# Recursive Neural-Symbolic Architectures for Deep Content Retrieval: Synergizing Recursive Language Models, Latent Hierarchical Indexing, and Reasoning-Based Search

## 1. Introduction: The Epistemological Crisis of Retrieval-Augmented Generation

The contemporary landscape of Large Language Model (LLM) deployment is dominated by the Retrieval-Augmented Generation (RAG) paradigm, a methodology that—while revolutionary—is beginning to exhibit fundamental limitations when applied to "Deep Research" tasks. As we transition from simple query-response interactions to complex, multi-hop reasoning over massive, unstructured corpora, the architectural deficiencies of standard RAG become increasingly obstructive. The prevailing approach treats documents as flat, homogenous sequences of text, arbitrarily fragmented into "chunks" and indexed via dense vector embeddings. This reductionist strategy, while computationally convenient, strips information of its structural skeleton—the hierarchical context, the semantic containment, and the logical flow that human authors painstakingly engineer into documents.

The user's inquiry identifies a critical intersection of three nascent technologies that promises to resolve these limitations: **Recursive Language Models (RLMs)**, which utilize Read-Eval-Print Loop (REPL) environments to decouple reasoning from context length; **PageIndex**, a framework proposing reasoning-based navigation over hierarchical tree structures rather than probabilistic vector similarity; and **Latent Hierarchical Indexing**, a necessary innovation for applying tree-based retrieval to documents lacking explicit Table of Contents (TOC) metadata.

This report proposes a unified, highly specialized content retrieval architecture—the **Recursive Neural-Symbolic Retriever (RNSR)**. This system does not merely "search" text; it constructs a traversable, symbolic environment from unstructured data and deploys an LLM agent to recursively explore, reason, and extract information within a persistent REPL environment. By synthesizing the "Prompt-as-Environment" paradigm of RLMs with the "Tree-Search" logic of PageIndex, and resolving the "No-TOC" constraint through advanced visual-semantic layout analysis, we articulate a robust methodology for the next generation of information retrieval systems.

### 1.1 The Limitations of Flat Retrieval in Deep Context

To understand the necessity of this new architecture, one must first deconstruct the failures of the current standard. "Flat" retrieval assumes that the semantic meaning of a text segment is self-contained. However, in complex documents—such as legal discovery bundles, financial audit trails, or technical manuals—meaning is inherently relational. A paragraph discussing "Termination Clauses" in a sub-contract is meaningless without the context of the "Master Service Agreement" it belongs to.

When standard RAG systems chunk this document into 512-token segments, they sever these relational links. A vector search might retrieve the "Termination Clause" based on keyword similarity but fail to retrieve the parent clause that dictates its jurisdiction. This phenomenon, often exacerbated by the "Lost-in-the-Middle" effect where LLMs neglect information in the center of long contexts, necessitates a shift from *probabilistic* retrieval (what looks similar?) to *deterministic* navigation (where does this concept live in the document structure?).

### 1.2 The Convergence of Recursive Methodologies

The proposed architecture leverages a "Recursive Turn" in AI research. We observe this recursion operating at three distinct levels:

1. **Recursive Data Structures:** The transformation of flat text into hierarchical trees (Latent TOCs).
2. **Recursive Inference:** The ability of the model to call itself (RLM) to solve sub-problems.
3. **Recursive Search:** The traversal algorithms (Depth-First, Breadth-First) used to navigate the data structure.

By integrating these levels, we move beyond simple "Question Answering" to "Reasoning via Planning" (RAP), where the system formulates a retrieval plan, executes it against a structured environment, and refines its understanding iteratively.[1]

---

## 2. The Recursive Language Model (RLM): The Engine of Infinite Context

The core engine of our proposed architecture is the Recursive Language Model (RLM). Traditional LLMs are constrained by their context window—the fixed buffer of tokens they can attend to at any given moment. While context windows are expanding (e.g., to 1M tokens), the quadratic complexity of attention mechanisms ($O(n^2)$) and the degradation of reasoning quality over long contexts remain unsolved physical constraints. The RLM paradigm, as detailed in recent literature [3], offers a radical alternative: treating the context not as input to the model, but as an *environment* the model interacts with.

## 2.1 The Prompt-as-Environment Abstraction

In the standard LLM paradigm, the prompt is ephemeral—it passes through the model layers and vanishes. In the RLM paradigm, the prompt is reified as a persistent variable within a Python REPL (Read-Eval-Print Loop) environment.

