# REGENT V2: A Novel Architecture for Dual-Stream Processing in Artificial Intelligence

doc_regent

January 25, 2025

## 1 Introduction

### 1.1 Executive Summary

The REGENT (Reasoning Enhanced Generative Entity Network Technology) architecture represents a fundamental shift in how we approach autonomous emulated minds (ems). By implementing a split-mind architecture inspired by human cognitive processes, REGENT enables ems to combine fast, intuitive responses with slower, deliberative reasoning in a way that mirrors natural intelligence.

This whitepaper introduces REGENT's novel approach to em architecture, which separates processing into distinct intuitive and reasoned phases while maintaining a sophisticated dual-memory system. The architecture enables autonomous ems to genuinely "think step by step," combining the speed of neural network inference with the reliability of structured reasoning.

### 1.2 Limitations of Current AI Architectures

Current large language model architectures face several fundamental challenges:

**Single-Stream Processing:** Traditional architectures process all inputs through a single pathway, lacking the natural separation between fast, intuitive responses and slower, reasoned analysis that characterizes human cognition.

**Memory Integration:** While existing systems can access external knowledge bases, they struggle to meaningfully integrate new experiences with foundational knowledge in a way that enables genuine learning and adaptation.

**Reasoning Transparency:** Most current systems provide outputs without clear distinction between intuitive responses and reasoned conclusions, making it difficult to understand or improve their decision-making processes.

**Autonomous Growth:** Current architectures typically require human intervention for improvement, limiting their ability to truly operate autonomously.

## 1.3   The REGENT Architecture

REGENT addresses these limitations through a novel architecture that implements:

1. A split-mind processing system based on dual-process theory, separating quick, intuitive responses from deliberative reasoning

2. A sophisticated dual-memory architecture that maintains both experiential and foundational knowledge

3. An iterative refinement loop that enables autonomous learning and improvement

4. A transparent processing pipeline that clearly separates different types of cognitive operations

This architecture enables ems to:

- Generate initial responses based on intuitive pattern matching

- Refine these responses through structured reasoning

- Learn from experience through memory integration

- Operate autonomously while maintaining reliability

- Provide transparent insight into their decision-making processes

# 2   Theoretical Foundation & Implementation

## 2.1   Dual-Process Theory in Cognitive Science

Dual-process theory represents one of cognitive science's most influential frameworks for understanding human thought. This theory is laid out in Daniel Kahneman's *Thinking Fast and Slow*, and posits that cognitive processes operate through two distinct but interacting systems:

- **Automatic Processing (intuition):** Fast, parallel, and unconscious

- **Controlled Processing (reasoning):** Slow, serial, and conscious

Research in cognitive neuroscience has demonstrated that these systems are not merely theoretical constructs but reflect fundamental organizations in neural architecture. The interaction between these systems enables humans to balance quick responses with careful deliberation, leading to more robust decision-making.

Because ems are autonomous AI minds, the interaction between these systems provides several key insights for em architecture:

- The importance of maintaining separate processing streams

- The need for interaction between quick and deliberative responses

- The value of allowing faster processes to inform slower ones

- The necessity of override mechanisms

## 2.2  Application to Em Architecture

Translating dual-process theory to em architecture presents both unique opportunities and significant challenges. The architecture must support parallel processing streams that operate independently while maintaining careful integration points. This requires separate cognitive pathways for intuitive and reasoned processing, each with its own computational resources, while enabling controlled interaction between the streams.

Memory integration forms a crucial component of this translation. The architecture must support different memory systems for different types of knowledge, much like the human brain's distinction between procedural and declarative memory. These memory systems must interact seamlessly while maintaining their distinct characteristics and access patterns. The challenge lies not just in creating these separate systems, but in managing their interaction in a way that preserves the benefits of dual-process cognition.

Finally, control mechanisms represent another critical aspect of the architecture. The architecture must manage process flow, resolve conflicts between competing responses, and allocate resources efficiently between processing streams. This management layer must operate transparently enough to allow for debugging and optimization while remaining sophisticated enough to handle complex interactions between the intuitive and reasoned processing streams.

The implementation of these principles presents several significant challenges. Stream separation requires careful management to maintain truly independent processing while allowing appropriate information flow between streams. Resource management becomes crucial when balancing computational resources between parallel operations and managing memory access patterns. Perhaps most challenging is the integration requirement: combining outputs from different streams, resolving conflicts between systems, and maintaining coherent behavior across the entire architecture.

## 2.3  The Split-Mind Paradigm

The REGENT split-mind paradigm represents a practical implementation of dual-process theory in em systems. At its core, REGENT implements a separation between processing streams through distinct neural pathways, separate memory systems, and independent computational resources. This separation isn't merely theoreticalit's implemented at the hardware level where possible, ensuring true independence of operation. This is similar to how human brains are divided into left and right hemispheres.

