# Report on the Article Retrieval System

## System Design

The system is designed to provide answers to user queries based on a given dataset using a combination of vector-based semantic search and keyword-based search. It employs a custom retriever that integrates both methods for more accurate and comprehensive retrieval.

The architecture includes the following components:
1. **Data Loading**: The system loads data from a provided CSV file (`medium.csv`) using a `SimpleDirectoryReader`.
2. **Node Extraction**: Nodes representing documents are extracted from the loaded data using a node parser.
3. **Storage Context Initialization**: A storage context is initialized to store the extracted nodes.
4. **Retrieval Mechanisms**:
     - Vector Index: A vector index is created to enable efficient retrieval of documents based on vector representations.
     - Keyword Index: A keyword index is created to facilitate retrieval of documents based on keyword matches.
     - Custom Retriever: This component combines both the vector and keyword retrievers to perform either an AND or OR search mode based on user preference.
5. **Quantization configuration**: Used to reduce the memory footprint and increase inference speed of neural networks.
6. **The HuggingFaceLLM with specified parameters**: "Llama-2-7b-chat-hf" for natural language processing tasks.
7. **Global service context:** This component allows the query engine to use the HuggingFaceLLM for more human-like responses.
8. **Query Engine:** The query engine processes user queries using the custom retriever and retrieves relevant documents.

## Technologies Selected

The system leverages several technologies and libraries in Python:
1. **Google Colab Pro:** Used as an environment for code development.
2. **Python:** version 3.10.
3. **LLama Index:** Used for managing data storage, retrieval, and query processing.
4. **Transformers Library**: Utilized for handling natural language processing tasks, including tokenization and model inference.
5. **PyTorch:** A deep learning framework used for model deployment and inference.
6. **Pandas:** Used for data manipulation and analysis, particularly for loading and preprocessing the dataset.
7. **NumPy:** Used internally for numerical computations and array manipulations.
8. **Accelerate:** Used for accelerating PyTorch computations, especially on GPU devices.
9. **Hugging Face Transformers:** Utilized for model deployment, specifically for the "Llama-2-7b-chat-hf" model.
10. **BitAndBytes**: Employed for model quantization to reduce memory footprint and increase inference speed.

## Challenges Encountered

While the system appears well-designed and implemented, several challenges have been encountered during its development and deployment:

1. **Data Preprocessing:** Ensuring the dataset (`medium.csv`) is properly formatted and preprocessed to extract meaningful information and avoid noise.
2. **Model Integration:** Integrating the "Llama-2-7b-chat-hf" model into the system and configuring it for efficient inference with custom prompts and quantization.
3. **Retrieval Optimization:** Optimizing the retrieval mechanisms, including vector indexing and keyword matching, to improve search efficiency and accuracy.

## Potential Areas for Future Development

Several areas can be explored for enhancing the capabilities and functionalities of the system:

1. **More Advanced Retrieval Techniques:** Experimenting with more advanced retrieval techniques, such as deep learning-based embeddings or graph-based methods, to further improve search accuracy.
2. **User Interface:** Developing a user-friendly interface to interact with the system, allowing users to input queries easily and visualize search results.
3. **Feedback Mechanism:** Implementing a feedback mechanism where users can provide feedback on retrieved answers, which can be used to continuously improve the system's performance.
4. **Multi-language Support:** Extending the system to support multiple languages by integrating language-specific models and datasets.
5. **Scalability Enhancements:** Optimizing the system architecture for scalability to handle larger datasets and increased user traffic efficiently.
6. **Error Handling:** Enhancing error handling mechanisms and ensuring the system is robust to handle unexpected inputs or errors gracefully.
7. **Integration with External Systems:** Integrating the system with external databases, knowledge bases, or APIs to enrich the dataset and provide more comprehensive answers to user queries.

By addressing these areas for future development, the system can evolve into a more robust, reliable, and feature-rich solution for information retrieval tasks.