As described in *Zhang et al. (2025)*, the RLM initializes a computational environment where the massive input context (the "Long Prompt") is loaded into a variable, typically denoted as P or DOC_VAR.[3] The Neural Network itself does not ingest P. Instead, it ingests a system instruction that informs it of P's existence and provides a set of tools to interact with it.

This abstraction allows the model to perform "symbolic introspection." The model can write Python code to:

- Query the length of the document: len(DOC_VAR).
- Sample specific segments: print(DOC_VAR[0:1000]).
- Decompose the document: segments = split_by_regex(DOC_VAR, pattern).

This decoupling of **Memory** (the REPL state) from **Processing** (the LLM inference) allows the system to process effectively unbounded context lengths.[5] The model acts as a controller, paging data in and out of its immediate view as necessary, much like a CPU manages data between RAM and Cache.

## 2.2 The Recursive Loop and Variable Stitching

The "Recursive" nomenclature refers to the model's capability to divide a complex reasoning task into smaller, manageable sub-tasks and invoke instances of itself to solve them. This mirrors the divide-and-conquer algorithms in computer science.

In practice, this manifests as a "Recursive REPL." When the RLM encounters a task that exceeds its immediate reasoning capacity (e.g., "Summarize the changes in liability clauses across these 50 contracts"), it does not attempt to generate the answer in one pass. Instead, it generates code to:

1. Iterate through the contracts (sub-contexts).
2. For each contract, call the LLM API (invoking itself) with a specific sub-prompt: "Extract liability clause from this text."
3. Assign the output of each sub-call to a variable in the REPL: liability_1 = sub_llm(contract_1).
4. Synthesize the results stored in these variables into a final answer: final_summary = compare(liability_1, liability_2,...).[4]

This mechanism, known as **Variable Stitching**, is critical for preventing "Context Rot." In standard long-context generation, early details are often forgotten as the generation proceeds. In an RLM, intermediate results are "frozen" into variables. The model can retrieve

liability_1 with perfect fidelity hours later in the computation, ensuring that the final synthesis is based on exact data rather than fading activations.[4]

## 2.3 Optimization via Batching

A significant challenge with recursive calls is latency—sequential calls to an LLM can be slow. The RLM framework addresses this through batching strategies. The system prompt explicitly instructs the model to "batch as much information as reasonably possible into each call".[3]

For instance, instead of making 1,000 individual calls to summarize 1,000 paragraphs, the RLM writes code to group paragraphs into chunks of 5 and processes them in parallel threads (or batched API calls), aiming for an optimal token density (e.g., 200k characters per call). This balances the granularity of recursion with the throughput efficiency of the underlying hardware.

---

# 3. PageIndex and Reasoning-Based Retrieval: The Structural Paradigm

While RLMs provide the *process* for infinite context interaction, they require a *structure* to navigate. Navigating a 10,000-page unstructured string via print(DOC_VAR[i:j]) is inefficient; it is akin to reading a book by randomly opening pages. To enable effective navigation, we must organize the unstructured text into a structured index. This brings us to the **PageIndex** methodology.

## 3.1 From Vector Similarity to Tree Search

Current RAG systems rely almost exclusively on **Vector Similarity Search**. This involves

embedding text chunks into high-dimensional vectors and retrieving the top-$k$ chunks closest to the query vector. While effective for simple pattern matching, vector search suffers from a fundamental "semantic gap."

- **The "Vibe" Problem:** Vector search retrieves text that "sounds like" the query. If a user asks, "What are the limitations of the proposed method?", vector search might retrieve a section titled "Limitations" from a completely different context (e.g., limitations of the *dataset* rather than the *algorithm*).
- **Lack of Global Awareness:** Vector search is local. It cannot answer queries that require understanding the document's overall argument structure.

**PageIndex** [6] proposes a paradigm shift: **Reasoning-Based Retrieval**. Instead of matching embeddings, the system performs a **Tree Search** over a hierarchical index of the document.

- **The Tree Index:** The document is represented as a tree where the Root is the document

title, branches are sections (e.g., "1.0 Introduction"), and leaves are paragraphs or data tables.

- **The Agentic Walk:** Retrieval becomes a decision process. The agent starts at the Root. It inspects the children nodes (the high-level Table of Contents). It reasons: "Given the user's query about 'Financial Risk,' the answer is likely within the 'Risk Factors' branch, not the 'Executive Summary' branch." It then descends into the 'Risk Factors' node and repeats the process.

