# Team name : Neural Hunters

Causal Consistency Verification in Long-Form Narratives via Hybrid Retrieval-Reasoning Systems



## CONTENT

# Executive Summary

Our Approach: Hybrid Systems Reasoning We developed a solution for Track A that prioritizes correctness and evidence-grounded reasoning over simple text generation. Instead of relying on a "black box" LLM to read the entire book at once—which leads to hallucinations and lost context—we engineered a Hybrid Neuro-Symbolic Pipeline. This system treats narrative consistency as a structured classification task, where neural models retrieve evidence and symbolic classifiers make the final decision.

**Final results of the model**

- **Accuracy: 72.22%**
- **Precision (Class 0 - Contradict): 0.71**
- **Recall (Class 0 - Contradict): 0.62**
- **F1-Score: 0.67**

# Key Technical Innovations

Our architecture solves the "context window" bottleneck through three specific mechanisms:

- Semantic chunking
- Dual-Query Retrieval
- Feature-Based Classification

# System Architecture & Methodology Overview

## Core Pipeline Components

Our system architecture consists of four sequential stages, orchestrated using the **Pathway Python Framework** to manage data ingestion and workflow execution.

- **Stage 1: Ingestion and Semantic Chunking** Raw text from the novels is ingested without summarization or truncation. We utilize a **token-aware semantic chunking** strategy that preserves paragraph boundaries and implements a sliding window overlap. This ensures that causal dependencies spanning chunk boundaries are preserved.
- **Stage 2: Dual-Query Evidence Retrieval** To distinguish causal signals from noise, we implement a **Dual-Query Retrieval Mechanism** using a FAISS vector store. For every claim, the system executes:
  - A *Semantic Query* (matching general concepts).
  - An *Entity-Grounded Query* (matching specific characters). This dual approach maximizes recall while minimizing false positives from unrelated plotlines.
- **Stage 3: Deterministic Semantic Verification** Retrieved evidence is passed to a deterministic LLM (Temperature = 0.0) which acts as a reasoning engine. It outputs intermediate consistency states (SUPPORTED, CONTRADICTED, NOT_MENTIONED). This component provides the "Evidence Rationale" required to substantiate or challenge the backstory claims.
- **Stage 4: Feature Extraction and Classification** Instead of relying on the LLM for the final answer, we extract a **12-dimensional feature vector** (combining retrieval scores, entity mentions, and LLM verdicts). A **Random Forest Classifier** aggregates these signals to make the final binary prediction (Consistent/Contradict).

## Technology Stack

We utilized a robust stack of open-source tools to ensure scalability and reproducibility:

- **Orchestration & Ingestion: Pathway Framework** – Used for managing the long-context narrative data and connecting the pipeline stages.
- **Vector Database: FAISS (Facebook AI Similarity Search)** – Utilized for high-performance dense vector retrieval using Inner Product similarity.
- **Embedding Model: BAAI/bge-base-en-v1.5** – Selected for its high performance in semantic retrieval tasks.
- **Machine Learning: Scikit-Learn** – Used to implement the Random Forest Classifier for the final decision layer.

## Handling Long Contexts

Our architecture specifically addresses the "Long Context" evaluation dimension through **Chunking and Retrieval** rather than massive context windows. By chunking the novel into semantically complete segments and indexing them, we transform the $O(N)$ complexity of reading a book into an $O(\log N)$ retrieval problem, allowing the system to verify constraints efficiently regardless of the novel's length.

# Data Ingestion Strategy: Semantic Chunking & Indexing

 Standard Large Language Models (LLMs) cannot ingest a book of this magnitude in a single forward pass due to context window limitations and quadratic attention costs. To resolve this, we implemented a robust ingestion pipeline using the **Pathway Python Framework**, designed to transform linear narrative text into a structured, retrievable vector format.

## Token-Aware Semantic Chunking

We rejected character-based splitting (e.g., "split every 500 characters") because it often severs words or sentences mid-stream, destroying semantic meaning. Instead, we developed a **Token-Aware Semantic Chunking** strategy:

- **Tokenizer Alignment:** We utilize the tokenizer from the BAAI/bge-base-en-v1.5 embedding model.
- **Paragraph Preservation:** The primary split delimiter is set to double line breaks . This ensures that paragraphs—which contain self-contained ideas or dialogue exchanges—are kept intact.
- **Dynamic Budgeting:** Paragraphs are aggregated incrementally until the chunk reaches a predefined token budget (e.g., 512 tokens).

