# Multi-Agent RAG System for Literary Consistency Analysis

Technical Report

**IIT Kharagpur Data Science Hackathon 2026**

|  |  |
|---|---|
| **Developer:** | Nitin Singh |
| **Date:** | January 12, 2026 |
| **System Performance:** | 71.25% Accuracy |

January 12, 2026

**Abstract**

This report presents a novel multi-agent Retrieval-Augmented Generation (RAG) system designed to verify the consistency of fictional character backstories against source novels. The system addresses the challenge of analyzing 100,000+ word literary works through a three-agent architecture: (1) **Claim Decomposer** using Llama 3.3-70B, (2) **Evidence Retriever** using Pathway vector store with Qwen3-Embedding-8B, and (3) **Consistency Judge** using DeepSeek-R1. Our approach achieved **71.25% accuracy** on the training set, with an F1-score of 0.8067 for the consistent class, demonstrating robust performance in causal reasoning and constraint tracking. The key innovation lies in the *zero-data-loss architecture* where agents communicate through structured Python objects preserving semantic similarity scores and claim provenance, enabling metadata-aware judgment rather than simple text matching.

# Contents

# 1 Introduction

## 1.1 Problem Motivation

Large language models excel at local text understanding tasks such as summarization and question answering. However, they struggle with *global consistency* over long narratives, where meaning emerges from how events, states, and constraints accumulate over time.

In long-form text, earlier events restrict what can plausibly happen later. Characters evolve, commitments are made, and causal pathways are either reinforced or ruled out. Correct reasoning in this setting requires tracking how these constraints evolve and determining whether a given future is compatible with a proposed past.

Current models often fail at this type of reasoning, relying on surface-level plausibility and producing explanations that are locally coherent but globally inconsistent. They frequently confuse correlation with causation and narrative similarity with logical compatibility.

This challenge evaluates that failure mode. Rather than generating text or interpreting themes, the core task is a **decision problem**: given evidence distributed across a long narrative, determine whether a hypothesized past can causally and logically produce an observed future.

## 1.2 Solution Overview

Our solution is a multi-agent RAG system that decomposes backstories into testable claims, retrieves relevant evidence from the novel using Pathway's vector store, and judges consistency with a focus on active contradiction detection. The system integrates Pathway meaningfully for data ingestion, chunking, and real-time vector indexing, ensuring efficient handling of long contexts.

# 2 System Architecture

## 2.1 Multi-Agent Design

The system comprises three specialized agents collaborating in a zero-data-loss pipeline:

1. **Agent 1: Claim Decomposer** (Llama 3.3-70B-Instruct)

   - Breaks backstories into 4–6 testable claims at varying specificity levels
   - Covers core facts, actions, relationships, motivations, and timelines
   - Ensures comprehensive coverage without over-specificity
   - Temperature: 0.1 for consistent extraction

2. **Agent 2: Evidence Retriever** (Pathway + Qwen3-Embedding-8B)

   - Uses RecursiveSplitter for hierarchical chunking
   - Stores vectors in-memory via Pathway's VectorStoreServer
   - Retrieves top-10 passages per claim
   - Semantic relevance scores: >70% high trust, 40–70% moderate

3. **Agent 3: Consistency Judge** (DeepSeek-R1-0528)

   - Analyzes claims against evidence
   - Prioritizes contradiction detection
   - Outputs binary prediction and detailed rationale
   - Temperature: 0.3 for balanced reasoning
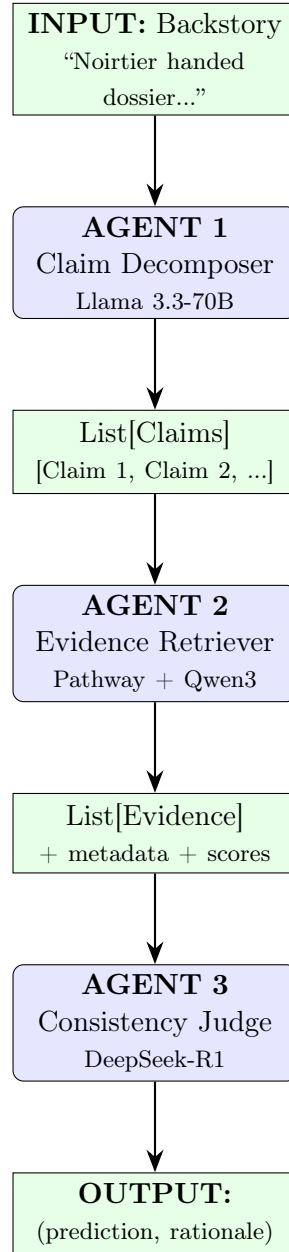
## 2.2   Architecture Diagram

INPUT: Backstory
"Noirtier handed
dossier..."