This approach transforms retrieval complexity from $O(N)$ (scanning all chunks) to $O(\log N)$ (traversing the tree depth). More importantly, it ensures **Contextual Integrity**. By navigating down the tree, the agent implicitly understands that the paragraph it eventually retrieves is part of "Section 3.2," which is part of "Chapter 3," providing the necessary scope for accurate interpretation.

## 3.2 Traceability and Explainability

A critical requirement for specialized retrieval tools (especially in regulated industries) is explainability. Vector RAG is opaque; explaining why a specific vector was closest to another is mathematically abstract and unintuitive.

In contrast, PageIndex enables **Traceable Retrieval**.[8] The path taken by the agent (Root $\rightarrow$ Node A $\rightarrow$ Node A.1 $\rightarrow$ Leaf) serves as an explicit audit trail. We can reconstruct the "Chain of Reasoning" that led the system to a specific piece of information: "The system selected 'Section 4' because the query pertained to 'Methodology,' then selected 'Subsection 4.2' because it specifically mentioned 'Data Preprocessing'." This transparency allows for debugging and verification, essential for high-stakes "Deep Research."

# 4. The "No-TOC" Challenge: Generating Latent Hierarchies

The user's request introduces a significant constraint: *"using sections... but not using Table of contents... As if you have a sufficiently large document that is composed of sub documents all without headings"*.

This is a common scenario in real-world data processing. Merged PDFs, legacy scans, and concatenated reports often lack explicit metadata (bookmarks) or consistent formatting that a simple parser could recognize as a Table of Contents. If PageIndex relies on a tree, and the document has no explicit tree, the system fails.

To bridge this gap, we must engineer a **Latent Hierarchy Generator**. We must synthesize a

tree structure from the raw signals available in the document: **Visual Geometry** (Layout) and **Semantic Cohesion** (Topic).

## 4.1 Visual-Geometric Analysis: The "Latent TOC"

Humans do not need an explicit Table of Contents to recognize the structure of a document. We rely on visual cues: a line of text that is larger, bolder, and centered is intuitively recognized as a header. A block of text indented from the left margin is recognized as a quote or list. We can codify these intuitions using computer vision and geometric analysis algorithms.

### 4.1.1 The Recursive XY-Cut Algorithm

The **Recursive XY-Cut (RXYC)** algorithm [9] is a top-down page segmentation technique that is particularly effective for discovering document structure without relying on text content.

- **Mechanism:** The algorithm views the document page as a binary image. It calculates the "projection profile"—the sum of black pixels—along the horizontal (X) and vertical (Y) axes.
- **Valleys as Separators:** Significant drops (valleys) in the projection profile correspond to whitespace channels. A wide horizontal white gap usually separates a header from the body text or one section from another. A vertical white gap separates columns.
- **Recursion:** The algorithm cuts the page at the widest valleys, dividing it into smaller rectangles. It then recursively applies the same logic to each sub-rectangle.
- **Tree Generation:** This process naturally produces a tree of bounding boxes. The largest bounding boxes at the top of the hierarchy (detected early in the cut process) often correspond to major structural elements (Titles, Headers), while smaller, deeply nested boxes correspond to paragraphs or table cells.

By analyzing the nesting relationship of these boxes, we can infer a **Geometric Hierarchy**. If Box A is physically above Box B, and Box A spans the full width while Box B is narrow (columnar), Box A is likely the parent/header of Box B.[11]

### 4.1.2 LayoutLM: Multimodal Structure Recognition

To refine the geometric cuts, we employ **LayoutLM** (specifically LayoutLMv3), a multimodal Transformer model pre-trained on document layouts.[12] LayoutLM ingests three modalities simultaneously:

1. **Text:** The token sequence (from OCR).
2. **Layout:** The 2D bounding box coordinates of each token.
3. **Image:** The visual features of the document patches.

Unlike pure NLP models (BERT) which only see text, LayoutLM "sees" that a specific token is bold, 24pt font, and centered. We can fine-tune LayoutLM on token classification tasks to label each geometric block identified by XY-Cut as Header, Footer, Title, Body, Table, or

Caption.

**Constructing the Tree:**

Once every text block is classified, we execute a linear scan to build the tree:

1. Identify all blocks labeled Header.
2. Analyze their font sizes (height of bounding box).
3. Establish a hierarchy based on size: The largest headers are Level 1 nodes (Chapters). Slightly smaller headers are Level 2 nodes (Sections).
4. Assign all Body text blocks to the nearest preceding Header of the appropriate level.

