003

005

006

800

009

010

012

013

014

015

016

017

018

019

020

021

024

025

026

027

028

029

030

033

034

035

036

037

038

041

042

043

046

047

049

050

051

053

060

061

064

074

075

087

RAG in the Aerospace Domain: A Comprehensive Retrieval, Generation, and User Evaluation for NASA Documentation

Anonymous Full Paper Submission 53

Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities in Natural Language Understanding and text generation, but their application is often limited by hallucinations, outdated knowledge, and lack of evidence. Retrieval-Augmented Generation (RAG) addresses these fundamental LLM limitations by integrating external knowledge sources, thereby improving the factual accuracy and traceability while maintaining the text generative capabilities. This work presents the design and implementation of a web-based RAG system for the aerospace domain, leveraging more than 10,000 NASA technical documents and lessons-learned mission reports. The system integrates open-source LLaMA and closed-source OpenAI models and performs an extensive comparative analysis of their performance within the RAG framework. Evaluation through both automated metrics and user studies demonstrates the effectiveness of the RAG approach for both technical and non-technical users. The findings provide insights and establish a foundation for future advancements in AI-driven knowledge management for specialized fields¹.

1 Introduction

Large organizations face a critical challenge: leveraging years or decades of accumulated knowledge and lessons learned to inform new projects and decisions. This problem is particularly acute in high-risk, high-stakes environments. An example of such an environment is NASA, an organization with thousands of past projects that span decades. NASA engineers who work on new missions often struggle to find relevant historical information that could prevent costly mistakes or accelerate innovation. The core challenge lies in the complexity of discovering knowledge within vast and diverse available collections. NASA, as a large and long-established organization that works with multiple contractors, faces several barriers: vocabulary differences across time periods, varying terminology between contracting companies, and the sheer scale of documentation. This creates a scenario where critical lessons learned, such as

the infamous O-ring failure that led to the Challenger disaster in 1986, may be documented, but remain inaccessible to engineers working on similar components in new projects [1]. The result is a knowledge gap where valuable insights from past missions remain isolated and unused. This challenge is exemplified by NASA's efforts to develop risk digital assistants that can extract and leverage past project data for predictive decision-making².

This problem extends beyond aerospace to any large organization with extensive historical documentation: healthcare systems with decades of patient data, legal firms with case histories, or manufacturing companies with safety records. In these domains, the ability to quickly and accurately retrieve relevant historical information can significantly impact decision quality, risk assessment, and project outcomes.

RAG has emerged as a leading approach that enhances LLMs by seamlessly integrating external knowledge sources with text generation capabilities [2]. This hybrid approach allows systems to take advantage of both the generative strengths of LLMs and the precision of Information Retrieval (IR) techniques, producing responses drawn from a substantially broader knowledge base than what is encoded in the model parameters alone. While RAG systems have demonstrated effectiveness across various domains, the specific challenges of applying them to highly specialized technical documentation in aerospace environments present unique opportunities for research.

The core contribution of this work is addressing this knowledge discovery challenge through a specialized RAG system designed for NASA's documentation. By processing more than 10,000 NASA technical documents [3] and lessons-learned mission reports [4], this system allows engineers to quickly access verified information with direct source references, significantly reducing the time and effort required to find relevant historical data. The system integrates and compares open-source LLaMA and closed-source OpenAI models, addressing the challenge of making NASA's extensive technical documentation and lessons-learned databases more accessible and trustworthy.

This paper is organized as follows: next section open reviews existing literature in the domains of RAG, open control of the c

 $^{^{1}\}mathrm{Code}$ and datasets will be freely provided after acceptance.

²https://techport.nasa.gov/projects/117547

094

095

096

097

098

099

100

101

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

125

126

127

128

129

130

131

133

134

135

136

137

138

140

141

142

143

144

145

157

158

159

161

162

166

169

172

173

174

184

185

191

192

193

195

203

204

conversational agents, and evaluation methodologies, focusing on applications in specialized domains and space sector implementations. Our approach is described in Section 3 and the evaluation is discussed in Section 4. We conclude the paper in Section 5.

2 Related Work

The RAG architecture was first introduced by Lewis et al. [5], combining a retrieval component with a sequence-to-sequence generative model to enhance text generation with external knowledge. Unlike traditional LLMs that rely solely on their pre-trained knowledge, RAG allows models to generate high-quality text supported by relevant external information, making it particularly effective for knowledge-intensive tasks [5]. The effectiveness of retrieval-augmented approaches has been demonstrated in open-domain question answering systems [6, 7] and dialogue systems [8], where the integration of retrieval mechanisms has improved contextual understanding and response generation capabilities.

Conversational Agents in Space Sector. The development of conversational Question Answering systems has advanced significantly with the rise of LLMs, enabling models to answer questions based on given contexts, often involving RAG when documents surpass the language model's context window [9]. Studies focus primarily on text-only documents (such as regulations, manuals, and technical reports) from various domains [10–12] and documents combining plain text with tables (such as Wikipedia articles with tabular data, financial reports, and semi-structured knowledge sources) [13–15].

Several projects have explored the use of AIpowered virtual assistants to aid engineers during spacecraft mission design. The most notable ones are Daphne [16], the Design Engineering Assistant (DEA) [17], and SpaceQA [18]. DEA is a nonintrusive decision support tool that enhances expert perception of different design alternatives and past decision outcomes through Natural Language Processing, Machine Learning, Knowledge Management, and Human-Machine Interaction methods [17]. Daphne, evaluated at NASA's Jet Propulsion Laboratory (JPL), enhances design task performance through a microservices architecture featuring a web front-end, server (Daphne Brain), and software roles that interface with structured knowledge graphs for better design inputs [19]. SpaceQA, developed by the European Space Agency (ESA), is an opendomain QA system for space mission design, employing a dense retriever and neural reader similar to an RAG pipeline [19].