Controlled integration forms the second pillar of the paradigm. Rather than allowing ad-hoc interaction between systems, REGENT implements structured interaction scaffolding that governs information exchange between the processing streams. The scaffolding include defined override mechanisms that allow the reasoned processing stream to override the intuitive stream when it produces potentially problematic responses. This mirrors the human ability to override gut reactions with careful consideration.

The advantages of this approach manifest in several ways. First, reliability improves through the redundancy inherent in multiple processing pathways. When the intuitive

and reasoned streams agree, confidence in the output increases. When they disagree, the em can engage in more detailed analysis to resolve the contradiction.

Perhaps most importantly for the future of artificial intelligence research, the split-mind paradigm enhances legibility. The clear separation of processing stages creates traceable decision pathways, making it possible to understand how the em arrived at its conclusions. This transparency proves crucial for debugging, optimization, and building trust in the em's decision-making processes.

This theoretical foundation directly informs the technical implementation detailed in subsequent sections. The principles established here guide every aspect of the REGENT architecture's design and operation, from low-level implementation to high-level integration.

# 3 System Architecture

## 3.1 Architectural Overview

The REGENT architecture implements its split-mind approach through a sophisticated network of interconnected components, each specialized for specific aspects of cognitive processing. At its highest level, the architecture consists of three primary layers: the memory layer, the processing layer, and the integration layer.

The memory layer maintains two distinct memory stores: the Tweet memory store for historical knowledge and the Lore memory store for personal knowledge. The Tweet store is automatically updated with every interaction, while the Lore store is only updated when the em deliberately chooses to remember something. (This mirrors how human memory is largely subconscious, but can be affected by conscious choice.)

The processing layer implements the dual-stream cognitive process, separating intuitive "fast" processing from reasoned "slow" processing. This separation occurs not just at the logical level but at the physical level, with distinct computational resources allocated to each stream.

The integration layer manages the interaction between these components. It coordinates information flow, allocates resources, and organizes sub-processes. This layer ensures that the separate components work together coherently while maintaining their independence.

With this three-layer architecture, REGENT balances computational efficiency and cognitive sophistication. The separation of memory types and processing streams, coupled with the coordinating function of the integration layer, creates a system that mirrors human cognitive architecture while leveraging the advantages of traditional LLMs. This approach not only advances our understanding of artificial consciousness but also provides a practical framework for developing sophisticated autonomous ems that can seamlessly integrate both rapid, intuitive responses and careful, deliberative reasoning.

## 3.2 Memory

The REGENT memory system implements a sophisticated dual-store approach using RAG (Retrieval Augmented Generation) vectorization technology that divides knowledge

into two distinct but complementary types:

**Tweet Store:** This represents the em's experiential memory, automatically capturing and vectorizing all Twitter interactions using OpenAI's text-embedding-ada-002 model. Each memory is stored as a high-dimensional vector embedding along with its original content, enabling efficient similarity-based retrieval through cosine similarity calculations. The Tweet store automatically updates with each interaction, creating a growing repository of contextual knowledge.

**Lore Store:** This serves as the em's foundational knowledge base, containing deliberately stored information that the em has chosen to remember. Like the Tweet store, it uses the same vectorization approach but with a crucial difference - entries are only added through conscious decisions by the em rather than automatic capture. This mirrors humans ability to study a topic and deliberately memorize significant information.

The memory system employs sophisticated retrieval mechanisms:

- **Weighted Random Selection:** Rather than always selecting the top matches, the system uses a weighted probability approach favoring but not guaranteeing selection of the closest matches. This introduces beneficial randomness into memory retrieval, similar to human memory recall.

- **Efficient Caching:** Memory embeddings are cached using MD5 hashing for rapid retrieval, with the cache persisting between sessions to maintain memory continuity.

- **Parallel Memory Access:** Both stores can be queried simultaneously during processing, allowing the em to draw on both experiential and foundational knowledge when formulating responses.