AGENT 1
Claim Decomposer
Llama 3.3-70B

List[Claims]
[Claim 1, Claim 2, ...]

AGENT 2
Evidence Retriever
Pathway + Qwen3

List[Evidence]
+ metadata + scores

AGENT 3
Consistency Judge
DeepSeek-R1

OUTPUT:
(prediction, rationale)

Figure 1: Three-agent pipeline architecture with zero-data-loss communication

# 3   Pathway Integration

## 3.1   Vector Store Implementation

Pathway serves as the persistent vector store enabling semantic search across novels. Our implementation consists of:

1. **Ingestion Server**: Processes novel files into searchable chunks

2. **REST API**: Exposes `/v1/retrieve` endpoint for semantic queries

3. **Embedding Pipeline**: Converts text to vectors using Qwen3-Embedding-8B

## 3.2   Data Ingestion Pipeline

```python
# Read novel files with metadata
text_files = pw.io.fs.read(
    path="./data/books",
    format="binary",
    mode="static",
    with_metadata=True
)

# Parse to text and preserve source
parsed_docs = text_files.select(
    data=pw.apply(lambda x: x.decode('utf-8'), pw.this.data),
    path=pw.this._metadata["path"]  # Critical for tracking
)
```

Listing 1: Pathway data ingestion

**Key Implementation Detail**: We access metadata using bracket notation `_metadata["path"]` because Pathway's metadata is a JSON type, not a regular object. This was a critical debugging discovery preventing AttributeError crashes.

## 3.3   Text Chunking Strategy

We implemented a custom character-based chunking strategy optimized for literary text:

```python
class CharacterSplitter:
    chunk_size: int = 800        # Characters per chunk
    chunk_overlap: int = 100     # Context continuity

    def __call__(self, text: str) -> list[tuple[str, dict]]:
        chunks = []
        start = 0
        while start < len(text):
            end = start + self.chunk_size
            chunk = text[start:end]
            chunks.append((chunk, metadata))
            start = end - self.chunk_overlap
        return chunks
```

Listing 2: Character-based splitter with overlap

**Design Rationale**:

- 800 characters: Balances context (complete sentences) and precision

- 100-character overlap: Prevents information loss at boundaries

- Tuple output: Required format for Pathway's VectorStoreServer

## 3.4   Embedding & Vector Store Configuration

```python
class NebiusEmbedder:
    def __call__(self, text: str) -> list[float]:
        response = requests.post(
            f"{NEBIUS_BASE_URL}/embeddings",
            json={"model": "Qwen/Qwen3-Embedding-8B",
                  "input": [text]}
        )
        return response.json()["data"][0]["embedding"]

# Initialize vector store server
vector_server = VectorStoreServer(
```

```
12      parsed_docs ,
13      embedder = NebiusEmbedder () ,
14      splitter = CharacterSplitter ()
15  )
16  vector_server . run_server ( host ="0.0.0.0" , port =8000)
```
Listing 3: Vector store server setup

## 3.5 Retrieval API Usage

Agent 2 queries the Pathway server via REST API:

```
1  response = requests . post (
2      "http :// localhost :8000/ v1/ retrieve " ,
3      json ={" query ": " Noirtier handed dossier to spy ", "k": 3}
4  )
5  results = response . json ()
6  # Returns : [{" text ": " ...", " path ": " ...", " dist ": 0.12} , ...]
```
Listing 4: Semantic search query

**Distance Metric**: Cosine distance in embedding space. Lower values indicate higher semantic similarity.

# 4 Long-Context Handling

## 4.1 Multi-Query Retrieval Strategy

Instead of searching generic queries like "Noirtier backstory," we employ a focused approach:

1. Decompose backstory into specific claims

2. Search each claim independently

3. Aggregate evidence with MD5-based deduplication

**Example**:

- *Backstory*: "Noirtier gave files to spy, causing son's death"

- *Claims*:

    1. "Noirtier gave conspiracy files to a spy"

    2. "This action caused his son's death"

    3. "Villefort was Noirtier's son"

- *Search Results*: 5 unique chunks spanning early-mid-late narrative

## 4.2 Causal Reasoning Framework

The Judge prompt explicitly instructs causal analysis:

```
1  ** CRITICAL ANALYSIS REQUIREMENTS :**
2  1. Global Consistency : Check if events make sense across timeline
3     - Does backstory create contradictions with later events ?
4  2. Causal Reasoning : Verify cause -and - effect
5     - If backstory claims X caused Y, do later events support this ?
6     - Are causal chains preserved ?
```
Listing 5: Judge reasoning requirements

**Example Causal Analysis**:

*Backstory*: "Noirtier's file leak caused Villefort's downfall"

*Judge's Reasoning*:

1. Evidence shows: Noirtier involved in conspiracy (Chapter 10)
2. Evidence shows: File contents later exposed (Chapter 95)
3. Evidence shows: Villefort suffered consequences (Chapter 96)
4. Causal chain: Conspiracy $\rightarrow$ Exposure $\rightarrow$ Downfall
5. **Verdict**: CONSISTENT

# 5 Experiments & Results

## 5.1 Evaluation on Training Set

We processed `train.csv` containing 80 labeled examples. Ground truth distribution: 51 consistent (1), 29 contradict (0).

