

# Error Audit

Error 429: Too Many Requests	Users: Students, researchers, professionals conducting extensive literature searches
Error type: System limitation - The system can't provide answers due to external API constraints (NewsAPI, arXiv rate limits)	User stakes: High
Error: Irrelevant Research Results	Users: All user types (students, researchers, professionals)
Error type: Context Limitation - System is "working as intended" but keyword extraction failed to disambiguate context - Poor assumptions about user intent	User Stakes : High
Error: Hallucinated Citations	Users: All user types (students, researchers, professionals)
Error type: Background - Severe credibility damage - Academic integrity violations if citations are used	User Stakes : High
Error: Outdated Information Presented as Current	Users: Researchers, professionals requiring current information
Error type: Background - System isn't working correctly but doesn't register an error - Data pipeline hasn't updated recently enough	User stakes: High
Error: Ambiguous Query Interpretation	Users: Users with domain-specific terminology that has multiple meanings

Error type: Context Limitation - System lacks context to disambiguate query intent	User stakes: Medium
Error: Vector Database Retrieval Failure	Users: All users during peak usage or system maintenance
Error type: System limitation - Infrastructure failure in vector database layer	User stakes: High
Error: Slow Response Time (>5 seconds)	Users: All users, especially those on slower connections
Error type: System limitation - LLM inference latency, multiple retrieval calls, or network issues	User stakes: Low

## Error Sources Analysis

### Error 1: API Rate Limit Exceeded

#### Failure state:

- Is your feature unusable as the result of multiple errors?  
If multiple APIs fail simultaneously, our system could become unusable. We currently query two API's, NewsAPI and arXiv API, to retrieve data. While a failure in one does not necessarily impact the other, if both APIs fail at the same time, the system will be unable to function.

### Error 2: Irrelevant Research Results

#### Relevance error signals:

- Is the model lacking available data or requirements for prediction accuracy?  
Yes, there are two possible reasons this could occur:
  - The model may not have enough data available to answer the question.
  - The retrieval component of the RAG model may not be correctly fetching all the relevant document chunks.

### Error 3: Outdated Information

### **Relevance error signals:**

- Is the model lacking available data or requirements for prediction accuracy?  
Yes, there are two possible reasons this could occur:
  1. The model may not have enough data available to answer the question.
  2. The retrieval component of the RAG model may not be correctly fetching all the relevant document chunks.

## **Error 4: Hallucinated Citations**

### **Input error signals:**

- Did the model improperly weigh a user action or other signal? No - the LLM produced incorrect information because the RAG system retrieved inaccurate data from the input source.

### **Relevance error signals:**

- Is the model lacking available data or requirements for prediction accuracy?  
Yes, there are two possible reasons this could occur:
  1. The model may not have enough data available to answer the question.
  2. The retrieval component of the RAG model may not be correctly fetching all the relevant document chunks.

## **Error 5: Vector Database Failure**

### **System hierarchy error:**

- Is your user connecting your product to another system?  
Yes, since our system relies on the FAISS vector database to store word embeddings, any failure in this dependency would also cause our system to fail.

## **Error 6: Ambiguous query interpretation**

### **Input Error Signals**

- Did the user anticipate the auto-correction of their input into an AI system?  
No clear expectation. User typed "Apple research" without specifying domain context. The user likely didn't realize the term has multiple meanings in different contexts (technology company vs. fruit).  
The system attempted to interpret intent without explicit user guidance, leading to mixed results.
- Did the model improperly weigh a user action or other signal?  
YES - This is a context error. The keyword extraction module treated "Apple" as equally likely to refer to the company or the fruit, without considering:

- User's browsing history (if they typically search AI/ML topics)
- Session context (previous queries in the conversation)
- Domain-specific priors (in an AI research tool, "Apple" more likely refers to Apple Inc.)
- The model failed to use available contextual signals to disambiguate.

## Error 7: Slow response time

### Relevance Error Signals

- Is the model receiving unstable or noisy data?  
Potentially contributing factor:
  - If external APIs (NewsAPI, arXiv) are slow or timing out, this cascades to user experience
  - Network latency to vector database or blob storage
  - Unstable infrastructure (database connection pools exhausted, CPU throttling)

### System Hierarchy Error

- Is your user connecting your product to another system, and it isn't clear which system is in charge?  
YES - Multiple system dependencies create latency:
  - External APIs (NewsAPI, arXiv, Google Scholar) - each can add 1-3 seconds
  - Vector database (FAISS/Pinecone) queries - 500ms-2s depending on index size and load
  - LLM inference (OpenAI API or self-hosted) - 2-4 seconds for complex generations
  - Sequential dependencies: keyword extraction → retrieval → LLM generation
- The user experiences the sum of all system latencies, but doesn't know which component is slow.
- No transparency about what's happening during the wait.

