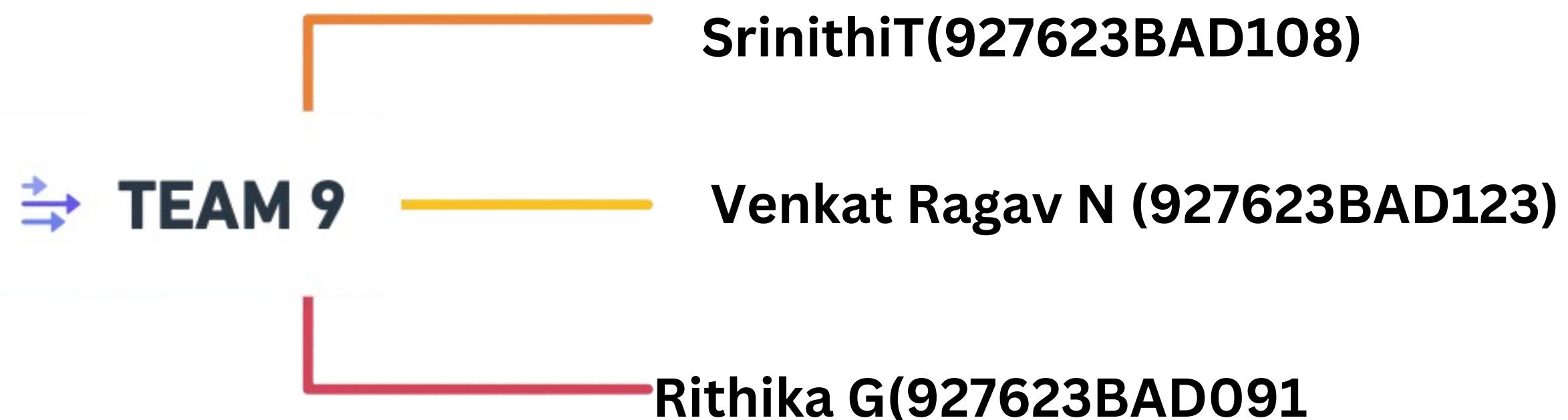


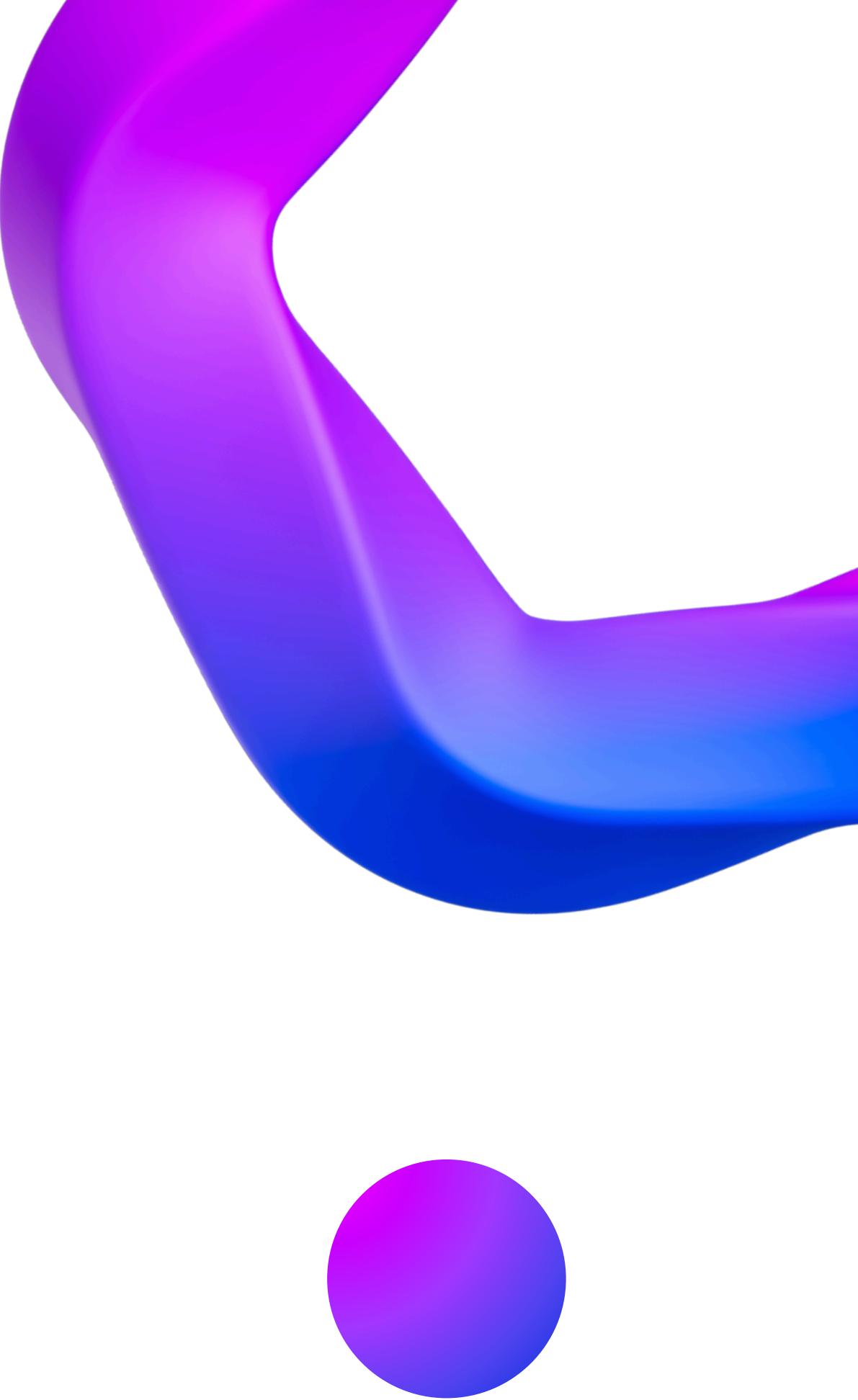


DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND DATA SCIENCE