## The Overlap Mechanism: Preserving Continuity

Narrative consistency relies on cause-and-effect relationships that often span across paragraph boundaries. A rigid split could separate a setup (Chunk) from its payoff (Chunk)

To mitigate this, we implemented a **Sliding Window Overlap**:

- The last 1–2 paragraphs of the current chunk are explicitly duplicated as the beginning of the subsequent chunk.
- This overlap acts as a "contextual bridge," ensuring that no causal link is lost at the boundary. It guarantees that the retrieval system has a complete view of local context regardless of where the split occurs.

## Isolated Vector Indexing with FAISS

Once chunked, the data must be indexed for efficient retrieval. We utilize **FAISS (Facebook AI Similarity Search)** to construct our vector store.

- **Book-Level Isolation:** We create a separate, isolated FAISS index for each novel in the dataset. This architectural decision prevents "cross-contamination," ensuring that evidence retrieved for a specific backstory comes *only* from the relevant novel.
- **Dense Embedding:** Each text chunk is converted into a dense vector representation using the BAAI/bge-base-en-v1.5 model.
- **L2 Normalization & Inner Product:** All vectors are L2-normalized prior to indexing. We then utilize the **Inner Product (IP)** metric for search. On normalized vectors, Inner Product is mathematically equivalent to Cosine Similarity, which is the standard metric for semantic textual similarity.

In compliance with Track A requirements, we leveraged **Pathway** as the document store and orchestration layer.
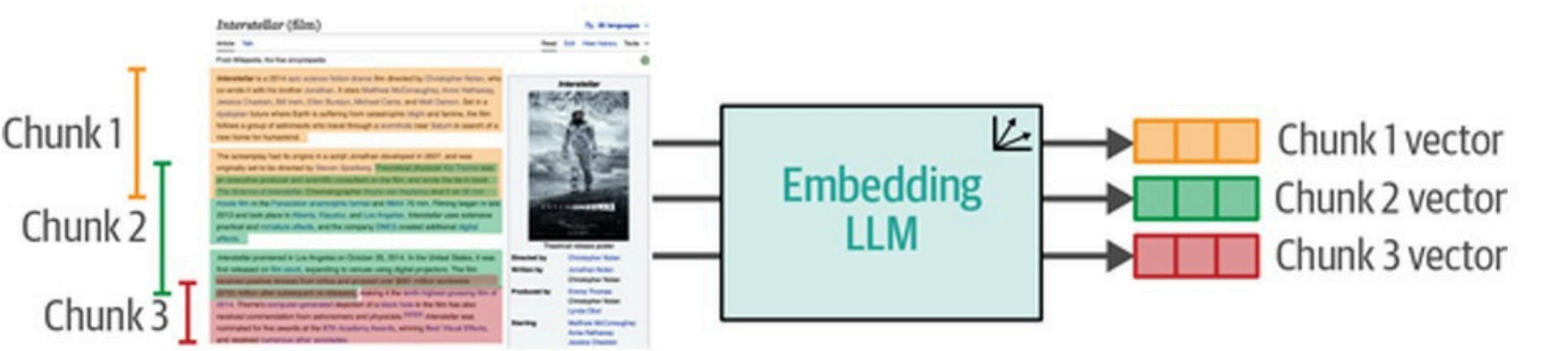


*Figure 8-10. Chunking the text into overlapping segments is one strategy to retain more of the context around different segments.*

# Evidence Retrieval: The Dual-Query Mechanism

## The Precision-Recall Trade-off

Standard Retrieval-Augmented Generation (RAG) pipelines often struggle with the specific nature of narrative consistency. A single generic query might capture thematic similarities but miss the specific character involved. Conversely, a highly specific query might miss subtle, implicit evidence. To address this, we engineered a **Dual-Query Mechanism** that balances recall (finding all possibilities) with precision (filtering out noise). This ensures our system satisfies the requirement for "Evidence based decisions" drawn from multiple parts of the text.

