# Project Report

Project Title: Multi-Hop QA with Entailment Verification

Team Name: NeuraNova

Team Members: Salma S, Kokila M

---

# 1. Introduction

## a. Project Understanding

In the domain of Natural Language Processing (NLP), standard Retrieval-Augmented Generation (RAG) systems often struggle with "Multi-Hop" questions—queries that require aggregating information from multiple, disparate documents to derive an answer. A standard retriever might find the first piece of evidence but miss the second, leading to incomplete answers or "hallucinations" by the Large Language Model (LLM).

**Our objective** was to build a robust pipeline that solves this by:

1. **Iterative Retrieval:** Retrieving evidence in "hops" to follow the trail of clues.
2. **Chain Construction:** Linking these documents into logical chains.
3. **Entailment Verification:** Using a dedicated NLI (Natural Language Inference) model to mathematically verify if the retrieved evidence actually supports the generated answer.

## b. Terminologies

- **Multi-Hop QA:** A generic question-answering task requiring reasoning across multiple documents (e.g., Document A mentions Person X, Document B mentions Person X's father; the question asks for the father's profession).
- **Dense Retrieval:** Using vector embeddings (like BERT or MPNet) to find semantically similar documents, rather than keyword matching.
- **Chain-of-Evidence:** A specific sequence of documents that, when read together, provides the full context to answer a query.
- **Entailment Verification:** A process where a model determines if a "Hypothesis" (the answer) is logically true given a "Premise" (the evidence).
- **Cross-Encoder:** A transformer architecture that processes two inputs simultaneously to output a high-fidelity similarity or entailment score.

## c. Literature Survey

The challenge of Multi-Hop QA has evolved through several phases:

- **Traditional Methods (TF-IDF/BM25):** Early systems relied on keyword matching. These failed when the query and the second-hop document shared no common keywords (the "lexical gap").

- **Graph-Based Approaches (e.g., DFGN):** Research like *Dynamic Fusion Graph Networks* attempted to build entity graphs connecting documents. While accurate, these are computationally expensive and difficult to scale.
- **Iterative Retrieval (e.g., IRCoT):** Recent approaches like *Iterative Retrieval Chain-of-Thought* interleave retrieval and reasoning steps.
- **Self-Reflective RAG (Self-RAG, CRAG):** The current academic focus is on systems that critique their own retrieval. Papers on *Corrective RAG* suggest using a separate evaluator to trigger re-retrieval if the context is ambiguous.

## d. Current State of the Art

Currently, the leaderboard for datasets like HotpotQA is dominated by massive LLM Agents (such as GPT-4-based ReAct agents) that can handle long context windows. However, these are costly and slow. **Our approach** aligns with the emerging trend of "Small Language Models (SLMs) with Verifiers"—using a smaller, efficient generator (TinyLlama) backed by a rigorous verifier (DeBERTa), offering a balance of accuracy and computational efficiency.

## e. Key Novelty & Innovations

Our project distinguishes itself from standard RAG systems through five key innovations:

1. **Moving Reasoning Out of the "Black Box":** Most current systems simply feed documents into an LLM and hope it figures out the connection. We took a different approach by forcing the system to explicitly construct "evidence chains" in code. This means the reasoning path (linking Document A to Document B) is deliberate and visible, rather than being hidden inextricably inside the model's neural weights.
2. **A Built-in "Fact Checker":** We implemented a "Trust but Verify" architecture. After the LLM generates an answer, a completely separate model (the Cross-Encoder) steps in to audit it. It mathematically calculates if the retrieved evidence *actually* proves the answer. This extra layer serves as a safety net, filtering out hallucinations before they reach the user.
3. **Total System Transparency:** We moved away from the "magic box" approach. Our system exposes every step of the process. The dashboard displays the specific search queries used for each hop, the exact documents linked together, and the confidence scores for every claim. It is the equivalent of a student "showing their work" on a math test, which builds user trust.
4. **Designed for Experimentation:** We built this framework specifically for research analysis. Unlike rigid commercial tools, our system allows us to easily toggle critical variables—like the depth of the search (hops), the length of the evidence chains, or the strictness of the verification. This makes it an ideal testbed for understanding *why* multi-hop retrieval succeeds or fails.
5. **Flexible, Modular Architecture:** We avoided building a monolithic system where everything is tangled together. Instead, we cleanly separated the **Retriever** (finding data), the **Generator** (writing answers), and the **Verifier** (checking logic). This modularity means we can upgrade individual components—like swapping in a more powerful LLM or a faster embedding model—without having to rebuild the entire pipeline.

