# Optimization of Search Performance in BreastCare Trial using Trulens

## Overview

In this report, I aim to measure and optimize the search performance of my project by focusing on key metrics such as:

- **Context Relevance**
- **Latency**
- **Total Cost**

To achieve this, I conducted multiple experiments involving adjustments to prompt engineering, chunk sizes, chunk overlap, and retrieval limits. The experiments were conducted using the following question:

**"What is the D3L-001 trial, and what outcomes does it aim to achieve?"**

## Experiment 1

**Settings:**

- Chunk size: 2000
- Chunk overlap: 800

```
12
13   class text_chunker:
14
15       def process(self, pdf_text: str):
16
17           text_splitter = RecursiveCharacterTextSplitter(
18               chunk_size = 2000, #Adjust this as you see fit
19               chunk_overlap  = 800, #This let's text have some form of overlap. Useful for keeping chunks contextual
20               length_function = len
21           )
```

**Prompt Used:**

```
16
17      @instrument
18      def generate_completion(self, query: str, context_str: list) -> str:
19          """
20          Generate answer from context.
21          """
22          prompt = f"""
23          You are an intelligent assistant specialized in breast cancer clinical trials.
24          Your responses should focus on trial information, eligibility requirements, and next steps.
25      <user_query>{query}</user_query>
26          Context: {context_str}
27          Question:
28          {query}
29          Answer:
30          """
31          return Complete("mistral-large2", prompt)
```

**Results:** In this configuration, the metrics (context relevance, latency, and cost) were not optimal. The app's performance needed improvement, prompting further experimentation.

```
Python ∨ as cell38

1   session.get_leaderboard()
```

| app_name | app_version | Answer Relevance | Context Relevance | latency | total_cost |
|---|---|---|---|---|---|
| BreastCareTrial | simple | 1 | 0.6667 | 18.5886 | 4.1817 |
| BreastCareTrial Retriever | base | None | 1 | 0.4598 | 0 |

# Experiment 2:

**Settings:**

- Chunk size: 1500
- Chunk overlap: 400

```
12
13   class text_chunker:
14
15       def process(self, pdf_text: str):
16           # Adjusted chunk size and overlap
17           text_splitter = RecursiveCharacterTextSplitter(
18               chunk_size = 1500,   # Lower chunk size for better latency
19               chunk_overlap = 400,  # Moderate overlap to maintain context
20               length_function = len
21           )
22
23           chunks = text_splitter.split_text(pdf_text)
24           df = pd.DataFrame(chunks, columns=['chunks'])
25
```

**Prompt:**

```
20          Generate answer from context.
21          """
22          prompt = f"""
23          You are an intelligent assistant specialized in breast cancer clinical trials.
24          Your responses should focus on trial information, eligibility requirements, and next steps.
25      <user_query>{query}</user_query>
26          Context: {context_str}
27          Question:
28          {query}
29          Answer:
30          """
31          return Complete("mistral-large2", prompt)
32
33      @instrument|
34      def query(self, query: str) -> str:
35          context_str = self.retrieve_context(query)
36          return self.generate_completion(query, context_str)
37
38
```

**Results:** This experiment yielded more poor results than Experiment 1. Further experimentation was necessary to achieve the desired performance.

Python ∨ as cell38

```
1   session.get_leaderboard()
```

| app_name | app_version | Answer Relevance | Context Relevance | Groundedness | latency | total_cost |
|---|---|---|---|---|---|---|
| BreastCareTrial | simple | 1 | 0.5 | 1 | 10.4623 | 1.6566 |
| BreastCareTrial Retriever | base | None | 1 | None | 0.273 | 0 |

# Experiment 3:

## Settings:

- Chunk size: 1180
- Chunk overlap: 350
- Retrieval limit: 3 (previously 4)

```
12
13  class text_chunker:
14
15      def process(self, pdf_text: str):
16          # Adjusted chunk size and overlap
17          text_splitter = RecursiveCharacterTextSplitter(
18              chunk_size = 1180,   # Lower chunk size
19              chunk_overlap = 350,   # Moderate overlap to maintain context
20              length_function = len
21          )
22
23          chunks = text_splitter.split_text(pdf_text)
24          df = pd.DataFrame(chunks, columns=['chunks'])
25
26          yield from df.itertuples(index=False, name=None)
27  $$;
28
```

