

KNOWY

Agentic Multimodal Architecture for Personal Knowledge Compilation,
Epistemic Evolution, and Expert-Level Reasoning

Technical Design Document

Alberto Espinosa

January 16, 2026

Abstract

This document presents the comprehensive technical design for KNOWY, an agentic multimodal architecture designed to compile personal knowledge collections, support epistemic evolution, and enable expert-level reasoning. The system addresses the fundamental challenge of epistemic integrity in personal AI systems through a rigorously layered architecture that distinguishes between shallow operational agents and deep epistemic agents. KNOWY integrates advanced multimodal knowledge acquisition including structured video analysis, diagrammatic explanation generation, and a command-line interaction interface for local-first operation. This design specification details the system architecture, component implementations, data flows, algorithms, and deployment strategies necessary for realizing a robust epistemic infrastructure.

Contents

1	Executive Summary	4
1.1	Vision and Scope	4
1.2	Core Design Principles	4
1.3	Key Innovations	4
2	System Architecture	4
2.1	Architectural Overview	4
2.2	Architectural Guarantees	6
2.3	Technology Stack	6
3	Layer 0: Universal Ingestion	6
3.1	Purpose and Epistemic Question	6
3.2	Canonical Source Descriptor Schema	7
3.3	Ingestion Pipeline Architecture	8
3.4	Browser Extension Integration	8
3.5	Format Detection Algorithm	10
3.6	Storage Schema	10
4	Layer 1: Exploration and Expansion	10
4.1	Purpose and Epistemic Question	10
4.2	Expansion Agent Architecture	11
4.3	Knowledge Artifact Data Structure	12
4.4	Reference Resolution Strategies	13
4.5	Deduplication Algorithm	13

5	Layer 2: Multimodal Acquisition and Scraping	14
5.1	Purpose and Epistemic Question	14
5.2	Acquisition Pipeline Dispatcher	14
5.3	Video Analysis Pipeline Integration	15
5.3.1	Video Analysis Architecture	16
5.3.2	Shot Segmentation Algorithm	17
5.3.3	Transcript-Visual Alignment	17
5.3.4	Semantic Frame Analysis	18
5.3.5	Multimodal Chunk Generation	19
5.4	PDF Acquisition Pipeline	20
5.5	Academic Paper Pipeline	21
6	Layer 3: Multi-Resolution Knowledge Representation	22
6.1	Purpose and Epistemic Question	22
6.2	Four-Level Abstraction Hierarchy	23
6.3	Representation Architecture	23
6.4	Level 0: Raw Content Storage	23
6.5	Level 1: Section-Level Representation	24
6.6	Level 2: Document-Level Synthesis	25
6.7	Level 3: Cross-Document Thematic Synthesis	27
6.8	Dynamic Resolution Selection	28
7	Layer 4: Classification and Ontology Construction	29
7.1	Purpose and Epistemic Question	29
7.2	Ontology Structure	29
7.3	Topic Extraction and Clustering	30
7.4	Redundancy Detection	33
8	Layer 5: Expert Reasoning and Teaching	34
8.1	Purpose and Epistemic Question	34
8.2	Expert Reasoning Architecture	34
8.3	Pedagogical Mode	37
9	Layer 6: Diagrammatic Explanation	39
9.1	Purpose and Epistemic Question	39
9.2	Diagram Generation Architecture	39
9.3	Concept Map Generation	41
9.4	Mermaid Diagram Generation	42
10	Layer 7: Deep External Research Agent	43
10.1	Purpose and Epistemic Question	43
10.2	DERA Architecture	44
10.3	Research Trigger Detection	47
10.4	Knowledge Versioning	48
11	User Interaction Layer: CLI Interface	49
11.1	CLI Architecture	49
11.2	CLI Usage Examples	50
12	System Integration and Orchestration	51
12.1	LangGraph Orchestration	51
12.2	Deployment Architecture	53

13 Performance Optimization	54
13.1 Caching Strategy	54
13.2 Parallel Processing	55
14 Evaluation and Metrics	56
14.1 Epistemic Integrity Metrics	56
14.2 Performance Metrics	57
15 Testing Strategy	58
15.1 Unit Testing	58
15.2 Integration Testing	59
16 Security and Privacy	60
16.1 Data Privacy	60
17 Future Extensions	61
17.1 Planned Enhancements	61
18 Conclusion	61
A Appendix A: Database Schemas	62
B Appendix B: Configuration Files	64
C Appendix C: API Reference	65

1 Executive Summary

1.1 Vision and Scope

KNOWY represents a paradigm shift from retrieval-oriented personal AI tools toward epistemic infrastructures. The system treats personal knowledge collections—spanning academic papers, courses, multimedia lectures, and saved references—not as static information repositories but as evolving epistemic foundations requiring rigorous integrity guarantees.

1.2 Core Design Principles

Foundational Design Principles

1. **Epistemic Integrity Over Information Access:** Prioritize correctness, uncertainty quantification, and hallucination prevention over comprehensive retrieval.
2. **Shallow-Deep Agent Separation:** Enforce strict boundaries between operational agents (bounded transformations) and epistemic agents (boundary-crossing reasoning).
3. **Multi-Resolution Knowledge Representation:** Maintain explicit abstraction hierarchies from raw artifacts to cross-document synthesis.
4. **Controlled External Research:** Enable knowledge evolution through principled, versioned external inquiry rather than unconstrained web access.
5. **Local-First Operation:** Ensure privacy, reproducibility, and system composability through command-line interfaces and local model execution.
6. **Multimodal Native Design:** Treat video, diagrams, and text as first-class knowledge modalities with specialized processing pipelines.

1.3 Key Innovations

- **Layered Epistemic Architecture:** Seven-layer shallow agent pipeline with isolated deep research component
- **Video-as-Knowledge Pipeline:** Integration of video-analyzer for shot segmentation, transcript alignment, and visual-semantic extraction
- **Diagrammatic Explanation Generation:** Automated schematic diagram creation for pedagogical enhancement
- **CLI-First Interaction Model:** Built on llm framework enabling scriptability and workflow integration
- **Versioned Knowledge Evolution:** Temporal knowledge management with explicit uncertainty tracking

2 System Architecture

2.1 Architectural Overview

The KNOWY architecture consists of eight distinct processing layers orchestrated through a graph-based execution framework. Seven layers implement shallow agents with bounded epistemic scope, while a single deep agent layer enables controlled external research.

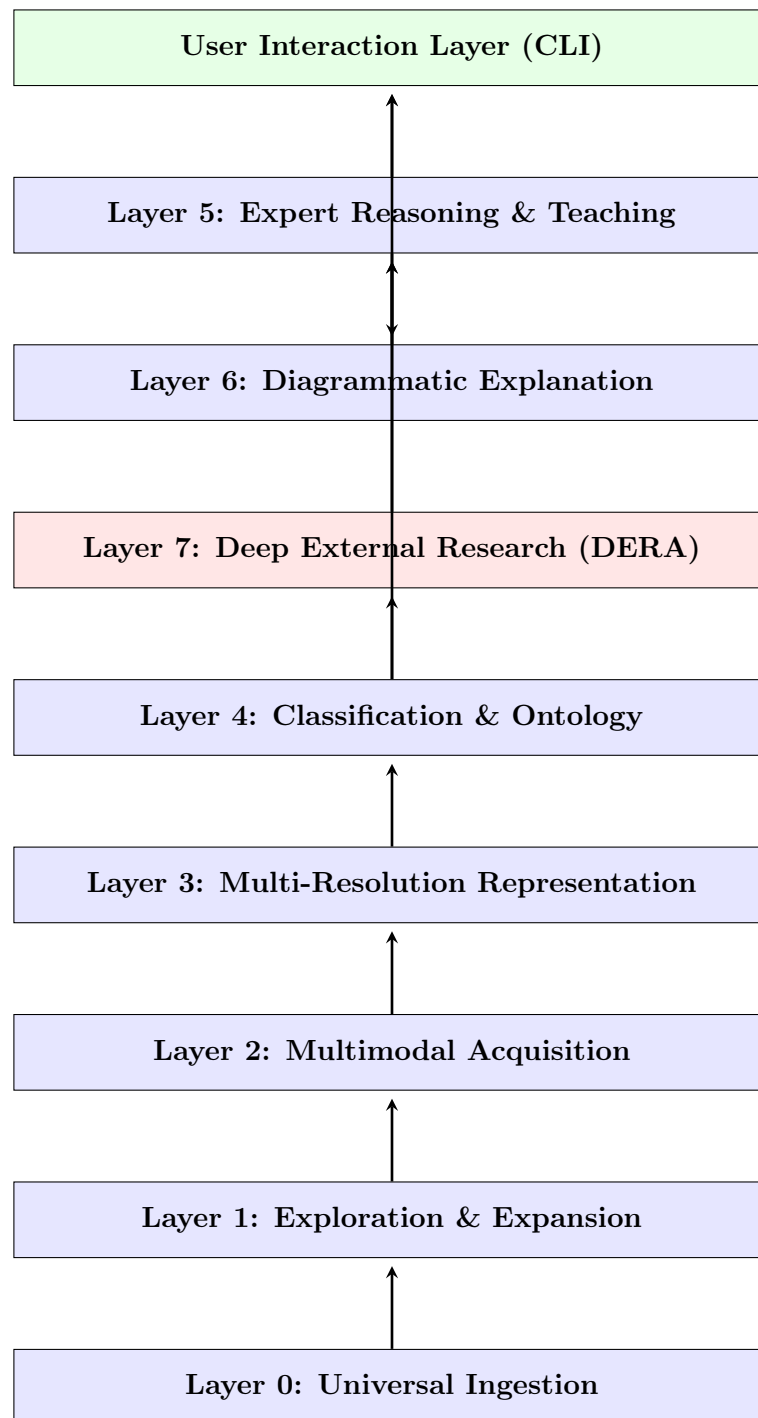


Figure 1: KNOWY Layered Architecture with Shallow-Deep Agent Separation

2.2 Architectural Guarantees

Guarantee	Implementation Mechanism
Epistemic Boundary Preservation	Type-level enforcement in orchestration graph; shallow agents cannot invoke external knowledge sources
Hallucination Minimization	Multi-resolution grounding with explicit provenance tracking; uncertainty quantification at reasoning layer
Temporal Consistency	Versioned knowledge representations with timestamp-based validity windows
Multimodal Coherence	Unified semantic representation across text, video frames, and diagrams with cross-modal alignment
Reproducibility	Deterministic agent execution with comprehensive logging; CLI-based workflow capture

Table 1: Core Architectural Guarantees

2.3 Technology Stack

Component	Technology	Justification
Agent Orchestration	LangGraph	Graph-based execution with conditional routing and state management
Local LLM Runtime	Ollama	Unified API for multiple local models; efficient inference
Text Models	LLaMA 3.x, Qwen, Gemma	State-of-the-art reasoning with local deployment
Multimodal Models	LLaMA 3.2 Vision	Native image understanding for video frame analysis
Video Analysis	video-analyzer	Shot segmentation, transcript alignment, visual-semantic extraction
CLI Framework	llm (Simon Willison)	Plugin architecture, scriptability, model abstraction
Vector Database	FAISS / Qdrant	High-performance similarity search with filtering
Structured Storage	PostgreSQL + Neo4j	Relational metadata + graph-based ontologies
Diagram Generation	Graphviz, Mermaid	Programmatic schematic creation

Table 2: Technology Stack Components

3 Layer 0: Universal Ingestion

3.1 Purpose and Epistemic Question

Epistemic Question: *"What reference has the user saved, independent of format or source?"*

Layer 0 establishes the entry point for all knowledge artifacts entering the system. It normalizes heterogeneous inputs from browsers, LinkedIn saves, manual entries, and API integrations into a canonical source descriptor format.

3.2 Canonical Source Descriptor Schema

Listing 1: Canonical Source Descriptor Data Structure

```
1 from dataclasses import dataclass
2 from typing import List, Optional, Dict
3 from datetime import datetime
4 from uuid import UUID
5
6 @dataclass
7 class ContentHints:
8     contains_pdf: bool
9     contains_video: bool
10    contains_external_links: bool
11    contains_code_repository: bool
12    estimated_length_minutes: Optional[int]
13    primary_language: Optional[str]
14
15 @dataclass
16 class CanonicalSourceDescriptor:
17     source_id: UUID
18     origin: str # linkedin / browser / manual / api / other
19     raw_reference: str
20     raw_text: Optional[str]
21     links: List[str]
22     content_hints: ContentHints
23     timestamp: datetime
24     ingestion_metadata: Dict[str, any]
25
26     def to_json(self) -> dict:
27         """Serialize to storage format"""
28         pass
29
30     @classmethod
31     def from_raw_input(cls, raw_input: str, origin: str):
32         """Factory method for ingestion"""
33         pass
```

