compression Framework

We introduce a multi-axis compression framework in which five orthogonal, independently operable axes compose multiplicatively. Critically, all five axes operate at the application layer without model weight modification, retraining, or architectural changes.

### 4.1 Axis Definitions

**Axis 1 -- Prompt Compression (3.33x):** Progressive skill withdrawal removes inference-time prompt scaffolding as model competence is demonstrated. Competence is tracked via exponential moving average of task success rates. A 1500-character prompt compresses to 450 characters.

**Axis 2 -- Output Compression (5.22x):** Bidirectional budget enforcement constrains output token generation via shrinking budgets with quality-driven expansion. Starting budget of 256 tokens converges to 32 tokens over 13 tightening cycles across 50 iterations, achieving 87.5% compression.

**Axis 3 -- Memory Compression (3.0x):** KV cache quantization reduces attention key and value tensor precision from full (f16) to reduced (q4\_0) format. Attention memory is reduced by at least 50% compared to full-precision storage, with negligible quality degradation.

**Axis 4 -- Context Compression (11.72x):** Streaming extractive memory scores input corpus chunks by information density and maintains compressed memory entries with importance decay and VQ-powered deduplication. Demonstrated: 24,000 tokens compressed to 2,048-token context windows while preserving instruction-bearing content.

**Axis 5 -- Compute Compression (Layer-Skip, 19% speedup):** Diffusion-guided layer-skip prediction classifies inputs by VQ fingerprint, retrieves or generates binary layer activation masks, and skips transformer layers whose learned importance falls below a configurable threshold. The first two (embedding) and last two (output) layers are protected from skipping.

### 4.2 Composite Compression

The four token/memory axes compose multiplicatively:  $3.33 \times 5.22 \times 3.0 \times 11.72 = \mathbf{612x}$  streaming compression, with layer-skip as an orthogonal 5th axis providing additional compute savings. This enables transformer models requiring N gigabytes of memory to operate within a fraction of N gigabytes through combined token, memory, and compute compression -- all without touching model weights.

**Table 2: Five-Axis Compression Framework Results**

Axis	Mechanism	Before	After	Ratio
1: Prompt	Skill withdrawal	1,500 chars	450 chars	3.33x
2: Output	Budget enforcement	256 tokens	32 tokens	5.22x (87.5%)
3: Memory	KV quantization	f16	q4_0	3.0x
4: Context	Extractive memory	24,000 tokens	2,048 tokens	11.72x
5: Compute	Layer-skip	32 layers	26 layers	19% speedup

**Composite (Axes 1-4):  $3.33 \times 5.22 \times 3.0 \times 11.72 = \mathbf{612x}$  streaming compression**

## 5. Domain-Agnostic Iterative Denoising Retrieval

### 5.1 Core Mechanism

We introduce a pattern recognition method based on iterative denoising toward codebook centroids. The algorithm proceeds as follows: (1) encode input data into a fixed-dimension pattern vector via a domain-specific feature encoder; (2) compute similarity between the input vector and codebook entries; (3) blend the input vector toward the highest-similarity entries with a decreasing blend weight per iteration; (4) re-normalize and repeat for N denoising steps; (5) return the nearest codebook entry as the identified pattern with confidence score.

The key insight is that this mechanism is completely domain-agnostic -- the same core algorithm applies to any domain where input data can be encoded into a fixed-dimension vector space and canonical patterns can be defined as codebook centroids. The only domain-specific component is the feature encoder; the denoising, blending, and ranking infrastructure is shared.

### 5.2 Multi-Domain Demonstration

We demonstrate the mechanism across 10+ signal domains, with verified results in 5 primary domains:

**Table 3: Iterative Denoising Retrieval Across Signal Domains**

Domain	Corpus Size	Patterns	Recall	Latency
NLP Document Retrieval	250M tokens (333K chunks)	20 needles	100% Rank-1	60ms avg
Financial Chart Patterns	OHLCV time-series	20 canonical	100% detect, 70% R1	<100ms
Medical Waveforms	ECG/EEG/CGM/CTG/etc.	92 across 8 sub-domains	64% (noisy synthetic)	<100ms
Response Caching	LLM outputs	VQ fingerprint	Sub-ms retrieval	<1ms
Layer-Skip Prediction	Per-layer importance	Binary masks	19% speedup	<1ms

### 5.3 Universal Signal Classification

The unified system maintains domain-specific codebooks while sharing all algorithmic infrastructure. Signal domain is determined via metadata or a lightweight classifier, the appropriate codebook is selected, and iterative denoising proceeds identically regardless of domain. Currently demonstrated codebook coverage includes: 20 NLP needles at 250M tokens, 20 financial chart patterns (head-and-shoulders, bull flag, double top, cup-and-handle, ascending triangle, etc.), 92 clinical patterns across 8 physiological domains (25 ECG, 15 EEG, 10 CGM, 12 CTG, 8 ventilator, 8 EMG, 8 ABP, 6 respiratory), seismological event patterns (tectonic, volcanic, induced, explosion signatures), and atmospheric patterns (cold front, mesocyclone, microburst, tropical cyclone, atmospheric river, derecho). Total: 132+

canonical patterns across 10+ signal domains.