## Error Resolution Plans

Error Rationale	Solution Type	Error Resolution
<b>Error 1: API Rate Limit Exceeded</b>  User expects immediate responses but external API limits prevent data retrieval	– User control - Display clear messaging about rate limits.  – Feedback - Allow users to retry after cooldown	<b>User Path:</b>  1. User submits query 2. System detects rate limit

		<ol style="list-style-type: none"> <li>3. System processes query when API available</li> <li>4. The user receives a response with notification: "Thank you for waiting!"</li> </ol> <p><b>Opportunity for Model Improvement:</b></p> <ul style="list-style-type: none"> <li>• Use multiple API keys</li> <li>• Log rate limit patterns to predict and prevent future occurrences</li> </ul>
<p><b>Error 2: Irrelevant Research Results</b></p> <p>User receives results from wrong domain because keyword extraction lacks context disambiguation</p>	<ul style="list-style-type: none"> <li>– Feedback - Allow users to refine queries and provide relevance feedback</li> <li>– User control - Offer domain filters (AI/ML, computer vision, NLP, etc.)</li> </ul>	<p><b>User Path:</b></p> <ol style="list-style-type: none"> <li>1. User submits ambiguous query</li> <li>2. System detects potential ambiguity</li> <li>3. System returns relevant filtered results</li> <li>4. User can provide feedback and ask the correct query if results are not accurate.</li> </ol> <p><b>Opportunity for Model Improvement:</b></p> <ul style="list-style-type: none"> <li>• Implement user history and preference learning</li> <li>• Use conversation context from chat history</li> <li>• Log ambiguous queries and user</li> </ul>

		selections for training data
<p><b>Error 3: Outdated Information</b></p> <p>User expects current information but data pipeline hasn't refreshed recently enough</p>	<p>Feedback</p> <ul style="list-style-type: none"> <li>• Display data freshness timestamps</li> <li>• Automated data pipeline monitoring and alerts</li> </ul>	<p><b>User Path:</b></p> <ol style="list-style-type: none"> <li>1. User submits query about current topic</li> <li>2. System retrieves best available data</li> <li>3. User sees response with timestamp: "Based on sources up to [date]. Last updated: 2 days ago"</li> <li>4. User queries for latest data</li> <li>5. System triggers priority data refresh for that topic</li> <li>6. User gets latest data</li> </ol> <p><b>Opportunity for Model Improvement:</b></p> <ul style="list-style-type: none"> <li>• Implement continuous data ingestion (near real-time)</li> <li>• Develop topic-based refresh prioritization (hot topics updated more frequently)</li> <li>• Monitor data staleness by topic area</li> <li>• Log user requests for fresh data to identify priority topics</li> <li>• Set up automated alerts when knowledge base falls behind threshold (e.g.,</li> </ul>

		>72 hours for AI news)
<b>Error 4: Hallucinated Citations</b>  LLM generates convincing but false citations, undermining system credibility	• Other: RAG grounding enforcement, citation verification system	<b>User Path:</b> <ol style="list-style-type: none"> <li>1. User submits query</li> <li>2. System retrieves relevant documents from vector DB</li> <li>3. LLM generates response ONLY from retrieved documents</li> <li>4. Every citation includes: source publication, date and other metadata</li> <li>5. User can verify each citation</li> </ol> <b>Opportunity for Model Improvement:</b> <ul style="list-style-type: none"> <li>• Enforce strict RAG-only responses (no knowledge outside retrieved docs) through prompt engineering</li> <li>• Add citation validator that cross-references with original sources</li> <li>• Log any detected hallucinations for model fine-tuning</li> <li>• Regular audits</li> </ul>
<b>Error 5: Vector Database Retrieval Failure</b>	• Other: Redundancy, monitoring, graceful degradation	<b>User Path:</b> <ol style="list-style-type: none"> <li>1. User submits query</li> <li>2. System detects vector DB unavailability</li> </ol>