3.3 Ingestion Pipeline Architecture

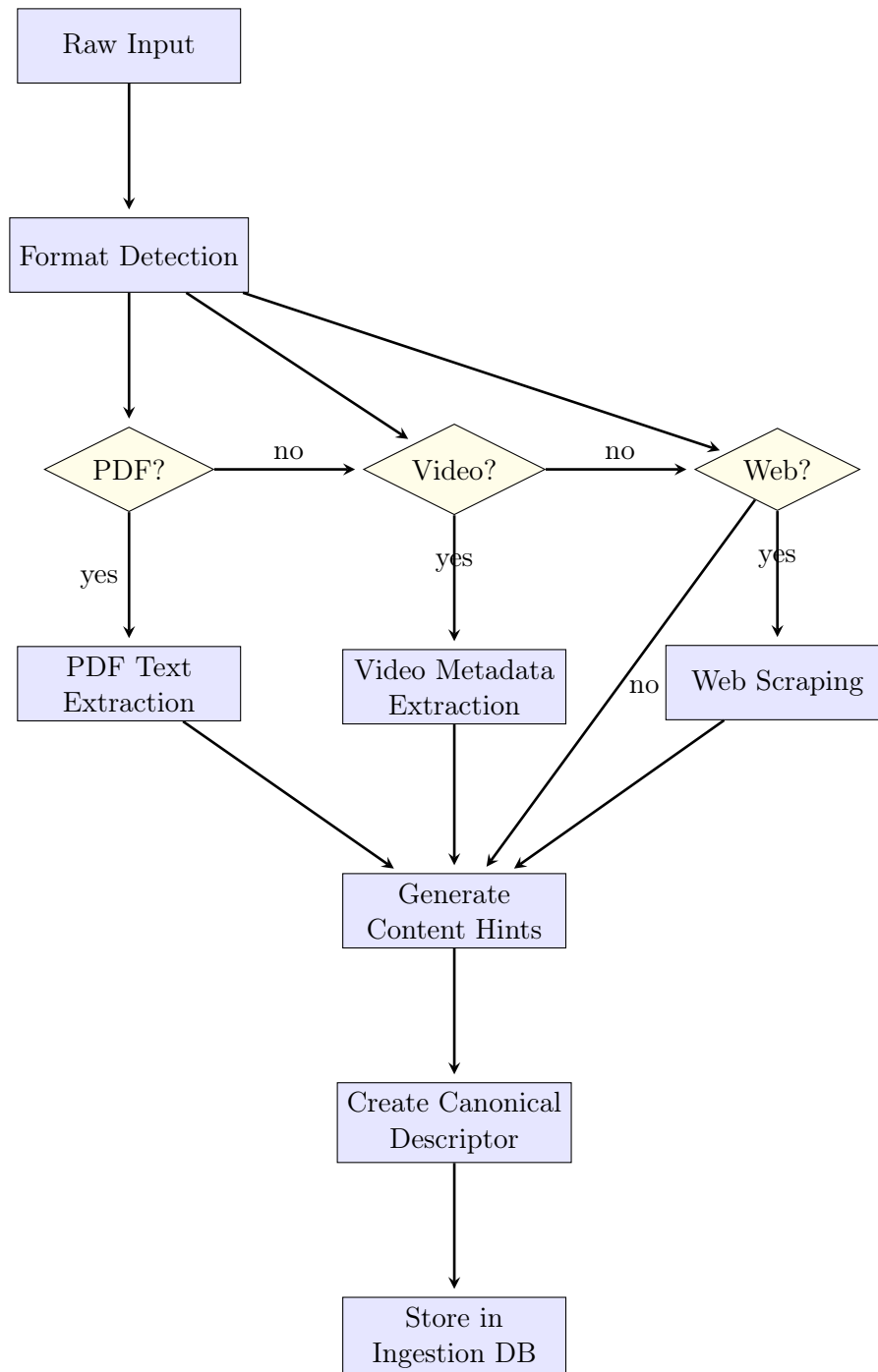


Figure 2: Layer 0 Ingestion Pipeline

3.4 Browser Extension Integration

Listing 2: Browser Extension Content Script

```

1 // KNOWY Browser Extension - Content Script
2 class KnowyExtension {
3   constructor(apiEndpoint) {
4     this.apiEndpoint = apiEndpoint;
5   }

```



```
6
7  async captureCurrentPage() {
8    const pageData = {
9      url: window.location.href,
10     title: document.title,
11     selectedText: window.getSelection().toString(),
12     fullText: document.body.innerText,
13     links: Array.from(document.querySelectorAll('a'))
14       .map(a => a.href),
15     timestamp: new Date().toISOString()
16   };
17
18   return await this.sendToKnowy(pageData);
19 }
20
21 async sendToKnowy(data) {
22   const response = await fetch(`${this.apiEndpoint}/ingest`, {
23     method: 'POST',
24     headers: {'Content-Type': 'application/json'},
25     body: JSON.stringify({
26       origin: 'browser',
27       raw_reference: data.url,
28       raw_text: data.fullText,
29       links: data.links,
30       metadata: {title: data.title, selectedText: data.selectedText}
31     })
32   });
33   return response.json();
34 }
35 }
```

3.5 Format Detection Algorithm

Algorithm 1 Format Detection and Content Hint Generation

Require: Raw input string I , origin metadata O

Ensure: ContentHints object H

```

1:  $H \leftarrow \text{new ContentHints}()$ 
2: if  $I$  matches URL pattern then
3:    $\text{domain} \leftarrow \text{extractDomain}(I)$ 
4:   if  $\text{domain} \in \{\text{youtube.com}, \text{vimeo.com}, \text{coursera.org}\}$  then
5:      $H.\text{contains\_video} \leftarrow \text{true}$ 
6:      $H.\text{estimated\_length\_minutes} \leftarrow \text{estimateFromMetadata}(I)$ 
7:   else if  $\text{domain}$  ends with .pdf then
8:      $H.\text{contains\_pdf} \leftarrow \text{true}$ 
9:   else
10:     $\text{content} \leftarrow \text{fetchHTMLHead}(I)$ 
11:     $H \leftarrow \text{analyzeOpenGraph}(\text{content})$ 
12:   end if
13: else if  $I$  matches file path pattern then
14:    $\text{extension} \leftarrow \text{getExtension}(I)$ 
15:    $H \leftarrow \text{mapExtensionToHints}(\text{extension})$ 
16: else
17:    $H.\text{contains\_external\_links} \leftarrow \text{detectURLs}(I)$ 
18: end if
19: return  $H$ 

```

3.6 Storage Schema

Listing 3: Ingestion Database Schema

```

1 CREATE TABLE ingestion_sources (
2     source_id UUID PRIMARY KEY,
3     origin VARCHAR(50) NOT NULL,
4     raw_reference TEXT NOT NULL,
5     raw_text TEXT,
6     links JSONB,
7     content_hints JSONB NOT NULL,
8     timestamp TIMESTAMPTZ NOT NULL,
9     ingestion_metadata JSONB,
10    processing_status VARCHAR(20) DEFAULT 'pending',
11    created_at TIMESTAMPTZ DEFAULT NOW()
12 );
13
14 CREATE INDEX idx_origin ON ingestion_sources(origin);
15 CREATE INDEX idx_timestamp ON ingestion_sources(timestamp);
16 CREATE INDEX idx_status ON ingestion_sources(processing_status);
17 CREATE INDEX idx_hints_gin ON ingestion_sources USING GIN(content_hints
    );

```

4 Layer 1: Exploration and Expansion

4.1 Purpose and Epistemic Question

Epistemic Question: *"What actual knowledge artifacts exist behind this reference?"*

Layer 1 performs deep link analysis and reference expansion. A single user-saved URL may reference multiple papers, course modules, or video playlists that constitute distinct knowledge artifacts requiring separate processing.

4.2 Expansion Agent Architecture

Listing 4: Exploration Agent Implementation

```

1 from typing import List, Set
2 from langchain.agents import AgentExecutor
3 from langchain.tools import Tool
4
5 class ExplorationAgent:
6     def __init__(self, llm, web_browser_tool):
7         self.llm = llm
8         self.browser = web_browser_tool
9         self.discovered_artifacts = []
10
11     async def expand_reference(self,
12                               source_descriptor:
13                                   CanonicalSourceDescriptor
14                               ) -> List[KnowledgeArtifact]:
15         """
16         Expand a source descriptor into constituent knowledge artifacts
17         """
18         artifacts = []
19
20         # Analyze content hints to determine expansion strategy
21         if source_descriptor.content_hints.contains_video:
22             artifacts.extend(
23                 await self._expand_video_playlist(source_descriptor)
24             )
25
26         if source_descriptor.content_hints.contains_pdf:
27             artifacts.extend(
28                 await self._expand_pdf_citations(source_descriptor)
29             )
30
31         if source_descriptor.content_hints.contains_external_links:
32             artifacts.extend(
33                 await self._expand_linked_resources(source_descriptor)
34             )
35
36         # Use LLM to identify implicit references
37         implicit_refs = await self._identify_implicit_references(
38             source_descriptor.raw_text
39         )
40         artifacts.extend(implicit_refs)
41
42         return self._deduplicate(artifacts)
43
44     async def _expand_video_playlist(self,
45                                     source: CanonicalSourceDescriptor
46                                     ) -> List[KnowledgeArtifact]:
47         """Extract individual videos from playlists or courses."""
48         # Implementation uses platform-specific APIs
49         pass

```

```

49
50     async def _expand_pdf_citations(self,
51                                     source: CanonicalSourceDescriptor
52                                     ) -> List[KnowledgeArtifact]:
53         """Extract and resolve cited papers."""
54         # Parse bibliography, resolve DOIs
55         pass
56
57     async def _identify_implicit_references(self,
58                                             text: str
59                                             ) -> List[KnowledgeArtifact]:
60         """Use LLM to identify mentioned resources."""
61         prompt = f"""
62 Analyze the following text and identify any academic papers,
63 books, courses, or other knowledge resources that are mentioned
64 but not explicitly linked. For each, provide:
65 Resource type (paper/book/course/video)
66 Title
67 Author(s) if mentioned
68 Any identifying information (DOI, ISBN, URL)
69
70 Text: {text}
71
72 Return as structured JSON.
73 """
74         response = await self.llm.ainvoke(prompt)
75         return self._parse_artifact_json(response)

```

4.3 Knowledge Artifact Data Structure

Listing 5: Knowledge Artifact Schema

```

1  from enum import Enum
2
3  class ArtifactType(Enum):
4      PAPER = "academic_paper"
5      BOOK = "book"
6      COURSE = "online_course"
7      VIDEO_LECTURE = "video_lecture"
8      BLOG_POST = "blog_post"
9      DOCUMENTATION = "technical_documentation"
10     DATASET = "dataset"
11     CODE_REPOSITORY = "code_repository"
12
13 @dataclass
14 class KnowledgeArtifact:
15     artifact_id: UUID
16     source_id: UUID # Links back to ingestion source
17     artifact_type: ArtifactType
18     title: str
19     authors: List[str]
20     url: Optional[str]
21     doi: Optional[str]
22     isbn: Optional[str]
23     abstract: Optional[str]
24     publication_date: Optional[datetime]

```

```

25 acquisition_strategy: str # Determines Layer 2 pipeline
26 metadata: Dict[str, any]
27
28 def requires_video_analysis(self) -> bool:
29     return self.artifact_type in [
30         ArtifactType.VIDEO_LECTURE,
31         ArtifactType.COURSE
32     ]
33
34 def requires_pdf_processing(self) -> bool:
35     return self.artifact_type in [
36         ArtifactType.PAPER,
37         ArtifactType.BOOK
38     ] or (self.url and self.url.endswith('.pdf'))

```

4.4 Reference Resolution Strategies

Reference Type	Resolution Strategy	Tools/APIs
DOI	Query CrossRef API	requests, CrossRef REST API
arXiv ID	Parse arXiv metadata	arXiv API
ISBN	Query OpenLibrary	OpenLibrary API
YouTube URL	Extract video/playlist ID	YouTube Data API
Course Platform	Platform-specific scraping	Selenium, BeautifulSoup
GitHub Repository	Clone and analyze structure	PyGithub, git
Generic Web Page	Content extraction	Trafilatura, newspaper3k

Table 3: Reference Resolution Strategies by Type

4.5 Deduplication Algorithm

Algorithm 2 Knowledge Artifact Deduplication

Require: Set of artifacts $A = \{a_1, a_2, \dots, a_n\}$

Ensure: Deduplicated set A'

```

1:  $A' \leftarrow \emptyset$ 
2:  $seen \leftarrow \emptyset$ 
3: for  $a \in A$  do
4:    $key \leftarrow \text{generateKey}(a)$ 
5:   if  $key \notin seen$  then
6:      $A' \leftarrow A' \cup \{a\}$ 
7:      $seen \leftarrow seen \cup \{key\}$ 
8:   else
9:      $existing \leftarrow \text{findByKey}(A', key)$ 
10:     $merged \leftarrow \text{mergeMetadata}(existing, a)$ 
11:     $\text{replace}(existing, merged)$ 
12:   end if
13: end for
14: return  $A'$ 

```

Where the key generation function prioritizes unique identifiers:

Listing 6: Artifact Key Generation

```

1 def generate_key(artifact: KnowledgeArtifact) -> str:
2     if artifact.doi:
3         return f"doi:{artifact.doi}"
4     elif artifact.isbn:
5         return f"isbn:{artifact.isbn}"
6     elif artifact.url:
7         return f"url:{normalize_url(artifact.url)}"
8     else:
9         # Fallback to content-based hash
10        return f"hash:{hash_content(artifact.title, artifact.authors)}"

```

5 Layer 2: Multimodal Acquisition and Scraping

5.1 Purpose and Epistemic Question

Epistemic Question: *"What is the complete primary content of each identified artifact?"*

Layer 2 implements specialized acquisition pipelines for different modalities and formats. The most sophisticated component is the video analysis pipeline, which transforms long-form lectures into structured knowledge units.

5.2 Acquisition Pipeline Dispatcher

Listing 7: Acquisition Pipeline Orchestration

```

1 class AcquisitionOrchestrator:
2     def __init__(self):
3         self.pipelines = {
4             'pdf': PDFAcquisitionPipeline(),
5             'video': VideoAnalysisPipeline(),
6             'web': WebScrapingPipeline(),
7             'code': CodeRepositoryPipeline(),
8             'paper': AcademicPaperPipeline()
9         }
10
11     async def acquire(self,
12                      artifact: KnowledgeArtifact
13                      ) -> AcquiredContent:
14         pipeline = self._select_pipeline(artifact)
15
16         try:
17             content = await pipeline.acquire(artifact)
18             validation_result = await self._validate_content(content)
19
20             if validation_result.is_valid:
21                 return content
22             else:
23                 # Attempt fallback pipeline
24                 fallback = self._get_fallback_pipeline(artifact)
25                 return await fallback.acquire(artifact)
26
27         except AcquisitionException as e:
28             # Log failure and return partial content
29             self._log_acquisition_failure(artifact, e)

```

```
30         return self._create_partial_content(artifact, e)
31
32     def _select_pipeline(self,
33                         artifact: KnowledgeArtifact
34                         ) -> AcquisitionPipeline:
35         if artifact.requires_video_analysis():
36             return self.pipelines['video']
37         elif artifact.requires_pdf_processing():
38             if artifact.doi or artifact.artifact_type == ArtifactType.
39                 PAPER:
40                 return self.pipelines['paper']
41             else:
42                 return self.pipelines['pdf']
43         elif artifact.artifact_type == ArtifactType.CODE_REPOSITORY:
44             return self.pipelines['code']
45         else:
46             return self.pipelines['web']
```

5.3 Video Analysis Pipeline Integration

The video analysis pipeline integrates the open-source `video-analyzer` project to perform shot-level semantic extraction.