This process essentially "hallucinates" a rigid Table of Contents for documents that never had one, creating the necessary scaffold for the PageIndex reasoning engine.

## 4.2 Semantic Boundary Detection: The "Topic" Approach

In cases where visual cues are absent (e.g., plain text files, flat OCR dumps), we must rely on **Semantic Boundary Detection**. This involves identifying points in the text stream where the topic shifts significantly.

### 4.2.1 Semantic Splitters and Embedding Analysis

We utilize **Semantic Splitters** (such as those implemented in LlamaIndex [15]) which operate on the "Manifold Hypothesis"—that semantically related text occupies a continuous region in the embedding vector space.

- **Sliding Window:** We define a sliding window (e.g., 3 sentences) that moves across the text.
- **Coherence Metric:** At each step, we calculate the cosine similarity between the embedding of the current window and the next window.
- **Breakpoint Detection:** A sharp drop in similarity (a local minimum) indicates a **Semantic Breakpoint**—a transition from one topic to another.

### 4.2.2 Hierarchical Clustering (H-SBM)

For more advanced segmentation, we can employ unsupervised hierarchical clustering techniques.[16] By clustering the sentence embeddings, we can discover "Latent Topics" at various resolutions.

- **Micro-Clusters:** Groups of 5-10 sentences forming a paragraph-level thought.
- **Macro-Clusters:** Groups of Micro-Clusters forming a chapter-level theme.

We can then structure these clusters into a tree. The "Title" of each node in this semantic tree is generated generatively: we feed the text of the cluster to a summarization LLM with the prompt *"Generate a concise 5-word header for this text section."* The resulting header becomes the handle for the PageIndex agent.

# 5. Architectural Synthesis: The Recursive Neural-Symbolic Retriever

We now synthesize these components—RLM/REPL, PageIndex, and Latent Hierarchy—into a unified system architecture: the **Recursive Neural-Symbolic Retriever (RNSR)**. This system is designed to ingest massive, unstructured/heterogeneous documents and provide a "Deep Research" interface.

## 5.1 System Architecture Diagram

The architecture consists of three distinct processing phases: **Ingestion (The Builder)**, **Indexing (The Map)**, and **Retrieval (The Navigator)**.

### Phase I: The Latent Structure Ingestion Engine

This phase is responsible for converting the raw "soup" of pixels and text into a structured JSON tree. It operates offline, before any query is received.

**Technology Stack:**

- **Unstructured.io:** For initial partitioning and OCR.[17]
- **LayoutLMv3:** For classifying layout elements.[12]
- **Recursive XY-Cut Algorithm:** For geometric segmentation.[9]

**Workflow:**

1. **Partitioning:** The document (PDF/Image) is passed to partition_pdf with strategy="hi_res". This invokes detection models to separate text, tables, and images.
2. **Element Classification:** The system iterates through the extracted elements. It checks metadata for font_size, font_weight, and coordinates.
3. **Latent Tree Construction Algorithm:**
   - *Input:* List of LayoutElement objects.
   - *State:* A Stack of active TreeNodes.
   - *Logic:*
     - Iterate through elements in reading order.
     - If element.type == "Title" or is_visually_prominent(element):
       - Determine hierarchy_level based on font size relative to the document global mode.
       - Pop items from Stack until Stack.top().level < element.level.
       - Create new Node, add as child to Stack.top(), push to Stack.
     - Else (Body Text):
       - Append text to Stack.top().content.
4. **Semantic Fallback:** If the Layout Analysis yields a flat tree (i.e., no headers detected),

the system triggers the **Semantic Splitter**. It partitions the text stream based on embedding similarity drops [15] and generates synthetic titles for each chunk using a fast LLM (e.g., GPT-3.5-Turbo).

## Phase II: The Hierarchical Index

Once the tree structure is defined, it must be indexed for efficient access by the LLM agent. We utilize **LlamaIndex** primitives, specifically the IndexNode.[19]

### The IndexNode Data Structure:

Each node in our tree is converted into an IndexNode.

- **index_id:** A unique UUID.
- **text:** A *summary* of the node's content (not the full text). This is critical. The agent navigates based on summaries to save context.
- **obj:** A pointer to the child nodes or the full-text payload.
- **metadata:** Includes page_number, bounding_box, hierarchy_level.