LLM-as-a-Judge Evaluation in RAG Systems. Traditional evaluation of the RAG systems require manually created ground truth data, which

are typically quite expensive to acquire. Traditional metrics like ROUGE and BLEU fail to capture the nuanced quality dimensions required for RAG evaluation, particularly factual correctness and contextual grounding. New evaluation methods, which employ LLMs, thus recently started to emerge in RAG evaluation. Wang et al. [20] demonstrate that ChatGPT can effectively evaluate text generation by providing scores on 0-100 or 1-5 star scales for aspects such as relevance, factual accuracy, and groundedness, achieving state-of-the-art correlation with human judgments across multiple NLG tasks. However, this approach exhibits sensitivity to prompt design and reduced effectiveness on datasets with strong lexical biases.

Muhamed [21] introduces the CCRS (Contextual Coherence and Relevance Score), employing LLaMA-70B as a zero-shot judge to evaluate RAG systems in five dimensions: contextual coherence, question relevance, information density, answer correctness, and information recall. Their evaluation on the BioASQ biomedical dataset shows promising results across multiple RAG configurations with different readers and retrievers, while highlighting the challenge of achieving perfect factual accuracy in complex biomedical domains.

Building on these evaluation foundations, broader investigations have explored the reliability and effectiveness of LLM judges across different contexts. Tseng et al. [22] conduct the first systematic evaluation of LLMs as expert-level data annotators across finance, biomedicine, and law domains, finding that models average 35% behind human expert performance despite showing promise in general NLP tasks. Ashktorab et al. [23] develop EvalAssist, comparing direct assessment versus pairwise comparison strategies, and demonstrate that practitioners prefer direct assessment for clarity while using pairwise comparison for subjective evaluations. Bavaresco et al. [24] present JUDGE-BENCH, a comprehensive benchmark evaluating 11 LLMs across 20 datasets, revealing substantial variance in model performance and better alignment with non-expert versus expert human judges. Thus, although the LLM-as-a-Judge Evaluation still suffers from multiple issues, due to its low cost, it is now a standard for RAG evaluation, often accompanying expensive user studies.

User Evaluation Studies in RAG Systems. While technical metrics provide important insights into RAG performance, user-centered evaluation remains critical for understanding real-world effectiveness and adoption. Hasan et al. [25] present a comprehensive study of five domain-specific RAG applications deployed across governance, cybersecurity, agriculture, industrial research, and medical diagnostics. Their approach combined Likert-scale surveys with open-ended qualitative feedback to capture both measurable insights and descriptive

208

209

211

213

214

215

216

217

218

220

221

222

224

225

226

227

228

229

231

232

233

234

235

236

239

240

242

243

244

246

247

248

249

250

252

253

257

258

260

261

262

263

264

265

266

273

274

275

279

280

281

283

285

286

user experiences, ultimately documenting twelve key lessons learned from user feedback to guide future RAG system development and deployment practices.

3 Method

RAG Framework 3.1

RAG is a hybrid approach that combines two essential components: (i) a retrieval system that pulls documents from an external knowledge base, and (ii) a generation component that uses this information to create natural, human-like text [26, 27]. By blending these capabilities, RAG models can produce coherent and fluent responses while anchoring their output in current real-world information. The workflow is illustrated in Figure 1.

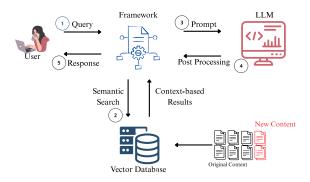


Figure 1. RAG Architecture workflow: (1) User submits query to the system, (2) Query is vectorized and semantic search is performed against the document vector database, (3) Top-k relevant documents are retrieved and passed as context, (4) LLM generates response using both the original query and retrieved context, and (5) presents it to the user.

The generation component leverages the contextual information to produce responses that maintain the natural language capabilities of LLMs while being factually anchored. The system can be configured to prioritize retrieved context over pre-trained knowledge, ensuring that responses are based primarily on the retrieved NASA documents rather than the model's training data [28].

3.2**Dataset Preparation**

To be able to use the data, which would support our use case of vast knowledge of a single organization, we decided to create a novel collection consisting of freely available NASA technical resources. We used web scraping techniques to collect information from NASA Technical Report Server³, which contains papers, patents, reports and other technical materials created or funded by NASA. The server

contains thousands of documents, collected since 1915, though our system focuses on documents from the 2000s onwards to ensure relevance to current aerospace practices. In addition to this, we used documents available at the NASA Lessons Learned⁴ which contains reviewed lessons learned from NASA programs and projects. Together, we collected 1.859 records from the Lessons Learned database, and 8,143 PDF files from the NASA Technical Report Server. The collection contains various aerospacerelated materials, on topics such as spacecraft design, propulsion systems, and mission operations.

The Selenium library was employed to automate browser interactions and handle JavaScript-rendered content, while BeautifulSoup was utilized for parsing HTML and navigating through hundreds of pages efficiently. This combination was particularly needed for NASA's websites, which required navigating through complex page hierarchies and extracting data from dynamically loaded content. The web scraping process required careful pacing to avoid triggering rate-limiting mechanisms, with individual PDF files often exceeding 30MB due to extensive technical content and images.

3.3 RAG System Implementation

The core pipeline consists of six key steps: data gathering, chunking, vectorization, storage, retrieval, and response generation.

3.3.1 **Document Processing and Chunking**

Chunking Strategy Analysis. Three primary chunking approaches were evaluated for this domainspecific application: (i) recursive chunking, which splits text based on hierarchical separators (paragraphs, sentences, words), (ii) semantic chunking, which groups text based on semantic similarity using embedding models to identify natural breakpoints, and (iii) section-based chunking, which leverages document structure to maintain semantic boundaries.

Recursive chunking, while computationally efficient, was unsuitable for NASA documentation as it often split technical procedures across multiple chunks. Semantic chunking required significantly higher processing power than section chunking, and the chunks were too large to be useful for a system with limited context window.

Section-Based Chunking Implementation. 282 The implemented approach utilizes section-based chunking through the pymupdf4llm⁵ library, which employs a multi-stage algorithm specifically designed to extract PDF content in formats optimized for LLM and RAG environments. The library's

³https://ntrs.nasa.gov/search

⁴https://llis.nasa.gov/

 $^{^5}$ https://pypi.org/project/pymupdf4llm/

algorithm first converts PDF pages to GitHub-compatible markdown format while preserving document structure through font analysis, header detection, and formatting preservation. Subsequently, it identifies section boundaries using font size variations, header patterns (both # markdown syntax and bold formatting), and hierarchical document structure.