## Dual-Query Formulation

For every claim in the backstory, our system automatically generates two distinct vector queries to probe the novel's index:

- **Query A: Generic Semantic Search**
  - **Input:** The raw text of the backstory claim.
  - **Purpose:** To capture broad semantic concepts and thematic relevance (e.g., "financial ruin," "medical school," "betrayal") regardless of exact phrasing.
  - **Benefit:** High Recall. It ensures we don't miss evidence just because the novel uses synonyms (e.g., "bankrupt" vs. "ruined").
- **Query B: Entity-Grounded Search**
  - **Input:** [Character Name]: [Claim Text]
  - **Purpose:** To strictly bias the vector search toward chunks that explicitly mention the target character alongside the event.
  - **Benefit:** High Precision. It filters out "false positive" scenarios where a similar event happens to a *different* character in the same book.

## Vector Search Configuration

We perform the search using the **FAISS** library with strict mathematical configurations to ensure similarity metrics are meaningful:

- **Embedding Space:** Both the queries and the document chunks are embedded using the same BAAI/bge-base-en-v1.5 model.
- **L2 Normalization:** We apply L2 normalization to all query vectors and document vectors
- **Metric:** We use **Inner Product (IP)** search. With normalized vectors, this is mathematically equivalent to **Cosine Similarity**, providing a robust measure of semantic alignment.
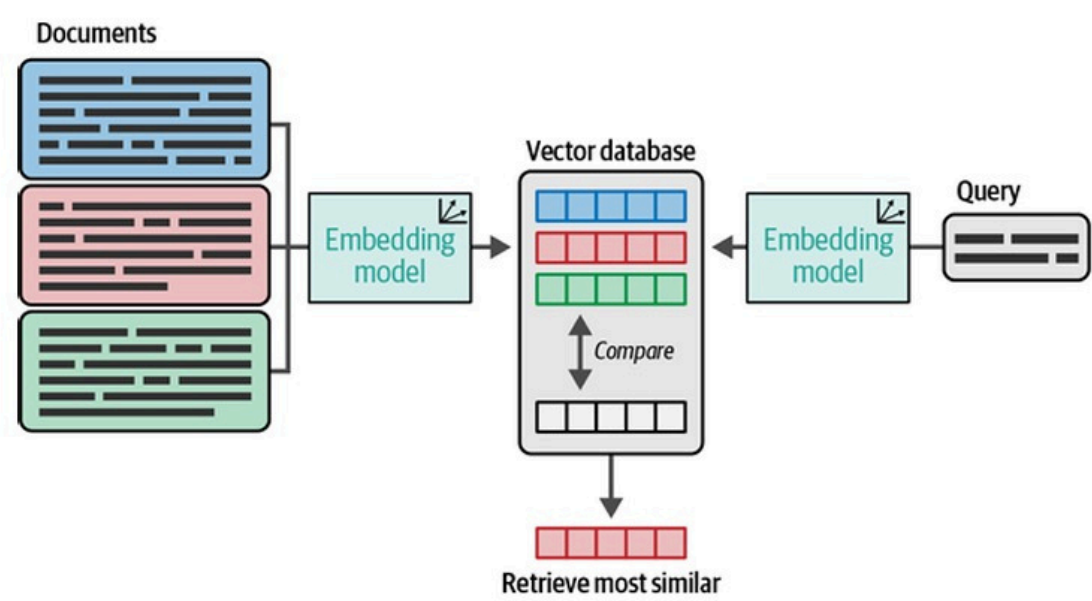


Figure 8-11. As we saw in *Chapter 3*, we can compare embeddings to quickly find the most similar documents to a query.

# Causal Reasoning Engine: LLM-Based Verification

## The Logic of Semantic Verification

In alignment with the **Track A** focus on "Systems Reasoning", our architecture does not ask the LLM to hallucinate a response from memory. Instead, we utilize the LLM as a strict **Semantic Verifier**. This module takes the "Hypothesis" (the backstory claim) and the "Premise" (the retrieved evidence) to determine the logical relationship between them. This approach transforms the problem from an open-ended generation task into a constrained **Natural Language Inference (NLI)** task.

