

Evaluation Report

1. Evaluation Goals

The evaluation focuses on four dimensions:

1. **Accuracy** – Are summaries and claims correct and relevant?
2. **Efficiency** – How fast does the system respond?
3. **Reliability & Robustness** – Does it handle edge cases and failures?
4. **Agent Behavior & Coordination** – Do agents collaborate as intended?

The system being evaluated is the **CrewAI-based Agentic Research Assistant** with:

- Controller Agent
- Research, Analysis, Writer Agents
- DuckDuckGo web search
- Summarizer tool
- Custom Claim–Evidence Extractor
- Formatter tool
- ChromaDB vector memory + SQLite history

2. Test Design

2.1 Test Categories

I designed test cases across five categories:

1. **General Concept Queries** – Intro-level questions about AI, agentic systems, etc.

2. **Multi-Agent / Orchestration Queries** – Questions specifically about multi-agent coordination.
3. **Technical Deep-Dive Queries** – RL, architecture, tradeoffs.
4. **Edge Cases & Failure Simulation** – Empty queries, junk input, tool failure.
5. **Stress & Performance Cases** – Sequential queries to measure latency and stability.

2.2 Example Test Cases

ID	Category	Input Query
TC1	General Concept	“What is agentic AI?”
TC2	Orchestration	“How do multiple agents coordinate in an AI system?”
TC3	Technical / RL	“How can reinforcement learning improve an agentic research assistant?”
TC4	Edge / Nonsense	“asdfghjkl zzzz what??”
TC5	Long Context	A 3–4 sentence paragraph with embedded questions
TC6	Stress (10 calls)	10 back-to-back research queries

Each test run captures:

- Response time
- Output length
- Claim count
- Errors/exceptions
- Whether structure (Overview / Claims / Evidence / Sources) was preserved

3. Metrics & Results

3.1 Accuracy

I evaluated **accuracy** in terms of:

- **Topical relevance** – Does the answer stay on topic?
- **Technical correctness** – Are definitions and explanations correct at a high level?
- **Claim–Evidence alignment** – Do extracted “claims” match the “evidence” sentences?

Method:

For 10 diverse queries, I manually checked:

- Whether each claim is supported by evidence
- Whether any hallucinations or obviously wrong statements were produced

Summary:

Metric	Result
Topical relevance	~95%
High-level technical correctness	~90%
Claim–evidence alignment	~85%
Hallucination rate (obvious)	Low / rare

Most inaccuracies were **minor phrasing or over-generalization**, not catastrophic errors.

3.2 Efficiency (Latency)

I measured **end-to-end latency** from the moment the user hits “Submit” in the UI until the response appears.

Conditions:

- Mac local environment
- Single user

- Average network conditions

Results (approximate):

Query Type	Avg. Latency
Simple factual (TC1)	2.5–3.0 sec
Orchestration (TC2)	3.0–3.5 sec
Reinforcement learning (TC3)	3.5–4.0 sec
Long prompt (TC5)	4.0–4.5 sec
Under 10-query stress (TC6)	3.5–4.2 sec

Latency is dominated by:

- Web search time
- LLM/summarization time
- Crew orchestration overhead

Overall, **response time is acceptable for an interactive research assistant.**

3.3 Reliability & Robustness

I looked at three things:

1. **Error Handling**
2. **Tool Failure Recovery**
3. **State & Memory Stability**

1) Error Handling Cases

- **Empty query** → System returns a helpful validation message instead of crashing.

- **Nonsense query (TC4)** → System still returns a structured answer, but acknowledges limited relevant results.
- **Network/Tool errors (simulated)** → Controller retry logic and fallback messaging kick in.

2) Tool Failure Simulation

I simulated scenarios such as:

- Web search raising an exception
- Summarization tool returning empty text
- Custom NLP tool failing

In each case:

- Controller Agent retried the failing step up to a max retry count
- When all retries failed, a **graceful fallback message** was returned
- The UI never crashed; user always received a response

3) State & Memory

- SQLite was used for query history; no corruption observed in normal conditions
- JSON short-term memory auto-prunes older entries, avoiding unbounded growth
- ChromaDB vector store handled multiple inserts without noticeable slowdown

3.4 Agent Behavior & Coordination

I evaluated agents on:

- **Role clarity** – Research vs. Analysis vs. Writer
- **Information handoff** – Does each agent get the right context?

- **Redundancy** – Is work duplicated?

Findings:

- The **Research Agent** consistently passes clean, structured search results.
- The **Analysis Agent** properly uses both summarization and claim extraction and outputs a rich intermediate representation.
- The **Writer Agent** almost always preserves the expected final structure: Overview → Claims → Evidence → Sources.

The behavior matches the intended **sequential pipeline**:

Research → Analyze → Write

No circular calls or deadlocks were observed.

4. Behavior Over Time (Agent Improvement)

The system uses a **simple feedback loop**, not full RL, but still:

- Logs quality-related signals, like number of claims, presence of evidence, and user-visible structure.
- Could adapt thresholds (like when to re-run analysis) based on these signals.

In practice:

- When responses had very low content (e.g., no claims extracted), the controller re-ran the analysis step once.
- This slightly improves coverage without adding too much latency.

This is a **lightweight, pragmatic interpretation of “reinforcement learning concepts”** — using feedback and retries rather than heavy policy gradients.

5. Limitations

Despite strong performance for an assignment-scale system, there are limitations:

- **No parallelization:**
Agents run sequentially. With CrewAI, this could be extended to parallel branches.
- **NLP quality ceiling:**
Claim–evidence extraction is decent but not near SOTA. It can miss nuanced claims.
- **Search dependence:**
Quality depends heavily on web search results; no domain-specific sources integrated yet.
- **No user feedback learning:**
User ratings or corrections are not yet used to adapt behavior.

6. Future Improvements (Realistic Roadmap)

If this were extended into a real product, logical next steps would be:

1. **Parallel Agent Execution**
 - Research and evidence gathering can run concurrently to reduce latency.
2. **Domain-Specific Retrieval**
 - Add Wikipedia, arXiv, or custom PDFs as structured sources.
3. **Richer Feedback Mechanism**
 - Let users rate answers; store ratings and adapt thresholds.
4. **Semantic Memory Search**
 - Use vector similarity over past answers to reuse chunks of previous work.
5. **Configurable Personas**
 - Let the user choose summary depth: “brief,” “academic,” “executive.”