For a document D with n sections, the chunking function $\phi: D \to C$ produces chunks $C = \{c_1, c_2, \ldots, c_n\}$ where each chunk c_i maintains semantic coherence and includes metadata (chunk ID, title, content, page number, section number, source URL, and statistics) enabling full traceability to original NASA materials.

Dataset Structure Analysis. As illustrated in Figure 2, most Technical Reports contain between 5 and 20 chunks, as the NASA documents vary widely from brief single-page reports to extensive technical manuals exceeding 50 pages⁶. This distribution shows consistent chunking patterns across the collection. The figure shows a distribution capped at 80 chunks per document. We opted for this to limit the processing of the outlier PDFs with complex tables and figures.

Distribution of Documents by Chunk Count

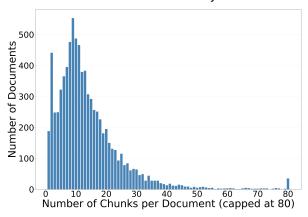


Figure 2. Distribution of NASA Technical Reports by chunk count showing the frequency of documents containing different numbers of chunks. Documents are capped at 80 chunks per document due to processing limitations that come from complex table formatting and images in certain PDFs.

3.3.2 Vectorization and Retrieval

For vectorization, the all-MiniLM-L6-v2⁷ transformer model is utilized, which is a lightweight but effective embedding model designed for semantic search tasks. The vectors are stored in a FAISS⁸

database for efficient similarity search. For a user query q, its embedding \mathbf{v}_q is computed and the top-k most similar chunks are retrieved using cosine similarity. The retrieval process returns chunks $R_k = \{c_{r_1}, c_{r_2}, \ldots, c_{r_k}\}$ where r_j represents the rank of chunk c_{r_j} based on similarity scores. We limit the system to retrieve top 3 documents per query, as retrieving a small number of top documents (commonly between 3 and 5) is standard practice, though the optimal k depends on the application and data characteristics. Given the relatively small response window in our system, 3 retrieved documents per query provided a good balance between context richness and computational efficiency.

Retrieval quality is evaluated using standard metrics including precision, recall, Mean Reciprocal Rank (MRR), Hit Rate, and Normalized Discounted Cumulative Gain (NDCG).

3.3.3 Language Model Integration

The system integrates both open-source and closed-source language models to enable a comprehensive comparison between paid and freely available options. The GPT models (GPT-40-mini and GPT-3.5-turbo) were selected for their cost-effectiveness while maintaining strong performance capabilities, making them accessible for budget-focused research applications. For open-source alternatives, LLaMA-3.3-70B-Instruct-Turbo and LLaMA-3.2-1B-Vision-Instruct represent the only freely available models on Together.ai's platform⁹, ensuring consistent cloud-based resource allocation critical for fair model comparison.

The system uses ChatPromptTemplate¹⁰ from LangChain, a library designed for creating structured and concise prompt instructions. This template formats the user's question with the retrieved context, creating a structured input that guides the language model to provide accurate, contextually grounded answers based on the NASA documentation rather than relying solely on its pre-trained knowledge.

For evaluation purposes, all four models are assessed in the automated generation evaluation to provide a comprehensive performance comparison. However, given the limited participant count in the user study, only GPT-40-mini and LLaMA-3.3-70B-Instruct-Turbo are included in the user evaluation to ensure sufficient statistical power while representing both commercial and open-source model categories.

⁶The Lessons Learned are typically single-page reports that convert entirely into individual chunks and are thus not included in the distribution.

⁷https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

⁸https://faiss.ai/index.html

⁹https://www.together.ai/models

¹⁰https://python.langchain.com/api_reference/
core/prompts/langchain_core.prompts.chat.
ChatPromptTemplate.html

369

370

371

375

376

377

378

379

380

382

383

384

385

386

387

388

389

391

392

393

394

395

396

397

398

399

400

401

402

403

404

406

407

415

417

418

419

421

426

428

431

432

433

439

440

443

446

447

448

450

451

456

459

3.4 Evaluation Framework

The evaluation approach is built upon a our constructed test dataset consisting of 50 records that are used for both retrieval and generation evaluation. Each record contains a manually crafted question created from randomly selected document chunks, paired with expected answers derived from the chunk content. For answer preparation, either the complete chunk content was used or GPT-4 was employed to extract only the parts relevant to the specific problem, ensuring that evaluation answers remained strictly grounded in the original NASA documentation while avoiding potential human error in test dataset creation. The expected answers are standardized to 150-200 tokens in length, with each record referencing 1-3 relevant document chunks. Table 1 illustrates the structure of a typical evaluation record.

Table 1. Example evaluation record from the test dataset.

Question	What issue was discovered with the ILT radiator during the CALIPSO satellite-level thermal vacuum test?
Expected Answer	During the satellite-level thermal vacuum test of CALIPSO, which aimed to demonstrate positive thermal control of all payload components within their required temperature limits under
Relevant Chunks	20070021525.pdf_11, 20070021525.pdf_12

3.4.1 Automated Evaluation

Generation quality is evaluated using a multifaceted approach that uses ROUGE scores, BERTScore for estimating semantic similarity, and LLM-as-a-Judge evaluation.

BERTS core computes semantic similarity by summing cosine similarities between token embeddings. For reference sentence x and predicted sentence y, the F1 measure is calculated as:

$$F_{BERT} = 2\frac{P \times R}{P + R} \tag{1}$$

where P and R represent precision and recall based on maximum cosine similarities between tokens of the generated and ground truth answer [29].

ROUGE metrics evaluate n-gram overlap between generated and reference texts. The ROUGE-N score is calculated as:

ROUGE-N =
$$\frac{\sum_{S \in \mathcal{R}} \sum_{g \in S} \text{Count}_{match}(g)}{\sum_{S \in \mathcal{R}} \sum_{g \in S} \text{Count}(g)}$$
(2)

where \mathcal{R} represents reference summaries, g represents n-grams, and $\operatorname{Count}_{match}(g)$ is the maximum number of matching n-grams of the generated and ground truth answer [30].

