# FA 25 6513-C BIG DATA FINAL PRESENTATION

NYU

LLM DataPrep: Scalable Pipeline for Biomedical Fine-Tuning

Professor: Amit Patel

# PROBLEM STATEMENT:

- Biomedical text (PubMed: **2.2M entries**) is **noisy, duplicated, and unstructured**

- Raw data **cannot** be used directly for LLM fine-tuning

- Need a **scalable pipeline** to transform big biomedical datasets into high-quality training data



**NYU**

# DATA OVERVIEW

- Data Source Name: PubMed

- **Data Source Link:** https://huggingface.co/datasets/MedRAG/pubmed

- Dataset File Size: 2.78 GB (Parquet) & ~70GB(JSONL)

- Approximate Number of Records: 2.21 Million

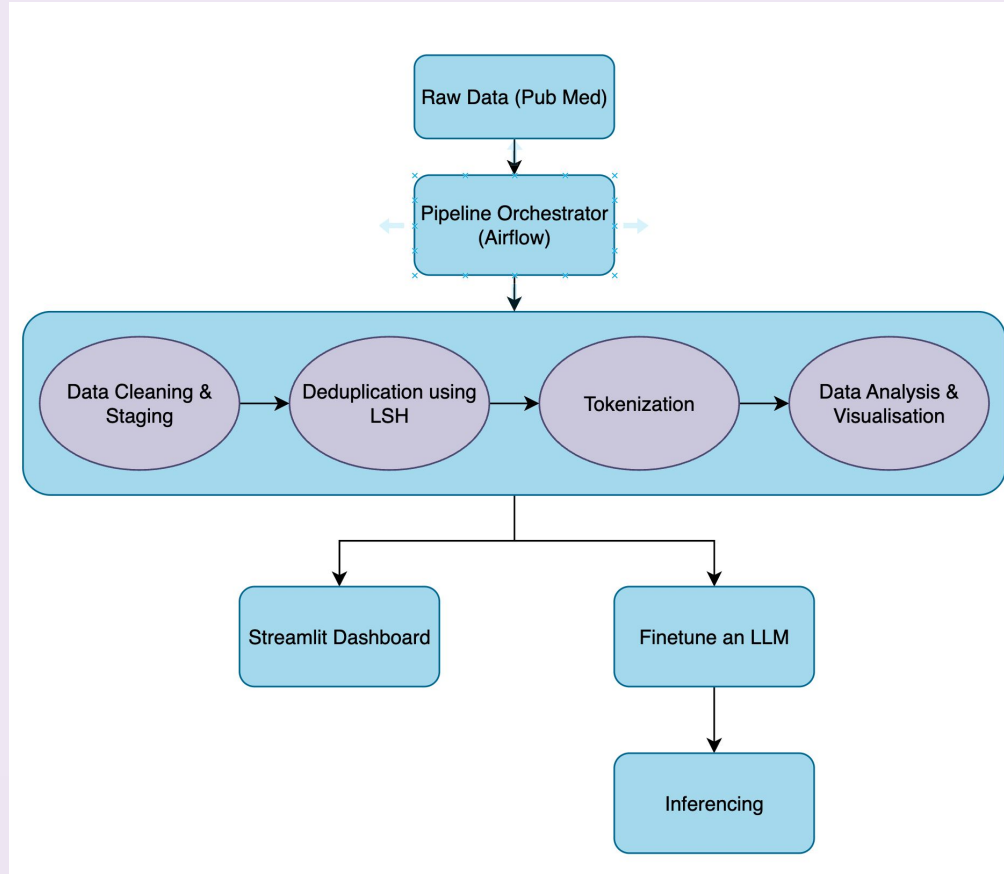- **Columns Include:** id, title, content, contents, PMID

**NYU**

# WHAT TECH STACK WE USED

# PROPOSED SOLUTION

- Build a scalable end-to-end ETL pipeline to process millions of biomedical text records efficiently.

- Automatically clean, normalize, deduplicate, and tokenize PubMed data to produce high-quality training inputs.

- Orchestrate and automate all pipeline stages using Apache Airflow for reliability and repeatability.

- Generate a structured, LLM-ready dataset for domain-specific fine-tuning.

