# Multi-Source RAG for Technical Support
## TaskFlow Lite (Fictional Product)

Srutileka Suresh

February 2026

## 1 Problem Statement

Build a Retrieval-Augmented QA system that answers customer support questions about a fictional software product (TaskFlow Lite) by retrieving evidence from **three distinct sources**: product documentation, customer forums, and technical blog posts. The system must implement source-specific chunking, multi-source retrieval with weighting, reranking, contradiction handling, and logging.

## 2 Product Overview: TaskFlow Lite

TaskFlow Lite is a lightweight project/task management tool. Users can manage projects (with a task limit), attach files to tasks (with file type + size limits), and use an API (with key expiry and rate limits). The knowledge sources intentionally include cases where community forums or blogs conflict with official docs.

## 3 Dataset (Custom made)

- **Documentation (.md):** Authoritative product rules such as attachment limits, export formats, API key expiry, and error codes.

- **Customer Forums(.jsonl):** User questions and replies; may be noisy or contradictory.

- **Technical Blogs(.md):** How-to posts and best practices; sometimes outdated.

## 4 Chunking Strategy

Different sources have different structure, so chunking is tuned per source:

- **Docs:** Larger, structured chunks to preserve bullet lists (limits, error codes) as a single coherent unit.

- **Blogs:** Medium chunks to keep narrative explanations together.

- **Forums:** Small chunks per question or per reply to keep distinct opinions separated (and to support contradiction detection).

# 5 Retrieval (Multi-Source + Weighting)

## 5.1 Why separate retrieval per source?

Three separate FAISS indexes are built (`docs`, `blogs`, `forums`). At query time, the system retrieves top-$k$ from each index to ensure every source can contribute evidence.

## 5.2 Source weighting

After retrieving from each source, the system applies source weights (docs prioritized):

$$\texttt{weighted\_score} = \texttt{similarity\_score} \times w_{source}$$

Weights used in the run: `docs=1.2`, `blogs=1.0`, `forums=0.8` :contentReferenceindex=2.

# 6 Reranking (Cross-Encoder)

The system uses a two-stage ranking pipeline:

1. **Bi-encoder retrieval** (fast): retrieve a shortlist of candidates via embedding similarity (FAISS).

2. **Cross-encoder reranking** (precise): rerank only the top candidates

This design balances speed (vector search) and relevance (token-level interaction in cross-encoder).

# 7 Contradiction Handling

Contradictions can occur between sources (e.g., size limits, expiry durations, rate limits). The system:

- extracts numeric claims (e.g., `25MB`, `60 days`, `300/hour`) from top evidence,

- flags conflicts when multiple values appear for the same category,

- resolves by a reliability policy: **Docs > Blogs > Forums**.

# 8 Logging

Each query is logged to a JSONL file including:

- query string,

- top candidates pre- and post-rerank,

- sources used and scores,

- detected contradictions and the selected reliability policy.

# 9 Performance Analysis

We compare retrieval quality before reranking (`pre`) vs after reranking (`post`) using file-level "gold" targets for the synthetic dataset :contentReferenceindex=4.

## 9.1 Summary Metrics

From `metrics.json` :contentReferenceindex=5:

| Metric | Pre-rerank | Post-rerank |
|--------|------------|-------------|
| Hit@5  | 0.80       | 1.00        |
| Hit@1  | 0.70       | 0.40        |

Table 1: Retrieval performance before vs after reranking (10 queries).

**Interpretation:**

- Hit@5 improves from 0.80 to 1.00, meaning reranking + multi-source merging increased the chance that the correct file appears in the top-5.

- Hit@1 decreases from 0.70 to 0.40 in this run. This happens because the cross-encoder sometimes ranked short, direct forum/blog snippets above the authoritative doc chunk (even though the doc chunk remained in the top results). The system mitigates this via a **source reliability policy** during answer selection.

# 10 Query Result Screenshots



(a) Q1 Output

(b) Q2 Output

(c) Q3 Output

(d) Q4 Output

(e) Q5 Output

(f) Q6 Output

(g) Q7 Output

(h) Q8 Output

(i) Q9 Output

(j) Q10 Output

Figure 1: Placeholder grid for 10 screenshot outputs (replace images in `figs/`).

# 11 Limitations and Future Work

- The contradiction detector focuses on numeric claims (MB/days/hour). More general contradiction detection would require deeper semantic modeling.

- Hit@1 decreased after reranking in this run :contentReferenceindex=17; a refinement is to incorporate source priors directly into reranking or to select the highest reranked chunk within the preferred source.

- The answer synthesis is deterministic (no external LLM). Adding a small LLM could improve readability, but would introduce nondeterminism and potential hallucinations.