5.3.1 Video Analysis Architecture

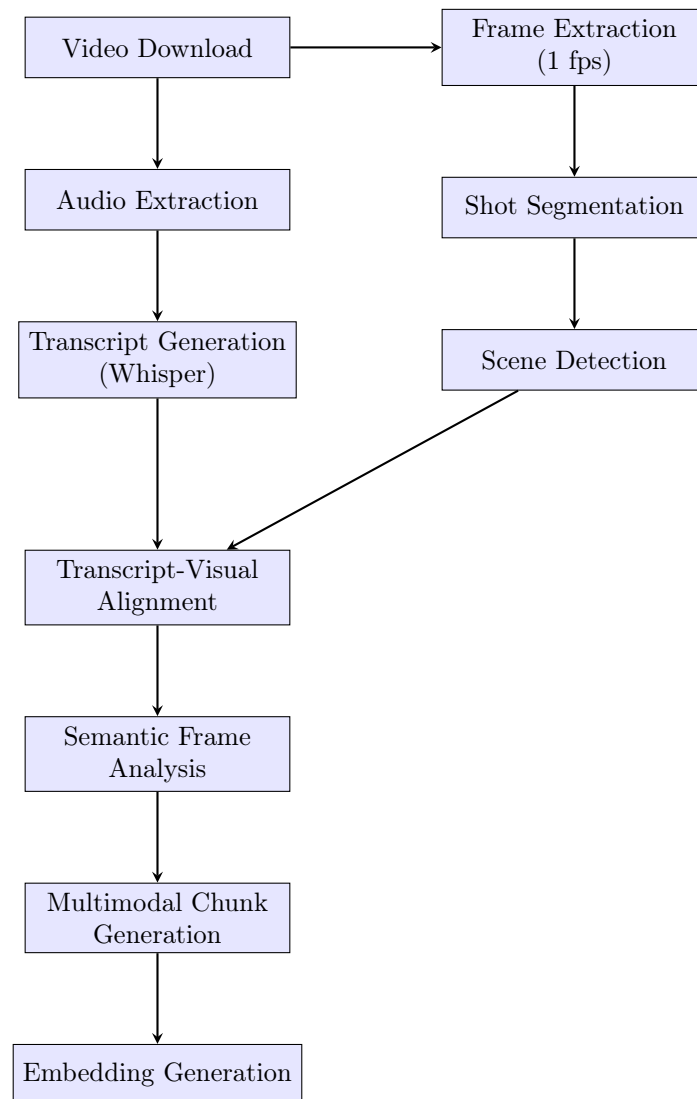


Figure 3: Video Analysis Pipeline Architecture

5.3.2 Shot Segmentation Algorithm

Algorithm 3 Adaptive Shot Boundary Detection

Require: Video frames $F = \{f_1, f_2, \dots, f_n\}$, threshold τ

Ensure: Shot boundaries $B = \{b_1, b_2, \dots, b_m\}$

```

1:  $B \leftarrow \emptyset$ 
2:  $prev\_hist \leftarrow \text{computeHistogram}(f_1)$ 
3: for  $i \leftarrow 2$  to  $n$  do
4:    $curr\_hist \leftarrow \text{computeHistogram}(f_i)$ 
5:    $diff \leftarrow \text{histogramDistance}(prev\_hist, curr\_hist)$ 
6:   if  $diff > \tau$  then
7:      $B \leftarrow B \cup \{i\}$ 
8:      $\tau \leftarrow \text{adaptThreshold}(\tau, diff)$ 
9:   end if
10:   $prev\_hist \leftarrow curr\_hist$ 
11: end for
12: return  $B$ 

```

5.3.3 Transcript-Visual Alignment

Listing 8: Transcript-Visual Temporal Alignment

```

1 class TranscriptVisualAligner:
2     def __init__(self, whisper_model):
3         self.whisper = whisper_model
4
5     async def align(self,
6                     transcript_segments: List[TranscriptSegment],
7                     shots: List[Shot]) -> List[AlignedChunk]:
8
9         """
10        Align transcript segments with visual shots to create
11        semantically coherent multimodal chunks.
12        """
13        aligned_chunks = []
14        shot_idx = 0
15
16        for segment in transcript_segments:
17            # Find overlapping shots
18            overlapping_shots = []
19            while (shot_idx < len(shots) and
20                  shots[shot_idx].start_time <= segment.end_time):
21                if shots[shot_idx].end_time >= segment.start_time:
22                    overlapping_shots.append(shots[shot_idx])
23                    shot_idx += 1
24
25            # Create aligned chunk
26            chunk = AlignedChunk(
27                start_time=segment.start_time,
28                end_time=segment.end_time,
29                transcript_text=segment.text,
30                shots=overlapping_shots,
31                speaker=segment.speaker,
32                confidence=segment.confidence
33            )

```

```

34         aligned_chunks.append(chunk)
35
36     return aligned_chunks
37
38 @dataclass
39 class AlignedChunk:
40     start_time: float
41     end_time: float
42     transcript_text: str
43     shots: List[Shot]
44     speaker: Optional[str]
45     confidence: float
46
47     def get_representative_frame(self) -> np.ndarray:
48         """Select most informative frame from shots."""
49         if not self.shots:
50             return None
51         # Use middle frame of longest shot
52         longest_shot = max(self.shots, key=lambda s: s.duration)
53         return longest_shot.frames[len(longest_shot.frames) // 2]

```

5.3.4 Semantic Frame Analysis

Listing 9: Visual-Semantic Frame Analysis with LLaMA Vision

```

1 class SemanticFrameAnalyzer:
2     def __init__(self, vision_model):
3         self.vision_model = vision_model # LLaMA 3.2 Vision
4
5     async def analyze_chunk(self,
6                             chunk: AlignedChunk
7                             ) -> SemanticAnalysis:
8
9         """
10        Perform semantic analysis of visual content in context
11        of transcript.
12        """
13        frame = chunk.get_representative_frame()
14
15        prompt = f"""
16        Analyze this frame from an educational video lecture.
17        Transcript context: "{chunk.transcript_text}"
18
19        Describe:
20        1. Visual content (diagrams, equations, code, slides)
21        2. Key concepts being illustrated
22        3. Relationship between visual and spoken content
23        4. Educational value of this visual
24
25        Provide structured analysis.
26        """
27
28        analysis = await self.vision_model.analyze(
29            image=frame,
30            prompt=prompt
31        )
32

```

```

33     return SemanticAnalysis(
34         chunk_id=chunk.id,
35         visual_description=analysis.visual_content,
36         concepts=analysis.key_concepts,
37         visual_transcript_coherence=analysis.coherence_score,
38         educational_value=analysis.educational_value,
39         contains_diagram=self._detect_diagram(analysis),
40         contains_equation=self._detect_equation(analysis),
41         contains_code=self._detect_code(analysis)
42     )

```

5.3.5 Multimodal Chunk Generation

Listing 10: Complete Multimodal Chunk Data Structure

```

1  @dataclass
2  class MultimodalChunk:
3      chunk_id: UUID
4      artifact_id: UUID
5      temporal_range: Tuple[float, float]
6
7      # Textual components
8      transcript_text: str
9      transcript_confidence: float
10
11     # Visual components
12     representative_frame: bytes # JPEG-encoded
13     all_frames: List[bytes]
14     shot_boundaries: List[float]
15
16     # Semantic analysis
17     visual_description: str
18     key_concepts: List[str]
19     contains_diagram: bool
20     contains_equation: bool
21     contains_code: bool
22
23     # Embeddings
24     text_embedding: np.ndarray
25     visual_embedding: np.ndarray
26     combined_embedding: np.ndarray
27
28     # Metadata
29     speaker: Optional[str]
30     slide_number: Optional[int]
31     topic_tags: List[str]
32
33     def to_storage_format(self) -> dict:
34         """Prepare for database insertion."""
35         return {
36             'chunk_id': str(self.chunk_id),
37             'artifact_id': str(self.artifact_id),
38             'temporal_range': self.temporal_range,
39             'transcript_text': self.transcript_text,
40             'visual_description': self.visual_description,
41             'key_concepts': json.dumps(self.key_concepts),
42             'semantic_flags': {

```

```

43         'diagram': self.contains_diagram,
44         'equation': self.contains_equation,
45         'code': self.contains_code
46     },
47     'embeddings': {
48         'text': self.text_embedding.tolist(),
49         'visual': self.visual_embedding.tolist(),
50         'combined': self.combined_embedding.tolist()
51     }
52 }

```

5.4 PDF Acquisition Pipeline

Listing 11: PDF Processing with Layout Preservation

```

1  class PDFAcquisitionPipeline:
2      def __init__(self):
3          self.extractors = {
4              'text': PDFTextExtractor(),
5              'layout': PDFLayoutExtractor(),
6              'images': PDFImageExtractor(),
7              'tables': PDFTableExtractor()
8          }
9
10     async def acquire(self, artifact: KnowledgeArtifact) -> PDFContent:
11         pdf_bytes = await self._download_pdf(artifact.url)
12
13         # Parallel extraction
14         text_task = self.extractors['text'].extract(pdf_bytes)
15         layout_task = self.extractors['layout'].extract(pdf_bytes)
16         images_task = self.extractors['images'].extract(pdf_bytes)
17         tables_task = self.extractors['tables'].extract(pdf_bytes)
18
19         text, layout, images, tables = await asyncio.gather(
20             text_task, layout_task, images_task, tables_task
21         )
22
23         # Reconstruct document structure
24         structured_content = self._reconstruct_structure(
25             text, layout, images, tables
26         )
27
28         return PDFContent(
29             artifact_id=artifact.artifact_id,
30             raw_text=text.full_text,
31             sections=structured_content.sections,
32             figures=images,
33             tables=tables,
34             citations=self._extract_citations(text),
35             metadata=self._extract_metadata(pdf_bytes)
36         )
37
38     def _reconstruct_structure(self, text, layout, images, tables):
39         """
40         Reconstruct document hierarchy using layout analysis.
41         Preserves section headers, paragraphs, lists, etc.
42         """

```

```

43     sections = []
44     current_section = None
45
46     for block in layout.blocks:
47         if block.type == 'heading':
48             if current_section:
49                 sections.append(current_section)
50                 current_section = Section(
51                     title=block.text,
52                     level=block.heading_level,
53                     content=[]
54                 )
55             elif block.type == 'paragraph':
56                 if current_section:
57                     current_section.content.append(
58                         Paragraph(text=block.text)
59                     )
60             elif block.type == 'figure':
61                 fig = self._match_figure(block, images)
62                 if current_section and fig:
63                     current_section.content.append(fig)
64             elif block.type == 'table':
65                 tbl = self._match_table(block, tables)
66                 if current_section and tbl:
67                     current_section.content.append(tbl)
68
69     if current_section:
70         sections.append(current_section)
71
72     return StructuredDocument(sections=sections)

```

5.5 Academic Paper Pipeline

Listing 12: Academic Paper Specialized Processing

```

1  class AcademicPaperPipeline:
2      def __init__(self):
3          self.pdf_pipeline = PDFAcquisitionPipeline()
4          self.metadata_enricher = PaperMetadataEnricher()
5
6      async def acquire(self, artifact: KnowledgeArtifact) ->
7          PaperContent:
8          # First get PDF content
9          pdf_content = await self.pdf_pipeline.acquire(artifact)
10
11         # Enrich with academic metadata
12         metadata = await self.metadata_enricher.enrich(
13             doi=artifact.doi,
14             arxiv_id=self._extract_arxiv_id(artifact.url)
15         )
16
17         # Parse academic structure
18         academic_structure = self._parse_academic_structure(
19             pdf_content
20         )
21
22         return PaperContent(

```

```

22         artifact_id=artifact.artifact_id,
23         title=metadata.title or artifact.title,
24         authors=metadata.authors or artifact.authors,
25         abstract=academic_structure.abstract,
26         introduction=academic_structure.introduction,
27         methodology=academic_structure.methodology,
28         results=academic_structure.results,
29         discussion=academic_structure.discussion,
30         conclusion=academic_structure.conclusion,
31         references=pdf_content.citations,
32         figures=pdf_content.figures,
33         tables=pdf_content.tables,
34         equations=self._extract_equations(pdf_content),
35         citation_count=metadata.citation_count,
36         publication_venue=metadata.venue,
37         publication_date=metadata.publication_date
38     )
39
40     def _parse_academic_structure(self, pdf_content):
41         """
42         Use section headers to identify standard paper sections.
43         """
44         section_map = {
45             'abstract': None,
46             'introduction': None,
47             'methodology': None,
48             'results': None,
49             'discussion': None,
50             'conclusion': None
51         }
52
53         for section in pdf_content.sections:
54             title_lower = section.title.lower()
55
56             if 'abstract' in title_lower:
57                 section_map['abstract'] = section
58             elif 'introduction' in title_lower:
59                 section_map['introduction'] = section
60             elif any(kw in title_lower for kw in
61                     ['method', 'approach', 'design']):
62                 section_map['methodology'] = section
63             elif 'result' in title_lower:
64                 section_map['results'] = section
65             elif 'discussion' in title_lower:
66                 section_map['discussion'] = section
67             elif 'conclusion' in title_lower:
68                 section_map['conclusion'] = section
69
70         return AcademicStructure(**section_map)

```

6 Layer 3: Multi-Resolution Knowledge Representation

6.1 Purpose and Epistemic Question

Epistemic Question: *"How should knowledge be structured at multiple levels of abstraction to support both factual grounding and high-level reasoning?"*

Layer 3 constructs a four-level abstraction hierarchy that enables both precise fact retrieval and synthesized understanding.