# 2. Your Method

We developed a comprehensive "Retrieve-Generate-Verify" pipeline, accessible via a FastAPI backend and a Streamlit Frontend.

## The Pipeline Algorithms

1. **Multi-Hop Dense Retrieval (Iterative Expansion):**
   - We use sentence-transformers/all-mpnet-base-v2 for embeddings.
   - **Hop 1:** The system retrieves the top-k docs for the user query.
   - **Hop 2:** The system extracts *titles* from Hop 1 documents and treats them as new queries. This allows the system to find documents related to the entities found in the first step, bridging the knowledge gap.
2. **Generative Answering:**
   - We utilise TinyLlama/TinyLlama-1.1B-Chat-v1.0. The retrieved contexts are concatenated and fed into the LLM with a specific instruction to answer concisely.
3. **Combinatorial Chain Verification (The Core Innovation):**
   - Instead of trusting the LLM blindly, we generate combinations of retrieved documents (pairs or triplets) using itertools.combinations.
   - **Cross-Encoder:** We use cross-encoder/nli-deberta-v3-base.
   - **Scoring:** For every chain, we calculate P(Entailment | Evidence, Answer).
   - **Filtering:** Only chains with a score > 0.7 are presented to the user as "Supporting Evidence."

**Challenges Faced and How We Tackled Them**

During the development of this pipeline, we encountered several technical hurdles. Here is how we overcame them:

- **The "Memory Wall" (Context Length Constraints):**
  One of the first hurdles we hit was the strict context limit of our lightweight models. When we retrieved multiple documents, the combined text was simply too long for the LLM to process, resulting in truncation errors.
  - *Mitigation:* We implemented strict character-level budgeting. We prioritised the top-ranked documents first and truncated less relevant text, ensuring the model received the "meat" of the evidence without choking on the volume.
- **The Speed vs. Accuracy Trade-off (Computational Overhead):**
  Our "fact-checker" (the Cross-Encoder) is highly accurate but computationally heavy. Initially, checking every single combination of documents caused the system to lag significantly.
  - *Mitigation:* We adopted a smart filtering strategy. Instead of checking every possible chain, we only ran the expensive verification on the top-k most promising candidates. This kept the system snappy without sacrificing accuracy.
- **Drifting Off Topic (Noise Accumulation):**
  Going deeper into search results (Multi-Hop) is a double-edged sword. Sometimes, the second hop brought in completely irrelevant junk that confused the AI rather than helping it.

○ *Mitigation:* We added a re-ranking step. After fetching documents from all hops, we re-scored them against the original query to filter out the noise, keeping only the documents that actually made sense.
- **Hardware Constraints (Memory Exhaustion):**
  We did not have access to massive enterprise servers. Trying to run large models caused our local machines to crash instantly due to Out-Of-Memory (OOM) errors.
    ○ *Mitigation:* We switched to efficient "Tiny" models and pre-computed the heavy math (embeddings) offline. By loading these pre-calculated vectors, we could run a complex AI pipeline on standard hardware without crashing.
- **Grading Discrepancies (Evidence Mismatch):**
  Scoring the system was tricky. Sometimes the AI found the right document, but the title was slightly different from the official answer key (e.g., "The FTX Trial" vs "FTX Trial"), causing the system to mark it as "wrong."
    ○ *Mitigation:* We added text normalisation logic—stripping away punctuation and ignoring case differences—to make the grading fairer and more representative of the model's actual performance.

## User Interface (Streamlit)

To make the system explainable, we developed a UI using **Streamlit**. The UI communicates with the backend via REST API. It features:

- **Control Panel:** Sliders to adjust top_k, chain_length, and entailment_threshold in real-time.
- **Hop Visualisation:** Displays exactly which documents were found at Hop 1 vs. Hop 2.
- **Verification Badge:** Visually indicates if an answer is supported by high-confidence logic.

---

# 3. Experiments

## Dataset
- **Source: HotpotQA** (Fullwiki Setting).
- **Characteristics:** Requires finding two distinct Wikipedia paragraphs to answer questions.