**NYU**

# High Level Pipeline Architecture

# Step 1 - Data Cleaning

```
--- Output of STEP 1: Data Cleaning & Staging ---
STEP 1: Starting Data Cleaning and Staging...
    Initial record count: 2209839
    Cleaned records saved: 2209839
    Stage 1 data written to hdfs:///user/gg3039_nyu_edu/LLM_DataPrep/cleaned/stage1_data
```

- **Loads the raw PubMed text and builds a clean main_text field by merging and normalizing the content.**
- **Removes noise such as citations, extra spaces, and short/invalid entries.**
- **Saves a cleaned, standardized dataset to HDFS for the next stage.**

**NYU**

# Step 2: Deduplication Using MinHash LSH

```
--- Output of STEP 2: Near-Deduplication (LSH) ---
STEP 2: Starting Near-Deduplication (LSH)...
   Records removed by deduplication: 186606
   Stage 2 data written to hdfs:///user/gg3039_nyu_edu/LLM_DataPrep/cleaned/stage2_data
```
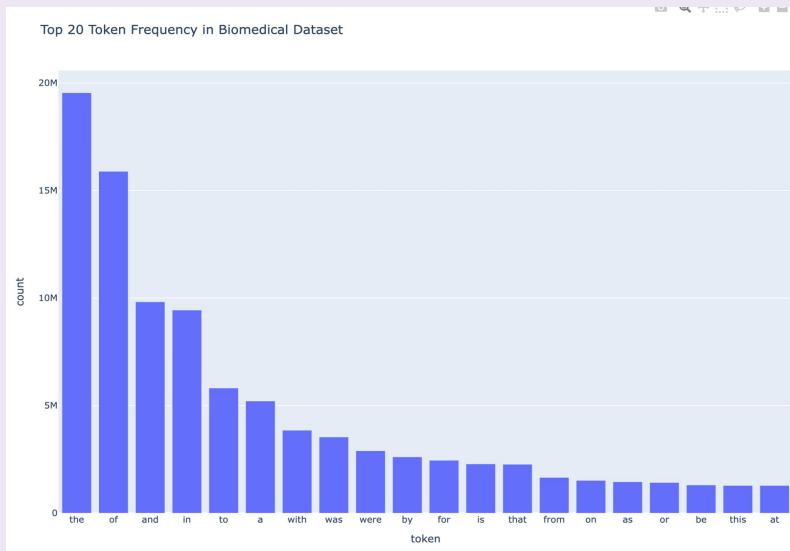
- **Tokenizes each document and converts it into a hashed feature vector using HashingTF, enabling efficient similarity checks.**
- **Uses MinHash LSH to group together texts that are highly similar and flags them as potential duplicates.**
- **Removes these near-duplicate records and saves a cleaned, deduplicated dataset to HDFS while reporting how many were dropped.**

NYU

# Step 3: Distributed Tokenization

```
--- Output of STEP 3: Tokenization & Final Save ---
STEP 3: Starting Tokenization and Final Save...
   Final records saved: 2023233
   Final LLM dataset written to hdfs:///user/gg3039_nyu_edu/LLM_DataPrep/final/llm_dataset
```

- **Converts main_text and title into numeric token sequences that LLMs use for learning**
- **Standardizes the dataset to keep only the fields required for downstream LLM training: id, PMID, title, main_text, tokens**
- **Processes tokenization in parallel across the Spark cluster, enabling millions of records to be handled efficiently**
- **Writes the final structured dataset as Parquet to HDFS for fast retrieval during model fine-tuning**
- **Produces a terminal summary with the total number of tokenized records**

**NYU**

# Step 4: Data Analysis & Visualization



Top 20 Token Frequency in Biomedical Dataset

- **Analyzes the cleaned and tokenized dataset to compute basic statistics like total records and average text length.**
- **Identifies the most common tokens in the corpus to understand vocabulary patterns.**
- **Generates an interactive Plotly chart and saves it for visualization and reporting.**

NYU

# Data Analysis - Hive Queries

```
        VERTICES        MODE        STATUS   TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
Map 1 ......... container    SUCCEEDED     52        52        0        0        0       0
Reducer 2 ...... container    SUCCEEDED    414       414        0        0        0       0
Reducer 3 ...... container    SUCCEEDED    207       207        0        0        0       0
Reducer 4 ...... container    SUCCEEDED      1         1        0        0        0       0
VERTICES: 04/04  [==========================>>] 100%  ELAPSED TIME: 41.06 s
INFO  : Completed executing command(queryId=hive_20251210044813_2ccb8ee7-e83d-4f7b-a150-bf1b4f2ea9b0); Time taken: 41.358 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
+------------------------------------------+-----------------+--------------+
|             lc.length_bin                | lc.record_count |  percentage  |
+------------------------------------------+-----------------+--------------+
| A. Very Short (<100 Chars)               | 326             | 0.02         |
| B. Standard Abstract (100-500 Chars)     | 280822          | 13.88        |
| C. Medium (501-1000 Chars)               | 786033          | 38.85        |
| D. Long (>1000 Chars)                    | 956052          | 47.25        |
+------------------------------------------+-----------------+--------------+
4 rows selected (41.501 seconds)
0: jdbc:hive2://localhost:10000>
```