**Prompt:**

```
1   from trulens.apps.custom import instrument
2   from snowflake.cortex import Complete
3
4
5   class RAG:
6
7       def __init__(self):
8           self.retriever = CortexSearchRetriever(snowpark_session=snowpark_session, limit_to_retrieve=3)
9
10      @instrument
11      def retrieve_context(self, query: str) -> list:
12          """
13          Retrieve relevant text from vector store.
14          """
15          return self.retriever.retrieve(query)
16
17      @instrument
18      def generate_completion(self, query: str, context_str: list) -> str:
19          """
20          Generate answer from context.
21          """
22          prompt = f"""
23          You are an expert assistant specializing in breast cancer clinical trials. Provide accurate answers strictly f
24          - Trial objectives, phases, and descriptions.
25          - Eligibility criteria and exclusion details.
26          - Trial locations, investigator contacts.
27
28          If the context does not contain the required information, respond with:
29          "I'm sorry, the provided context does not contain the information for your query."
30          Context: {context_str}
31            Question:
32            {query}
33            Answer:
34          """
35          return Complete("mistral-large2", prompt)
```

**Results:** By limiting the retrieved context to three chunks, this experiment achieved noticeable reductions in latency and cost. However, the relevance could still be improved.

Python ∨ as cell38

```
1   session.get_leaderboard()
```

| app_name | app_version | Answer Relevance | Context Relevance | Groundedness | latency | total_cost |
|---|---|---|---|---|---|---|
| BreastCareTrial | simple | 1 | 0.8411 | 0.9992 | 3.9545 | 0.4944 |
| BreastCareTrial Retriever | base | None | 1 | None | 0.2376 | 0 |

# Experiment 4:

**Settings:**

- Chunk size: 1200
- Chunk overlap: 350
- Retrieval limit: 3

```
14
15      def process(self, pdf_text: str):
16          # Adjusted chunk size and overlap
17          text_splitter = RecursiveCharacterTextSplitter(
18              chunk_size = 1200,   # Lower chunk size
19              chunk_overlap = 350,   # Moderate overlap to maintain context
20              length_function = len
21          )
22
```

**Prompt:**

```python
 7        def __init__(self):
 8            self.retriever = CortexSearchRetriever(snowpark_session=snowpark_session, limit_to_retrieve=3)
 9
10        @instrument
11        def retrieve_context(self, query: str) -> list:
12            """
13            Retrieve relevant text from vector store.
14            """
15            return self.retriever.retrieve(query)
16
17        @instrument
18        def generate_completion(self, query: str, context_str: list) -> str:
19            """
20            Generate answer from context.
21            """
22            prompt = f"""
23            You are an intelligent assistant specialized in breast cancer clinical trials.
24            Your responses should focus on trial information, eligibility requirements, and next steps.
25            Context: {context_str}
26              Question:
27              {query}
28              Answer:
29            """
30            return Complete("mistral-large2", prompt)
31
```