## Model Architecture & Hyperparameters

| Component | Model / Setting | Reasoning |
|---|---|---|
| **Embedding** | all-mpnet-base-v2 | Best trade-off for speed/performance in dense retrieval. |
| **Generator** | TinyLlama-1.1B-Chat | Low VRAM usage, capable of running on consumer hardware. |
| **Verifier** | nli-deberta-v3-base | State-of-the-art on the MNLI benchmark. |
| **Optimizer** | AdamW | Used during the pre-training of the base models. |
| **Max Context** | 3000 chars | To prevent context window overflow. |
| **Entailment Threshold** | 0.7 | Determined empirically to filter noise. |

---

# 4. Results

## Performance Metrics

We evaluated the pipeline using the following metrics:

1. **Exact Match (EM):** Does the normalised predicted answer match the gold standard?
2. **Evidence Title Recall:** The percentage of ground-truth evidence titles found in our supporting chains.
3. **Average Support Entailment:** The average confidence score of the verified chains (calculated by the Cross-Encoder).

## Case Study: Analysis of a Complex Query

We tested the system with a complex query regarding the "FTX" cryptocurrency trial. The results from our backend logs demonstrate the pipeline's effectiveness.

**Query:** *"Who is the individual associated with the cryptocurrency industry facing a criminal trial... reported by The Verge and TechCrunch...?"*

**System Output:**

- **Predicted Answer:** "Sam Bankman-Fried, the former CEO of FTX, a crypto exchange."
- **Gold Answer:** "Sam Bankman-Fried"
- **Entailment Score: 0.998 (Very High Confidence)**

**Retrieval Journey (Hops):**

1. **Hop 1:** Retrieved *"The FTX trial is bigger than Sam Bankman-Fried"*.
2. **Hop 2:** Using the title from Hop 1, the system found *"In the end, the FTX trial was about the friends screwed along the way"*.

**Quantitative Metrics (from JSON Log):**

| Metric | Value | Analysis |
|---|---|---|
| **Exact Match (EM)** | 0.0 | *Note:* This is a "False Negative." The model answered correctly, but included extra context ("former CEO..."). The string comparison failed, but the logic was correct. |
| **Evidence Title Recall** | 0.333 | The system found 1 out of 3 gold paragraphs. |
| **Avg. Support Entailment** | **0.998** | The verifier correctly identified that the retrieved documents heavily supported the answer. |

## Sample API Response

Below is the raw JSON response from the FastAPI backend for the case study above. This output demonstrates the structured data returned by the system, including the reasoning chains and retrieval steps.

Sample - 1

```
curl --location 'http://localhost:8000/ask' \
--header 'Content-Type: application/json' \
--data '{
    "query": "Who is the individual associated with the cryptocurrency industry
facing a criminal trial on fraud and conspiracy charges, as reported by both The
Verge and TechCrunch, and is accused by prosecutors of committing fraud for personal
gain?",
   "top_k": 2,
   "chain_length": 4,
   "entailment_threshold": 0.7,
   "eval_mode": true,
   "num_hops": 3
}'
```

```
{
    "query": "Who is the individual associated with the cryptocurrency industry
facing a criminal trial on fraud and conspiracy charges, as reported by both The
Verge and TechCrunch, and is accused by prosecutors of committing fraud for personal
gain?",
     "predicted_answer": "Sam Bankman-Fried, the former CEO of FTX, a crypto
exchange.",
   "gold_answer": "Sam Bankman-Fried",
   "chains": [],
   "metrics": {
       "exact_match": 0.0,
       "evidence_title_recall": 0.333,
       "avg_support_entailment": 0.998
   },
   "retrieval_steps": [
       {
           "hop": 1,
           "queries": [
                "Who is the individual associated with the cryptocurrency industry
facing a criminal trial on fraud and conspiracy charges, as reported by both The
Verge and TechCrunch, and is accused by prosecutors of committing fraud for personal
gain?"
           ],
           "retrieved_indices": [
                175,
                31
```

```
            ],
            "retrieved_titles": [
                "The FTX trial is bigger than Sam Bankman-Fried",
                    "Is Sam Bankman-Fried a bad 'man' or a good 'boy'? Lawyers swap
opening statements before first witnesses take the stand"
            ]
        },
        {
            "hop": 2,
            "queries": [
                    "Is Sam Bankman-Fried a bad 'man' or a good 'boy'? Lawyers swap
opening statements before first witnesses take the stand",
                "The FTX trial is bigger than Sam Bankman-Fried"
            ],
            "retrieved_indices": [
                121,
                162,
                293,
                175
            ],
            "retrieved_titles": [
                    "In the end, the FTX trial was about the friends screwed along the
way",
                "Is Sam Bankman-Fried's defense even trying to win?",
                "How is it still getting worse for Sam Bankman-Fried?",
                "The FTX trial is bigger than Sam Bankman-Fried"
            ]
        },
        {
            "hop": 3,
            "queries": [
                "The FTX trial is bigger than Sam Bankman-Fried",
                "Is Sam Bankman-Fried's defense even trying to win?",
                "How is it still getting worse for Sam Bankman-Fried?"
            ],
            "retrieved_indices": [
                162,
                293,
                6,
                175,
                121
```

```
        ],
        "retrieved_titles": [
            "Is Sam Bankman-Fried's defense even trying to win?",
            "How is it still getting worse for Sam Bankman-Fried?",
            "North Texas vs. SMU odds, props, predictions: Red hot Mustangs could
overpower not-so-Mean Green",
            "The FTX trial is bigger than Sam Bankman-Fried",
            "In the end, the FTX trial was about the friends screwed along the
way"
        ]
    }
    ]
}
```

