# Result Management System ( Large-Scale University)

**1.Introduction:** One of the most important tasks in major universities is effectively controlling and analysing student results. Using cutting-edge computational methods, the "Result Management System" seeks to automate the creation, processing, and analysis of student performance data. As part of this project, student profiles will be created, subjects will be assigned, marks will be processed using MapReduce, and a dashboard will show statistical insights.

Ten thousand student records in six common subjects can be handled by the system. Its use of big data processing techniques guarantees efficiency and scalability. Large datasets can be processed effectively in a distributed setting by utilising MapReduce, which enables quick calculations of pass/fail rates, average marks, and greatest and lowest scores. Administrators may more easily evaluate patterns in student performance because to the system's integration of data visualisation techniques, which provide insightful information in graphical form. Additionally, this research shows how AI/ML approaches may be used in the real world for processing data and making decisions in education.

## 2.Dataset Creation:

1.**Generating 10,000 Student Profiles:** We generated 10,000 distinct student profiles to mimic an actual campus setting. Every profile includes:

- *Student ID*: A special number assigned to every student.
- *Name*: names that are produced at random.
- *Department*: Assigned according to predetermined classifications.
- *Enrolment Year*: Allotted at random from a range.

My dataset is called **"student_results_final.csv"** and I created it with the assistance of chat gpt.

2.**Generating 6 General Subject**: DSA, Electronics, Programming, Database, Data science, Mathematics.

3.**Generating Marks for 10,000 Students**: For these subjects, grades were given to each student. In order to replicate actual scoring patterns, marks were created at random using a normal distribution within the range of "0 to 100".

## 3.MapReduce Data Processing

1. **Concept of MapReduce and Its Usage**: An Overview and Application large datasets may be processed in parallel and distributed using the MapReduce programming technique. In this project, student grades were effectively aggregated and analysed using "PySpark". Utilising Spark's distributed computing technology, we made sure that calculations were carried out effectively over a sizable dataset.

### 2.Student Mark Processing :

- *Data Loading and Preprocessing:*
  - A "Spark DataFrame" was created for distributed processing once the dataset was loaded using "Pandas".
  - Values that were missing were examined and dealt with properly.
  - Consistency was ensured by mapping gender values to numerical representations.

- *Map Phase:*
  - The grades of every student were mapped into key-value pairs, with the "subject name" serving as the key and the "student's marks" as the value.
  - Examples of key-value pairs are `(DSA, 76)`, `(Mathematics, 85)`, and so on.
- *Reduce Phase*: "Spark SQL functions" were used to aggregate operations such "average, highest, lowest grades, and pass/fail counts".

  **computed statistics were:**
  Average Marks: subject-wise average marks are calculated using the `avg()` function.
  Highest and Lowest Marks: Peak scores were determined using the `max()` and `min()` methods.

  Pass/Fail Counts: A **passing threshold of 40%** was established, and pass/fail ratios were calculated using count aggregation.

# 4.Basic Analytics and Statistics

- **Data insights derived from student grades**:
  1.**Average grades for each subject**: were calculated using the `avg () ` function in Spark SQL.

  2. **Highest and Lowest Marks**: Identified across all students using max () and min () functions.

  3. **Pass/Fail Statistics**: Students who received scores over "40%" were classified as "pass", while those who had scores below 40% were classified as "fail".

  4.**Subject-Wise Performance Trends**: To identify the topics with the greatest and lowest average scores, trend analysis was conducted.

  **Additional Statistical Analysis:**

  1.**Mark Distribution:** Histograms were made to illustrate how marks varied across participant.

  2.**Department-wise Performance**: Students' performance in several departments was examined.

  **correlation analysis**: Used to find relationships between subjects (e.g., correlation between **Mathematics and Electronic** scores), KMeans clustering was applied to group students based on their performance patterns.

# 5. Visualisation and Dashboard

- **Graphical Representation of Results**: To graphically depict student performance data, a \*\*dashboard\*\* was made with "Matplotlib, Seaborn, and Plot". The following graphic aids were used
  - **Bar Charts:** Displayed **average marks per subject**.
  - **Violin Plots:** Combined aspects of box plots and density plots to visualize the **distribution of marks** while also displaying data density.
  - **Box Plots:** Showed the **spread and outliers** in student marks for each subject

# 6. Difficulties and Resolutions

**Challenges Faced:**

- **Processing large-scale data efficiently**: Solved using Spark and MapReduce, which optimized parallel execution.
- **Ensuring realistic data generation**: Used normal distributions to simulate real-world student performance.
- **Creating meaningful visualizations**: Selected appropriate graph types to maximize insights.
- **Handling missing data**: Implemented preprocessing techniques to clean and standardize the dataset

# 7. Conclusion and Future scope

### Summary of Project Outcomes

i) produced and processed marks for 10,000 students with success.

ii) PySpark was used to implement MapReduce for effective data management.

iii) carried out comprehensive analytics, such as clustering, subject trends, and pass/fail analysis.

iv) created a decision-supporting dashboard with interactive visualisations.

### Future Improvements

(1) **Integration with a Database**: Store findings for future retrieval and scalability.

(2) **Web-Based Interface**: Enable students to view their own outcomes by logging in.

(3) **Predictive Analysis**: Use machine learning models to forecast patterns in student performance.

(4) **Automated Reporting System**: Produce PDF reports that provide an overview of student performance.