Distribution of text lengths

```
        VERTICES        MODE        STATUS   TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
Map 1 ......... container    SUCCEEDED     52        52        0        0        0       0
Reducer 2 ...... container    SUCCEEDED    414       414        0        0        0       0
Reducer 3 ...... container    SUCCEEDED      1         1        0        0        0       0
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 34.88 s
INFO  : Completed executing command(queryId=hive_20251210045636_dfd796f2-087a-4251-afdb-e5c93ca7b6be); Time taken: 35.002 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
+--------------------------------+------------+
|            title               | duplicates |
+--------------------------------+------------+
| Hypertension in the elderly.   | 15         |
| Laparoscopic cholecystectomy.  | 13         |
| Malignant hyperthermia.        | 12         |
| Lyme disease.                  | 12         |
| Neuroleptic malignant syndrome.| 11         |
+--------------------------------+------------+
5 rows selected (35.439 seconds)
0: jdbc:hive2://localhost:10000>
```
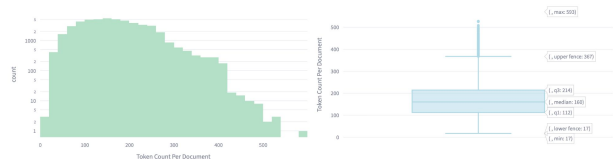
Top 5 Most repeated titles

```
LATERAL VIEW EXPLODE(tokens) exploded_table AS token_word
WHERE token_word NOT IN ('the', 'is', 'that', 'from', 'on', 'as', 'or', 'at', 'be', 'this')
GROUP BY token_word
ORDER BY token_count DESC
LIMIT 10
INFO  : Query ID = hive_20251210044446_311cebac-0df0-4094-b726-1ac7e8594de6
INFO  : Total jobs = 1
INFO  : Launching Job 1 out of 1
INFO  : Starting task [Stage-1:MAPRED] in serial mode
INFO  : Subscribed to counters: [] for queryId: hive_20251210044446_311cebac-0df0-4094-b726-1ac7e8594de6
INFO  : Session is already open
INFO  : Dag name: SELECT
token_word,
COUNT(1) AS token_co...10 (Stage-1)
INFO  : Status: Running (Executing on YARN cluster with App id application_1756163132607_32901)

        VERTICES        MODE        STATUS   TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
Map 1 ......... container    SUCCEEDED     52        52        0        0        0       0
Reducer 2 ...... container    SUCCEEDED    621       621        0        0        0       0
Reducer 3 ...... container    SUCCEEDED      1         1        0        0        0       0
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 77.05 s
INFO  : Completed executing command(queryId=hive_20251210044446_311cebac-0df0-4094-b726-1ac7e8594de6); Time taken: 77.259 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
+------------+-------------+
| token_word | token_count |
+------------+-------------+
| of         | 15325178    |
| and        | 9521786     |
| in         | 9090934     |
| to         | 5621679     |
| a          | 5025575     |
| with       | 3729574     |
| was        | 3428281     |
| were       | 2834538     |
| by         | 2505923     |
| for        | 2383933     |
+------------+-------------+
10 rows selected (77.978 seconds)
```

Most Common Non-stop word tokens

NYU

# Visualization - Streamlit

# Pipeline orchestration through Airflow

# Big Data Aspects We Focused On

- **Distributed processing:** Spark jobs across Dataproc cluster
- **Scalability:** Pipeline works on millions of biomedical records
- **Fault tolerance:** HDFS storage + Spark retry mechanisms
- **Optimized storage:** Parquet format for compression & speed
- **Parallel ETL stages:** Cleaning, dedupe, tokenization handled at scale
- **Workflow Orchestration at Scale (Airflow DAGs)** : Automated all pipeline stages using Apache Airflow, ensuring reproducibility, scheduling, and monitoring of large-scale ETL workflows.
- **Data Analysis:** Large-scale SQL analytics performed through Hive for dataset insights

# Challenges faced

- **Setting up Airflow in NYU Dataproc**
- **Unavailability of GPU resources**
- **Preparing LLM fine-tuning-ready format**

# Lessons Learned

- **Preprocessing quality drives LLM quality**
- **Distributed processing is essential at this scale**
- **Modular ETL design reduces debugging time**
- **Parquet + Spark = fast and scalable**

**NYU**

# THANK YOU

**Student Names and Net IDs:**

Geetha Krishna Guruju (Net Id: gg3039)

Tejdeep Chippa (Net Id: tc4263)

Isha Kookanapalli Math (Net Id: ik2688)

**NYU**

**ANY QUESTIONS**