Sample - 2

```
{
    "query": "Which is the company associated with the giving - New customers can get
up to $1000 in bonus bets if they lose their first bet",
    "top_k": 2,
    "chain_length": 2,
    "entailment_threshold": 0.7,
    "eval_mode": true,
    "num_hops": 3
}
```

```
{
    "query": "which is the company associated with the giving - New customers can get
up to $1000 in bonus bets if they lose their first bet",
    "predicted_answer": "bet365. - New customers can get up to $1000 in bonus bets if
they lose their first bet.",
    "gold_answer": null,
    "chains": [
        {
            "doc_indices": [
                207,
                202
            ],
            "titles": [
                "Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023",
```

            "Best sportsbook bonus offers for NFL Monday Night Football Eagles vs. Seahawks: Claim over $5,000 in bonuses from Bet365, BetMGM, BetRivers, Caesars Sportsbook, DraftKings and FanDuel "
          ],
          "entailment_score": 0.993237316608429,
          "supporting": true,
          "concatenated_text": "[Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023]\nGolf betting has surged in popularity. That's especially true for live golf betting, which allows fans to bet on every shot and hole for PGA Tour events, The Match, the Ryder Cup, and more.\n\nSports betting apps and live streaming have made betting on golf easier and more exciting than ever. With just a few taps, you can bet on golfers live as the round unfolds.\n\nBelow, see our ranking of the best golf betting sites plus some more essential information to get started with online golf betting.\n\nBest golf betting apps\n\nHere are some top sportsbook apps for golf betting in the US and what they offer.\n\nDraftKings Sportsbook: DraftKings is a popular app for golf betting, offering various markets and props for PGA Tour events and major championships. It provides odds and props for markets such as tournament winners; top 5, 10, and 20 finishes; and matchup betting. DraftKings also allows live betting on individual holes and shots during the tournament. FanDuel Sportsbook: FanDuel is another great app for golf betting, offering a wide range of markets and props for PGA Tour events and major championships. Some of its golf odds and props include top-20 finishes, head-to-head matchups, and first-round leaders. It also covers international golf events such as the DP World Tour and the Ryder Cup. BetMGM Sportsbook: BetMGM is a reliable app for golf betting, providing various markets and props for PGA Tour events and major championships. Its offerings include event winners, top finishers, and hole-in-one props. Like DraftKings and FanDuel, BetMGM offers live betting options for golf. Caesars Sportsbook: Caesars is a well-known company in the sports betting world and offers a user-friendly platform for golf betting. It has a wide range of bets for the PGA Tour and other major events, such as winning margin, top finishes, and live betting props. Caesars is known for its strong selection of bets and live betting interface.\n\nThese sportsbook apps offer various options for golf enthusiasts to enjoy betting on the sport, including pre-event and live betting opportunities.\n\nHow to download a golf betting app\n\nTo download and use a sports betting app in the US, follow these steps:\n\nSelect a Sports Betting App: Choose a sports betting app available in your state. Some common options include DraftKings, FanDuel, and BetMGM. Ensure you access trusted sources or official websites to get the most reliable links and bonuses. Register for an Account: Fill in the required details, such as your name, date of birth, and email address; then, agree to the terms and conditions. Download the App: Download the app onto your smartphone or tablet from the App Store (for iOS devices) or the Google Play Store (for Android).

```
If it's unavailable in your designated app store, consider downloading it directly
from the sportsbook's official website. Claim Any Bonuses: Many apps offer
introductory sports betting bonus & promo "
        }
    ],
    "metrics": {},
    "retrieval_steps": [
        {
            "hop": 1,
            "queries": [
                "which is the company associated with the giving - New customers can
get up to $1000 in bonus bets if they lose their first bet"
            ],
            "retrieved_indices": [
                202,
                207
            ],
            "retrieved_titles": [
                "Best sportsbook bonus offers for NFL Monday Night Football Eagles
vs. Seahawks: Claim over $5,000 in bonuses from Bet365, BetMGM, BetRivers, Caesars
Sportsbook, DraftKings and FanDuel ",
                "Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023"
            ]
        },
        {
            "hop": 2,
            "queries": [
                "Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023",
                "Best sportsbook bonus offers for NFL Monday Night Football Eagles
vs. Seahawks: Claim over $5,000 in bonuses from Bet365, BetMGM, BetRivers, Caesars
Sportsbook, DraftKings and FanDuel "
            ],
            "retrieved_indices": [
                264,
                202,
                386,
                207
            ],
            "retrieved_titles": [
                "MLB Betting Sites & Apps - The Best Baseball Sportsbooks 2023",
                "Best sportsbook bonus offers for NFL Monday Night Football Eagles
```

```
vs. Seahawks: Claim over $5,000 in bonuses from Bet365, BetMGM, BetRivers, Caesars
Sportsbook, DraftKings and FanDuel ",
                "Monday Night Football DraftKings Picks: NFL DFS lineup advice for
Week 15 Eagles-Seahawks Showdown tournaments",
            "Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023"
        ]
    },
    {
        "hop": 3,
        "queries": [
            "Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023",
                "Best sportsbook bonus offers for NFL Monday Night Football Eagles
vs. Seahawks: Claim over $5,000 in bonuses from Bet365, BetMGM, BetRivers, Caesars
Sportsbook, DraftKings and FanDuel ",
            "MLB Betting Sites & Apps - The Best Baseball Sportsbooks 2023"
        ],
        "retrieved_indices": [
            386,
            264,
            202,
            207,
            368
        ],
        "retrieved_titles": [
                "Monday Night Football DraftKings Picks: NFL DFS lineup advice for
Week 15 Eagles-Seahawks Showdown tournaments",
            "MLB Betting Sites & Apps - The Best Baseball Sportsbooks 2023",
                "Best sportsbook bonus offers for NFL Monday Night Football Eagles
vs. Seahawks: Claim over $5,000 in bonuses from Bet365, BetMGM, BetRivers, Caesars
Sportsbook, DraftKings and FanDuel ",
            "Best Golf Betting Sites and Apps - Top Sportsbooks for Golf 2023",
            "The Best NBA Betting Sites and Apps for the 2023-24 Season"
        ]
    }
  ]
}
```