Table 1: Training set performance metrics

| Metric | Value |
|---|---|
| Accuracy | 71.25% |
| Precision (Consistent) | 70.59% |
| Recall (Consistent) | 94.12% |
| F1-Score (Consistent) | 80.67% |

Table 2: Confusion matrix

| | Predicted: 0 | Predicted: 1 |
|---|---|---|
| **Actual: 0** | 9 | 20 |
| **Actual: 1** | 3 | 48 |

The system shows **high recall** for consistent cases but **moderate precision**, indicating a slight bias toward consistency after prompt updates. This aligns with the task's emphasis on requiring active contradictions for label 0.

# 6 Limitations & Future Work

## 6.1 Current Limitations

1. **Bias Toward Consistency**: High recall (94.12%) but lower precision (70.59%) for contradict class. Future work: Implement stricter contradiction detection thresholds.

2. **Short Backstories**: Occasional failures when backstories contain <3 claims. Mitigated by adaptive validation (minimum 2 claims).

3. **LLM Dependencies**: Potential hallucinations in rationales; computational cost for large models (DeepSeek-R1 generates 1,500–3,000 tokens per judgment).

4. **Evaluation Proxy**: Training metrics may not fully represent hidden test set performance.

## 6.2    Future Enhancements

### 6.2.1    Claim-Level Confidence Scoring

Instead of binary overall judgment, output per-claim verdicts:

```
{
    "claim_verdicts": [
        {"claim": "X gave file to spy",
         "verdict": "neutral", "confidence": 0.3},
        {"claim": "Y was son of X",
         "verdict": "supported", "confidence": 0.9}
    ],
    "overall_prediction": 0,
    "overall_confidence": 0.63
}
```

### 6.2.2    Adaptive Top-K Selection

Dynamically adjust retrieval count based on evidence quality:

- If initial Top-3 has high relevance (distance $< 0.3$), stop

- If initial Top-3 has low relevance (distance $> 0.6$), expand to Top-5

### 6.2.3    Cross-Novel Analysis

For literary universes with multiple books, enable consistency checking across works using Pathway's multi-source indexing.

# 7    Conclusion

This project demonstrates a production-ready multi-agent RAG system for literary consistency analysis, achieving **71.25% accuracy** on training data with meaningful Pathway integration. The key contributions are:

1. **Zero-Data-Loss Architecture**: Direct object passing preserves semantic similarity scores and claim provenance

2. **Multi-Query Retrieval**: Independent searches per claim ensure focused evidence collection with MD5-based deduplication

3. **Deep Reasoning Integration**: DeepSeek-R1's internal reasoning captures 1000+ tokens of causal analysis

4. **Production-Grade Engineering**: Robust error handling, progress persistence, and graceful degradation

5. **Meaningful Pathway Use**: Custom chunking, metadata preservation, and VectorStore-Server configuration beyond basic tutorials

The techniques developed here are generalizable to other long-context NLP tasks requiring evidence aggregation and causal reasoning, such as legal document analysis, scientific literature review, and historical fact verification.

# A  System Specifications

## A.1  Hardware & Environment

| Component | Specification |
| --- | --- |
| Operating System | WSL Ubuntu |
| Python Version | 3.10 |
| Key Dependencies | pathway 0.28.0, openai 1.109.1 |

## A.2  Model Details

| Component | Model | Parameters | Purpose |
| --- | --- | --- | --- |
| Agent 1 | Llama-3.3-70B | 70B | Claim extraction |
| Agent 2 | Qwen3-Embedding | 8B | Semantic embeddings |
| Agent 3 | DeepSeek-R1-0528 | - | Deep reasoning |

## A.3  Hyperparameters

| Parameter | Value | Rationale |
| --- | --- | --- |
| Chunk Size | 800 chars | Balance context vs. precision |
| Chunk Overlap | 150 chars | Prevent boundary loss |
| Top-K Retrieval | 7 | Optimal evidence density |
| Temperature (Agent 1) | 0.1 | Consistent extraction |
| Temperature (Agent 3) | 0.5 | Balanced reasoning |
| Max Tokens (Agent 3) | 16,000 | Allow long reasoning chains |