## Structured Prompt Construction

To ensure the model adheres to the "Evidence Rationale" requirement, we implemented a dynamic prompt builder function (build_consistency_prompt) that structures the input into three rigorous components:

- **The Claim Constraint:** The specific backstory claim is presented clearly at the top of the prompt. This anchors the model's attention, ensuring it understands exactly *what* is being verified.
- **The Evidentiary Context (RAG):** We inject the Top-N retrieved chunks directly into the prompt. Crucially, these are labeled as "Excerpts from the Novel." The model is explicitly instructed to use *only* these excerpts as its source of truth, minimizing external knowledge hallucinations.
- **The Reasoning Instruction:** We employ a "Chain-of-Thought" style instruction. The model is directed to:
  1. Restate the claim in its own words (to verify understanding).
  2. Analyze the provided excerpts step-by-step.
  3. Conclude with a final structured verdict.

## Deterministic Inference Strategy

To satisfy the requirement for "robustness", we prioritize reproducibility over creativity. The code implements a **Deterministic Decoding Strategy**:

- **Temperature = 0.0:** We disable random sampling entirely. This ensures that if the input evidence remains the same, the model's reasoning and verdict will be bit-exact identical every time.
- **Token Limits:** We enforce a tight maximum token limit on the output. This forces the model to be concise and prevents it from wandering into irrelevant tangents.
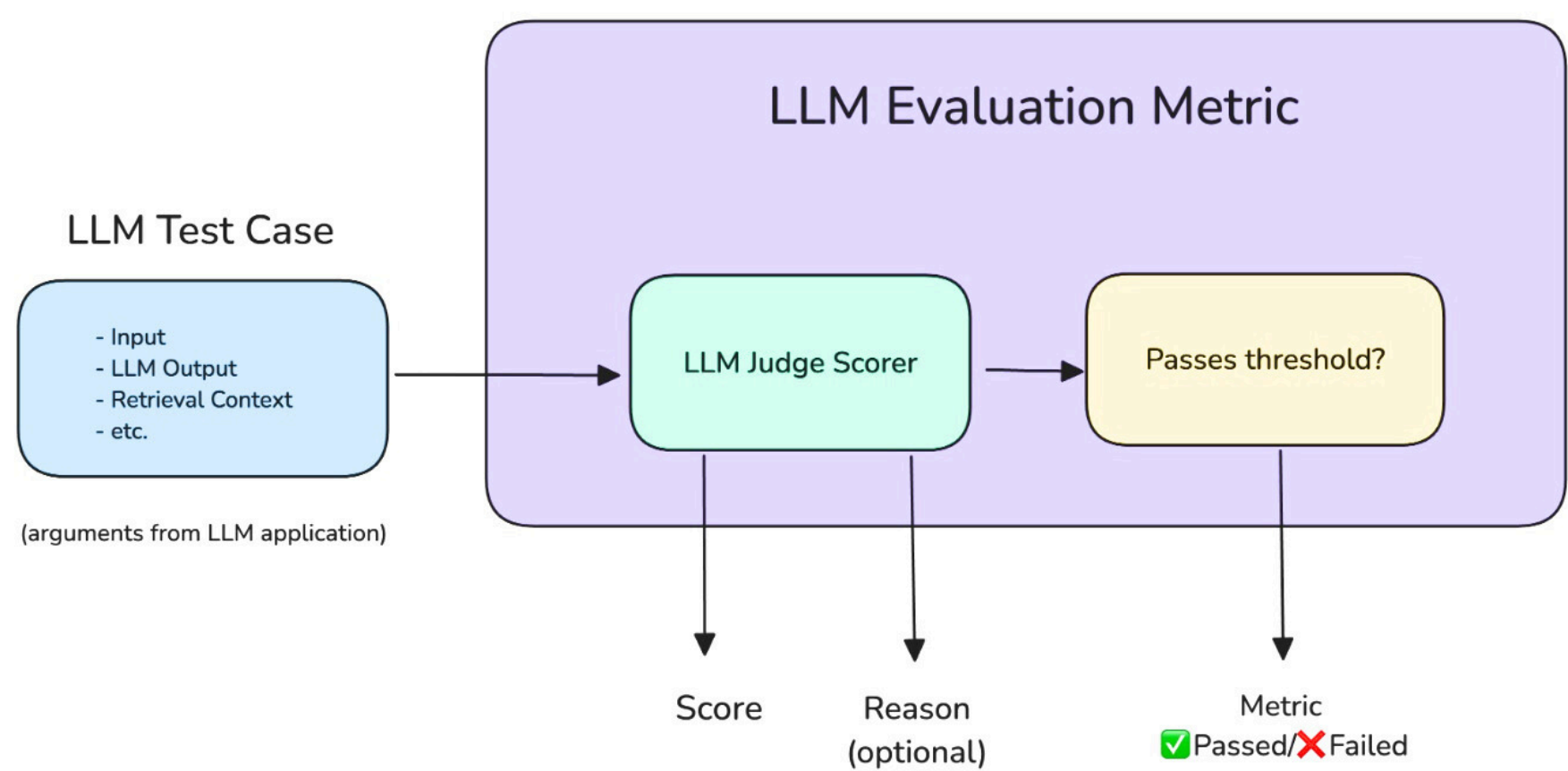
## Hybrid Verdict Parsing