3.4.2 LLM-as-a-Judge Methodology

The LLM-as-a-Judge approach has emerged as a promising method to replace traditional statistical

metrics and human evaluation with LLMs for assessment tasks [31]. Compared to traditional evaluation methods, LLM judges demonstrate significant advantages: they can adjust evaluation criteria based on specific task context rather than relying on fixed metrics, generate interpretive evaluations that offer comprehensive feedback on model performance, and provide a scalable and reproducible alternative to human evaluation while significantly reducing associated costs and time [31].

However, LLM-as-a-Judge approaches face several critical challenges. Evaluation results are often influenced by prompt templates, which can lead to biased or inconsistent assessments [32]. Additionally, LLMs may inherit implicit biases from their training data, impacting the fairness and reliability of their evaluations, while distinct tasks and domains require specific evaluation criteria that make dynamic adaptation challenging [33].

This work utilizes two LLM-as-a-Judge methods: the first uses OpenAI's GPT-4o-mini model for evaluation, while the second implements G-Eval [34], a framework that uses a chain-of-thought (CoT) approach to evaluate the quality of generated text through structured evaluation forms. The main LLM-as-a-Judge scores are Correctness, Relevance, Accuracy, and Groundedness. Correctness measures factual correctness and completeness, Relevance assesses how well the answer addresses the question, Accuracy assesses technical detail alignment, and Groundedness verifies that the context is indeed obtained from the source documents.

The custom GPT-4o-mini evaluation uses a single prompt¹¹ requesting scores on a 0-10 scale for three criteria (relevance, factual accuracy, groundedness) with direct numerical output, while G-Eval employs a more structured CoT approach that decomposes evaluation into explicit reasoning steps for four criteria (correctness, relevance, accuracy, groundedness). This dual approach enables evaluation using both single-prompt and multi-prompt methodologies for comprehensive RAG system assessment.

3.4.3 User Study Design

The user study was conducted with 20 participants from diverse non-technical backgrounds, with only 3 participants having IT experience, ensuring evaluation from the general audience perspective. Each session lasted 15-25 minutes, during which each participant was asked to complete five distinct engineering tasks designed to simulate realistic scenarios requiring specialized aerospace knowledge retrieval.

Task Design and Implementation. The five tasks covered critical aerospace domains: (1) Engine Rollback Investigation, requiring analysis of Propul-

 $^{^{11}\}mathrm{GPT\text{-}4o\text{-}mini}$ uses custom Chat PromptTemplate prompts.

465

466

467

468

469

472

473

474

475

476

479

480

481

482

483

484

485

486

488

489

490

491

492

493

494

499

501

502

503

505

509

513

514

516

521

523

524

525

526

527

530

531

535

sion Systems Laboratory data on ice particle effects; (2) Satellite Thermal System Evaluation, focusing on CALIPSO payload thermal performance; (3) Aircraft Noise Profile Assessment, examining noise components across different flight conditions; (4) Critical Power System Design, investigating Uninterruptible Power Supply applications; and (5) Electronics System Safety Review, identifying analytical methods for circuit problem detection. Tasks were created by randomly choosing document chunks, manually creating problems from those chunks. GPT model was only used for acquiring supplementary information such as creating scenarios to assist users in understanding the problem. Each task required participants to find at least three of seven predefined keywords to demonstrate successful knowledge acquisition. Figure 3 illustrates the structure and complexity of a typical user study task.

Task 1: Engine Rollback Investigation

You're an aerospace engineer analyzing engine performance in icing conditions. Your team needs to understand what particle characteristics lead to engine rollback events. Research the Propulsion Systems Laboratory (PSL) test data findings to determine critical ice particle sizes and temperature conditions that contribute to these events.

Expected findings should include:

- Analysis of Propulsion Systems Laboratory (PSL) data points on the LF11 engine model
- Critical particle size requirements for engine rollback (when thrust unexpectedly decreases)
- Relevant wet bulb temperature range in the Low Pressure Compressor (LPC) region

Figure 3. Example of Task 1 in User Evaluation Study: Engine Rollback Investigation interface showing the structured engineering scenario presented to participants.

Study Protocol and Bias Mitigation. Participants received authentication credentials and accessed a web application that provided clear instructions and system descriptions. To prevent cheating and ensure authentic interaction with the RAG system, copy functionality was disabled for task descriptions, requiring users to reformulate queries in their own words. The system randomly assigned different models (LLaMA-3.3-70B-Instruct-Turbo or GPT-40-mini) to each task per user, ensuring balanced model comparison. Upon completion, participants provided feedback through the System Usability Scale¹² (SUS) questionnaire supplemented with custom questions addressing task difficulty, AI assistant

helpfulness, and system improvement suggestions.

4 Results

Retrieval Performance. The retrieval system demonstrates strong capability in addressing the core knowledge discovery challenge presented in the introduction. Table 2 presents the evaluation results, showing that the system successfully retrieves relevant NASA documentation with a recall of 0.66 and hit rate of 0.68, meaning it finds the majority of expected documents and successfully locates at least one relevant document for most queries.

Table 2. Retrieval evaluation results across 50 test questions, where each question has one or several expected answer chunks from potentially different source

Metric	Score		
Precision	0.23		
Recall	0.66		
MRR	0.64		
Hit Rate	0.68		
NDCG	0.64		

The MRR of 0.64 indicates that relevant documents consistently appear in top positions, crucial for engineers who need quick access to historical information. While the 0.23 precision score appears low, this may reflect the system's ability to find additional relevant documents beyond those manually marked in the test dataset. Across all test questions, the system retrieved 150 documents (we use top 3 retrieved documents only) with 53 identified as relevant, demonstrating its effectiveness in navigating NASA's extensive 10,000+ document collection.

Generation Capabilities. The generation evaluation demonstrates that the RAG system successfully addresses the core challenge of providing accurate, contextually grounded responses from NASA documentation. Table 3 presents the results using only queries for which at least one relevant document was retrieved (34 of 50 questions), corresponding to the 68% hit rate. Performance across all four models is remarkably consistent, underscoring the critical role of high-quality retrieval in RAG systems.