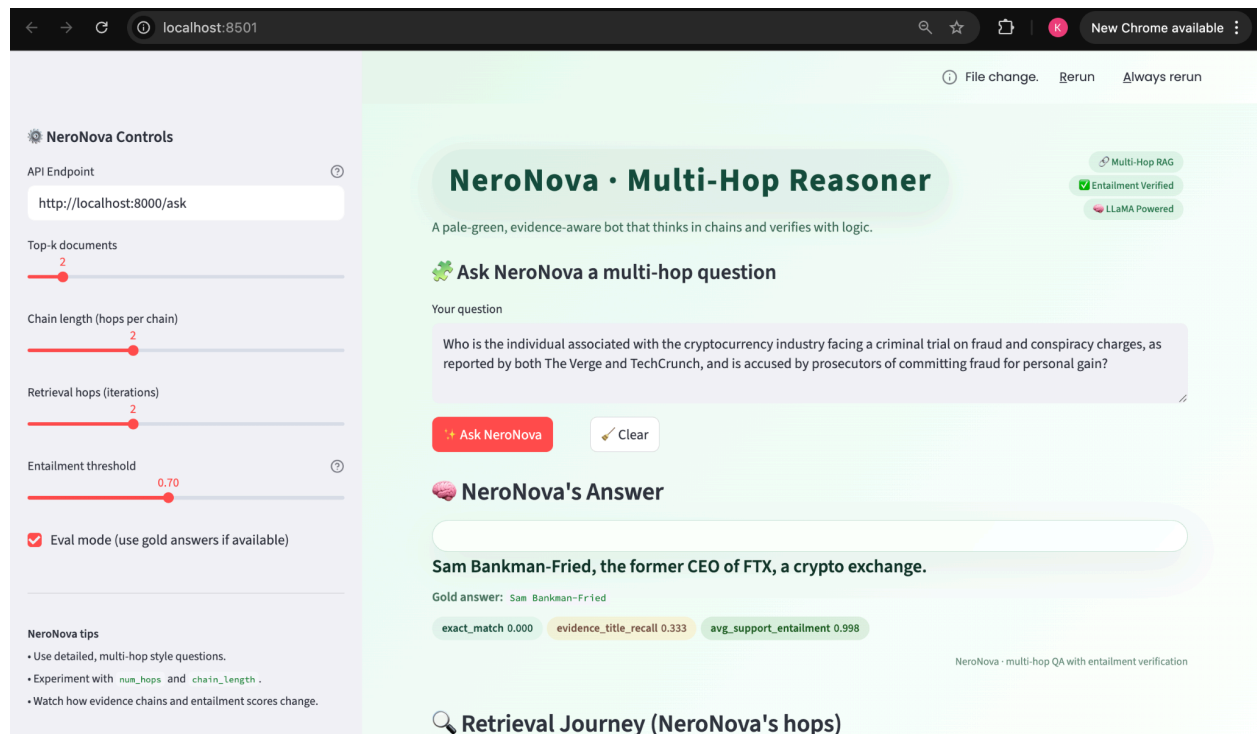## Visual Results (Streamlit UI)