This creates a "Skeleton Index." The full content (the meat) is stored in a separate Key-Value store, retrieved only when the agent explicitly requests to "read" a leaf node.

## Phase III: The Recursive REPL Agent (The Navigator)

This is the runtime engine where the RLM comes into play. We implement this using **LangGraph** [21] to manage the state of the recursive agent.

### The REPL Environment:

The agent does not interact with the index via vector search. It interacts via a Python REPL pre-loaded with a Navigator API.

- list_children(node_id): Returns the summaries of immediate children.
- read_node(node_id): Returns the full text content of a node.
- search_index(query): A keyword/hybrid search to find initial entry points in the tree.

**The Recursive Logic (Reasoning via Planning - RAP):** When a user asks a complex question (e.g., *"Compare the indemnification clauses in the 2023 and 2024 agreements"*), the Agent executes a **Reasoning via Planning (RAP)** loop [1]:

1. **Planning:** The agent calls list_children('root'). It sees nodes for "2023 Agreement" and "2024 Agreement."
2. **Decomposition:** It decomposes the query into two sub-tasks: "Find indemnification in 2023" and "Find indemnification in 2024."
3. **Recursive Execution:**
   - *Sub-task A:* The agent enters the "2023 Agreement" node. It recursively lists children until it finds a node likely to contain "Indemnification" (perhaps under "Legal Terms"). It calls read_node() on that leaf.

- *Sub-task B:* Repeats the process for 2024.
4. **Variable Accumulation:** The text found in Step 3 is stored in REPL variables (clause_2023, clause_2024).
5. **Synthesis:** The agent uses the LLM to compare the content of these two variables and generates the final response.

---

# 6. Detailed Technical Implementation: Building the "Latent Table of Contents"

The success of this architecture hinges on the robustness of the **Latent Hierarchy Generator**. If the inferred tree is garbage, the agent will get lost. Here, we detail the specific algorithmic logic for reconstructing structure from unstructured PDFs.

## 6.1 Font-Size Histogram Analysis

The most reliable proxy for document hierarchy in the absence of tags is font size. However, simply picking "large text" is insufficient. We need a statistical approach.

**Algorithm:**

1. **Extraction:** Extract all text elements $E = \{e_1, e_2, ...\}$ with attributes $\{text, size, weight, x, y\}$.
2. **Histogramming:** Create a frequency histogram of size values, rounded to the nearest integer.
3. **Mode Detection:** Identify the global mode $M_{size}$. This represents the "Body Text."
4. **Outlier Detection:** Identify all elements where $e.size > M_{size} + \sigma$ (where $\sigma$ is a significance threshold, typically 2pt). These are candidate headers.
5. **Clustering:** Cluster the candidate headers into levels using 1D k-means clustering.
   - *Cluster 1 (Largest):* Level 1 Headers (Titles).
   - *Cluster 2 (Medium):* Level 2 Headers (Section Heads).
   - *Cluster 3 (Small):* Level 3 Headers (Subsection Heads).
6. **Tree Assembly:** Iterate through the document. Use a stack-based parser to reconstruct the tree. When a Level 1 header is found, push it to the stack. When a Level 2 header is found, append it as a child of the current Level 1 header. If a new Level 1 header is found, pop the previous one and push the new one.

## 6.2 Handling "L-Shaped" and Columnar Layouts

A major failure mode of simple parsing is complex layouts (e.g., a figure spanning two

columns, or an L-shaped text wrap). The **Recursive XY-Cut** handles this.[9]

**Implementation:**

We implement a Python class RecursiveXYCutter.

Python

```python
def recursive_xy_cut(image_pixels, bounds):
    # Calculate projection profiles
    x_proj = np.sum(image_pixels[bounds], axis=0)
    y_proj = np.sum(image_pixels[bounds], axis=1)

    # Find valleys (gaps of white space)
    y_gaps = find_valleys(y_proj, threshold)

    if y_gaps:
        # Split horizontally (Top/Bottom sections)
        for gap in y_gaps:
            recursive_xy_cut(image_pixels, gap_bounds)
    else:
        # Try splitting vertically (Columns)
        x_gaps = find_valleys(x_proj, threshold)
        if x_gaps:
            for gap in x_gaps:
                recursive_xy_cut(image_pixels, gap_bounds)
        else:
            # Leaf node (Text Block)
            register_leaf_node(bounds)
```

This recursive splitting creates a tree of bounding boxes. We then map the OCR text elements into these boxes. If a box contains a LayoutLM-detected "Header" and subsequent boxes contain "Body Text," we link them structurally.