The results show strong semantic performance with BERTScore F1 values between 0.78-0.80, indicating that generated answers are semantically very close to ground truth. ROUGE-1 scores of 0.46-0.48 demonstrate strong unigram overlap, considered very good for domain-specific datasets. Most importantly, the LLM-as-a-judge metrics reveal that all models generate highly relevant and accurate answers, with factual accuracy and groundedness scores of 0.87-0.90, demonstrating that responses

¹²https://www.surveylab.com/blog/
system-usability-scale-sus/

538

539

540

541

542

543

544

545

546

547

549

550

551

552

553

554

555

557

558

559

560

561

563

564

565

566

567

568 569 575

576

579

Table 3. Generation evaluation results on the retrieved documents (34 questions).

Metric	GPT4o-mini	GPT3.5-turbo	Llama-70b	Llama-vision
Semantic Metrics				
BERTScore Precision	0.74	0.78	0.75	0.75
BERTScore Recall	0.83	0.82	0.83	0.83
BERTScore F1	0.78	0.80	0.79	0.78
Semantic Similarity	0.85	0.85	0.85	0.82
ROUGE Metrics				
ROUGE-1	0.46	0.48	0.46	0.47
ROUGE-2	0.26	0.29	0.26	0.29
ROUGE-L	0.45	0.48	0.45	0.47
LLM-as-judge Metrics				
Answer Relevance	0.86	0.85	0.87	0.84
Factual Accuracy	0.89	0.88	0.90	0.89
Groundedness	0.88	0.87	0.88	0.88
GEval Metrics				
GEval Correctness	0.79	0.75	0.79	0.74
GEval Relevance	0.98	0.97	0.97	0.96
GEval Accuracy	0.85	0.81	0.84	0.81
GEval Groundedness	0.87	0.86	0.87	0.87

are not only correct but well-supported by retrieved context.

The comparative analysis between custom LLMas-a-Judge and G-Eval metrics reveals important methodological insights for RAG evaluation. While the custom GPT-40-mini approach yields higher factual accuracy scores (0.88-0.90) compared to G-Eval accuracy scores (0.81-0.85), this difference may reflect the inherent limitations of single-prompt evaluation. The custom approach's reliance on a single, comprehensive prompt could potentially make it more susceptible to prompt-specific biases and may oversimplify complex evaluation criteria into direct numerical outputs. In contrast, G-Eval's multiprompt CoT methodology decomposes evaluation into explicit reasoning steps, possibly providing more nuanced and reliable assessments despite lower absolute scores. The remarkably high G-Eval relevance scores (0.96-0.98) compared to custom relevance scores (0.84-0.87) further suggest that structured multi-prompt approaches might offer more granular analysis capabilities, while single-prompt evaluation may conflate different evaluation dimensions.

This performance directly addresses the "knowledge gap" problem from the introduction by enabling engineers to access verified information with direct source references. The high groundedness scores (0.87-0.88) confirm that the system successfully leverages NASA documentation rather than relying on pre-trained knowledge, ensuring factual accuracy and traceability. The consistency across models suggests that when provided with high-quality retrieved context, all evaluated models are capable of producing accurate and relevant answers to tech-

nical questions, validating the RAG approach for specialized aerospace documentation.

User Evaluation Results. The user evaluation provides crucial validation that the RAG system successfully addresses the challenge of knowledge discovery from the introduction. Figure 4 shows the completion times of tasks in different models, while Figure 5 shows the number of attempts by users to complete various tasks.

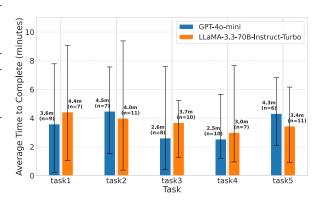


Figure 4. Comparison of GPT-40-mini and LLaMA-3.3-70B-Instruct-Turbo models showing average task completion times in minutes. Each bar represents the mean completion time, with n indicating the number of users who completed each task. Vertical black lines show the minimum and maximum completion times for each model and task combination.

Most notably, the low attempt counts reveal that users with no aerospace domain experience successfully solved relatively complex aerospace engineering scenarios in just a few tries. The first three tasks

585

586

588

589

591

592

593

595

596

597

598

599

600

601

602

603

604

605

606

607

608

615

622

623

628

630

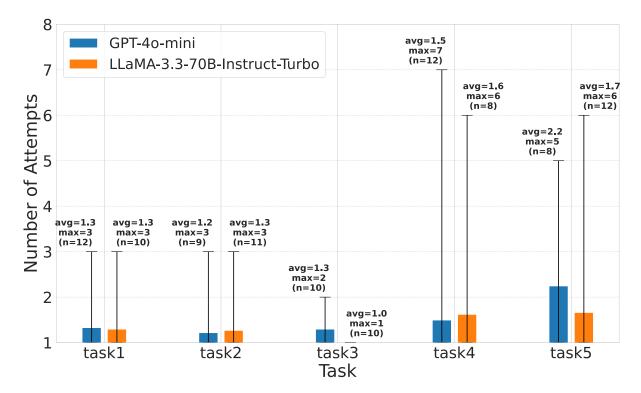


Figure 5. Comparison of GPT-4o-mini and LLaMA-3.3-70B-Instruct-Turbo models showing the average number of attempts users required to complete each task. Each bar displays the mean number of attempts, with n indicating the number of users who completed each task. Vertical black lines show the minimum and maximum number of attempts for each model and task combination.

required in average only 1.3 attempts to complete, demonstrating that the system effectively bridges the knowledge gap where valuable insights from past missions previously remained isolated.

These user evaluation results correspond closely with the generation evaluation findings, where all models performed very similarly when provided with high-quality retrieved context. This consistency validates the RAG approach: engineers can now quickly access verified information with direct source references. The evaluation with 20 participants completing realistic engineering scenarios provides strong evidence of the system's practical utility for NASA engineers.