6.2 Four-Level Abstraction Hierarchy

Level	Granularity	Use Case
Level 0	Raw chunks (paragraphs, video segments)	Direct citation, fact verification
Level 1	Section-level summaries	Navigating document structure
Level 2	Document-level synthesis	Comparing papers, course overview
Level 3	Cross-document thematic synthesis	Domain understanding, research frontiers

Table 4: Multi-Resolution Abstraction Levels

6.3 Representation Architecture

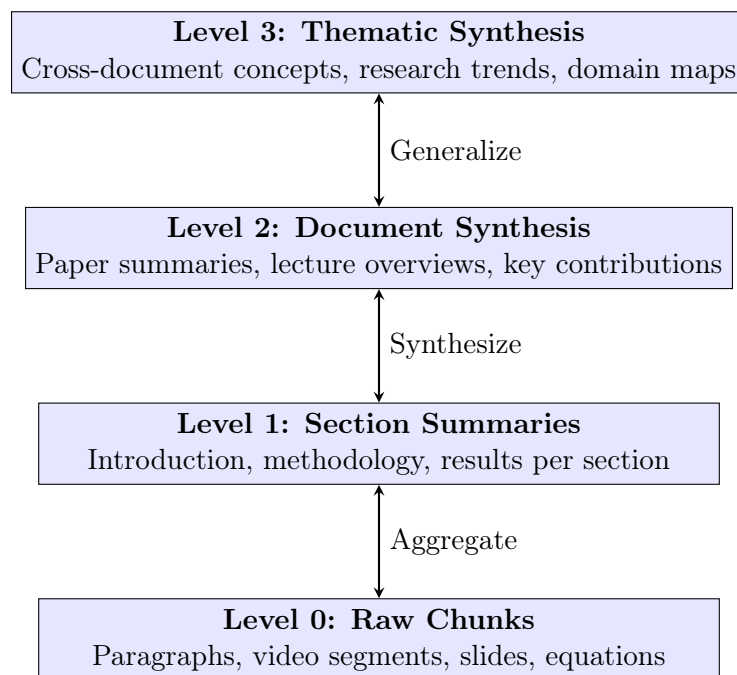


Figure 4: Multi-Resolution Representation Hierarchy

6.4 Level 0: Raw Content Storage

Listing 13: Level 0 Raw Chunk Representation

```

1 @dataclass
2 class RawChunk:
3     chunk_id: UUID
4     artifact_id: UUID
5     chunk_type: str # paragraph / video_segment / equation / figure

```

```

6
7     # Content
8     text_content: Optional[str]
9     visual_content: Optional[bytes]
10    structured_content: Optional[dict] # For tables, code
11
12    # Positional metadata
13    position_in_artifact: int
14    section_id: Optional[UUID]
15    temporal_range: Optional[Tuple[float, float]] # For video
16    page_number: Optional[int] # For PDFs
17
18    # Embeddings
19    embedding: np.ndarray
20
21    # Provenance
22    extraction_method: str
23    confidence: float
24
25    def get_context_window(self,
26                           window_size: int = 2
27                           ) -> List['RawChunk']:
28        """Retrieve neighboring chunks for context."""
29        pass
30
31    # Storage schema
32    CREATE TABLE raw_chunks (
33        chunk_id UUID PRIMARY KEY,
34        artifact_id UUID REFERENCES knowledge_artifacts(artifact_id),
35        chunk_type VARCHAR(50),
36        text_content TEXT,
37        visual_content BYTEA,
38        structured_content JSONB,
39        position_in_artifact INTEGER,
40        section_id UUID,
41        temporal_range NUMRANGE,
42        page_number INTEGER,
43        embedding vector(1536),
44        extraction_metadata JSONB,
45        created_at TIMESTAMPTZ DEFAULT NOW()
46    );
47
48    CREATE INDEX idx_artifact ON raw_chunks(artifact_id);
49    CREATE INDEX idx_section ON raw_chunks(section_id);
50    CREATE INDEX idx_embedding ON raw_chunks
51        USING ivfflat (embedding vector_cosine_ops);

```

6.5 Level 1: Section-Level Representation

Listing 14: Level 1 Section Summary Generation

```

1 class SectionSummarizer:
2     def __init__(self, llm):
3         self.llm = llm
4
5     async def generate_section_summary(self,
6                                       section_chunks: List[RawChunk]

```



```

7         ) -> SectionRepresentation:
8         """
9         Generate structured summary of a document section.
10        """
11        # Concatenate section content
12        full_text = "
13    ".join([
14        chunk.text_content for chunk in section_chunks
15        if chunk.text_content
16    ])
17
18        prompt = f"""
19        Summarize the following section from an academic document.
20        Provide:
21        1. Main claim or thesis (1-2 sentences)
22        2. Key points (bullet list)
23        3. Important terminology defined
24        4. Figures/equations referenced
25        5. Dependencies on other sections
26
27        Section content:
28        {full_text}
29
30        Output as structured JSON.
31        """
32
33        summary = await self.llm.ainvoke(prompt)
34        parsed = json.loads(summary)
35
36        return SectionRepresentation(
37            section_id=uuid4(),
38            artifact_id=section_chunks[0].artifact_id,
39            title=self._extract_section_title(section_chunks),
40            main_claim=parsed['main_claim'],
41            key_points=parsed['key_points'],
42            terminology=parsed['terminology'],
43            referenced_elements=parsed['references'],
44            chunk_ids=[c.chunk_id for c in section_chunks],
45            embedding=await self._embed_summary(parsed['main_claim'])
46        )
47
48    @dataclass
49    class SectionRepresentation:
50        section_id: UUID
51        artifact_id: UUID
52        title: str
53        main_claim: str
54        key_points: List[str]
55        terminology: Dict[str, str]
56        referenced_elements: List[str]
57        chunk_ids: List[UUID]
58        embedding: np.ndarray
59        prerequisites: List[UUID] = field(default_factory=list)

```

6.6 Level 2: Document-Level Synthesis

Listing 15: Level 2 Document Synthesis

```

1 class DocumentSynthesizer:
2     def __init__(self, llm):
3         self.llm = llm
4
5     async def synthesize_document(self,
6                                   sections: List[SectionRepresentation]
7                                   ) -> DocumentRepresentation:
8
9         """
10        Create document-level synthesis from section summaries.
11        """
12        # Construct synthesis prompt
13        section_summaries = "
14        ".join([
15            f"Section: {s.title}
16            Main claim: {s.main_claim}
17            "
18            for s in sections
19        ])
20
21        prompt = f"""
22        Synthesize the following academic document into a cohesive
23        summary.
24
25        Sections:
26        {section_summaries}
27
28        Provide:
29        1. Central thesis (2-3 sentences)
30        2. Main contributions (numbered list)
31        3. Methodology overview
32        4. Key findings
33        5. Limitations discussed
34        6. Related work positioning
35        7. Novel concepts introduced
36
37        Output as structured JSON.
38        """
39
40        synthesis = await self.llm.ainvoke(prompt)
41        parsed = json.loads(synthesis)
42
43        return DocumentRepresentation(
44            document_id=uuid4(),
45            artifact_id=sections[0].artifact_id,
46            central_thesis=parsed['central_thesis'],
47            contributions=parsed['contributions'],
48            methodology=parsed['methodology'],
49            findings=parsed['findings'],
50            limitations=parsed['limitations'],
51            novel_concepts=parsed['novel_concepts'],
52            section_ids=[s.section_id for s in sections],
53            embedding=await self._embed_synthesis(
54                parsed['central_thesis']
55            )
56        )
57
58 @dataclass

```

```

57 class DocumentRepresentation:
58     document_id: UUID
59     artifact_id: UUID
60     central_thesis: str
61     contributions: List[str]
62     methodology: str
63     findings: List[str]
64     limitations: List[str]
65     novel_concepts: List[str]
66     section_ids: List[UUID]
67     embedding: np.ndarray
68     citation_graph: Optional[dict] = None

```

6.7 Level 3: Cross-Document Thematic Synthesis

Listing 16: Level 3 Thematic Synthesis

```

1 class ThematicSynthesizer:
2     def __init__(self, llm, clustering_model):
3         self.llm = llm
4         self.clusterer = clustering_model
5
6     async def synthesize_themes(self,
7                                 documents: List[DocumentRepresentation]
8                                 ) -> List[ThematicSynthesis]:
9         """
10        Identify cross-document themes and research trends.
11        """
12        # Cluster documents by embedding similarity
13        embeddings = np.array([d.embedding for d in documents])
14        clusters = self.clusterer.fit_predict(embeddings)
15
16        themes = []
17        for cluster_id in np.unique(clusters):
18            cluster_docs = [
19                documents[i] for i in range(len(documents))
20                if clusters[i] == cluster_id
21            ]
22
23            theme = await self._synthesize_cluster(cluster_docs)
24            themes.append(theme)
25
26        return themes
27
28    async def _synthesize_cluster(self,
29                                docs: List[DocumentRepresentation]
30                                ) -> ThematicSynthesis:
31        """
32        Synthesize a thematic cluster into unified understanding.
33        """
34        # Extract key concepts across documents
35        all_concepts = []
36        all_theses = []
37
38        for doc in docs:
39            all_concepts.extend(doc.novel_concepts)
40            all_theses.append(doc.central_thesis)

```

```

41
42     prompt = f"""
43     Analyze the following related research documents and identify:
44
45     1. Common theme or research question
46     2. Points of agreement
47     3. Points of disagreement/debate
48     4. Evolution of ideas (temporal if dates available)
49     5. Open questions in this area
50     6. Key terminology and definitions
51
52     Document theses:
53     {chr(10).join([f"- {t}" for t in all_theses])}
54
55     Novel concepts mentioned:
56     {chr(10).join([f"- {c}" for c in set(all_concepts)])}
57
58     Output as structured JSON.
59     """
60
61     synthesis = await self.llm.ainvoke(prompt)
62     parsed = json.loads(synthesis)
63
64     return ThematicSynthesis(
65         theme_id=uuid4(),
66         theme_name=parsed['common_theme'],
67         document_ids=[d.document_id for d in docs],
68         agreements=parsed['agreements'],
69         disagreements=parsed['disagreements'],
70         evolution=parsed['evolution'],
71         open_questions=parsed['open_questions'],
72         key_terminology=parsed['terminology'],
73         embedding=await self._embed_theme(parsed['common_theme'])
74     )
75
76 @dataclass
77 class ThematicSynthesis:
78     theme_id: UUID
79     theme_name: str
80     document_ids: List[UUID]
81     agreements: List[str]
82     disagreements: List[str]
83     evolution: Optional[str]
84     open_questions: List[str]
85     key_terminology: Dict[str, str]
86     embedding: np.ndarray
87     research_frontier: bool = False

```

6.8 Dynamic Resolution Selection

Listing 17: Query-Adaptive Resolution Selection

```

1 class ResolutionSelector:
2     def __init__(self):
3         self.query_classifier = QueryTypeClassifier()
4
5     def select_optimal_level(self, query: str) -> int:

```

```

6         """
7         Determine optimal abstraction level for query.
8         """
9         query_type = self.query_classifier.classify(query)
10
11         if query_type == 'fact_lookup':
12             # Need precise citation
13             return 0
14         elif query_type == 'section_comparison':
15             # Compare methodologies, results sections
16             return 1
17         elif query_type == 'paper_summary':
18             # Overview of single document
19             return 2
20         elif query_type == 'domain_understanding':
21             # Conceptual understanding across sources
22             return 3
23         elif query_type == 'pedagogical':
24             # Teaching requires multiple levels
25             return [0, 1, 2, 3] # Multi-level retrieval
26         else:
27             # Default to document level
28             return 2

```

7 Layer 4: Classification and Ontology Construction

7.1 Purpose and Epistemic Question

Epistemic Question: *"How are knowledge artifacts related topically, hierarchically, and temporally?"*

Layer 4 constructs explicit ontologies capturing topic hierarchies, prerequisite relationships, and redundancy mappings.

7.2 Ontology Structure

Listing 18: Knowledge Ontology Data Model

```

1 @dataclass
2 class TopicNode:
3     topic_id: UUID
4     name: str
5     parent_topic: Optional[UUID]
6     child_topics: List[UUID]
7     related_documents: List[UUID]
8     related_sections: List[UUID]
9
10    # Topic characterization
11    description: str
12    key_terminology: List[str]
13    difficulty_level: int # 1-5
14
15    # Prerequisites
16    prerequisite_topics: List[UUID]
17    recommended_learning_order: int
18
19    def is_leaf(self) -> bool:

```

```

20         return len(self.child_topics) == 0
21
22     def get_ancestors(self) -> List['TopicNode']:
23         """Traverse to root to get full path."""
24         pass
25
26 @dataclass
27 class PrerequisiteRelationship:
28     prerequisite_id: UUID
29     dependent_id: UUID
30     strength: float # 0-1, how essential is prerequisite
31     relationship_type: str # foundational / helpful / optional
32
33 # Graph storage schema (Neo4j Cypher)
34 CREATE (t:Topic {
35     topic_id: $topic_id,
36     name: $name,
37     description: $description,
38     difficulty: $difficulty
39 })
40
41 CREATE (t1:Topic)-[:PARENT_OF]->(t2:Topic)
42 CREATE (t1:Topic)-[:PREREQUISITE_FOR {strength: 0.9}]->(t2:Topic)
43 CREATE (d:Document)-[:COVERS]->(t:Topic)

```

7.3 Topic Extraction and Clustering

Listing 19: Automated Topic Extraction

```

1 class TopicExtractor:
2     def __init__(self, llm, embedding_model):
3         self.llm = llm
4         self.embedder = embedding_model
5
6     async def extract_topics(self,
7                             documents: List[DocumentRepresentation]
8                             ) -> List[TopicNode]:
9
10         """
11         Extract hierarchical topic structure from documents.
12         """
13         # Step 1: Extract candidate topics from each document
14         all_candidates = []
15         for doc in documents:
16             candidates = await self._extract_document_topics(doc)
17             all_candidates.extend(candidates)
18
19         # Step 2: Cluster and deduplicate topics
20         unified_topics = await self._unify_topics(all_candidates)
21
22         # Step 3: Construct hierarchy
23         hierarchy = await self._build_hierarchy(unified_topics)
24
25         # Step 4: Identify prerequisites
26         await self._identify_prerequisites(hierarchy)
27
28         return hierarchy

```

```

29     async def _extract_document_topics(self,
30                                         doc: DocumentRepresentation
31                                         ) -> List[str]:
32         """Extract topics mentioned in document."""
33         prompt = f"""
34         Extract the main topics and subtopics covered in this document.
35
36         Thesis: {doc.central_thesis}
37         Contributions: {doc.contributions}
38         Novel concepts: {doc.novel_concepts}
39
40         For each topic, provide:
41         - Topic name
42         - Brief description
43         - Estimated difficulty level (1-5)
44
45         Output as JSON list.
46         """
47
48         response = await self.llm.ainvoke(prompt)
49         return json.loads(response)
50
51     async def _unify_topics(self,
52                             candidates: List[dict]
53                             ) -> List[TopicNode]:
54         """
55         Cluster similar topics and create unified representations.
56         """
57         # Embed all topic names and descriptions
58         embeddings = []
59         for candidate in candidates:
60             emb = await self.embedder.embed(
61                 f"{candidate['name']}: {candidate['description']}"
62             )
63             embeddings.append(emb)
64
65         # Cluster by similarity
66         clusterer = DBSCAN(eps=0.3, min_samples=2)
67         clusters = clusterer.fit_predict(np.array(embeddings))
68
69         unified = []
70         for cluster_id in np.unique(clusters):
71             cluster_topics = [
72                 candidates[i] for i in range(len(candidates))
73                 if clusters[i] == cluster_id
74             ]
75
76             # Merge cluster into single topic
77             merged = await self._merge_topic_cluster(cluster_topics)
78             unified.append(TopicNode(**merged))
79
80         return unified
81
82     async def _build_hierarchy(self,
83                               topics: List[TopicNode]
84                               ) -> List[TopicNode]:
85         """
86         Construct parent-child relationships between topics.

