# Retrieval-Augmented Generation System for EPL Match Queries

Ahmed Hani (ID: 7387)

Ahmed Younis (ID: 7745)

Mohamed Hazem (ID: 7729)

May 2025

# Contents

# 1    Introduction

The English Premier League (EPL) is one of the most popular and data-rich sports leagues globally, generating vast amounts of match data, including scores, player statistics, and performance metrics. This wealth of data offers significant opportunities for enhancing fan engagement, supporting sports analytics, and enabling data-driven decision-making for teams and analysts. However, accessing specific information from this data, such as match results or performance details, often requires manual navigation of large datasets or reliance on predefined dashboards, which can be inefficient for both casual fans and professional analysts.

Our project addresses this challenge by developing a Retrieval-Augmented Generation (RAG) system that allows users to query EPL match data using natural language. The system retrieves relevant match summaries from a dataset spanning the 2019–2023 seasons and generates accurate, context-aware responses. This approach bridges the gap between complex datasets and user-friendly interfaces, enabling seamless access to insights for diverse audiences, from fans seeking match scores to analysts exploring performance trends.

The significance of this problem lies in its potential to democratize access to sports data. Traditional methods for querying sports data often require technical expertise or familiarity with specific tools, limiting their accessibility. By leveraging advanced natural language processing (NLP) techniques, including vector search and large language models (LLMs), our RAG system provides an intuitive interface that interprets user queries and delivers precise, human-readable answers. For example, a user can ask, "What was the score of Manchester City vs. Arsenal in 2022?" and receive a concise response based on relevant match data.

This project builds on recent advancements in RAG architectures, which combine information retrieval with text generation to provide contextually grounded answers. By integrating semantic search with a fine-tuned language model, the system ensures both relevance and accuracy in its responses. The following sections describe the dataset, tools, experimental methodology, results, and future directions for this project, demonstrating its potential as a scalable framework for sports analytics.

# 2    Dataset

The dataset used in this project comprises comprehensive match data from the English Premier League for the 2019–2023 seasons, sourced from publicly available sports data repositories. It contains 23 columns capturing detailed match-specific information for all 20 EPL teams across multiple seasons. The dataset is structured to include both numerical and categorical variables, making it suitable for both statistical analysis and natural language processing tasks. Key columns include:

- **date**: The date of the match (e.g., 2022-08-07).

- **team**: The name of the home team (e.g., Manchester United).

- **opponent**: The name of the opposing team (e.g., Liverpool).

- **gf**: Goals scored by the home team.

- **ga**: Goals conceded by the home team.

- **venue**: Indicates whether the match was played at home or away.

- **result**: The match outcome (W for Win, D for Draw, L for Loss).

- **xg**, **xga**: Expected goals for and against, reflecting the quality of scoring opportunities.

- **poss**: Possession percentage during the match.

- **sh**, **sot**: Total shots and shots on target.

- **attendance**: Number of spectators, with missing values for 2020 matches due to COVID-19 restrictions.

## 2.1 Exploratory Data Analysis

Exploratory data analysis (EDA) was conducted to understand the dataset's structure and quality. The dataset includes over 3,000 match records, covering 38 match weeks per season for each team. Key findings from the EDA include:

- **Data Consistency**: The dataset maintains consistent column definitions across seasons, facilitating longitudinal analysis.

- **Team Representation**: Top-performing teams like Manchester City, Liverpool, and Arsenal appear frequently, with Manchester City averaging 2.1 expected goals (`xg`) per match, compared to lower-performing teams like Norwich City (mean `xg` of 0.9).

- **Missing Values**: Missing data was minimal, primarily in the `attendance` column for the 2020–2021 season due to pandemic-related restrictions. These were handled by excluding attendance from match summaries.

- **Statistical Insights**: The average goals scored (`gf`) and conceded (`ga`) across all matches were 1.5 and 1.3, respectively. Possession (`poss`) showed a slight correlation with match outcomes, with winning teams averaging 55% possession.

## 2.2 Preprocessing

To prepare the dataset for the RAG system, we filtered it to include only home matches, reducing the dataset size by approximately 50% to focus on a consistent perspective (home team performance). Each match was transformed into a concise summary, combining key fields (e.g., date, teams, score, result) into a single text string suitable for embedding and retrieval. For example, a summary might read: "On August 7, 2022, Manchester City defeated West Ham 2-0 at home." This preprocessing step ensured that the RAG system could efficiently retrieve and process relevant match information.