Conclusion 5

In this study, we introduced a specialized RAG system designed to address the critical knowledge discovery challenge faced by large organizations with extensive historical documentation. Our approach processes over 10,000 NASA technical documents and lessons-learned reports, offering engineers quick access to verified information with direct source references, addressing a significant need for efficient knowledge retrieval in high-stakes aerospace envi-

The system demonstrated effectiveness across mul-

tiple evaluation dimensions, achieving a recall of 0.66 and hit rate of 0.68 across the vast document collection. All four evaluated models achieved consistent generation quality (BERTScore F1 0.78-0.80, groundedness scores 0.87-0.88) when provided with high-quality retrieved context, suggesting that retrieval quality, rather than model choice, deter- 616 mines RAG performance. The 20-participant user study validated practical utility, showing that nontechnical users can solve complex aerospace engineering problems with minimal attempts, indicating the system's potential to bridge knowledge gaps where valuable historical insights previously remained isolated.

Future work should test different chunking approaches such as recursive and semantic chunking, which were barely explored in this paper, and investigate whether LLMs have already been trained on NASA documentation data, as this may affect evaluation validity.

References

K. Jenab, S. Khoury, F. Troy, and S. 631 Moslehpour. "Cause-Consequence Analysis for 632 NASA's Space Transportation System (STS)-Solid Rocket Booster (SRB)". In: International Journal of Business and Management

696

699

701

702

710

713

716

723

724

726

727

737

- 10 (July 2015), pp. 23-23. DOI: 10.5539/ijbm. 636 v10n8p23. 637
- S. Jafari and S. A. Olah. "Development and 638 Implementation of a Retrieval-Augmented 639 Generation System: The Lookinglass". In: 640 641 2023, pp. 18-19. URL: https://hdl.handle. net/20.500.12608/66789. 642
- NASA. NASA Technical Report Server. Ac-643 cessed: 2025-01-27. URL: https://ntrs.nasa. 644 645 gov/search.
- NASA. NASA Lessons Learned. Accessed: 646 2025-01-27. URL: https://llis.nasa.gov/. 647
- P. Lewis, E. Perez, A. Piktus, F. Petroni, V. 648 Karpukhin, N. Goyal, H. Küttler, M. Lewis, 649 W.-t. Yih, T. Rocktäschel, S. Riedel, and D. 650 Kiela. Retrieval-Augmented Generation for 651 Knowledge-Intensive NLP Tasks. 2021. arXiv: 2005.11401 [cs.CL]. 653
- H. Trivedi, N. Balasubramanian, T. Khot, and A. Sabharwal. "Interleaving Retrieval with 655 Chain-of-Thought Reasoning for Knowledge-656 Intensive Multi-Step Questions". In: (2023). 657 arXiv: 2212.10509 [cs.CL]. URL: https:// 658 arxiv.org/abs/2212.10509. 659
- J. Chen, H. Lin, X. Han, and L. Sun. "Bench-660 marking large language models in retrieval-661 augmented generation". In: Proceedings of 662 the AAAI Conference on Artificial Intelli-663 gence. Vol. 38. 16. 2024, pp. 17754–17762. URL: 664 665 https://arxiv.org/abs/2309.01431.
- H. Wang, W. Huang, Y. Deng, R. Wang, Z. 666 Wang, Y. Wang, F. Mi, J. Z. Pan, and K.-F. 667 Wong. "UniMS-RAG: A Unified Multi-source 668 Retrieval-Augmented Generation for Person-669 670 alized Dialogue Systems". In: (2024). arXiv: 2401.13256 [cs.CL]. URL: https://arxiv. 671 org/abs/2401.13256. 672
- Z. Liu, W. Ping, R. Roy, P. Xu, C. Lee, M. 673 Shoeybi, and B. Catanzaro. "ChatQA: Sur-674 passing GPT-4 on Conversational QA and 675 RAG". In: 2024. arXiv: 2401.10225 [cs.CL]. 676 URL: https://arxiv.org/abs/2401.10225.
- [10]M. Saeidi, M. Bartolo, P. Lewis, S. Singh, 678 T. Rocktäschel, M. Sheldon, G. Bouchard, 679 and S. Riedel. Interpretation of Natural Lan-680 guage Rules in Conversational Machine Read-681 ing. 2018. arXiv: 1809.01494 [cs.CL]. 682
- S. Feng, H. Wan, C. Gunasekara, S. Pa-|11|683 tel, S. Joshi, and L. Lastras. "doc2dial: A 684 685 Goal-Oriented Document-Grounded Dialogue Dataset". In: Proceedings of the 2020 Confer-686 ence on Empirical Methods in Natural Lan-687 guage Processing (EMNLP). Ed. by B. Web-688 ber, T. Cohn, Y. He, and Y. Liu. Online: Association for Computational Linguistics, Nov. 690