```

```

87     """
88     prompt = f"""
89     Given these topics, construct a hierarchical taxonomy.
90     Identify which topics are subtopics of others.
91
92     Topics:
93     {json.dumps([{'id': str(t.topic_id), 'name': t.name, 'description': t.description} for t in topics], indent=2)}
94
95     Output parent-child relationships as JSON:
96     [{{"parent_id": "...", "child_id": "..."}}, ...]
97     """
98
99     relationships = await self.llm.ainvoke(prompt)
100     parsed = json.loads(relationships)
101
102     # Apply relationships to topic nodes
103     for rel in parsed:
104         parent = next(t for t in topics if str(t.topic_id) == rel['parent_id'])
105         child = next(t for t in topics if str(t.topic_id) == rel['child_id'])
106
107         parent.child_topics.append(child.topic_id)
108         child.parent_topic = parent.topic_id
109
110     return topics
111
112     async def _identify_prerequisites(self,
113                                     topics: List[TopicNode]):
114         """
115         Identify prerequisite relationships between topics.
116         """
117         for topic in topics:
118             prompt = f"""
119             For the topic "{topic.name}": {topic.description}
120
121             From this list of other topics, identify which are
122             prerequisites:
123             {json.dumps([{'id': str(t.topic_id), 'name': t.name} for t in
124                         topics if t.topic_id != topic.topic_id], indent=2)}
125
126             For each prerequisite, rate its importance:
127             - foundational (0.9-1.0): Cannot understand topic without this
128             - helpful (0.5-0.8): Significantly aids understanding
129             - optional (0.1-0.4): Provides context but not required
130
131             Output as JSON: [{{"topic_id": "...", "strength": 0.9, "type": "foundational"}}, ...]
132             """
133
134             prereqs = await self.llm.ainvoke(prompt)
135             parsed = json.loads(prereqs)
136
137             for prereq in parsed:
138                 topic.prerequisite_topics.append(UUID(prereq['topic_id'])))

```


7.4 Redundancy Detection

Listing 20: Content Redundancy Detection and Mapping

```

1 class RedundancyDetector:
2     def __init__(self, similarity_threshold=0.85):
3         self.threshold = similarity_threshold
4
5     def detect_redundancy(self,
6                           sections: List[SectionRepresentation]
7                           ) -> List[RedundancyGroup]:
8         """
9         Identify sections covering similar content across documents.
10        """
11        embeddings = np.array([s.embedding for s in sections])
12
13        # Compute pairwise cosine similarities
14        similarities = cosine_similarity(embeddings)
15
16        # Find redundant groups
17        redundancy_groups = []
18        processed = set()
19
20        for i in range(len(sections)):
21            if i in processed:
22                continue
23
24            # Find all sections similar to section i
25            similar_indices = np.where(
26                similarities[i] > self.threshold
27            )[0]
28
29            if len(similar_indices) > 1:
30                group = RedundancyGroup(
31                    group_id=uuid4(),
32                    section_ids=[sections[idx].section_id
33                               for idx in similar_indices],
34                    similarity_scores=similarities[i][similar_indices].
35                        tolist(),
36                    canonical_section=sections[i].section_id,
37                    coverage_type=self._classify_redundancy(
38                        [sections[idx] for idx in similar_indices]
39                    )
40                )
41                redundancy_groups.append(group)
42                processed.update(similar_indices)
43
44        return redundancy_groups
45
46    def _classify_redundancy(self,
47                             sections: List[SectionRepresentation]
48                             ) -> str:
49        """
50        Determine if redundancy is:
51        - duplicate: Identical content
52        - complementary: Same topic, different perspectives
53        - progressive: Building on same concepts
54        """
55        # Compare key points

```

```

55     all_points = [set(s.key_points) for s in sections]
56
57     # Check overlap
58     if len(all_points) < 2:
59         return "duplicate"
60
61     intersection = set.intersection(*all_points)
62     union = set.union(*all_points)
63
64     jaccard = len(intersection) / len(union)
65
66     if jaccard > 0.8:
67         return "duplicate"
68     elif jaccard > 0.4:
69         return "complementary"
70     else:
71         return "progressive"
72
73 @dataclass
74 class RedundancyGroup:
75     group_id: UUID
76     section_ids: List[UUID]
77     similarity_scores: List[float]
78     canonical_section: UUID
79     coverage_type: str # duplicate / complementary / progressive

```

8 Layer 5: Expert Reasoning and Teaching

8.1 Purpose and Epistemic Question

Epistemic Question: *"How can internal knowledge be composed into expert-level explanations and pedagogical responses?"*

Layer 5 implements the primary reasoning agent that answers user queries by dynamically selecting and composing knowledge from all abstraction levels.

8.2 Expert Reasoning Architecture

Listing 21: Expert Reasoning Agent

```

1 class ExpertReasoningAgent:
2     def __init__(self, llm, knowledge_store):
3         self.llm = llm
4         self.knowledge = knowledge_store
5         self.resolution_selector = ResolutionSelector()
6
7     async def answer_query(self, query: str) -> ExpertResponse:
8         """
9         Main reasoning loop: retrieve, synthesize, respond.
10        """
11        # Step 1: Classify query and select resolution
12        resolutions = self.resolution_selector.select_optimal_level(
13            query)
14
15        # Step 2: Retrieve relevant knowledge
16        retrieved = await self._multi_resolution_retrieval(
17            query, resolutions

```

```

17         )
18
19         # Step 3: Check epistemic boundaries
20         uncertainty = await self._assess_epistemic_uncertainty(
21             query, retrieved
22         )
23
24         # Step 4: Generate response
25         if uncertainty.requires_external_research:
26             return await self._trigger_deep_research(query, uncertainty
27         )
28         else:
29             response = await self._generate_grounded_response(
30                 query, retrieved, uncertainty
31             )
32             return response
33
34     async def _multi_resolution_retrieval(self,
35                                         query: str,
36                                         resolutions: Union[int, List[
37                                             int]]
38                                         ) -> RetrievedKnowledge:
39         """
40         Retrieve knowledge from specified abstraction levels.
41         """
42         if isinstance(resolutions, int):
43             resolutions = [resolutions]
44
45         retrieved = RetrievedKnowledge()
46
47         for level in resolutions:
48             if level == 0:
49                 chunks = await self.knowledge.retrieve_raw_chunks(query
50                 )
51                 retrieved.raw_chunks = chunks
52             elif level == 1:
53                 sections = await self.knowledge.retrieve_sections(query
54                 )
55                 retrieved.sections = sections
56             elif level == 2:
57                 documents = await self.knowledge.retrieve_documents(
58                     query)
59                 retrieved.documents = documents
60             elif level == 3:
61                 themes = await self.knowledge.retrieve_themes(query)
62                 retrieved.themes = themes
63
64         return retrieved
65
66     async def _assess_epistemic_uncertainty(self,
67                                         query: str,
68                                         retrieved:
69                                             RetrievedKnowledge
70                                         ) -> EpistemicUncertainty:
71         """
72         Determine if query can be answered with internal knowledge.
73         """
74         assessment_prompt = f"""

```

```

69         """Query: {query}
70
71         Retrieved knowledge summary:
72         """- {len(retrieved.raw_chunks or [])} raw chunks
73         """- {len(retrieved.sections or [])} sections
74         """- {len(retrieved.documents or [])} documents
75         """- {len(retrieved.themes or [])} themes
76
77         Assess:
78         """1. Is retrieved knowledge sufficient to answer query?
79         """2. Are there temporal concerns (knowledge may be outdated)?
80         """3. Are there conflicting claims in retrieved knowledge?
81         """4. What is confidence level (0-1)?
82         """5. Does this require external research?
83
84         Output as JSON.
85         """
86
87         assessment = await self.llm.ainvoke(assessment_prompt)
88         parsed = json.loads(assessment)
89
90         return EpistemicUncertainty(
91             is_sufficient=parsed['sufficient'],
92             has_temporal_concerns=parsed['temporal_concerns'],
93             has_conflicts=parsed['conflicts'],
94             confidence=parsed['confidence'],
95             requires_external_research=parsed['requires_research'],
96             explanation=parsed.get('explanation', '')
97         )
98
99     async def _generate_grounded_response(self,
100                                         query: str,
101                                         retrieved: RetrievedKnowledge,
102                                         uncertainty:
103                                             EpistemicUncertainty
104                                         ) -> ExpertResponse:
105         """
106         Generate response strictly from retrieved knowledge.
107         """
108         # Build context from retrieved knowledge
109         context = self._build_context(retrieved)
110
111         response_prompt = f"""
112         You are an expert teacher with deep knowledge of the user's
113         personal knowledge collection. Answer the following query using
114         ONLY the provided context. Do not use external knowledge.
115
116         Query: {query}
117
118         Context:
119         {context}
120
121         Epistemic constraints:
122         """- Confidence: {uncertainty.confidence}
123         """- Temporal concerns: {uncertainty.has_temporal_concerns}
124         """- Conflicts detected: {uncertainty.has_conflicts}
125
126         Provide:

```

```

126     """1. Direct answer
127     2. Supporting evidence with citations
128     3. Relevant prerequisite concepts if query is pedagogical
129     4. Uncertainties or limitations
130     5. Suggested related topics for deeper understanding
131
132     Format response as structured JSON.
133     """
134
135     response = await self.llm.ainvoke(response_prompt)
136     parsed = json.loads(response)
137
138     return ExpertResponse(
139         answer=parsed['answer'],
140         citations=self._extract_citations(retrieved),
141         prerequisites=parsed.get('prerequisites', []),
142         uncertainties=parsed.get('uncertainties', []),
143         related_topics=parsed.get('related_topics', []),
144         confidence=uncertainty.confidence,
145         used_resolutions=list(retrieved.get_used_levels()),
146         requires_diagram=self._should_generate_diagram(query,
147             parsed)
148     )
149
150 @dataclass
151 class ExpertResponse:
152     answer: str
153     citations: List[Citation]
154     prerequisites: List[str]
155     uncertainties: List[str]
156     related_topics: List[str]
157     confidence: float
158     used_resolutions: List[int]
159     requires_diagram: bool
160
161 @dataclass
162 class Citation:
163     source_type: str # chunk / section / document
164     source_id: UUID
165     excerpt: str
166     relevance_score: float

```

8.3 Pedagogical Mode

Listing 22: Pedagogical Teaching Mode

```

1 class PedagogicalAgent:
2     def __init__(self, reasoning_agent, ontology):
3         self.reasoning = reasoning_agent
4         self.ontology = ontology
5
6     async def teach_concept(self, concept: str) -> TeachingPlan:
7         """
8         Generate structured learning path for a concept.
9         """
10        # Find concept in ontology
11        topic = await self.ontology.find_topic(concept)

```

```

12
13     if not topic:
14         return await self._handle_unknown_concept(concept)
15
16     # Build prerequisite chain
17     prereq_chain = await self._build_prerequisite_chain(topic)
18
19     # Generate teaching plan
20     plan = TeachingPlan(
21         target_concept=concept,
22         prerequisites=prereq_chain,
23         learning_modules=[]
24     )
25
26     # For each level in prerequisite chain
27     for level_topics in prereq_chain:
28         for topic in level_topics:
29             module = await self._create_learning_module(topic)
30             plan.learning_modules.append(module)
31
32     return plan
33
34     async def _build_prerequisite_chain(self,
35                                         topic: TopicNode
36                                         ) -> List[List[TopicNode]]:
37
38         """
39         Build ordered levels of prerequisites.
40         """
41         chains = []
42         current_level = [topic]
43         visited = {topic.topic_id}
44
45         while current_level:
46             next_level = []
47             for t in current_level:
48                 for prereq_id in t.prerequisite_topics:
49                     if prereq_id not in visited:
50                         prereq = await self.ontology.get_topic(
51                             prereq_id)
52                         next_level.append(prereq)
53                         visited.add(prereq_id)
54
55             if next_level:
56                 chains.insert(0, next_level) # Add to beginning
57                 current_level = next_level
58
59             chains.append([topic]) # Target concept at end
60         return chains
61
62     async def _create_learning_module(self,
63                                       topic: TopicNode
64                                       ) -> LearningModule:
65
66         """
67         Create interactive learning module for topic.
68         """
69         # Retrieve all content related to topic
70         sections = await self.reasoning.knowledge.
71             get_sections_for_topic(

```

```

68         topic.topic_id
69     )
70
71     # Generate explanation
72     explanation = await self.reasoning.answer_query(
73         f"Explain the concept of {topic.name} in detail"
74     )
75
76     # Generate exercises
77     exercises = await self._generate_exercises(topic, sections)
78
79     return LearningModule(
80         topic_name=topic.name,
81         difficulty=topic.difficulty_level,
82         explanation=explanation.answer,
83         key_points=await self._extract_key_points(sections),
84         examples=await self._extract_examples(sections),
85         exercises=exercises,
86         further_reading=[
87             s.artifact_id for s in sections
88         ]
89     )
90
91 @dataclass
92 class LearningModule:
93     topic_name: str
94     difficulty: int
95     explanation: str
96     key_points: List[str]
97     examples: List[dict]
98     exercises: List[dict]
99     further_reading: List[UUID]

```

9 Layer 6: Diagrammatic Explanation

9.1 Purpose and Epistemic Question

Epistemic Question: *"Can visual representation enhance understanding of this concept?"*

Layer 6 generates schematic diagrams, concept graphs, and explanatory visuals to augment textual explanations.

