# Strategy to Handle Large JSON Files:

1. **Read Line by Line**: Instead of loading the entire JSON file into memory, the code reads the file line by line. This approach ensures that only a single line of the file is loaded into memory at a time, significantly reducing memory usage.
2. **Incremental Processing**: Each line of the JSON file is processed individually. After reading a line, it is immediately parsed as JSON and inserted into the MongoDB collection. This allows for incremental processing of the data without the need to store the entire dataset in memory.
3. **Batch Insertion**: To optimize database performance and minimize overhead, the data insertion is performed in batches. This means that multiple documents are inserted into the MongoDB collection at once, rather than inserting them one by one. Batch insertion reduces the number of database transactions and can significantly improve overall throughput.
4. **Adjustable Batch Size**: The code allows for the adjustment of the batch size based on the available memory and system resources. By experimenting with different batch sizes, you can find the optimal value that balances memory usage and insertion speed.
5. **Robust Error Handling**: Error handling mechanisms are implemented to gracefully handle any potential issues during data loading. For example, connection errors to the MongoDB server are caught and appropriate error messages are displayed to the user.
6. **Feedback and Progress Reporting**: Progress messages are printed during the data loading process to provide feedback to the user. This includes notifications about the number of messages inserted and the current progress of the operation. This helps users monitor the progress of the data loading process and ensures transparency.

By employing these strategies, the code effectively handles large JSON files and loads data into MongoDB efficiently, while minimizing memory usage and ensuring robustness against potential errors.

## Effect of Indexing on Query Performance

---

## Introduction:

In this report, we analyze the impact of indexing on the runtime of four different queries executed on a MongoDB database. The queries involve searching for messages containing specific text, finding the sender with the most messages, counting the number of messages where the sender's credit is 0, and updating the credits for senders with a credit less than 100.

---

## 1. No Index:

### Query 1:

- **Description:** Searches for messages containing the text 'ant' in their body.
- **Runtime:** 2.002 seconds

- **Explanation:** Without an index, MongoDB performs a collection scan to find documents matching the query criteria. This involves examining every document in the collection, resulting in longer runtime, especially with a large number of documents.

## Query 2:

- **Description:** Finds the sender with the most messages.
- **Runtime:** 0.952 seconds
- **Explanation:** Similar to Query 1, Query 2 also requires scanning the entire collection to compute the sender with the most messages, resulting in a moderate runtime.

## Query 3:

- **Description:** Counts the number of messages where the sender's credit is 0.
- **Runtime:** 0.589 seconds
- **Explanation:** This query involves filtering documents based on the credit field in the senders collection. Without an index on the credit field, MongoDB performs a collection scan to identify documents matching the criteria, resulting in a relatively shorter runtime compared to Queries 1 and 2.

## Query 4:

- **Description:** Doubles the credits for senders with a credit less than 100.
- **Runtime:** 0.013 seconds
- **Explanation:** As this query involves updating documents based on a condition, it tends to have a shorter runtime compared to read-heavy queries. Without an index, MongoDB still needs to scan the collection to identify documents matching the condition, leading to a slightly longer runtime.

---

## 2. After Indexing:

## Query 1:

- **Runtime:** 1.883 seconds
- **Explanation:** Despite indexing, the runtime of Query 1 remains relatively similar. This suggests that the query primarily benefits from the indexing when executing more complex queries.

## Query 2:

- **Runtime:** 1.055 seconds
- **Explanation:** Similarly, Query 2's runtime remains comparable even after indexing. This indicates that the nature of the query may not be significantly affected by indexing.

## Query 3:

- **Runtime:** 0.021 seconds
- **Explanation:** With indexing on the credit field, Query 3's runtime drastically decreases. This is because indexing allows MongoDB to efficiently locate documents based on the indexed field, resulting in a significant performance improvement.

**Query 4:**

- **Runtime:** 0.010 seconds
- **Explanation:** Query 4's runtime also slightly improves after indexing, although the impact is not as pronounced as in Query 3. This suggests that while indexing may offer some performance benefits for update operations, the improvement may not be as significant as for read operations.

---

**Conclusion:**

In summary, indexing plays a crucial role in improving query performance, particularly for read-heavy queries involving filtering or sorting based on indexed fields. While indexing may not significantly impact the runtime of all queries, it can lead to substantial improvements in scenarios where queries heavily rely on indexed fields for document retrieval. Therefore, careful consideration and optimization of indexes are essential for achieving optimal performance in MongoDB databases.

**Analysis of Normalized vs Embedded Database Models**

---

**Introduction:**

In this report, we compare the performance of normalized and embedded database models using MongoDB for a messaging system. The normalized model stores data across multiple collections, while the embedded model embeds related data within a single collection.

---

**1. Normalized Database Model:**

**Query 1:**

- **Runtime:** 4.221 seconds
- **Explanation:** Querying messages with specific text involves scanning the entire messages collection. Since the text field is indexed, the query runtime is relatively high due to the large volume of documents to scan.

**Query 2:**

- **Runtime:** 1.288 seconds

- **Explanation:** Finding the sender with the most messages requires scanning the messages collection and grouping by sender. Although the sender field is indexed, the aggregation operation still incurs overhead due to document retrieval and grouping.

**Query 3:**

- **Runtime:** 0.051 seconds
- **Explanation:** Counting messages with sender's credit as 0 involves a simple lookup on the senders collection. Since the credit field is indexed, the query completes quickly.

**Query 4:**

- **Runtime:** 0.055 seconds
- **Explanation:** Updating credits for senders with credit less than 100 involves scanning the senders collection and updating documents. Although indexing helps, the update operation still incurs overhead due to document modification.

**2. Embedded Database Model:**

**Query 1:**

- **Runtime:** 2.714 seconds
- **Explanation:** Searching for messages with specific text involves scanning the entire messages collection. Despite embedding sender information, the runtime is similar to the normalized model due to document scanning.

**Query 2:**

- **Runtime:** 1.333 seconds
- **Explanation:** Finding the sender with the most messages involves scanning the entire messages collection and extracting sender information. Despite embedding sender information, the runtime is similar to the normalized model.

**Query 3:**

- **Runtime:** 0.825 seconds
- **Explanation:** Counting messages with sender's credit as 0 involves scanning the entire messages collection and filtering based on embedded sender information. Despite embedding sender information, the runtime is similar to the normalized model.

**Query 4:**

- **Runtime:** 0.002 seconds
- **Explanation:** Checking for senders with credit less than 100 involves scanning the entire messages collection and filtering based on embedded sender information. Despite embedding sender information, the runtime is similar to the normalized model.

**Conclusion:**

1. **Performance Difference:** The performance between the normalized and embedded models varies across different queries. Queries involving simple lookups or filtering based on indexed fields tend to perform similarly in both models. However, queries involving scanning entire collections or complex aggregation operations show minimal performance difference between the two models.
2. **Better Choice:** The choice between normalized and embedded models depends on the specific use case and query patterns. If the application requires complex queries involving multiple collections or relationships, a normalized model may offer better flexibility and maintainability. However, if the application frequently accesses related data together and requires high read performance, an embedded model may be a better choice due to reduced query complexity and fewer database operations.