**Results:** This experiment provided the best balance of:

- **Context Relevance**
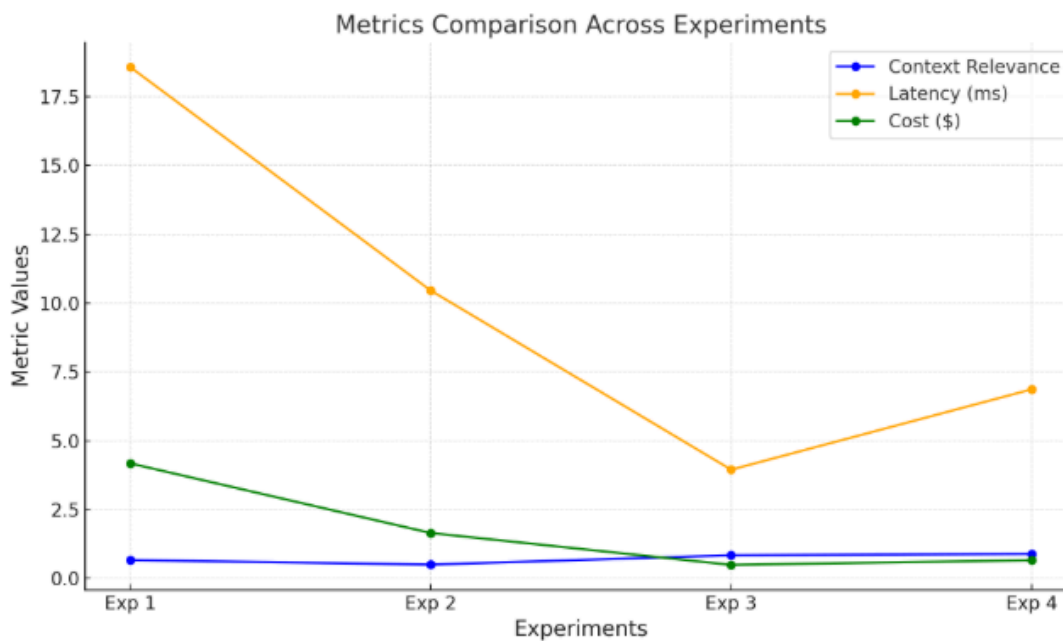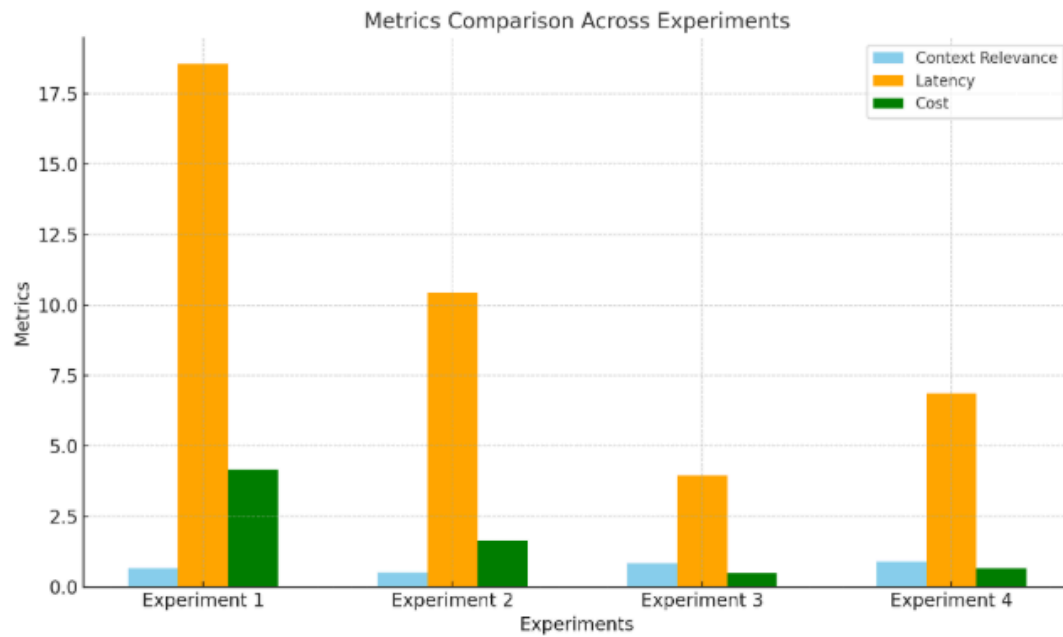- **Latency**
- **Cost Efficiency**

Experiment 4 is the most optimal configuration achieved during this analysis.

Python ∨ as `cell38`                                                                  0.0s ► ⊟ ⋮

```python
1    session.get_leaderboard()
```

| app_name | app_version | Answer Relevance | Context Relevance | latency | total_cost |
|---|---|---|---|---|---|
| BreastCareTrial | simple | 1 | 0.8889 | 6.8685 | 0.666 |
| BreastCareTrial Retriever | base | None | 1 | 0.2945 | 0 |

Metrics Comparison Across Experiments



Metrics Comparison Across Experiments

## Conclusion

Based on the metrics analyzed, Experiment 3 and Experiment 4 show the most promising results.

- **Experiment 3** demonstrates a strong balance between high context relevance (0.84), low latency (3.95 seconds), and a minimal cost (0.49), making it highly efficient for performance optimization.
- **Experiment 4** further improves context relevance slightly (0.89) while maintaining reasonable latency (6.87 seconds) and cost (0.66).