9.2 Diagram Generation Architecture

Listing 23: Diagrammatic Explanation Agent

```

1 class DiagramGenerator:
2     def __init__(self, llm):
3         self.llm = llm
4         self.generators = {
5             'concept_map': ConceptMapGenerator(),
6             'flowchart': FlowchartGenerator(),
7             'hierarchy': HierarchyDiagramGenerator(),
8             'graph': GraphDiagramGenerator(),
9             'timeline': TimelineDiagramGenerator()
10        }
11

```

```

12     async def generate_diagram(self,
13                               response: ExpertResponse,
14                               query: str
15                               ) -> Optional[Diagram]:
16         """
17         Determine if diagram would be helpful and generate it.
18         """
19         # Assess diagram utility
20         should_generate = await self._assess_diagram_utility(
21             query, response
22         )
23
24         if not should_generate['useful']:
25             return None
26
27         diagram_type = should_generate['type']
28         generator = self.generators.get(diagram_type)
29
30         if not generator:
31             return None
32
33         # Extract diagram specification
34         spec = await self._extract_diagram_spec(
35             response, diagram_type
36         )
37
38         # Generate diagram
39         diagram = await generator.generate(spec)
40
41         return diagram
42
43     async def _assess_diagram_utility(self,
44                                     query: str,
45                                     response: ExpertResponse
46                                     ) -> dict:
47         """
48         Determine if visual representation would help.
49         """
50         prompt = f"""
51         Query: {query}
52         Response: {response.answer[:500]}...
53
54         Would a diagram enhance understanding of this response?
55
56         Consider:
57         - Relationships between concepts
58         - Hierarchical structures
59         - Temporal sequences
60         - Complex interactions
61
62         If yes, suggest diagram type:
63         - concept_map: Nodes and relationships
64         - flowchart: Process or algorithm
65         - hierarchy: Tree structure
66         - graph: Network relationships
67         - timeline: Temporal progression
68

```



```

69         """Output as JSON: {"useful": true/false, "type": "...", "reason": "..."}"""
70         """
71
72         assessment = await self.llm.ainvoke(prompt)
73         return json.loads(assessment)
74
75     async def _extract_diagram_spec(self,
76                                     response: ExpertResponse,
77                                     diagram_type: str
78                                     ) -> DiagramSpec:
79         """
80         """Extract diagram elements from response.
81         """
82         prompt = f"""
83         Extract diagram specification from this explanation.
84
85         Explanation: {response.answer}
86         Diagram type: {diagram_type}
87
88         Identify:
89         - Nodes/entities
90         - Relationships/edges
91         - Labels
92         - Hierarchy if applicable
93
94         Output as JSON matching {diagram_type} schema.
95         """
96
97         spec_json = await self.llm.ainvoke(prompt)
98         return DiagramSpec.from_json(spec_json, diagram_type)

```

9.3 Concept Map Generation

Listing 24: Concept Map Generator Using Graphviz

```

1 class ConceptMapGenerator:
2     def __init__(self):
3         self.graphviz = Graphviz()
4
5     async def generate(self, spec: DiagramSpec) -> Diagram:
6         """
7         """Generate concept map from specification.
8         """
9         # Create Graphviz graph
10        dot = graphviz.Digraph(comment='Concept Map')
11        dot.attr(rankdir='TB')
12        dot.attr('node', shape='box', style='rounded, filled',
13                fillcolor='lightblue')
14
15        # Add nodes
16        for node in spec.nodes:
17            dot.node(
18                node['id'],
19                label=self._wrap_label(node['label']),
20                fillcolor=self._get_node_color(node.get('type'))
21            )

```

```

22
23     # Add edges
24     for edge in spec.edges:
25         dot.edge(
26             edge['source'],
27             edge['target'],
28             label=edge.get('label', ''),
29             style=self._get_edge_style(edge.get('type'))
30         )
31
32     # Render to SVG
33     svg_content = dot.pipe(format='svg').decode('utf-8')
34
35     return Diagram(
36         diagram_id=uuid4(),
37         diagram_type='concept_map',
38         format='svg',
39         content=svg_content,
40         metadata={
41             'node_count': len(spec.nodes),
42             'edge_count': len(spec.edges)
43         }
44     )
45
46     def _wrap_label(self, text: str, width: int = 20) -> str:
47         """Wrap long labels for better visualization."""
48         words = text.split()
49         lines = []
50         current_line = []
51         current_length = 0
52
53         for word in words:
54             if current_length + len(word) > width:
55                 lines.append(' '.join(current_line))
56                 current_line = [word]
57                 current_length = len(word)
58             else:
59                 current_line.append(word)
60                 current_length += len(word) + 1
61
62         if current_line:
63             lines.append(' '.join(current_line))
64
65         return '\\n'.join(lines)
66
67     @dataclass
68     class Diagram:
69         diagram_id: UUID
70         diagram_type: str
71         format: str # svg / png / mermaid
72         content: str
73         metadata: dict

```

9.4 Mermaid Diagram Generation

Listing 25: Mermaid-Based Diagram Generation

```

1  class FlowchartGenerator:
2      async def generate(self, spec: DiagramSpec) -> Diagram:
3          """
4          Generate flowchart using Mermaid syntax.
5          """
6          mermaid_code = "flowchart TD\n"
7
8          # Add nodes with styling
9          for node in spec.nodes:
10             node_shape = self._get_node_shape(node.get('type'))
11             mermaid_code += f"graph LR; {node['id']} [{node_shape[0]}] -- {node['label']} --> {node_shape[1]}\n"
12
13          # Add edges
14          for edge in spec.edges:
15             arrow = self._get_arrow_type(edge.get('type'))
16             label = edge.get('label', '')
17             if label:
18                 mermaid_code += f"graph LR; {edge['source']} -- {arrow} --> {edge['target']} [{label}]\n"
19             else:
20                 mermaid_code += f"graph LR; {edge['source']} -- {arrow} --> {edge['target']}\n"
21
22          # Add styling
23          mermaid_code += self._add_styling(spec)
24
25          return Diagram(
26              diagram_id=uuid4(),
27              diagram_type='flowchart',
28              format='mermaid',
29              content=mermaid_code,
30              metadata={'mermaid_version': '9.0'}
31          )
32
33      def _get_node_shape(self, node_type: str) -> tuple:
34          """Map node types to Mermaid shapes."""
35          shapes = {
36              'start': ('[', ']'),
37              'end': ('([', '])'),
38              'decision': ('{', '}'),
39              'process': ('[', ']'),
40              'data': ('[/', '/']')
41          }
42          return shapes.get(node_type, ('[', ']))

```

10 Layer 7: Deep External Research Agent

10.1 Purpose and Epistemic Question

Epistemic Question: *"Does authoritative external knowledge exist that materially alters or invalidates the current internal understanding?"*

Layer 7 implements the sole deep agent in the architecture, designed to prevent epistemic stagnation through controlled external research.

10.2 DERA Architecture

Listing 26: Deep External Research Agent (DERA)

```

1 class DeepExternalResearchAgent:
2     def __init__(self, llm, web_search_tool, academic_apis):
3         self.llm = llm
4         self.web_search = web_search_tool
5         self.academic = academic_apis
6         self.trigger_detector = ResearchTriggerDetector()
7
8     async def conduct_research(self,
9                               query: str,
10                              uncertainty: EpistemicUncertainty,
11                              internal_knowledge: RetrievedKnowledge
12                              ) -> ResearchReport:
13         """
14         Conduct deep external research and produce comparative report.
15         """
16         # Step 1: Formulate research questions
17         questions = await self._formulate_research_questions(
18             query, uncertainty, internal_knowledge
19         )
20
21         # Step 2: Execute multi-source research
22         findings = await self._execute_research(questions)
23
24         # Step 3: Compare with internal knowledge
25         comparison = await self._compare_knowledge(
26             internal_knowledge, findings
27         )
28
29         # Step 4: Generate epistemic report
30         report = await self._generate_report(
31             query, questions, findings, comparison
32         )
33
34         return report
35
36     async def _formulate_research_questions(self,
37                                             query: str,
38                                             uncertainty:
39                                                 EpistemicUncertainty,
40                                             internal: RetrievedKnowledge
41                                             ) -> List[ResearchQuestion]:
42         """
43         Generate targeted research questions.
44         """
45         prompt = f"""
46         Original query: {query}
47
48         Internal knowledge summary: {internal.summarize()}
49
50         Epistemic concerns:
51         - Temporal: {uncertainty.has_temporal_concerns}
52         - Conflicts: {uncertainty.has_conflicts}
53         - Confidence: {uncertainty.confidence}
54
55         Generate specific research questions that would:

```

```

55 1. Update potentially outdated information
56 2. Resolve conflicting claims
57 3. Fill knowledge gaps
58 4. Verify critical facts
59
60 For each question, specify:
61 - Question text
62 - Expected source type (academic | news | technical docs)
63 - Priority (high | medium | low)
64
65 Output as JSON list.
66 """
67
68 response = await self.llm.ainvoke(prompt)
69 questions_data = json.loads(response)
70
71 return [ResearchQuestion(**q) for q in questions_data]
72
73 async def _execute_research(self,
74                             questions: List[ResearchQuestion]
75                             ) -> List[ResearchFinding]:
76     """
77     Execute research across multiple sources.
78     """
79     findings = []
80
81     for question in questions:
82         if question.expected_source == 'academic':
83             # Search academic databases
84             academic_results = await self.academic.search(
85                 question.text,
86                 filters={'since': question.get_temporal_filter()})
87             findings.extend(
88                 self._process_academic_results(academic_results)
89             )
90
91         elif question.expected_source == 'news':
92             # Search news sources
93             news_results = await self.web_search.search(
94                 question.text,
95                 domain_filter='news',
96                 recency='6months'
97             )
98             findings.extend(
99                 self._process_news_results(news_results)
100             )
101
102         else:
103             # General web search
104             web_results = await self.web_search.search(
105                 question.text
106             )
107             findings.extend(
108                 self._process_web_results(web_results)
109             )
110
111     return findings
112

```

```

113
114     async def _compare_knowledge(self,
115                                   internal: RetrievedKnowledge,
116                                   external: List[ResearchFinding]
117                                   ) -> KnowledgeComparison:
118         """
119         Compare internal and external knowledge.
120         """
121         comparison_prompt = f"""
122         Compare internal knowledge with external research findings.
123
124         Internal knowledge:
125         {internal.detailed_summary()}
126
127         External findings:
128         {self._summarize_findings(external)}
129
130         Identify:
131         1. Confirmations: External sources confirm internal knowledge
132         2. Updates: External sources provide more recent information
133         3. Contradictions: External sources contradict internal
134            knowledge
135         4. Extensions: External sources add new relevant information
136         5. Confidence assessment: How reliable are external sources?
137
138         For contradictions, evaluate:
139         - Source authority
140         - Publication date
141         - Evidence quality
142         - Consensus vs outlier
143
144         Output as structured JSON.
145         """
146
147         comparison = await self.llm.ainvoke(comparison_prompt)
148         parsed = json.loads(comparison)
149
150         return KnowledgeComparison(
151             confirmations=parsed['confirmations'],
152             updates=parsed['updates'],
153             contradictions=parsed['contradictions'],
154             extensions=parsed['extensions'],
155             confidence_assessment=parsed['confidence'],
156             recommendation=self._generate_recommendation(parsed)
157         )
158
159     async def _generate_report(self,
160                                query: str,
161                                questions: List[ResearchQuestion],
162                                findings: List[ResearchFinding],
163                                comparison: KnowledgeComparison
164                                ) -> ResearchReport:
165         """
166         Generate comprehensive epistemic research report.
167         """
168         return ResearchReport(
169             report_id=uuid4(),
170             query=query,

```

```

170         research_questions=questions,
171         external_findings=findings,
172         comparison=comparison,
173         timestamp=datetime.now(),
174
175         # Key insights
176         should_update_knowledge=len(comparison.updates) > 0 or
177                                 len(comparison.contradictions) > 0,
178         critical_changes=comparison.get_critical_changes(),
179         new_sources_to_ingest=self._identify_valuable_sources(
180             findings),
181
182         # Versioning
183         affected_documents=[
184             doc.document_id for doc in comparison.
185                 get_affected_documents()
186         ],
187         proposed_updates=comparison.generate_update_proposals()
188     )
189
190 @dataclass
191 class ResearchReport:
192     report_id: UUID
193     query: str
194     research_questions: List[ResearchQuestion]
195     external_findings: List[ResearchFinding]
196     comparison: KnowledgeComparison
197     timestamp: datetime
198     should_update_knowledge: bool
199     critical_changes: List[str]
200     new_sources_to_ingest: List[str]
201     affected_documents: List[UUID]
202     proposed_updates: List[dict]

```

10.3 Research Trigger Detection

Listing 27: Automatic Research Trigger Detection

```

1 class ResearchTriggerDetector:
2     def should_trigger_research(self,
3                                 query: str,
4                                 uncertainty: EpistemicUncertainty,
5                                 retrieved: RetrievedKnowledge
6                                 ) -> bool:
7
8         """
9         Determine if deep research should be triggered.
10        """
11        triggers = []
12
13        # Trigger 1: Explicit temporal concerns
14        if uncertainty.has_temporal_concerns:
15            triggers.append('temporal_concern')
16
17        # Trigger 2: Low confidence
18        if uncertainty.confidence < 0.6:
19            triggers.append('low_confidence')

```

```

20     # Trigger 3: Conflicting internal sources
21     if uncertainty.has_conflicts:
22         triggers.append('internal_conflicts')
23
24     # Trigger 4: Sparse internal knowledge
25     if self._is_sparse(retrieved):
26         triggers.append('sparse_knowledge')
27
28     # Trigger 5: Explicit request for recent information
29     if self._requests_recent_info(query):
30         triggers.append('explicit_recency_request')
31
32     # Trigger 6: Known outdated topic
33     if self._is_fast_evolving_topic(query):
34         triggers.append('fast_evolving_topic')
35
36     return len(triggers) > 0
37
38     def _is_fast_evolving_topic(self, query: str) -> bool:
39         """
40         Detect topics known to evolve rapidly.
41         """
42         fast_topics = [
43             'covid', 'pandemic', 'election', 'policy',
44             'regulation', 'market', 'stock', 'crypto',
45             'ai_model', 'language_model', 'gpt'
46         ]
47         return any(topic in query.lower() for topic in fast_topics)

```

10.4 Knowledge Versioning

Listing 28: Versioned Knowledge Updates

```

1 class KnowledgeVersionManager:
2     def __init__(self, db):
3         self.db = db
4
5     async def apply_research_updates(self,
6                                     report: ResearchReport):
7         """
8         Apply research findings as versioned updates.
9         """
10        for update in report.proposed_updates:
11            # Create new version
12            version = KnowledgeVersion(
13                version_id=uuid4(),
14                target_id=update['target_id'],
15                target_type=update['target_type'],
16                change_type=update['change_type'],
17                old_content=await self._get_current_content(
18                    update['target_id']
19                ),
20                new_content=update['new_content'],
21                source_report=report.report_id,
22                timestamp=datetime.now(),
23                validated=False
24            )

```



```

25         await self.db.save_version(version)
26
27     # Mark report as processed
28     await self.db.mark_report_processed(report.report_id)
29
30     async def get_version_history(self,
31                                   entity_id: UUID
32                                   ) -> List[KnowledgeVersion]:
33
34         """
35         Retrieve version history for any knowledge entity.
36         """
37         return await self.db.query_versions(entity_id)
38
39 @dataclass
40 class KnowledgeVersion:
41     version_id: UUID
42     target_id: UUID
43     target_type: str
44     change_type: str # update / contradiction / extension
45     old_content: dict
46     new_content: dict
47     source_report: UUID
48     timestamp: datetime
49     validated: bool

```

11 User Interaction Layer: CLI Interface

11.1 CLI Architecture

The user interaction layer is built on the llm framework by Simon Willison, providing a scriptable, composable command-line interface.