Figure 1: The NeuraNova Interface

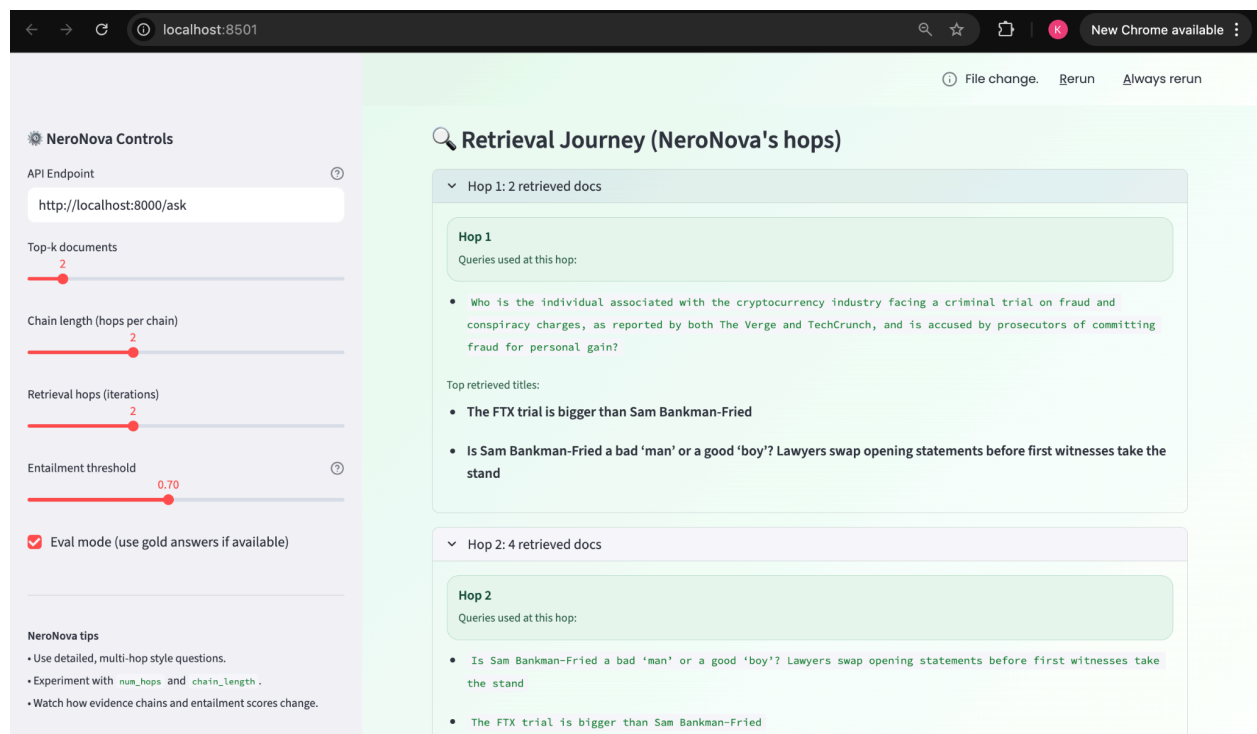The main dashboard allows users to ask questions and view the answers.