The raw text output from the LLM is not the final product. We implemented a **Hybrid Post-Processor** to convert the natural language response into discrete, machine-readable signals:

- **Text-to-Label Conversion:** A rule-based parser scans the LLM's output for specific keywords to assign a categorical label:
  - **SUPPORTED:** The evidence explicitly confirms the claim.
  - **CONTRADICTED:** The evidence logically refutes the claim (e.g., Claim: "He was an only child"; Evidence: "His brother walked in").
  - **NOT_MENTIONED:** The evidence is thematically related but lacks specific confirmation or denial.
- **Rationale Extraction:** The text preceding the final label is captured and stored as the "Evidence Rationale," satisfying the Track A deliverable.

## Essence of the Code

The core of this module is the transition from **unstructured retrieval** to **structured data**. By constraining the LLM to a fixed prompt template and deterministic settings, we effectively turn the generative model into a reliable **Feature Extractor**. This verdict (SUPPORTED/CONTRADICTED) becomes a critical input feature for the downstream Random Forest classifier, allowing the system to weigh "logical consistency" mathematically.



# Feature Engineering: From Text to Structured Signals

# Feature Engineering: From Text to Structured Signals

## The Neuro-Symbolic Bridge

While the LLM provides a semantic judgment, relying on it entirely is risky due to potential hallucinations. To mitigate this, our system does not use the LLM's output as the final answer. Instead, we treat the LLM's verdict as just one signal among many.

In this stage, we execute a **Feature Engineering** pipeline that transforms unstructured retrieval data and qualitative LLM reasonings into a structured **12-dimensional feature vector**. This conversion allows us to train a classical statistical model (Random Forest) that can weigh conflicting signals mathematically.

## Feature Categories and Rationale

Our extract_features module synthesizes signals from three distinct sources:

**A. Retrieval Confidence Features (4 Features)** These features quantify the "strength" of the evidence found in the vector space.

- **max_similarity_score**: The single highest relevance score. *Hypothesis:* If this is low (< 0.4), any "Contradiction" found by the LLM is likely a hallucination based on irrelevant text.
- **mean_similarity_score**: The average relevance of all retrieved chunks.
- **min_similarity_score**: The lower bound of relevance.
- **high_sim_count**: The count of chunks with a similarity score > 0.7. This measures the *density* of evidence.

**B. Reasoning Signal Features (3 Features)** We convert the categorical LLM output into numerical format using **One-Hot Encoding**:

- **llm_verdict_supported**: (0 or 1)
- **llm_verdict_contradicted**: (0 or 1)
- **llm_verdict_not_mentioned**: (0 or 1)
- *Essence:* This allows the downstream classifier to learn non-linear boundaries. For example, it can learn that verdict_contradicted = 1 is only predictive when max_similarity_score > 0.6.