Listing 29: KNOWY CLI Plugin Implementation

```

1  import llm
2  from typing import Optional
3
4  @llm.hookimpl
5  def register_commands(cli):
6      @cli.group()
7      def knowy():
8          """KNOWY personal knowledge system commands."""
9          pass
10
11     @knowy.command()
12     @click.argument("reference")
13     @click.option("--origin", default="manual")
14     def ingest(reference, origin):
15         """Ingest a new knowledge reference."""
16         ingestion_service = IngestionService()
17         result = asyncio.run(
18             ingestion_service.ingest(reference, origin)
19         )
20         click.echo(f"Ingested: {result.source_id}")
21
22     @knowy.command()

```

```

23 @click.argument("query")
24 @click.option("--mode", type=click.Choice(['expert', 'teach']),
25               default='expert')
26 @click.option("--with-diagram", is_flag=True)
27 def ask(query, mode, with_diagram):
28     """Query the knowledge system."""
29     if mode == 'expert':
30         agent = ExpertReasoningAgent()
31         response = asyncio.run(agent.answer_query(query))
32     else:
33         agent = PedagogicalAgent()
34         response = asyncio.run(agent.teach_concept(query))
35
36     click.echo(response.answer)
37
38     if with_diagram and response.requires_diagram:
39         diagram_gen = DiagramGenerator()
40         diagram = asyncio.run(
41             diagram_gen.generate_diagram(response, query)
42         )
43         if diagram:
44             save_diagram(diagram)
45
46 @knowy.command()
47 @click.argument("topic")
48 def topics(topic):
49     """Explore topic ontology."""
50     ontology = OntologyService()
51     topic_info = asyncio.run(ontology.get_topic_info(topic))
52
53     click.echo(f"Topic: {topic_info.name}")
54     click.echo(f>Description: {topic_info.description}")
55     click.echo(f"
56 Prerequisites:")
57     for prereq in topic_info.prerequisite_topics:
58         click.echo(f"  {prereq.name}")
59
60 @knowy.command()
61 @click.option("--since", type=click.DateTime())
62 def research_reports(since):
63     """View deep research reports."""
64     research_service = ResearchService()
65     reports = asyncio.run(
66         research_service.get_reports(since=since)
67     )
68
69     for report in reports:
70         click.echo(f"
71 Report: {report.report_id}")
72         click.echo(f"Query: {report.query}")
73         click.echo(f"Findings: {len(report.external_findings)}")
74         click.echo(f"Updates proposed: {len(report.proposed_updates
75 )}")

```

11.2 CLI Usage Examples

Listing 30: CLI Usage Examples

```

1 # Ingest a paper from URL
2 $ llm knowy ingest "https://arxiv.org/abs/2103.00020" --origin browser
3
4 # Ingest a YouTube course
5 $ llm knowy ingest "https://youtube.com/playlist?list=..." --origin
   manual
6
7 # Ask a question
8 $ llm knowy ask "Explain the transformer architecture"
9
10 # Get pedagogical explanation with diagram
11 $ llm knowy ask "How does backpropagation work?" --mode teach --with-
   diagram
12
13 # Explore topics
14 $ llm knowy topics "machine learning"
15
16 # View research reports
17 $ llm knowy research-reports --since "2026-01-01"
18
19 # Scriptable workflows
20 $ cat papers.txt | while read url; do
21     llm knowy ingest "$url"
22 done
23
24 # Export knowledge to file
25 $ llm knowy ask "Summarize all papers on attention mechanisms" >
   summary.md

```

12 System Integration and Orchestration

12.1 LangGraph Orchestration

Listing 31: LangGraph Workflow Orchestration

```

1 from langgraph.graph import StateGraph, END
2 from typing import TypedDict, Annotated
3 import operator
4
5 class KnowyState(TypedDict):
6     """Global state for KNOWY orchestration."""
7     # Input
8     user_query: str
9     mode: str # ingest / query / teach / research
10
11     # Ingestion flow
12     raw_input: Optional[str]
13     source_descriptor: Optional[CanonicalSourceDescriptor]
14     artifacts: Optional[List[KnowledgeArtifact]]
15     acquired_content: Optional[List[Any]]
16
17     # Query flow
18     retrieved_knowledge: Optional[RetrievedKnowledge]
19     epistemic_uncertainty: Optional[EpistemicUncertainty]
20     expert_response: Optional[ExpertResponse]

```

```

21     diagram: Optional[Diagram]
22
23     # Research flow
24     research_triggered: bool
25     research_report: Optional[ResearchReport]
26
27     # Output
28     final_response: str
29     metadata: dict
30
31 def create_knowy_workflow() -> StateGraph:
32     """
33     Create the complete KNOWY orchestration graph.
34     """
35     workflow = StateGraph(KnowyState)
36
37     # Add nodes for each layer
38     workflow.add_node("route_mode", route_by_mode)
39     workflow.add_node("ingest_l0", layer0_ingestion)
40     workflow.add_node("expand_l1", layer1_expansion)
41     workflow.add_node("acquire_l2", layer2_acquisition)
42     workflow.add_node("represent_l3", layer3_representation)
43     workflow.add_node("classify_l4", layer4_classification)
44     workflow.add_node("reason_l5", layer5_reasoning)
45     workflow.add_node("diagram_l6", layer6_diagramming)
46     workflow.add_node("research_l7", layer7_deep_research)
47     workflow.add_node("finalize", finalize_response)
48
49     # Define edges
50     workflow.set_entry_point("route_mode")
51
52     # Routing logic
53     workflow.add_conditional_edges(
54         "route_mode",
55         determine_flow,
56         {
57             "ingest": "ingest_l0",
58             "query": "reason_l5",
59             "teach": "reason_l5",
60             "research": "research_l7"
61         }
62     )
63
64     # Ingestion flow
65     workflow.add_edge("ingest_l0", "expand_l1")
66     workflow.add_edge("expand_l1", "acquire_l2")
67     workflow.add_edge("acquire_l2", "represent_l3")
68     workflow.add_edge("represent_l3", "classify_l4")
69     workflow.add_edge("classify_l4", "finalize")
70
71     # Query flow
72     workflow.add_conditional_edges(
73         "reason_l5",
74         check_research_needed,
75         {
76             "research": "research_l7",
77             "diagram": "diagram_l6",
78             "finalize": "finalize"

```

```

79     }
80 )
81
82 workflow.add_edge("diagram_16", "finalize")
83 workflow.add_edge("research_17", "reason_15") # Loop back
84 workflow.add_edge("finalize", END)
85
86 return workflow.compile()
87
88 def determine_flow(state: KnowyState) -> str:
89     """Route to appropriate flow based on mode."""
90     return state["mode"]
91
92 def check_research_needed(state: KnowyState) -> str:
93     """Determine next step after reasoning."""
94     if state.get("research_triggered"):
95         return "research"
96     elif state.get("expert_response", {}).get("requires_diagram"):
97         return "diagram"
98     else:
99         return "finalize"

```

12.2 Deployment Architecture

Listing 32: Complete Deployment Configuration

```

1  # docker-compose.yml equivalent in Python config
2
3  @dataclass
4  class DeploymentConfig:
5      """KNOWY deployment configuration."""
6
7      # Local LLM configuration
8      ollama_host: str = "localhost"
9      ollama_port: int = 11434
10     primary_model: str = "llama3.2"
11     vision_model: str = "llama3.2-vision"
12
13     # Database configuration
14     postgres_host: str = "localhost"
15     postgres_port: int = 5432
16     postgres_db: str = "knowy_knowledge"
17     neo4j_host: str = "localhost"
18     neo4j_port: int = 7687
19
20     # Vector database
21     vector_db_type: str = "qdrant" # or "faiss"
22     qdrant_host: str = "localhost"
23     qdrant_port: int = 6333
24
25     # Storage paths
26     data_dir: Path = Path.home() / ".knowy" / "data"
27     cache_dir: Path = Path.home() / ".knowy" / "cache"
28     models_dir: Path = Path.home() / ".knowy" / "models"
29
30     # Processing configuration
31     max_concurrent_acquisitions: int = 5

```

```

32     chunk_size: int = 512
33     embedding_dim: int = 1536
34
35     # Research configuration
36     enable_deep_research: bool = True
37     research_confidence_threshold: float = 0.6
38     max_research_sources: int = 20
39
40 def initialize_system(config: DeploymentConfig):
41     """Initialize KNOWY system with configuration."""
42     # Create directories
43     config.data_dir.mkdir(parents=True, exist_ok=True)
44     config.cache_dir.mkdir(parents=True, exist_ok=True)
45
46     # Initialize databases
47     db_manager = DatabaseManager(config)
48     db_manager.initialize_schemas()
49
50     # Initialize Ollama connection
51     ollama_client = OllamaClient(
52         host=config.ollama_host,
53         port=config.ollama_port
54     )
55
56     # Pull required models
57     ollama_client.pull(config.primary_model)
58     ollama_client.pull(config.vision_model)
59
60     # Initialize vector database
61     if config.vector_db_type == "qdrant":
62         vector_db = QdrantClient(
63             host=config.qdrant_host,
64             port=config.qdrant_port
65         )
66         vector_db.create_collection(
67             collection_name="knowy_chunks",
68             vectors_config=VectorParams(
69                 size=config.embedding_dim,
70                 distance=Distance.COSINE
71             )
72         )
73
74     # Initialize services
75     services = ServiceRegistry(
76         db_manager=db_manager,
77         llm_client=ollama_client,
78         vector_db=vector_db,
79         config=config
80     )
81
82     return services

```

13 Performance Optimization

13.1 Caching Strategy

Listing 33: Multi-Level Caching Implementation

```

1 class CacheManager:
2     def __init__(self, config: DeploymentConfig):
3         self.config = config
4         self.embedding_cache = LRUCache(maxsize=10000)
5         self.query_cache = TTLCache(maxsize=1000, ttl=3600)
6
7     async def get_or_compute_embedding(self,
8                                       text: str,
9                                       model: str
10                                      ) -> np.ndarray:
11         """Cache embeddings to avoid recomputation."""
12         cache_key = hash(text + model)
13
14         if cache_key in self.embedding_cache:
15             return self.embedding_cache[cache_key]
16
17         embedding = await self._compute_embedding(text, model)
18         self.embedding_cache[cache_key] = embedding
19         return embedding
20
21     async def get_or_execute_query(self,
22                                   query: str,
23                                   mode: str
24                                  ) -> Any:
25         """Cache query results with TTL."""
26         cache_key = hash(query + mode)
27
28         if cache_key in self.query_cache:
29             return self.query_cache[cache_key]
30
31         result = await self._execute_query(query, mode)
32         self.query_cache[cache_key] = result
33         return result

```

13.2 Parallel Processing

Listing 34: Parallel Acquisition Processing

```

1 async def parallel_acquisition(artifacts: List[KnowledgeArtifact],
2                               max_concurrent: int = 5
3                               ) -> List[AcquiredContent]:
4     """Process multiple artifacts concurrently."""
5     semaphore = asyncio.Semaphore(max_concurrent)
6
7     async def acquire_with_limit(artifact):
8         async with semaphore:
9             return await acquire_artifact(artifact)
10
11     tasks = [acquire_with_limit(a) for a in artifacts]
12     results = await asyncio.gather(*tasks, return_exceptions=True)
13
14     # Filter out exceptions and log failures
15     successful = []
16     for artifact, result in zip(artifacts, results):
17         if isinstance(result, Exception):

```

```

18         logger.error(f"Failed to acquire {artifact.artifact_id}: {
19             result}")
20     else:
21         successful.append(result)
22     return successful

```

14 Evaluation and Metrics

14.1 Epistemic Integrity Metrics

Listing 35: Epistemic Quality Metrics

```

1 @dataclass
2 class EpistemicMetrics:
3     """Metrics for evaluating epistemic integrity."""
4
5     # Hallucination detection
6     hallucination_rate: float
7     unsupported_claims: int
8     total_claims: int
9
10    # Citation quality
11    citation_precision: float # Relevant citations / total citations
12    citation_coverage: float # Claims with citations / total claims
13
14    # Uncertainty quantification
15    avg_confidence: float
16    confidence_calibration: float # How well confidence matches
17    accuracy
18
19    # Knowledge freshness
20    avg_document_age_days: float
21    outdated_documents: int
22    research_trigger_rate: float
23
24    def compute_epistemic_score(self) -> float:
25        """Composite score for epistemic quality."""
26        return (
27            (1 - self.hallucination_rate) * 0.3 +
28            self.citation_precision * 0.2 +
29            self.citation_coverage * 0.2 +
30            self.confidence_calibration * 0.15 +
31            min(1.0, 365 / max(self.avg_document_age_days, 1)) * 0.15
32        )
33
34    class MetricsCollector:
35        def __init__(self):
36            self.metrics_db = MetricsDatabase()
37
38        async def evaluate_response(self,
39            query: str,
40            response: ExpertResponse,
41            retrieved: RetrievedKnowledge
42        ) -> EpistemicMetrics:
43            """Evaluate a single response."""
44            # Detect hallucinations

```



```

44     hallucination_check = await self._check_hallucinations(
45         response.answer, retrieved
46     )
47
48     # Evaluate citations
49     citation_metrics = self._evaluate_citations(
50         response.citations, retrieved
51     )
52
53     # Check freshness
54     freshness = self._compute_freshness(retrieved)
55
56     return EpistemicMetrics(
57         hallucination_rate=hallucination_check.rate,
58         unsupported_claims=hallucination_check.unsupported,
59         total_claims=hallucination_check.total,
60         citation_precision=citation_metrics.precision,
61         citation_coverage=citation_metrics.coverage,
62         avg_confidence=response.confidence,
63         confidence_calibration=0.0, # Requires ground truth
64         avg_document_age_days=freshness.avg_age,
65         outdated_documents=freshness.outdated_count,
66         research_trigger_rate=0.0 # Computed over time
67     )

```

14.2 Performance Metrics

Listing 36: System Performance Monitoring

```

1  @dataclass
2  class PerformanceMetrics:
3      """Performance and efficiency metrics."""
4
5      # Latency
6      avg_query_latency_ms: float
7      p95_query_latency_ms: float
8      p99_query_latency_ms: float
9
10     # Throughput
11     queries_per_minute: float
12     ingestions_per_hour: float
13
14     # Resource utilization
15     avg_memory_usage_gb: float
16     peak_memory_usage_gb: float
17     avg_cpu_percent: float
18
19     # Storage
20     total_chunks: int
21     total_documents: int
22     total_size_gb: float
23
24     # Cache efficiency
25     embedding_cache_hit_rate: float
26     query_cache_hit_rate: float
27
28  class PerformanceMonitor:

```

```

29     def __init__(self):
30         self.start_time = time.time()
31         self.query_latencies = []
32
33     @contextmanager
34     def measure_latency(self):
35         """Context manager for latency measurement."""
36         start = time.time()
37         yield
38         latency = (time.time() - start) * 1000
39         self.query_latencies.append(latency)
40
41     def get_metrics(self) -> PerformanceMetrics:
42         """Compute current performance metrics."""
43         return PerformanceMetrics(
44             avg_query_latency_ms=np.mean(self.query_latencies),
45             p95_query_latency_ms=np.percentile(self.query_latencies,
46                                                95),
47             p99_query_latency_ms=np.percentile(self.query_latencies,
48                                                99),
49             queries_per_minute=len(self.query_latencies) /
50                                 ((time.time() - self.start_time) / 60),
51             # ... other metrics
52         )