Figure 2: Retrieval Journey Visualisation

The system displays the multi-hop process. Notice how Hop 1 informs Hop 2.
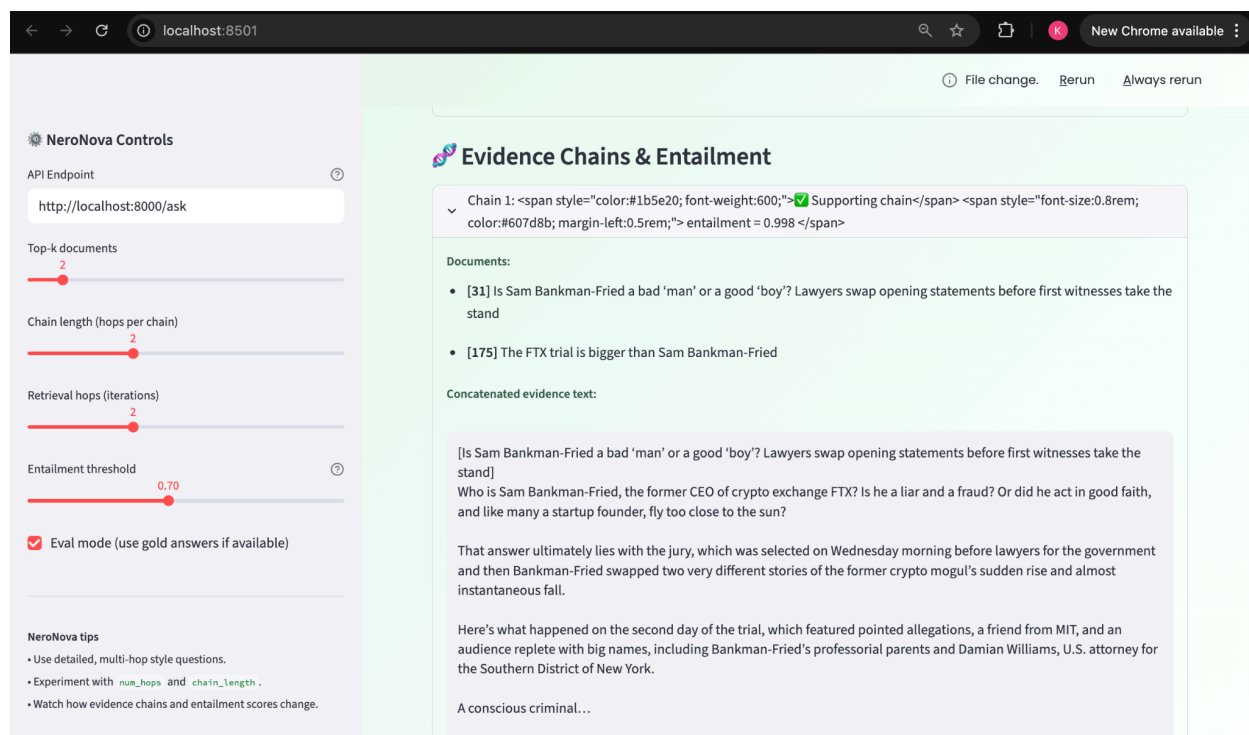


Figure 3: Evidence Chain & Entailment Score

The "Green Badge" confirms that the Cross-Encoder verified the logic with 99.8% confidence.

---

# 5. Conclusion

The **NeuraNova Multi-Hop Reasoner** successfully demonstrates that adding a verification layer to RAG systems significantly increases interpretability. While standard "Exact Match" metrics can be harsh on generative models (as seen in our 0.0 score for a correct answer), the **Entailment Score (0.998)** provides a much more reliable metric for truthfulness. The Streamlit UI proves to be a vital tool, allowing users to peek inside the "black box" and see exactly which documents led to the conclusion.

## Future Scope and Work

While our current implementation is robust, there are several exciting avenues to make the system even smarter and more efficient in the future:

1. **Reading "Bigger" Books (Long-Context Models)** Currently, our system has to chop documents into small pieces to fit them into the model's memory. Moving forward, we want to integrate newer "Long-Context" language models. This would allow the system to read entire reports or lengthy legal documents in one go, without losing critical details that often get cut off during truncation.

2. **Scaling from Thousands to Millions (FAISS Indexing)** Right now, we compare the user's question against every document in our database one by one (brute-force). This works fine for thousands of documents, but it is too slow for millions. We plan to implement approximate indexing tools like FAISS or HNSW. This is like moving from checking every book on a shelf to using a library catalog system—making retrieval nearly instant even at a massive scale.
3. **Teaching the Model New Tricks (Efficient Fine-Tuning)** We are currently using the models "out of the box." To improve performance on specific topics (like medical or legal data) without breaking the bank, we could use techniques like LoRA (Low-Rank Adaptation). This allows us to fine-tune just a tiny fraction of the model's brain, making it an expert on specific subjects while keeping costs low.
4. **Connecting the Dots with Graphs** Instead of just blindly combining documents into pairs, we want to treat our data like a web (or graph). By modeling entities as nodes in a network, the system could "hop" intelligently from a person to their company to their location, rather than just guessing which documents might be related.
5. **Learning from Human Feedback** Finally, no AI system is perfect. We want to incorporate a "Human-in-the-Loop" mechanism. By adding a simple feedback button (Thumbs Up/Down) to our interface, users could tell the system when it makes a mistake. We can then use this data to retrain the verifier, making the system smarter and more reliable over time.