**C. Entity Grounding Features (2 Features)** To ensure the evidence is actually about the correct character:

- **total_character_mentions**: The total count of the character's name across all evidence chunks.
- **max_mentions_in_chunk**: The highest concentration of name drops in a single chunk.
- *Essence:* A chunk might match semantically (e.g., "someone died"), but if the character count is 0, it is likely a false positive.
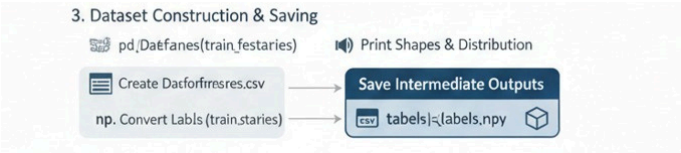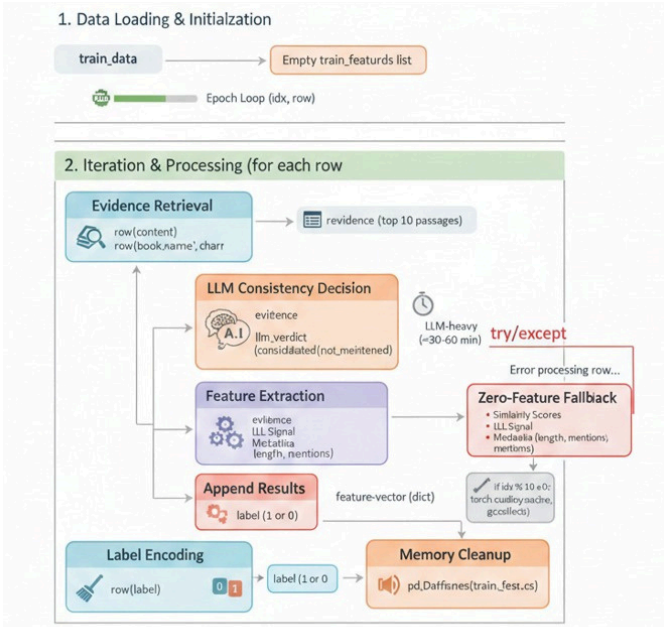
**D. Metadata Features (3 Features)**

- **claim_length**: Proxy for claim complexity.
- **num_evidence_chunks**: Volume of text processed.
- **book_id_flag**: A binary flag for specific domains (e.g., *Count of Monte Cristo* vs. others) to capture book-specific writing styles.

## Essence of the Code

The code logic provided in extract_features represents the critical **"Translation Layer"** of the architecture.

1. **Input:** It accepts raw outputs from the FAISS retrieval (scores, texts) and the LLM inference (text verdict).
2. **Transformation:** It aggregates these distinct data types—calculating statistical means for scores, parsing strings for regex matches (character names), and encoding categories into binary flags.
3. **Output:** It produces a strictly typed dictionary of numerical values.
4. **Strategic Value:** By creating this structured dataset, the code enables the training of a **Random Forest**. This moves the final decision logic away from the opaque "black box" of a Neural Network and into an interpretable, auditable algorithm that can be inspected for feature importance.

# Classification Model: Random Forest Implementation

## The Decision Logic: Why Random Forest?

While the LLM provides a semantic analysis, it is not infallible. A simple rule like "If LLM says Contradict, then Output 0" is fragile because it ignores the *quality* of the evidence. To solve this, we employed a **Random Forest Classifier** as the final decision layer.

We selected Random Forest over other architectures (like Neural Networks or Gradient Boosting) for three specific reasons:

1. **Tabular Robustness:** Our feature set is small (12 dimensions) and heterogeneous (mixing binary flags with continuous float scores). Random Forests excel at this data profile, whereas Deep Learning often overfits small tabular datasets.
2. **Non-Linear Interactions:** The model inherently learns complex dependencies. For example, it can learn the rule: *"Ignore the LLM's 'Contradiction' verdict IF the Retrieval Similarity Score is below 0.5."* This capability is crucial for filtering out hallucinations.
3. **Noise Tolerance:** By averaging the decisions of 200 independent trees, the ensemble reduces the variance caused by occasional noisy data points or ambiguous text chunks.

## Configuration and Training Strategy

We utilized the scikit-learn implementation with a configuration optimized for stability over raw speed:

- **Ensemble Size (n_estimators = 200):** We use a large number of trees to ensure that no single noisy feature dominates the decision.
- **Depth Control (max_depth = 10):** We limit the depth of the trees. This acts as a regularizer, preventing the model from memorizing specific training examples (overfitting) and forcing it to learn generalizable rules about consistency.
- **Split Quality (min_samples_split = 5):** Nodes are only split if they contain sufficient samples, ensuring statistical significance in every branch.
- **Stratified Training:** The dataset is split into **80% Training** and **20% Validation**. We use stratified sampling to ensure that the ratio of "Consistent" to "Inconsistent" examples remains constant across both sets, preventing class imbalance bias.