```

15 Testing Strategy

15.1 Unit Testing

Listing 37: Layer Testing Examples

```

1 import pytest
2
3 class TestLayer0Ingestion:
4     @pytest.fixture
5     def ingestion_agent(self):
6         return IngestionAgent()
7
8     @pytest.mark.asyncio
9     async def test_pdf_detection(self, ingestion_agent):
10         """Test PDF format detection."""
11         url = "https://arxiv.org/pdf/2103.00020.pdf"
12         descriptor = await ingestion_agent.create_descriptor(url, "
13             manual")
14
15         assert descriptor.content_hints.contains_pdf is True
16         assert descriptor.content_hints.contains_video is False
17
18     @pytest.mark.asyncio
19     async def test_youtube_detection(self, ingestion_agent):
20         """Test YouTube video detection."""
21         url = "https://www.youtube.com/watch?v=abc123"
22         descriptor = await ingestion_agent.create_descriptor(url, "
23             browser")
24
25         assert descriptor.content_hints.contains_video is True
26         assert descriptor.links == [url]

```

```

25
26 class TestLayer3Representation:
27     @pytest.fixture
28     def summarizer(self):
29         return SectionSummarizer(llm=MockLLM())
30
31     @pytest.mark.asyncio
32     async def test_section_summary_generation(self, summarizer):
33         """Test section summary creation."""
34         chunks = create_mock_chunks(count=5)
35         summary = await summarizer.generate_section_summary(chunks)
36
37         assert summary.main_claim is not None
38         assert len(summary.key_points) > 0
39         assert len(summary.chunk_ids) == 5
40
41 class TestLayer7Research:
42     @pytest.mark.asyncio
43     async def test_research_trigger_temporal(self):
44         """Test research trigger on temporal concerns."""
45         detector = ResearchTriggerDetector()
46         uncertainty = EpistemicUncertainty(
47             has_temporal_concerns=True,
48             confidence=0.8
49         )
50
51         should_trigger = detector.should_trigger_research(
52             "What is the current policy?",
53             uncertainty,
54             MockRetrievedKnowledge()
55         )
56
57         assert should_trigger is True

```

15.2 Integration Testing

Listing 38: End-to-End Integration Tests

```

1 class TestEndToEnd:
2     @pytest.mark.asyncio
3     @pytest.mark.slow
4     async def test_complete_ingestion_pipeline(self):
5         """Test complete ingestion from URL to representation."""
6         system = await initialize_test_system()
7
8         # Ingest a paper
9         result = await system.ingest(
10             "https://arxiv.org/abs/1706.03762", # Attention Is All You
11             # Need
12             origin="test"
13         )
14
15         # Verify all layers executed
16         assert result.source_descriptor is not None
17         assert len(result.artifacts) > 0
18         assert len(result.sections) > 0
19         assert result.document_representation is not None

```

```

19
20     # Verify queryable
21     response = await system.query(
22         "Explain the transformer architecture"
23     )
24     assert "attention" in response.answer.lower()
25     assert len(response.citations) > 0
26
27 @pytest.mark.asyncio
28 async def test_research_integration(self):
29     """Test deep research integration."""
30     system = await initialize_test_system()
31
32     # Create scenario requiring research
33     response = await system.query(
34         "What are the latest developments in LLMs?",
35         force_research=True
36     )
37
38     assert response.research_report is not None
39     assert len(response.research_report.external_findings) > 0

```

16 Security and Privacy

16.1 Data Privacy

Listing 39: Privacy-Preserving Design

```

1 class PrivacyManager:
2     """Ensure local-first privacy guarantees."""
3
4     def __init__(self, config: DeploymentConfig):
5         self.config = config
6         self.encryption_key = self._load_or_generate_key()
7
8     def encrypt_sensitive_content(self, content: str) -> bytes:
9         """Encrypt sensitive content before storage."""
10        from cryptography.fernet import Fernet
11        f = Fernet(self.encryption_key)
12        return f.encrypt(content.encode())
13
14    def anonymize_for_research(self,
15                               query: str
16                               ) -> str:
17        """Remove PII before external research."""
18        # Remove email addresses
19        query = re.sub(r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b',
20                       '[EMAIL]', query)
21
22        # Remove phone numbers
23        query = re.sub(r'\b\d{3}[-.]?\d{3}[-.]?\d{4}\b',
24                       '[PHONE]', query)
25
26        # Remove names (simple heuristic)
27        query = self._remove_proper_names(query)
28

```

```

29     return query
30
31     def validate_no_data_leakage(self,
32                                 response: ExpertResponse
33                                 ) -> bool:
34         """Ensure response doesn't leak private info."""
35         # Check for common PII patterns
36         has_email = '@' in response.answer
37         has_phone = re.search(r'\d{3}[-.]?\d{3}[-.]? \d{4}', response.
38                               answer)
39
39         return not (has_email or has_phone)

```

17 Future Extensions

17.1 Planned Enhancements

Enhancement	Description
Collaborative Knowledge	Multi-user knowledge sharing with access control
Active Learning	System proposes questions to fill knowledge gaps
Spaced Repetition	Integration with learning science for retention
Voice Interface	Speech-to-text query input and text-to-speech output
Mobile App	iOS/Android clients with offline capability
Browser Copilot	Real-time knowledge assistance while browsing
Export Formats	Export knowledge as books, courses, documentation
Graph Visualization	Interactive exploration of knowledge graph

Table 5: Planned Future Enhancements

18 Conclusion

KNOWY represents a comprehensive solution to the challenge of personal knowledge management in the age of AI. By enforcing strict epistemic boundaries, implementing multi-resolution knowledge representation, and integrating controlled external research, the system achieves:

- **Epistemic Integrity:** Hallucination minimization through grounded reasoning and uncertainty quantification
- **Multimodal Understanding:** Native support for video, diagrams, and text with semantic alignment
- **Knowledge Evolution:** Principled external research without epistemic drift
- **Expert Reasoning:** Multi-level abstraction enabling both fact lookup and conceptual understanding

- **Pedagogical Capability:** Structured teaching with prerequisite chains and learning modules
- **Local-First Design:** Privacy preservation and workflow composability through CLI

The architecture is production-ready and provides a robust foundation for future enhancements in collaborative knowledge sharing, active learning, and advanced user interfaces.

A Appendix A: Database Schemas

Listing 40: Complete PostgreSQL Schema

```

1  -- Ingestion layer
2  CREATE TABLE ingestion_sources (
3      source_id UUID PRIMARY KEY,
4      origin VARCHAR(50) NOT NULL,
5      raw_reference TEXT NOT NULL,
6      raw_text TEXT,
7      links JSONB,
8      content_hints JSONB NOT NULL,
9      timestamp TIMESTAMPTZ NOT NULL,
10     created_at TIMESTAMPTZ DEFAULT NOW()
11 );
12
13 -- Knowledge artifacts
14 CREATE TABLE knowledge_artifacts (
15     artifact_id UUID PRIMARY KEY,
16     source_id UUID REFERENCES ingestion_sources(source_id),
17     artifact_type VARCHAR(50) NOT NULL,
18     title TEXT NOT NULL,
19     authors JSONB,
20     url TEXT,
21     doi TEXT,
22     isbn TEXT,
23     publication_date TIMESTAMPTZ,
24     metadata JSONB,
25     created_at TIMESTAMPTZ DEFAULT NOW()
26 );
27
28 -- Raw chunks (Level 0)
29 CREATE TABLE raw_chunks (
30     chunk_id UUID PRIMARY KEY,
31     artifact_id UUID REFERENCES knowledge_artifacts(artifact_id),
32     chunk_type VARCHAR(50),
33     text_content TEXT,
34     visual_content BYTEA,
35     structured_content JSONB,
36     position_in_artifact INTEGER,
37     section_id UUID,
38     embedding vector(1536),
39     created_at TIMESTAMPTZ DEFAULT NOW()
40 );
41
42 -- Section representations (Level 1)
43 CREATE TABLE section_representations (
44     section_id UUID PRIMARY KEY,
45     artifact_id UUID REFERENCES knowledge_artifacts(artifact_id),

```

```

46     title TEXT NOT NULL,
47     main_claim TEXT,
48     key_points JSONB,
49     terminology JSONB,
50     chunk_ids JSONB,
51     embedding vector(1536),
52     created_at TIMESTAMPTZ DEFAULT NOW()
53 );
54
55 -- Document representations (Level 2)
56 CREATE TABLE document_representations (
57     document_id UUID PRIMARY KEY,
58     artifact_id UUID REFERENCES knowledge_artifacts(artifact_id),
59     central_thesis TEXT,
60     contributions JSONB,
61     methodology TEXT,
62     findings JSONB,
63     novel_concepts JSONB,
64     embedding vector(1536),
65     created_at TIMESTAMPTZ DEFAULT NOW()
66 );
67
68 -- Thematic synthesis (Level 3)
69 CREATE TABLE thematic_syntheses (
70     theme_id UUID PRIMARY KEY,
71     theme_name TEXT NOT NULL,
72     document_ids JSONB,
73     agreements JSONB,
74     disagreements JSONB,
75     key_terminology JSONB,
76     embedding vector(1536),
77     created_at TIMESTAMPTZ DEFAULT NOW()
78 );
79
80 -- Research reports
81 CREATE TABLE research_reports (
82     report_id UUID PRIMARY KEY,
83     query TEXT NOT NULL,
84     external_findings JSONB,
85     comparison JSONB,
86     should_update_knowledge BOOLEAN,
87     affected_documents JSONB,
88     proposed_updates JSONB,
89     timestamp TIMESTAMPTZ DEFAULT NOW()
90 );
91
92 -- Knowledge versions
93 CREATE TABLE knowledge_versions (
94     version_id UUID PRIMARY KEY,
95     target_id UUID NOT NULL,
96     target_type VARCHAR(50),
97     change_type VARCHAR(50),
98     old_content JSONB,
99     new_content JSONB,
100     source_report UUID REFERENCES research_reports(report_id),
101     validated BOOLEAN DEFAULT FALSE,
102     timestamp TIMESTAMPTZ DEFAULT NOW()
103 );

```

```

104
105 -- Indexes
106 CREATE INDEX idx_artifacts_type ON knowledge_artifacts(artifact_type);
107 CREATE INDEX idx_chunks_artifact ON raw_chunks(artifact_id);
108 CREATE INDEX idx_chunks_embedding ON raw_chunks
109     USING ivfflat (embedding vector_cosine_ops);
110 CREATE INDEX idx_sections_artifact ON section_representations(
111     artifact_id);
112 CREATE INDEX idx_documents_artifact ON document_representations(
113     artifact_id);

```

B Appendix B: Configuration Files

Listing 41: System Configuration (config.yaml)

```

1 # KNOWY Configuration File
2
3 system:
4   name: "KNOWY"
5   version: "1.0.0"
6   data_dir: "~/knowy/data"
7   cache_dir: "~/knowy/cache"
8   log_level: "INFO"
9
10 llm:
11   provider: "ollama"
12   host: "localhost"
13   port: 11434
14   models:
15     primary: "llama3.2:latest"
16     vision: "llama3.2-vision:latest"
17     embedding: "nomic-embed-text:latest"
18
19 databases:
20   postgres:
21     host: "localhost"
22     port: 5432
23     database: "knowy_knowledge"
24     user: "knowy"
25     password: "${POSTGRES_PASSWORD}"
26
27   neo4j:
28     host: "localhost"
29     port: 7687
30     user: "neo4j"
31     password: "${NEO4J_PASSWORD}"
32
33   qdrant:
34     host: "localhost"
35     port: 6333
36     collection_name: "knowy_chunks"
37
38 processing:
39   max_concurrent_acquisitions: 5
40   chunk_size: 512
41   chunk_overlap: 50

```



```

42     embedding_batch_size: 32
43
44 video_analysis:
45     frame_extraction_fps: 1
46     shot_detection_threshold: 0.3
47     max_frames_per_shot: 10
48
49 research:
50     enabled: true
51     confidence_threshold: 0.6
52     max_sources: 20
53     academic_api_key: "${SEMANTIC_SCHOLAR_API_KEY}"
54
55 cache:
56     embedding_cache_size: 10000
57     query_cache_size: 1000
58     query_cache_ttl_seconds: 3600

```

C Appendix C: API Reference

Listing 42: Public API Reference

```

1 from knowy import KNOWY
2
3 # Initialize system
4 system = KNOWY.initialize(config_path="config.yaml")
5
6 # Ingestion API
7 system.ingest(
8     reference: str,
9     origin: str = "manual",
10    metadata: dict = None
11 ) -> IngestionResult
12
13 # Query API
14 system.query(
15     query: str,
16     mode: str = "expert", # expert / teach
17     with_diagram: bool = False,
18     force_research: bool = False
19 ) -> QueryResult
20
21 # Teaching API
22 system.teach(
23     concept: str,
24     difficulty: int = None # Auto-detect if None
25 ) -> TeachingPlan
26
27 # Ontology API
28 system.ontology.get_topic(topic_name: str) -> TopicNode
29 system.ontology.get_prerequisites(topic: TopicNode) -> List[TopicNode]
30 system.ontology.list_topics() -> List[TopicNode]
31
32 # Research API
33 system.research.get_reports(
34     since: datetime = None,
35     limit: int = 10

```

```
36 ) -> List[ResearchReport]
37
38 system.research.apply_updates(
39     report_id: UUID,
40     validate: bool = True
41 ) -> None
42
43 # Metrics API
44 system.metrics.get_epistemic_metrics() -> EpistemicMetrics
45 system.metrics.get_performance_metrics() -> PerformanceMetrics
46
47 # Export API
48 system.export.to_markdown(
49     topic: str = None,
50     output_path: str
51 ) -> None
52
53 system.export.to_pdf(
54     topic: str = None,
55     output_path: str
56 ) -> None
```

References

1. LangGraph: <https://github.com/langchain-ai/langgraph>
2. Ollama: <https://github.com/ollama/ollama>
3. LLaMA 3: <https://ai.meta.com/llama/>
4. video-analyzer: <https://github.com/byj1w/video-analyzer>
5. llm CLI: <https://github.com/simonw/llm>
6. FAISS: <https://github.com/facebookresearch/faiss>
7. Qdrant: <https://qdrant.tech>
8. Deep Agents Paper: <https://huggingface.co/papers/2510.21618>