## 6.3 Semantic Refinement for "Flat" Documents

For documents that are essentially flat text (e.g., plain text legal dumps), visual analysis fails. Here, we use **Semantic Boundary Detection**.

**Technique:** We integrate the **LlamaIndex SemanticSplitterNodeParser**.[15]

1. **Embedding:** We use a high-throughput embedding model (e.g., text-embedding-3-small or bge-m3).
2. **Breakpoint Calculation:** We calculate the cosine distance between sliding windows (window size=3 sentences).
3. **Thresholding:** We set a breakpoint_percentile_threshold (e.g., 95th percentile). Any distance above this percentile is marked as a section break.
4. **Synthetic Header Generation:**
   - For each identified section, we execute an LLM call:*Prompt:* "Read the following text segment. Generate a descriptive, hierarchical title for it (e.g., 'Section 3: Liability Limitations'). Return ONLY the title."
   - This synthetic title acts as the "Virtual Header" in our PageIndex tree, giving the agent a semantic handle to grasp during traversal.

---

# 7. The Recursive Agent Implementation Strategy

The implementation of the agentic navigator requires sophisticated state management. We utilize **LangGraph** to model the recursive control flow.

## 7.1 State Definition

The agent's state is a typed dictionary that persists across the recursive steps:

Python

```python
class AgentState(TypedDict):
    query: str
    current_node_id: str
    path_history: List[str]
    accumulated_context: Dict[str, str]  # The REPL variables
    plan: List[str]
    final_answer: Optional[str]
```

## 7.2 The "Tree of Thoughts" (ToT) Prompting Strategy

We do not use standard Chain-of-Thought (CoT). We use **Tree of Thoughts (ToT)** [23], which explicitly encourages the model to explore multiple branches of the index before committing to a path.

**System Prompt:**

"You are a Deep Research Agent navigating a document tree.

You are currently at Node: {current_node_summary}.

Children Nodes: {children_summaries}.

Your Goal: {query}.

1. **Evaluate:** For each child node, estimate the probability that it contains relevant information.
2. **Plan:** Select the top-k most promising nodes.
3. **Recursion:** Generate a plan to visit these nodes.
4. **Backtrack:** If a node yields no useful info, report 'Dead End' and return to parent."

## 7.3 Integration with LlamaIndex RecursiveRetriever

LlamaIndex provides a RecursiveRetriever class [19], but it is typically used as a black-box query engine. For our architecture, we "crack open" this class.

- We use the IndexNode logic to link summaries to raw text.
- However, instead of using the automated query_engine.query(), we expose the graph traversal methods (get_children, get_node_content) as **Tools** to the LangGraph agent.
- This gives the agent **agency**. It allows the agent to decide *not* to read a node, or to read *all* children, or to jump back to the root—behaviors that standard RAG pipelines cannot execute.

---

# 8. Comparative Analysis: Vector RAG vs. PageIndex vs. Recursive Neural-Symbolic

To demonstrate the efficacy of this proposed architecture, we compare it against existing paradigms.

| Feature | Standard Vector RAG | Standard PageIndex | Recursive Neural-Symbolic (Proposed) |
|---|---|---|---|
| **Ingestion** | Flat chunking (fixed size). | Requires explicit TOC. | **Latent Hierarchy Generation** (Visual/Semantic). |
| **Context** | Limited by context | Limited by traversal | **Infinite (REPL** |

| | window. | depth. | **Variable Stitching).** |
|---|---|---|---|
| **Retrieval** | Similarity Search (Probabilistic). | Tree Search (Deterministic). | **Agentic Reasoning via Planning (RAP).** |
| **Handling No-TOC** | Excellent (ignores structure). | Fails. | **Excellent (reconstructs structure).** |
| **Latency** | Low (milliseconds). | Medium (sequential). | **High** (requires multiple agent steps). |
| **Explainability** | Low (black box). | High (traceable path). | **Maximum** (full execution trace). |
| **Use Case** | Factoid QA. | Structured Doc Search. | **Deep Research / Complex Synthesis.** |

## 8.1 The Latency-Accuracy Trade-off

It is crucial to acknowledge that the proposed RNSR system is slower than vector RAG. Vector search takes milliseconds. An agentic walk through a tree, involving multiple LLM inference steps to evaluate node summaries and plan the next move, may take 10-60 seconds.