- 2020, pp. 8118-8128. DOI: 10.18653/v1/2020. 691 emnlp-main.652.
- S. Vakulenko, S. Longpre, Z. Tu, and R. Anan- 693 tha. Question Rewriting for Conversational Question Answering. 2020. arXiv: 2004.14652 [cs.IR].
- P. Pasupat and P. Liang. Compositional Semantic Parsing on Semi-Structured Tables. 2015. arXiv: 1508.00305 [cs.CL].
- Z. Chen, S. Li, C. Smiley, Z. Ma, S. Shah, 700 and W. Y. Wang. "ConvFinQA: Exploring the Chain of Numerical Reasoning in Conversational Finance Question Answering". In: 703 Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Ed. by Y. Goldberg, Z. Kozareva, and Y. 706 Zhang. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 708 2022, pp. 6279–6292. DOI: 10.18653/v1/2022. emnlp-main.421.
- K. Nakamura, S. Levy, Y.-L. Tuan, W. 711 [15]Chen, and W. Y. Wang. HybriDialogue: An Information-Seeking Dialogue Dataset Grounded on Tabular and Textual Data. 2022. 714 arXiv: 2204.13243 [cs.CL].
- A. V. i. Martin and D. Selva. "Daphne: A Virtual Assistant for Designing Earth Observation Distributed Spacecraft Missions". In: 718 IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 13 (2020), pp. 30-48. DOI: 10.1109/JSTARS. 721 2019.2948921.
- F. M. et al. "Artificial intelligence for early design of space mis- sions in support of concurrent engineering session". In: 8th Inter- national Systems Concurrent Engineering for Space Applications Conference (2018).
- A. García-Silva, C. Berrío, J. M. Gómez-Pérez, J. A. Martínez-Heras, A. Donati, and I. Roma. "SpaceQA: Answering Questions about the Design of Space Missions and Space Craft Con- 731 cepts". In: 2022. arXiv: 2210.03422 [cs.CL]. 732
- E. F. Luca Cagliero. "Optimizing Retrieval- 733 Augmented Generation for Space Mission 734 Design via Multi-Task Learning". In: 2024, 735 pp. 11-12. URL: http://webthesis.biblio. 736 polito.it/id/eprint/33181.
- J. Wang, Y. Liang, F. Meng, Z. Sun, H. Shi, 738 |20|Z. Li, J. Xu, J. Qu, and J. Zhou. *Is Chat-* 739 GPT a Good NLG Evaluator? A Preliminary 740 Study. 2023. arXiv: 2303.04048 [cs.CL]. URL: 741 https://arxiv.org/abs/2303.04048.
- A. Muhamed. CCRS: A Zero-Shot LLM-asa-Judge Framework for Comprehensive RAG Evaluation. 2025. arXiv: 2506.20128 [cs.CL]. 745 URL: https://arxiv.org/abs/2506.20128.

809

813

814

815

819

- Y.-M. Tseng, W.-L. Chen, C.-C. Chen, and 747 H.-H. Chen. Are Expert-Level Language Mod-748 els Expert-Level Annotators? 2024. arXiv: 749 2410.03254 [cs.CL]. URL: https://arxiv. 750 org/abs/2410.03254. 751
- [23]Z. Ashktorab, M. Desmond, Q. Pan, J. M. 752 Johnson, M. S. Cooper, E. M. Daly, R. Nair, 753 T. Pedapati, H. J. Do, and W. Geyer. Align-754 ing Human and LLM Judgments: Insights 755 from EvalAssist on Task-Specific Evaluations 756 757 and AI-assisted Assessment Strategy Preferences. 2025. arXiv: 2410.00873 [cs.HC]. URL: 758 https://arxiv.org/abs/2410.00873. 759
- A. Bavaresco, R. Bernardi, L. Bertolazzi, D. [24]760 Elliott, R. Fernández, A. Gatt, E. Ghaleb, M. 761 Giulianelli, M. Hanna, A. Koller, A. F. T. Mar-762 tins, P. Mondorf, V. Neplenbroek, S. Pezzelle, 763 B. Plank, D. Schlangen, A. Suglia, A. K. 764 Surikuchi, E. Takmaz, and A. Testoni. LLMs 765 instead of Human Judges? A Large Scale Empirical Study across 20 NLP Evaluation 767 Tasks. 2025. arXiv: 2406.18403 [cs.CL]. URL: 768 https://arxiv.org/abs/2406.18403. 769
- [25]M. T. Hasan, M. Waseem, K.-K. Kemell, A. A. 770 Khan, M. Saari, and P. Abrahamsson. En-771 gineering RAG Systems for Real-World Ap-772 plications: Design, Development, and Evalua-773 tion. 2025. arXiv: 2506.20869 [cs.SE]. URL: 774 https://arxiv.org/abs/2506.20869. 775
- [26] V. Karpukhin, B. Oğuz, S. Min, P. Lewis, L. 776 Wu, S. Edunov, D. Chen, and W.-t. Yih. Dense 777 Passage Retrieval for Open-Domain Question Answering. 2020. arXiv: 2004.04906 [cs.CL]. 779
- [27]S. Gupta, R. Ranjan, and S. N. Singh. A 780 781 Comprehensive Survey of Retrieval-Augmented Generation (RAG): Evolution, Current Land-782 scape and Future Directions. 2024. arXiv: 2410. 783 12837 [cs.CL]. 784
- Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, [28]785 Y. Bi, Y. Dai, J. Sun, M. Wang, and H. 786 Wang. Retrieval-Augmented Generation for 787 Large Language Models: A Survey. 2024. arXiv: 788 2312.10997 [cs.CL]. 789
- [29] T. Zhang, V. Kishore, F. Wu, K. Q. Wein-790 berger, and Y. Artzi. BERTScore: Evaluat-791 ing Text Generation with BERT. 2020. arXiv: 792 1904.09675 [cs.CL].
- C.-Y. Lin. "ROUGE: A Package for Automatic [30]794 Evaluation of Summaries". In: Text Sum-795 marization Branches Out. Barcelona, Spain: 796 Association for Computational Linguistics, 797 July 2004, pp. 74-81. URL: https:// 798 aclanthology.org/W04-1013/.

- H. Li, Q. Dong, J. Chen, H. Su, Y. Zhou, Q. Ai, 800 Z. Ye, and Y. Liu. LLMs-as-Judges: A Comprehensive Survey on LLM-based Evaluation Methods. 2024. arXiv: 2412.05579 [cs.CL]. 803 URL: https://arxiv.org/abs/2412.05579.
- X. Xu, K. Kong, N. Liu, L. Cui, D. Wang, 805 J. Zhang, and M. Kankanhalli. An LLM can Fool Itself: A Prompt-Based Adversarial Attack. 2023. arXiv: 2310.13345 [cs.CR]. URL: 808 https://arxiv.org/abs/2310.13345.
- [33] J. Ye, Y. Wang, Y. Huang, D. Chen, Q. Zhang, 810 N. Moniz, T. Gao, W. Geyer, C. Huang, P.-Y. 811 Chen, N. V. Chawla, and X. Zhang. "Justice or Prejudice? Quantifying Biases in LLMas-a-Judge". In: (2024). arXiv: 2410.02736 [cs.CL].
- Y. Liu, D. Iter, Y. Xu, S. Wang, R. Xu, and C. 816 Zhu. G-Eval: NLG Evaluation using GPT-4 817 with Better Human Alignment. 2023. arXiv: 818 2303.16634 [cs.CL].