## 6. Quantum-Equivalent Classical Search

The composed compression pipeline achieves a search space reduction from  $N > 10^8$  to a single result through sequential compression stages:

**Table 4: Composed Search Space Reduction Pipeline**

Stage	Input	Output	Reduction
Corpus compression	250M tokens	333,333 entries	750x
VQ + Index compression	333,333 entries	~500 candidates	666x
Iterative denoising (5 steps)	500 candidates	1 Rank-1 result	500x
TOTAL	250,000,000 tokens	1 result	$2.5 \times 10^8 : 1$

The iterative denoising stage implements classical analogues of three quantum computational properties. **Superposition:** the query vector exists as a weighted combination of all codebook centroids prior to final ranking, analogous to a qubit in superposition across computational basis states. **Constructive and destructive interference:** each denoising step increases the weight of high-similarity centroids (constructive) while decreasing the weight of low-similarity centroids (destructive), analogous to quantum amplitude amplification. **Measurement collapse:** the final ranking operation selects the single highest-similarity centroid, analogous to wavefunction collapse.

Empirically, the system achieves 100% retrieval recall (20/20 Rank-1) at 250 million tokens on classical CPU hardware in 60ms average latency -- a search space reduction of  $2.5 \times 10^8 : 1$ . For comparison, Grover's algorithm would require  $\sqrt{333,333} \sim 577$  quantum oracle queries for the same corpus, but requires quantum hardware, quantum error correction, and cryogenic cooling. Our approach achieves comparable reduction on commodity classical hardware.

## 7. Post-Quantum Cryptographic Security Application

We apply the quantum-equivalent search mechanism to post-quantum cryptographic security. A codebook of known vulnerability patterns is constructed including: deprecated algorithm signatures (RSA-1024, MD5, SHA-1, DES, 3DES), quantum-vulnerable key exchange protocols (ECDH with insufficient key length, static Diffie-Hellman), implementation weaknesses (timing side-channels, padding oracle patterns), and quantum attack probe signatures.

The composed compression pipeline scans large-scale network traffic captures, source code repositories, or certificate databases for entries matching vulnerability patterns. Output includes pattern name, confidence score, CVE reference where applicable, and remediation recommendation aligned with NIST post-quantum standards FIPS 203 (ML-KEM), FIPS 204 (ML-DSA), and FIPS 205 (SLH-DSA). The constant-time search capability enables scanning datasets exceeding  $10^8$  elements, supporting the post-quantum migration timeline mandated for completion before 2035.

## 8. Unified Algorithmic Infrastructure

A distinguishing characteristic of MightyRayn is that the cognitive architecture, compression framework, and retrieval system share core algorithmic components rather than being independent techniques assembled post-hoc. Specifically:

**Vector quantization (VQ)** is used for: episodic memory encoding (Section 3.5), response cache fingerprinting (Section 4, Axis 1), approximate cache matching (Section 4, Axis 4), layer-skip mask retrieval (Section 4, Axis 5), and codebook centroid matching in iterative denoising (Section 5). **Locality-sensitive hashing (LSH)** is used for:  $O(1)$  episodic recall in the cognitive loop, multi-resolution index lookup in the retrieval pipeline, and bucket-based candidate pruning. **Importance scoring via exponential moving average (EMA)** is used for: competence tracking in prompt skill withdrawal, information density scoring in context compression, and per-layer importance profiling for layer-skip. This shared infrastructure means that improvements to any core component propagate across all three systems simultaneously.

## 9. Empirical Results

### 9.1 Needle-in-Haystack at 250M Tokens

The flagship benchmark evaluates retrieval of 20 planted needles within a 250-million-token corpus (333,333 chunks). Results: all 20 needles retrieved at Rank-1 (100% Rank-1 recall) in 60ms average latency on CPU. For comparison, LSH alone achieves approximately 10% recall on the same data, demonstrating the critical contribution of iterative denoising refinement.

### 9.2 Diffusion Retrieval at Scale

**Table 5: Diffusion Retrieval Benchmark at 250M Tokens**

Metric	Result
Corpus size	250,000,000 tokens (333,333 chunks)
Needles planted	20
Needles found	20/20 (100%)
Rank-1 accuracy	20/20 (100%)
Average latency	60ms (CPU)
LSH-only baseline	~10% recall
Search space reduction	$2.5 \times 10^8 : 1$

### 9.3 Cognitive Architecture Processing

The bilateral processing loop completes in 0.1ms (routine/convergent) to 7.4ms (novel/maximum exploration) per frame. Episodic memory achieves 83% hit rate with 6 episodes stored. The full self-regulating feedback loop -- disagreement to attention increase to temperature raise to exploration to resolution to episode storage -- executes within these time bounds without external supervision or manual tuning.

### 9.4 Compression Framework

Output budget convergence demonstrates the bidirectional enforcement mechanism: starting from a 256-token budget, the system converges to 32 tokens over 13 tightening cycles across 50 iterations (87.5% compression), with context expanding from 2,048 to 3,200 tokens as the system learns which information to retain. The 4-axis composite of 612x represents the multiplicative combination of independently measured per-axis ratios.