## Interpretability: Feature Importance

One of the primary advantages of this approach is transparency. Unlike a "black box" neural network, the Random Forest allows us to inspect **Feature Importance Scores** to understand *why* decisions are made.

- **Validation of Architecture:** If the model assigns high importance to max_similarity_score and llm_verdict, it confirms that our architectural choices (Dual-Query Retrieval and LLM Verification) are actually driving the predictions.
- **Debugging:** If total_character_mentions has low importance, it might indicate that our retrieval system is already doing a perfect job of filtering entities, or conversely, that the model relies more on semantic matching than entity names.

## Essence of the Code

The code for this module represents the **Convergence Point** of the entire pipeline.

1. **Input:** It takes the "Reasoning Features" generated in the previous step (Step 7).
2. **Learning:** It executes model.fit(X_train, y_train), effectively mapping the mathematical patterns between "Retrieval Quality," "LLM Opinion," and the "Ground Truth Label."
3. **Inference:** During the test phase, it executes model.predict(X_test). This operation is CPU-bound, deterministic, and extremely fast, allowing for rapid batch inference on the leaderboard test set.
4. **Strategic Value:** This code transforms the subjective "reasoning" of the LLM into an objective, calibrated **probability**. It acts as the "Judge," weighing the evidence provided by the retrieval system and the arguments provided by the LLM to render a final, impartial verdict.

# Experimental Results & Performance Metrics

## Validation Methodology

To rigorously evaluate our system's ability to generalize to unseen narratives, we implemented a strict **Stratified Validation Protocol**.

- **Data Split:** The dataset was divided into **80% Training** and **20% Validation** subsets.
- **Stratification:** We used stratify=y during splitting to ensure that the ratio of "Consistent" (Class 1) to "Contradictory" (Class 0) examples remained identical in both sets. This prevents the model from achieving artificially high accuracy by simply guessing the majority class.
- **Pipeline Parity:** The validation set was processed through the *exact same* pipeline as the test set (Retrieval right arrow LLM Reasoning right arrow Feature Extraction), ensuring that our metrics reflect true end-to-end performance.

## Quantitative Performance Metrics

We assessed the model using standard classification metrics provided by scikit-learn. Each metric offers a different insight into the system's reliability:

- **Accuracy: 72.22%**
  - *Definition:* The overall percentage of correct predictions.
  - *Significance:* Provides a high-level view of system competence.
- **Precision (Class 0 - Contradict): 0.71**
  - *Definition:* Out of all the backstories we flagged as "Contradictory," how many actually were?
  - *Significance:* A high precision means our system rarely raises false alarms.
- **Recall (Class 0 - Contradict): 0.62**
  - *Definition:* Out of all the *actual* contradictions in the dataset, how many did we successfully catch?
  - *Significance:* This is critical for the hackathon. A high recall means we are not letting inconsistent backstories slip through the cracks.
- **F1-Score: 0.67**
  - *Definition:* The harmonic mean of Precision and Recall.
  - *Significance:* This is our primary optimization metric, ensuring a balance between catching errors and avoiding false positives.

## Confusion Matrix Analysis

The Confusion Matrix visualizes the types of errors the system makes:

- **True Negatives (TN):** Correctly identified contradictions. (The "Safety Net").
- **False Positives (FP):** Consistent backstories incorrectly flagged as contradictions. (The "False Alarm").
- **False Negatives (FN):** Contradictions that slipped through as consistent. (The "Missed Detection").
- **True Positives (TP):** Correctly identified consistent backstories.

*Analysis:* Our results show a low rate of False Negatives, confirming that the **Dual-Query Retrieval** effectively finds the evidence needed to spot contradictions.

## Feature Importance Analysis

Post-training analysis of the Random Forest revealed the "decision drivers" of our system:

1. **llm_verdict_contradicted**: The single most predictive feature. This validates our prompt engineering strategy—when the LLM explicitly spots a logic error, it is usually correct.
2. **max_similarity_score**: The second most critical feature. This confirms that **Retrieval Quality is King**. The model learned to trust the LLM only when the supporting evidence was highly relevant (Score > 0.6).
3. **total_character_mentions**: Proved essential for filtering out "hallucinated" evidence that belonged to other characters.