# 3 Tools

The RAG system was built using a suite of open-source tools and libraries, each selected for its efficiency, compatibility, and ability to support specific components of the pipeline. Below is a detailed overview of the tools used:

- **Pandas**: A Python library for data manipulation, used to load the CSV dataset, filter home matches, handle missing values, and generate match summaries. Its DataFrame structure enabled efficient column-wise operations and data cleaning.

- **NumPy**: Supported numerical computations, including embedding normalization and statistical analysis of match data. NumPy's array operations were critical for processing embeddings generated by the SentenceTransformers model.

- **FAISS**: A library developed by Facebook AI for efficient similarity search and clustering of dense vectors. FAISS was used to create a vector index of match summary embeddings, enabling fast retrieval of the top-$k$ relevant summaries for a given query.

- **SentenceTransformers (all-MiniLM-L6-v2)**: A lightweight transformer model that generates 384-dimensional dense vector embeddings for text. It was used to encode both match summaries and user queries, enabling semantic similarity search.

- **Transformers (Hugging Face)**: Provided the `falcon-7b-instruct` language model and tokenizer for text generation. The `pipeline` API simplified inference, allowing the system to generate human-readable responses based on retrieved context.

- **PyTorch**: Served as the backend for tensor operations and GPU acceleration, supporting both the SentenceTransformers model and the `falcon-7b-instruct` model. PyTorch's dynamic computation graph facilitated rapid experimentation.

These tools were integrated into a cohesive pipeline, with Pandas and NumPy handling data preprocessing, FAISS and SentenceTransformers enabling retrieval, and Transformers and PyTorch powering the generation step. The combination ensured scalability and performance, even with large datasets and complex queries.

# 4 Experiments

The RAG system combines semantic search and text generation to process natural language queries. The pipeline consists of four main steps: 1. Preprocessing the dataset to create concise match summaries. 2. Generating embeddings for summaries using `all-MiniLM-L6-v2` and indexing them with FAISS. 3. Encoding user queries and retrieving the top-$k$ relevant summaries based on cosine similarity. 4. Passing the retrieved summaries as context to the `falcon-7b-instruct` model to generate a final answer.

## 4.1 Trials

Multiple trials were conducted to optimize the system's performance, exploring different retrieval strategies, model configurations, and prompt designs. Below are the detailed trials, including their methodologies, outcomes, and associated figures.

- **Trial 1: Baseline RAG with Default Parameters** The baseline system used `all-MiniLM-L6-v2` for embedding generation and `falcon-7b-instruct` with default parameters (`max_new_tokens=300`, `temperature=0.7`). Queries like "Who won Manchester City vs. Arsenal in 2022?" were tested. The system retrieved relevant match summaries but occasionally included irrelevant results due to overly broad semantic matches. For example, a query about a specific match sometimes retrieved summaries from unrelated matches involving the same teams.



Figure 1: Test run for Trial 1 showing match score query results.

- **Trial 2: Adjusted Retrieval (k=5)** To improve context coverage, we increased the number of retrieved summaries to $k = 5$. This trial aimed to capture more relevant matches for complex queries. However, the increased context introduced noise, as irrelevant summaries were sometimes included, reducing precision for specific queries like "Score of Tottenham vs. Chelsea in January 2023." The system struggled to prioritize the most relevant match, leading to answers that included extraneous details.

- **Trial 3: Fine-tuned Prompt** To address the issue of irrelevant responses, we modified the prompt to emphasize factual accuracy: "Provide a concise answer based only on the provided context. Avoid speculation and focus on the query." This improved response relevance, with 85% of queries returning accurate match details. However, some responses were overly brief, lacking additional context that users might expect (e.g., key player performances). Figure 2 shows a test run for a player-related query, illustrating the system's ability to handle such inputs.



Figure 2: Test run for Trial 3 showing player query results.

- **Trial 4: Alternative Embedding Model (distilbert-base-nli-stsb-mean-tokens)** We experimented with `distilbert-base-nli-stsb-mean-tokens` to improve semantic matching. This model produces 768-dimensional embeddings, compared to the 384-dimensional embeddings of `all-MiniLM-L6-v2`. While it slightly improved retrieval accuracy for some queries, it increased memory usage and retrieval time by 30%, making it less practical for real-time applications. We reverted to `all-MiniLM-L6-v2` for its balance of speed and accuracy.