## 10. Intellectual Property Structure

The system is protected by U.S. Provisional Patent Application No. 64/031,489 (filed April 7, 2026) comprising 31 claims organized into four divisional groups supporting independent continuation applications:

**Table 6: Patent Claim Structure and Divisional Strategy**

Division	Claims	Scope
A: Cognitive	1-20	Bilateral architecture, Kuramoto coupling, metacognition, episodic memory
B: Compression	21-23	Runtime compression proxy, layer-skip, 5-axis framework
C: Universal Pattern	24-29	Domain-agnostic denoising retrieval, 6 domain applications
D: Quantum-Equivalent	30-31	Classical quantum-equivalent search, post-quantum crypto

## 11. Discussion and Future Work

MightyRayn demonstrates that a unified algorithmic infrastructure based on vector quantization, locality-sensitive hashing, and importance scoring can simultaneously address cognitive architecture design, runtime model compression, and universal signal classification. The 612x compression ratio enables deployment of large transformer models on resource-constrained hardware, while the quantum-equivalent search mechanism provides an alternative path to large-scale search reduction without quantum hardware.

Several directions warrant further investigation. First, the medical waveform domain currently achieves 64% recall on noisy synthetic data; real clinical validation with annotated physiological datasets would establish clinical utility. Second, the quantum-equivalent search analogy invites formal complexity-theoretic analysis of the composed compression pipeline. Third, the post-quantum cryptographic application should be validated against real-world network traffic captures and CVE databases. Fourth, the bilateral cognitive architecture should be evaluated on established multi-agent benchmarks to compare against serial pipeline architectures like MAP.

A PCT international application is planned within the 12-month priority window (by April 2027), followed by non-provisional conversion with continuation applications along the four divisional strategies outlined in Section 10.

## 12. Conclusion

We have presented MightyRayn, a unified cognitive architecture integrating bilateral Kuramoto-coupled processing with metacognitive state transitions, 5-axis runtime compression achieving 612x composite reduction without weight modification, and domain-agnostic iterative denoising retrieval achieving 100% Rank-1 recall at 250 million tokens with  $2.5 \times 10^8:1$  search space reduction on classical hardware. The system has been reduced to practice with verified empirical results across 10+ signal domains and 132+ canonical patterns, protected by U.S. Provisional Patent Application No. 64/031,489. The shared algorithmic infrastructure demonstrates that cognitive architecture, model compression, and signal classification are not independent problems but manifestations of a unified framework built on vector quantization, locality-sensitive hashing, and importance-weighted iterative refinement.

## References

- [1] Nature Communications (2025). "Modular Autonomous Planner (MAP): A Brain-Inspired Modular Architecture for Multi-Agent Planning."
- [2] ICLR 2025. "AKOrN: Associative Kuramoto Oscillatory Recurrent Networks."
- [3] HoloBrain. "Kuramoto Coupling Models with Attention Mechanisms for Brain Oscillation Dynamics."
- [4] PhaseGPT. "Phase-Coupled Oscillator Attention within Transformer Architectures for Language Modeling."
- [5] Lewis, P. et al. (2020). "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks." NeurIPS.
- [6] MemInsight. "Memory-Augmented Retrieval for Enhanced LLM Outputs."
- [7] Indyk, P. and Motwani, R. (1998). "Approximate Nearest Neighbors: Towards Removing the Curse of Dimensionality." STOC.
- [8] Dettmers, T. et al. (2022). "GPT3.int8(): 8-bit Matrix Multiplication for Transformers at Scale." NeurIPS.
- [9] Frantar, E. and Alistarh, D. (2023). "SparseGPT: Massive Language Models Can Be Accurately Pruned in One-Shot." ICML.
- [10] Hinton, G. et al. (2015). "Distilling the Knowledge in a Neural Network." NeurIPS Workshop.
- [11] Leviathan, Y. et al. (2023). "Fast Inference from Transformers via Speculative Decoding." ICML.
- [12] Elhoushi, M. et al. (2024). "LayerSkip: Enabling Early-Exit Inference and Self-Speculative Decoding." arXiv:2404.16710.
- [13] DDCM (2025). "Denoising Diffusion with Codebook Models for Image Generation."
- [14] gDDCM. arXiv:2511.13387. "Generalized Denoising Diffusion Codebook Models."
- [15] DiffuRetrieval. ACM WWW 2024. "Diffusion-Based Sponsored Search Ranking via Chain-of-Thought Refinement."
- [16] DiffuGR. arXiv:2511.08150. "Discrete Diffusion for Generative Document Retrieval."
- [17] FTS-Diffusion. OpenReview 2024. "Financial Time Series Denoiser."
- [18] Distill-VQ. arXiv:2204.00185. "Vector Quantization for One-Shot Document Retrieval."
- [19] Beyveinel, A. (2026). "U.S. Provisional Patent Application No. 64/031,489: Self-Regulating Bilateral Cognitive Architecture." USPTO.