<p>Critical infrastructure component fails, making feature completely unusable</p>		<ol style="list-style-type: none"> <li>3. System sends alert to ops team</li> <li>4. User receives degraded but functional response</li> <li>5. Notification when full service restored</li> </ol> <p><b>Opportunity for Model Improvement:</b></p> <ul style="list-style-type: none"> <li>• Implement vector DB replication and failover</li> <li>• Set up health monitoring with automatic alerts (&lt;99.5% uptime triggers page)</li> <li>• Create backup retrieval systems (PostgreSQL full-text search, Elasticsearch)</li> <li>• Implement connection pooling and retry logic with exponential backoff</li> <li>• Log all failures for infrastructure improvement</li> <li>• Regular load testing to identify capacity limits</li> </ul>
<p><b>Error 6: Ambiguous Query Interpretation</b></p> <p>User uses terminology with multiple meanings, system lacks context to choose correctly</p>	<p>Feedback:</p> <ul style="list-style-type: none"> <li>• Interactive clarification</li> </ul>	<p><b>User Path:</b></p> <ol style="list-style-type: none"> <li>1. User submits "Apple research"</li> <li>2. System detects ambiguity (confidence: low)</li> </ol>



		<ol style="list-style-type: none"> <li>If the results are undesirable, user queries with clarification again</li> <li>System remembers preference for session unless user states otherwise</li> </ol> <p><b>Opportunity for Model Improvement:</b></p> <ul style="list-style-type: none"> <li>Use conversation history for context</li> <li>Implement query expansion/refinement suggestions</li> <li>Track ambiguous term resolution patterns</li> <li>Fine-tune query understanding model on user selections</li> </ul>
<p><b>Error 7: Slow Response Time</b></p> <p>Users expect sub-5-second responses but system exceeds latency target</p>	<ul style="list-style-type: none"> <li>Feedback - Progress indicators with estimated time</li> <li>Other: Performance optimization, caching</li> </ul>	<p><b>Opportunity for Model Improvement:</b></p> <ul style="list-style-type: none"> <li>Implement query result caching</li> <li>Optimize vector similarity search (approximate nearest neighbors)</li> <li>Parallelize API calls and retrieval operations</li> <li>Pre-compute embeddings during ingestion, not at query time</li> <li>Implement streaming responses (show results as they arrive)</li> </ul>

		<ul style="list-style-type: none"> <li>• Use faster embedding models for initial retrieval, detailed models for re-ranking</li> <li>• Log slow queries for performance analysis</li> </ul>
--	--	--

## 2. Quality Assurance Plan

Goal: Ensure research sources are output with timestamps	<b>Review frequency:</b> Weekly - Manual spot checks of source dates
Method:  - User feedback: "Was this information current?"  - Comparison of system responses to latest publications	
Start date: Week 1 of deployment Review / End date: Ongoing	
Goal: <3 second average response time, 99.5% uptime	<b>Review frequency:</b> Daily - System performance metrics  - Monthly - Deep dive user research interviews
Method: In-product metrics: response time, error rate, query success rate  - In-product surveys: Post-interaction "How helpful was this?"  - Quarterly diary studies with 10 power users  - A/B testing of UX improvements	
Start date: Week 2 of deployment Review / End date: Ongoing	

Goal: 100% citation accuracy - every citation must be verifiable	<b>Review frequency:</b> Weekly - Manual verification of flagged citations
Method: User reporting: "Citation not found"  - Weekly manual audit: randomly select 50 responses and verify every citation  - Customer reports of citation issues	
Start date: Week 1 of deployment Review / End date: Ongoing	
Goal: Identify and fix edge cases before they impact >1% of users	<b>Review frequency:</b>  Weekly - Review error logs and user reports  Monthly - Dedicated edge case testing session
Method: Anomaly detection in query patterns  - Error log analysis (group by error type, frequency, user impact)  - User research interviews specifically about error experiences  - Proactive testing of identified edge cases	
Start date: Week 3 of deployment Review / End date: Ongoing	

## Error Reporting Sources (Comprehensive Monitoring)

We will monitor errors through:

1. **Customer Service Reports**
  - Escalation of critical issues
  - Weekly summary of top reported issues
2. **In-Product Metrics** (Automated Dashboards)
  - Error rate
  - Response time distribution
  - Query success/failure rates
  - Usage analytics
3. **In-Product Surveys**

- Post-query satisfaction ("Was this helpful?")
- Periodic NPS surveys
- Feature-specific feedback forms
- Exit surveys for churned users

#### 4. **User Research**

- Monthly user interviews (5-10 users)
- Quarterly diary studies with power users
- Usability testing of new features
- Academic partner feedback sessions

#### 5. **System Monitoring**

- Infrastructure health (uptime, latency, errors)
- Data pipeline status
- API rate limit tracking
- Vector DB performance metrics