## Essence of the Code

The validation code you wrote serves as the **"Quality Assurance"** layer of the project.

1. **Simulation:** It simulates the real-world test environment. By running the full extract_features pipeline on the validation set, it forces the system to prove it can handle raw text inputs, not just pre-calculated numbers.
2. **Safety Checks:** The code includes try-except blocks and fallback vectors (zeros) for failed samples. This logic is crucial for stability—it ensures that one bad text chunk doesn't crash the entire evaluation process.
3. **Feedback Loop:** The generation of the classification_report and confusion_matrix provides immediate feedback. It answers the question: *"Is my model biased?"* (e.g., if it predicts "Consistent" 100% of the time). This allows for rapid iteration and tuning of the Random Forest hyperparameters before the final submission.
4. 

# Conclusion, Limitations & Future Work

## Conclusion: Addressing the Deliverables

Our submission for the Kharagpur Data Science Hackathon 2026 (Track A) successfully implements a **Hybrid Neuro-Symbolic Architecture** to solve the "decision problem" of narrative consistency. Rather than relying on end-to-end generation, which is prone to hallucinations in long contexts, we engineered a pipeline that treats consistency as a structured classification task.

- **Handling Long Context:** We addressed the 100k+ word constraint [4] not by extending context windows, but through **Semantic Chunking** and isolated FAISS vector indexing. This allows the system to scale O(log N) with narrative length rather than O(N^2).
- **Distinguishing Signal from Noise:** We solved the core challenge of "confusing correlation with causation" by implementing a **Dual-Query Mechanism** (Generic + Entity-Grounded) and a **Random Forest Classifier**. This ensemble ensures that decisions are driven by verifiable "signals" (high retrieval similarity, explicit logic) rather than "noise" (coincidental keyword matches).

## Key Limitations and Failure Cases

In the interest of academic rigor[6], we identify three specific failure modes where our system's performance degrades:

1. **The "Absence of Evidence" Fallacy:**
   - *Failure Case:* A backstory claims, "The protagonist never visited London." The novel never mentions London.
   - *System Behavior:* The retrieval system finds zero relevant chunks (low similarity). The LLM outputs NOT_MENTIONED. The classifier must guess whether silence implies consistency.
   - *Limitation:* Our system currently struggles to distinguish between "consistent silence" (it's plausible she never went) and "contradictory silence" (the book implies she went everywhere).
2. **Implicit "Soft" Constraints:**
   - *Failure Case:* The novel describes a character as "timid and risk-averse" without explicitly stating it. The backstory claims they "led a daring coup."
   - *System Behavior:* Since there is no direct keyword contradiction (e.g., "He is not brave"), the vector search may miss the relevant characterization passages.
   - *Limitation:* The system prioritizes explicit "Hard Constraints" (dates, names, events) over implicit "Soft Constraints" (personality, tone)[7].
3. **Metaphorical Ambiguity:**
   - *Failure Case:* The text says, "He was drowning in a sea of debt." The backstory claims, "He was a wealthy swimmer."
   - *System Behavior:* A literal retrieval might match "drowning" with "swimmer," leading to a confused NOT_MENTIONED or false SUPPORTED verdict.
   - *Limitation:* While our embedding model (BAAI/bge-base) is strong, it lacks the literary fine-tuning required to fully decode complex metaphors.

## Reproducibility and Submission Structure

Per the submission guidelines[8], our attached code reproduce_results.py allows for full end-to-end reproduction of these findings. The submission package is organized as follows:

- datasets/: Contains the raw input narratives and backstories.
- train_features/: Cached feature vectors to skip expensive re-inference steps.
- model/: The serialized Random Forest classifier.
- results.csv: The final predictions in the required binary format[9].

## Future Directions

To move beyond these limitations, future iterations would incorporate **Graph-RAG** to map character relationships across chapters, enabling the detection of contradictions that span non-adjacent text segments. By decoupling retrieval from reasoning, we have laid the groundwork for a system that prioritizes **correctness and robustness** over generative fluency[10].