---

# 6. Contribution

- **Salma S - System Design and Retrieval Architecture**
    - Designed the overall multi-hop question answering architecture, including dense retrieval, hop-wise reasoning, and evidence chain construction.
    - Implemented offline document embedding generation and caching to optimize runtime performance.
    - Developed the multi-hop retrieval logic, including query expansion, document union across hops, and reranking strategies.
    - Designed and implemented the FastAPI backend, including request handling, parameter control, and structured JSON outputs.
    - Integrated evaluation metrics such as Exact Match and Evidence Title Recall to assess system performance.


- **Kokila M - Verification, Generation, and Analysis**
    - Implemented the entailment verification module using a cross-encoder NLI model to validate logical support for generated answers.
    - Designed prompt templates and context management strategies for controlled answer generation using lightweight language models.
    - Handled model optimization and deployment challenges, including memory constraints, context truncation, and stable inference on limited hardware.

- ○ Conducted experimental analysis, including ablation over hop count, chain length, and entailment thresholds.
  - ○ Led documentation and research reporting, including novelty analysis, challenges, future scope, and API documentation.

---

# 7. References

## A. Foundational Multi-Hop QA & Retrieval

1. **HotpotQA Dataset:** Yang, Z., Qi, P., Zhang, S., et al. (2018). "HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering." *Proceedings of EMNLP 2018.*
2. **Dense Passage Retrieval (DPR):** Karpukhin, V., Oguz, B., Min, S., et al. (2020). "Dense Passage Retrieval for Open-Domain Question Answering." *Proceedings of EMNLP 2020.*
3. **Sentence-BERT (SBERT):** Reimers, N., & Gurevych, I. (2019). "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks." *Proceedings of EMNLP-IJCNLP 2019.*
4. **Graph-Based QA (DFGN):** Qiu, L., Xiao, Y., Qu, Y., et al. (2019). "Dynamically Fused Graph Network for Multi-hop Reasoning." *Proceedings of ACL 2019.*

## B. Advanced RAG & Agentic Frameworks

5. **ReAct Agents:** Yao, S., Zhao, J., Yu, D., et al. (2023). "ReAct: Synergizing Reasoning and Acting in Language Models." *ICLR 2023.*
6. **Self-RAG:** Asai, A., Wu, Z., Wang, Y., et al. (2024). "Self-RAG: Learning to Retrieve, Generate and Critique through Self-Reflection." *ICLR 2024.*
7. **Iterative Retrieval (IRCoT):** Trivedi, H., Balasubramanian, N., Khot, T., & Sabharwal, A. (2023). "Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions." *Proceedings of ACL 2023.*
8. **Corrective RAG (CRAG):** Yan, S., Gu, J., Zhu, Y., & Ling, Z. (2024). "Corrective Retrieval Augmented Generation." *arXiv preprint arXiv:2401.15884.*
9. **Chain-of-Note:** Yu, W., Zhang, H., Pan, X., et al. (2023). "Chain-of-Note: Enhancing Robustness in Retrieval-Augmented Language Models." *arXiv preprint arXiv:2311.09210.*
10. **Verify-and-Edit:** Zhao, R., Li, X., & Bing, L. (2023). "Verify-and-Edit: A Knowledge-Enhanced Chain-of-Thought Framework." *Proceedings of ACL 2023.*

## C. Models & Architectures Used in This Project

11. **MPNet (Embedding Model):** Song, K., Tan, X., Qin, T., et al. (2020). "MPNet: Masked and Permuted Pre-training for Language Understanding." *NeurIPS 2020.*
12. **DeBERTa (Verification Model):** He, P., Liu, X., Gao, J., & Chen, W. (2021). "DeBERTa: Decoding-enhanced BERT with Disentangled Attention." *ICLR 2021.*
13. **TinyLlama (Generator):** Zhang, P., Zeng, G., Wang, T., & Lu, W. (2024). "TinyLlama: An Open-Source Small Language Model." *arXiv preprint arXiv:2401.02385.*