- **Trial 5: Query Preprocessing** To handle ambiguous queries (e.g., "Who played well in 2022?"), we introduced a query preprocessing step to identify key entities (e.g., teams, dates) using a named entity recognition (NER) model from spaCy. This improved retrieval precision by filtering summaries based on identified entities before embedding. However, the NER model occasionally misclassified terms, leading to incomplete context for some queries.

- **Failed Trial: Using Smaller LLM (gpt2)** To reduce computational requirements, we tested `gpt2` as an alternative to `falcon-7b-instruct`. The smaller model struggled to generate coherent responses, often ignoring the provided context or producing factually incorrect answers. For example, when asked about a specific match, `gpt2` generated generic responses like "Manchester City played well." This trial was abandoned due to its 40% accuracy rate.

## 4.2 Results of Trials

The trials provided valuable insights into the system's performance:

- **Trial 1**: Achieved 80% accuracy but included irrelevant matches in 20% of cases, as shown in Figure 1.

- **Trial 2**: Accuracy dropped to 70% due to noisy context from retrieving $k = 5$ summaries.

- **Trial 3**: Improved accuracy to 85% with a fine-tuned prompt, though some responses lacked depth (see Figure 2).

- **Trial 4**: No significant accuracy improvement; processing time increased by 30%.

- **Trial 5**: Query preprocessing improved precision for specific queries but introduced errors for ambiguous ones.

- **Failed Trial**: `gpt2` achieved only 40% accuracy, deemed unsuitable.

## 5 Results

The final RAG system configuration used `all-MiniLM-L6-v2` for embeddings, FAISS with $k = 3$ for retrieval, and `falcon-7b-instruct` with a fine-tuned prompt. This setup balanced accuracy, speed, and response quality. Key results include:

- **Accuracy**: The system correctly answered 85% of test queries, accurately identifying match details such as scores, dates, and teams.

- **Response Time**: Average query response time was 2.5 seconds on GPU-enabled hardware, suitable for real-time applications.

- **Example Query**: For the query "Score of Manchester City vs. Arsenal in 2022?" the system responded: "On February 15, 2023, Manchester City won 3-1 against Arsenal." This matched the ground truth in the dataset.

- **Qualitative Performance**: The system excelled at specific queries but struggled with ambiguous or aggregate queries (e.g., "Who performed best in 2022?") due to the lack of aggregated statistics in the dataset.

## 5.1 Evaluation Metrics

To evaluate the system, we used a test set of 100 queries, manually annotated with expected answers. Metrics included:

- **Accuracy**: Percentage of queries with correct match details (85%).

- **Precision@3**: Proportion of relevant summaries in the top-3 retrieved results (90%).

- **Response Coherence**: Assessed by human evaluators, with 80% of responses rated as coherent and informative.

## 5.2 Limitations

The system faced challenges with ambiguous queries and queries requiring aggregated statistics (e.g., season totals). Additionally, the reliance on home match summaries limited its ability to provide comprehensive away match details. These limitations highlight the need for further dataset expansion and preprocessing enhancements.

# 6 Conclusion

The Retrieval-Augmented Generation system developed for EPL match queries demonstrates the power of combining semantic search with contextual text generation. Achieving 85% accuracy and a 2.5-second response time, the system provides an intuitive interface for accessing complex sports data. The use of `all-MiniLM-L6-v2` for embeddings and `falcon-7b-instruct` for generation ensures a balance of efficiency and performance, making the system suitable for both casual and professional use.

However, the system's limitations, particularly in handling ambiguous or aggregate queries, suggest several avenues for future work:

- **Enhanced Retrieval**: Fine-tuning the embedding model on sports-specific data could improve semantic matching for domain-specific terms.

- **Aggregate Query Handling**: Adding a preprocessing module to compute season-level statistics (e.g., total goals scored by a team) would enable the system to handle broader queries.

- **Model Optimization**: Exploring larger LLMs or quantization techniques could enhance response quality while maintaining computational efficiency.

- **Dataset Expansion**: Including away match summaries and additional seasons would provide a more comprehensive data source.

This project lays the foundation for scalable, user-friendly sports analytics tools, with potential applications beyond the EPL, such as other sports leagues or real-time match analysis. By addressing the identified limitations, the system could evolve into a robust platform for interactive data exploration in sports.