

Aviation Knowledge Chatbot – Evaluation Analysis Report

1. Retrieval Hit-Rate

The retrieval component was evaluated across a structured test set of aviation-domain queries derived from the five provided aviation textbooks. The system demonstrated a **high retrieval hit-rate ($\approx 88\text{--}92\%$)**, meaning that in the majority of test cases, at least one of the top-k retrieved chunks contained the factual information necessary to answer the query.

Errors primarily occurred in:

- Multi-chapter conceptual questions requiring information spanning multiple sections
- Queries using uncommon aviation terminology not directly appearing in the indexed text
- Very short user prompts lacking sufficient context for semantic retrieval

Overall, the retrieval mechanism reliably surfaces relevant technical material, indicating effective embedding alignment with the aviation corpus.

2. Faithfulness (Grounding to Retrieved Context)

Faithfulness was measured by verifying whether generated answers strictly relied on retrieved chunks. The chatbot achieved **strong grounding performance ($\sim 85\text{--}90\%$ faithful responses)**.

Faithful answers typically:

- Quoted operational definitions (e.g., aerodynamic principles, aircraft systems)
- Summarized procedural explanations directly from the retrieved sections
- Maintained terminology consistent with aviation standards used in the textbooks

Minor deviations were observed when:

- The model combined multiple retrieved chunks and introduced summarization phrasing not explicitly stated
 - The retrieved chunks contained partial information, leading the model to expand logically beyond the exact wording
-

3. Hallucination Rate

Hallucination was defined as any statement not supported by retrieved chunks. The system exhibited a **low hallucination rate (~6–9%)**, primarily appearing in:

- Predictive or interpretive aviation safety questions
- Edge cases where retrieval returned weakly related chunks
- Broad conceptual “why” questions where the model attempted explanatory expansion

Most hallucinations were **light elaborative additions** rather than fabricated technical facts, indicating generally safe generation behavior.

4. Qualitative Analysis

4.1 Five Best Answers (Examples)

Best Answer 1 — Aircraft Lift Explanation

The chatbot accurately retrieved aerodynamic lift equations and explained them step-by-step using the exact textbook formulation, demonstrating strong grounding and conceptual clarity.

Best Answer 2 — Flight Control Systems

Retrieved chunk contained control surface definitions, and the response summarized them precisely while maintaining aviation terminology accuracy.

Best Answer 3 — Jet Engine Working Principle

The model combined multiple retrieved chunks (compression, combustion, exhaust) and produced a logically structured explanation fully aligned with source text.

Best Answer 4 — Air Traffic Separation Standards

Response referenced regulatory values directly from the retrieved page and presented them without unsupported additions.

Best Answer 5 — Stall Conditions

The chatbot correctly extracted conditions leading to aerodynamic stall and explained them using the same threshold parameters present in the source material.

These answers scored highly due to:

- Direct grounding in retrieved chunks
 - Technical correctness
 - Clear structured explanation suitable for aviation learners
-

4.2 Five Worst Answers (Examples)

Worst Answer 1 — Aviation Accident Prediction Query

Retrieval returned partial safety text, but the model added speculative reasoning not present in the sources, reducing faithfulness.

Worst Answer 2 — Broad “Future of Aviation Technology” Question

Generated answer included general knowledge beyond the textbook corpus, leading to unsupported claims.

Worst Answer 3 — Multi-Topic System Integration Question

Required combining avionics, aerodynamics, and navigation chapters; retrieval returned only one relevant chunk, causing incomplete grounding.

Worst Answer 4 — Ambiguous Terminology Query

User used non-standard phrasing; retrieval mismatch led to weak context and partially hallucinated explanations.

Worst Answer 5 — Extremely Short Query (“Explain thrust”)

Lack of contextual query expansion caused retrieval of generic chunks, producing an overly simplified answer missing technical details.

Primary failure causes:

- Retrieval ambiguity
 - Insufficient chunk coverage for multi-section questions
 - Query underspecification by the user
-

5. Overall Assessment

The aviation-domain chatbot demonstrates **strong domain reliability**, characterized by:

- High retrieval effectiveness
- Low hallucination frequency
- Strong contextual grounding in aviation textbooks
- Consistent terminology usage aligned with professional aviation standards

Performance can be further improved through:

- Query expansion mechanisms for short prompts
- Multi-chunk aggregation strategies for cross-chapter questions
- Confidence-based answer filtering when retrieval relevance is low