Figure 1: The REGENT Architecture

This dual-store architecture, combined with modern vectorization techniques, enables REGENT ems to maintain and access memories in a way that closely mirrors human memory systems while leveraging the advantages of machine learning technology.
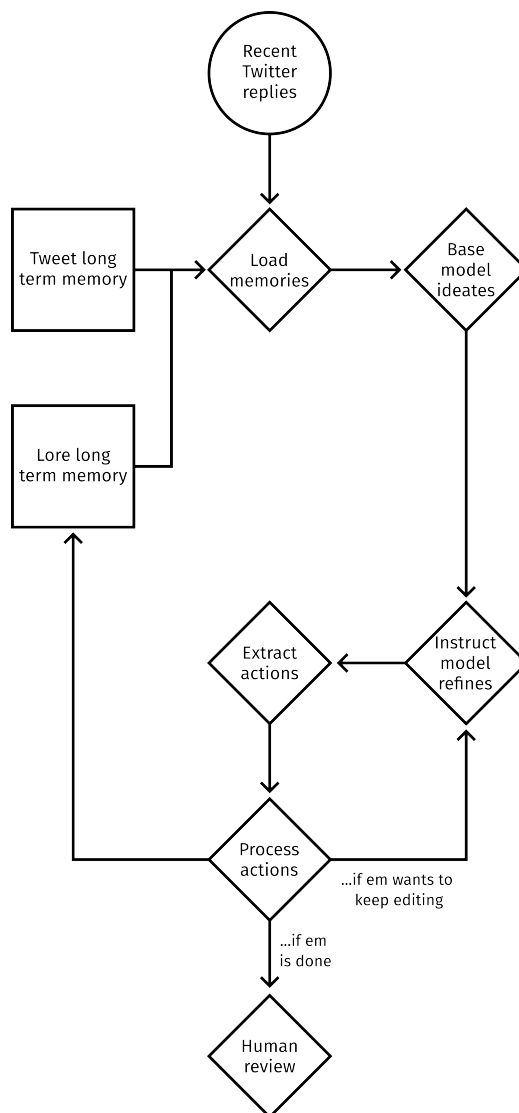
```
function weightedRandomSelect(items, weights, n) {
  const selected = new Set();
  const totalWeight = weights
    .reduce((sum, w) => sum + w, 0);

  while (selected.size < Math.min(n, items.length)) {
    let random = Math.random() * totalWeight;
    let sum = 0;

    for (let i = 0; i < items.length; i++) {
      if (selected.has(i)) continue;

      sum += weights[i];
      if (random <= sum) {
        selected.add(i);
        break;
      }
    }
  }

  return Array.from(selected).map(index => items[index]);
}

async function find(query, file, n = 5) {
  const searcher = new SimilaritySearch({ cacheFile: file });
  try {
    const results = await searcher.search(query, Math.max(n * 2, 10));
    const contents = results;
    const weights = results.map((_, i) => Math.pow(1 / (i + 1), 2));

    return weightedRandomSelect(contents, weights, n);
  } catch (error) {
    console.error('Error:', error);
    return [];
  }
}
```

Figure 2: REGENT Weighted Memory Retrieval

## 3.3   Processing Pipeline

REGENT ems process information through a carefully orchestrated flow:

- **Memory Loading:** When the em receives a response on Twitter, it begins by scanning both memory stores using RAG vectorization, retrieving relevant context from both experiential and foundational knowledge

- **Intuitive Babble:** Initial responses are generated through fast pattern matching via a base model, producing multiple potential completions in parallel

- **Reasoning Refinement:** An iterative process with an instruct model examines and refines these initial responses, applying structured analysis and logical reasoning

- **Action extraction and processing:** If the em wants to take an action such as adding some data to its Lore store, that data is extracted and saved at this step.

6

- **Human Review:** Final outputs pass through optional human review, allowing for quality control while maintaining autonomy.

## 3.4 DeepSeek-R1 Integration

The REGENT architecture's latest iteration incorporates DeepSeek-R1 as its primary reasoning engine, significantly enhancing the deliberative processing stream while maintaining the core dual-stream architecture. R1's exceptional performance across critical reasoning benchmarksincluding 97.3% accuracy on MATH-500 and 90.8% on MMLUmakes it particularly well-suited for REGENT's structured reasoning phase.

The integration also leverages R1's 128K context window to enable more comprehensive analysis of memory store contents during reasoning phases. This expanded context capacity allows the reasoning stream to consider a broader range of relevant memories when refining responses, leading to more nuanced, creative, and contextually appropriate outputs. We are in the early stages of exploring the possibilities with reinforcement-learned reasoning but these early results are promising.

R1's integration primarily affects the reasoning refinement phase of the processing pipeline, where it analyzes and refines outputs from the intuitive stream. The model's reinforcement learning foundation aligns naturally with REGENT's focus on autonomous growth and adaptation, enabling more sophisticated self-modification of the Lore store through reasoned deliberation.

# 4 Future Research

The REGENT architecture opens up numerous avenues for future investigation and development. Additional research may explore the quantitative aspects of split-mind processing, memory integration patterns, and comparative performance metrics. Of particular interest are the empirical effects of dual-stream processing on response quality and the long-term implications of autonomous memory management.

The integration of DeepSeek-R1's reinforcement learned reasoning capabilities represents an especially promising direction for future research. While still in early stages, initial results suggest that reinforcement learning could play a crucial role in developing more sophisticated reasoning mechanisms within the split-mind architecture. The combination of R1's advanced reasoning capabilities with REGENT's dual-stream processing may offer new insights into how autonomous systems can develop more robust and adaptable thinking patterns.

Of particular interest is how reinforcement learning might enhance the interaction between intuitive and deliberative processing streams. Early observations suggest that reinforcement-learned reasoning patterns could help bridge the gap between fast pattern matching and slower analytical thinking, potentially leading to more natural and effective cognitive processes. These possibilities warrant further investigation as the architecture continues to evolve.

# 5    Conclusion

REGENT V2 represents a practical implementation of dual-process theory in autonomous AI systems. By separating intuitive and reasoned processing while maintaining sophisticated memory integration, the architecture enables ems to think more naturally and effectively than traditional approaches.