However, for the intended use case—"Deep Research" over massive documents where accuracy is paramount (e.g., "Find all contradictions between these two 500-page mergers")—this latency is negligible compared to the time cost of human review. The system trades speed for **Recall** and **Contextual Understanding**.

# 9. Conclusion and Future Outlook

The user's query highlights a critical gap in current AI information retrieval: the inability to apply structured, reasoning-based retrieval methods to unstructured, "messy" documents that lack explicit organization. The proposed **Recursive Neural-Symbolic Retriever** fills this

gap.

By synergizing the **Recursive Language Model**'s ability to handle infinite context via REPL environments with **PageIndex**'s deterministic tree-search logic, and underpinning the entire system with **Latent Hierarchy Generation** derived from LayoutLM and Recursive XY-Cut algorithms, we create a tool that is not merely a search engine, but a digital researcher.

This system does not just find words; it understands structure. It does not just retrieve chunks; it navigates arguments. As LLMs continue to evolve from chatbots to agents, architectures like RNSR—which emphasize structure, agency, and recursion—will become the standard for interacting with the vast repository of human knowledge.

## 9.1 Recommendations for Development

1. **Prioritize the Latent Ingestion Pipeline:** The quality of the retrieval is strictly bounded by the quality of the constructed tree. Investing in fine-tuning LayoutLMv3 on domain-specific documents (e.g., finance, legal) to accurately detect headers is the highest-ROI activity.
2. **Adopt LangGraph for Control Flow:** The cyclic, stateful nature of recursive retrieval cannot be easily modeled in linear chains. LangGraph's support for cycles and state persistence is essential.[22]
3. **Hybridize Search:** Implement "Hybrid Search" at the agent level. Give the agent a vector_search tool as a shortcut. The agent can use vector search to find a starting node (e.g., "Jump to 'Risk' section") and then switch to tree traversal for local exploration. This balances the speed of RAG with the precision of PageIndex.

The future of retrieval is not flat; it is recursive, hierarchical, and agentic. The architecture detailed in this report provides the blueprint for building it.

---

**Citations:**

- [3]: "Recursive Language Models," arXiv:2512.24601 (RLM, REPL, Prompt-as-Variable).
- [7]: "PageIndex," VectifyAI (Reasoning-based RAG, Tree Search).
- [12]: "LayoutLM," Microsoft (Multimodal layout analysis).
- [9]: "Recursive XY-Cut" (Geometric segmentation).
- [15]: Semantic Boundary Detection & Splitters.
- [23]: "Tree of Thoughts" (Search strategy).
- [1]: "Reasoning via Planning" (RAP).
- [17]: "Unstructured.io" (PDF Partitioning).
- [21]: "LangGraph" (Agentic orchestration).
- [19]: "LlamaIndex" (IndexNode, RecursiveRetriever).

**Works cited**

1. Reasoning via Planning (RAP) - Emergent Mind, accessed on January 24, 2026, https://www.emergentmind.com/topics/reasoning-via-planning-rap
2. RAP-RAG: A Retrieval-Augmented Generation Framework with Adaptive Retrieval Task Planning - MDPI, accessed on January 24, 2026, https://www.mdpi.com/2079-9292/14/21/4269
3. Recursive Language Models - arXiv, accessed on January 24, 2026, https://arxiv.org/pdf/2512.24601
4. Recursive Language Models - RLM - arXiv, accessed on January 24, 2026, https://arxiv.org/html/2512.24601v1
5. Recursive Language Models - arXiv, accessed on January 24, 2026, https://arxiv.org/abs/2512.24601
6. VectifyAI/pageindex-mcp: MCP server for PageIndex. PageIndex is a vectorless reasoning-based RAG system which uses multi-step reasoning and tree search to retrieve information like a human expert would. - GitHub, accessed on January 24, 2026, https://github.com/VectifyAI/pageindex-mcp
7. VectifyAI/PageIndex: PageIndex: Document Index for ... - GitHub, accessed on January 24, 2026, https://github.com/VectifyAI/PageIndex
8. PageIndex, accessed on January 24, 2026, https://pageindex.ai/
9. A modified recursive x-y cut algorithm for solving block ordering problems - SciSpace, accessed on January 24, 2026, https://scispace.com/pdf/a-modified-recursive-x-y-cut-algorithm-for-solving-block-4z12yah9d3.pdf
10. Replicating the Recursive XY-Cut Algorithm for document layout to set a custom reading order template : r/computervision - Reddit, accessed on January 24, 2026, https://www.reddit.com/r/computervision/comments/e1wory/replicating_the_recursive_xycut_algorithm_for/
11. High Performance Document Layout Analysis, accessed on January 24, 2026, https://www.dfki.de/fileadmin/user_upload/import/2000_HighPerfDocLayoutAna.pdf
12. How Can LayoutLM Transform Text Extraction? - Nitor Infotech, accessed on January 24, 2026, https://www.nitorinfotech.com/blog/how-can-layoutlm-transform-text-extraction/
13. Information Extraction — Part 3. LayoutLM v3 — Explanation and... | by Tejpal Kumawat | Medium, accessed on January 24, 2026, https://medium.com/@tejpal.abhyuday/information-extraction-part-3-9c2487ec4930
14. layoutlmv2: multi-modal pre-training for visually-rich document understanding - Microsoft, accessed on January 24, 2026, https://www.microsoft.com/en-us/research/wp-content/uploads/2021/01/layoutlmv2.pdf
15. Node Parser Modules | LlamaIndex Python Documentation, accessed on January