A Appendix

820

A.1 System Interface Design

821

The RAG system was implemented as a web-based application to facilitate user evaluation and ensure consistent interaction experiences across all participants. Figure A.1 demonstrates the dual-panel interface design that integrates task management with conversational AI interaction.

Task 3: Aircraft Noise Profile Assessment You're working on noise reduction for a new aircraft design. Your task is to understand how engine power settings affect different noise components. Research how noise profiles vary between approach and takeoff conditions to inform your design recommendations. Expected findings should include: Inlet broadband component behavior (noise distributed across many frequencies) at low power Relationship between flight velocity and airframe noise (noise from the aircraft body) • Comparative noise levels at high takeoff power Submit your findings: components are at their highest levels due to high engine settings and flight velocities. The relationship between engine power settings, flight velocity, and noise components is crucial for understanding and addressing noise reduction in aircraft design Submit Answer Skip (3 more attempts) Task Progress Task 1: Completed successfully

- Hey, I am working on noise reduction for a new aircraft design. I would like to understand how engine power settings affect different noise components. Could you look into how noise profiles vary between approach and takeoff conditions? Try to include these main points: Inlet broadband component behavior, relationship between flight velocity and airframe noise, comparative noise levels at high takeoff power.
 - Brief Answer: Engine power settings significantly impact the noise profiles of an aircraft, with
 variations between approach and takeoff conditions. At low engine power settings, such as
 during approach, inlet broadband component levels and airframe noise are relatively low
 due to lower flight velocities. In contrast, high takeoff power results in higher levels of both
 inlet and airframe noise components.
 - 2. Key Points:
 - Inlet broadband component levels are relatively low at low engine power settings (approach) and higher at high engine power settings (takeoff).
 - Airframe noise is influenced by flight velocity, with lower velocities resulting in lower airframe noise levels.
 - At high takeoff power, both inlet and airframe noise components are at their highest levels due to high engine settings and flight velocities.
 - The relationship between engine power settings, flight velocity, and noise components is crucial for understanding and addressing noise reduction in aircraft design.
 - 4. Sources:
 - 20200002461.pdf: https://ntrs.nasa.gov/api/citations/20200002461/ downloads/20200002461.pdf?attachment=true
 - 20170008825.pdf: https://ntrs.nasa.gov/api/citations/20170008825/downloads/20170008825.pdf?attachment=true
 - 20190030361.pdf: https://ntrs.nasa.gov/api/citations/20190030361/downloads/20190030361.pdf?attachment=true

Figure A.1. Complete user interface of the RAG system showing the integrated evaluation environment. The left panel displays the current task description, answer submission area, and task progress tracker, while the right panel provides the interactive chat interface where users can query the AI assistant, receive contextually grounded responses, and navigate through retrieved source documents with direct links to original NASA materials.

A.2 User Evaluation Questionnaire Details

825

827

After users completed all tasks and gained familiarity with the system, a comprehensive questionnaire was administered consisting of System Usability Scale (SUS) questions and additional questions to assess AI performance and user experience.

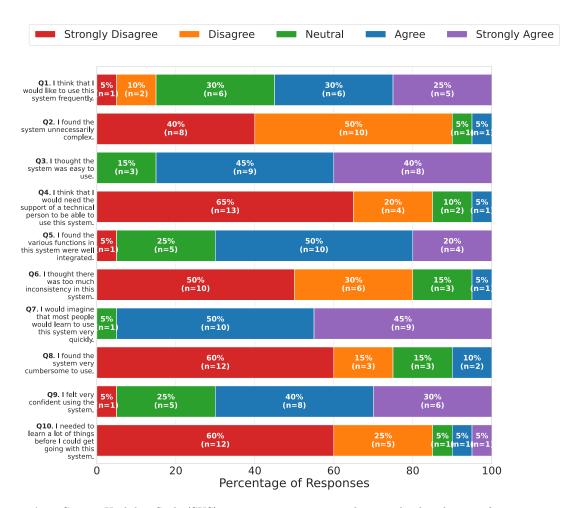


Figure A.2. System Usability Scale (SUS) questionnaire responses showing the distribution of user ratings across all ten SUS questions. Each horizontal bar represents one SUS question with response percentages and participant counts (n) for each rating level from Strongly Disagree to Strongly Agree.

User Feedback Analysis

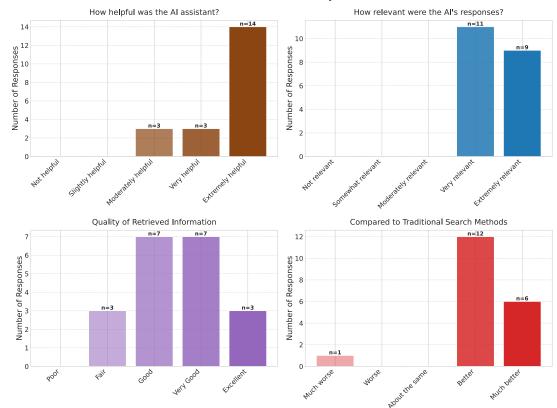


Figure A.3. User feedback analysis across four custom evaluation questions assessing AI system performance. The visualization shows response distributions for questions evaluating system accuracy, response quality, information usefulness, and overall satisfaction with percentage breakdowns and participant counts for each response category.

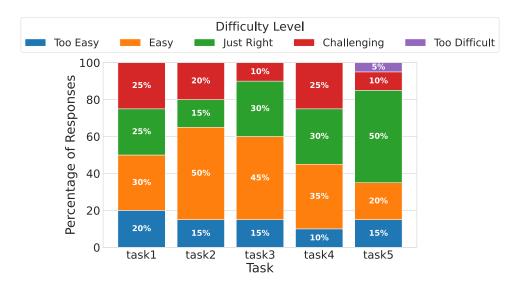


Figure A.4. Task difficulty assessment showing user-reported difficulty levels for each of the five evaluation tasks. The stacked bar chart displays the percentage distribution of difficulty ratings (Very Easy to Very Hard) with participant counts, providing insights into task complexity from the user perspective.