24, 2026,
https://developers.llamaindex.ai/python/framework/module_guides/loading/node_parsers/modules/

16. Hierarchy-Agnostic Unsupervised Segmentation: Parsing Semantic Image Structure - NIPS, accessed on January 24, 2026,
https://proceedings.neurips.cc/paper_files/paper/2024/file/b31c332c4cebcec31b788400b47c94b3-Paper-Conference.pdf

17. Partition configuration - Unstructured, accessed on January 24, 2026,
https://docs.unstructured.io/open-source/ingestion/ingest-configuration/partition-configuration

18. Partitioning - Unstructured, accessed on January 24, 2026,
https://docs.unstructured.io/open-source/core-functionality/partitioning

19. Recursive Retriever + Query Engine Demo | LlamaIndex Python Documentation, accessed on January 24, 2026,
https://developers.llamaindex.ai/python/examples/query_engine/pdf_tables/recursive_retriever/

20. A Cheat Sheet and Some Recipes For Building Advanced RAG - LlamaIndex, accessed on January 24, 2026,
https://www.llamaindex.ai/blog/a-cheat-sheet-and-some-recipes-for-building-advanced-rag-803a9d94c41b

21. LangGraph overview - Docs by LangChain, accessed on January 24, 2026,
https://langchain-ai.github.io/langgraph/

22. Designing Advanced RAG-Based Retrieval Systems with LangGraph: A Practical Guide, accessed on January 24, 2026,
https://mihirinamdar.medium.com/designing-advanced-rag-based-retrieval-systems-with-langgraph-a-practical-guide-b8ff7fec8800

23. What is Tree Of Thoughts Prompting? - IBM, accessed on January 24, 2026,
https://www.ibm.com/think/topics/tree-of-thoughts

24. Tree of Thoughts (ToT) - Prompt Engineering Guide, accessed on January 24, 2026, https://www.promptingguide.ai/techniques/tot

25. llama_index/llama-index-core/llama_index/core/retrievers/recursive_retriever.py at main · run-llama/llama_index - GitHub, accessed on January 24, 2026,
https://github.com/run-llama/llama_index/blob/main/llama-index-core/llama_index/core/retrievers/recursive_retriever.py

26. Recursive Language Models | Alex L. Zhang, accessed on January 24, 2026,
https://alexzhang13.github.io/blog/2025/rlm/

27. VectifyAI/pageindex-mcp · GitHub - Workflow runs, accessed on January 24, 2026, https://github.com/VectifyAI/pageindex-mcp/actions

28. LayoutLM: Pre-training of Text and Layout for Document Image Understanding - arXiv, accessed on January 24, 2026, https://arxiv.org/pdf/1912.13318

29. XY-Cut++: Advanced Layout Ordering via Hierarchical Mask Mechanism on a Novel Benchmark - arXiv, accessed on January 24, 2026,
https://arxiv.org/html/2504.10258v1

30. Semantic Boundary Detection for Improving RAG on Real-time Agents - inferable.ai, accessed on January 24, 2026,

https://www.inferable.ai/blog/posts/semantic-boundary-detection-rag-agents

31. Boosting Semantic Segmentation by Conditioning the Backbone with Semantic Boundaries - PMC - NIH, accessed on January 24, 2026, https://pmc.ncbi.nlm.nih.gov/articles/PMC10422643/