Streamlined Data Processing Pipeline

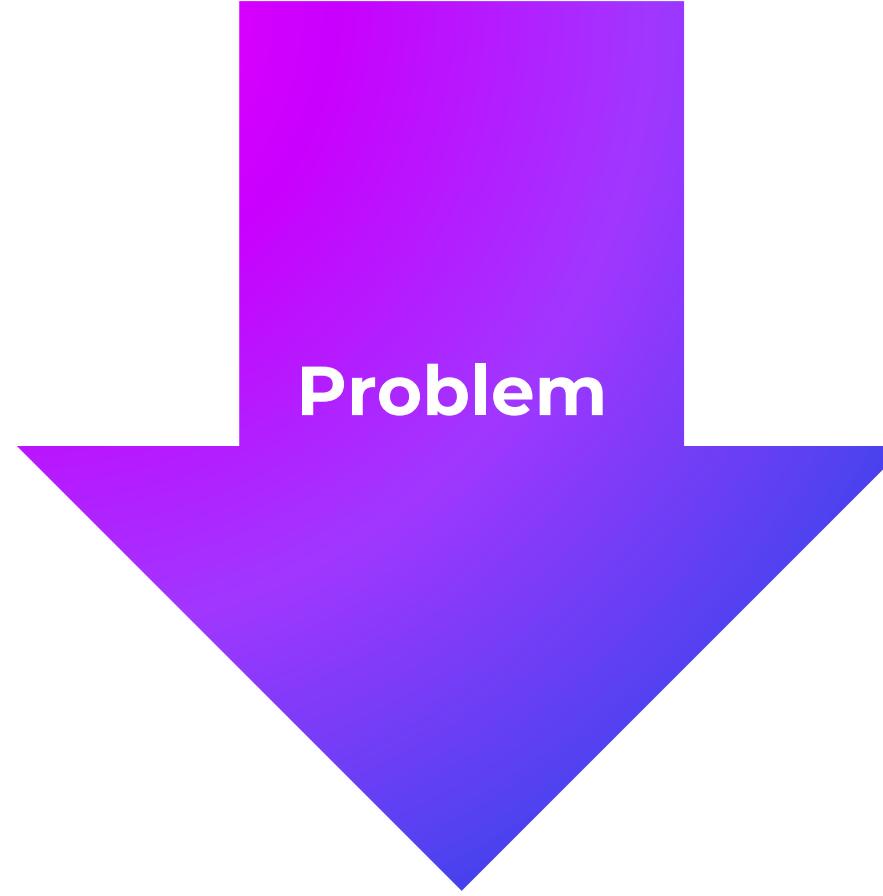


PandaPrep: Streamlined Data Processing Pipeline



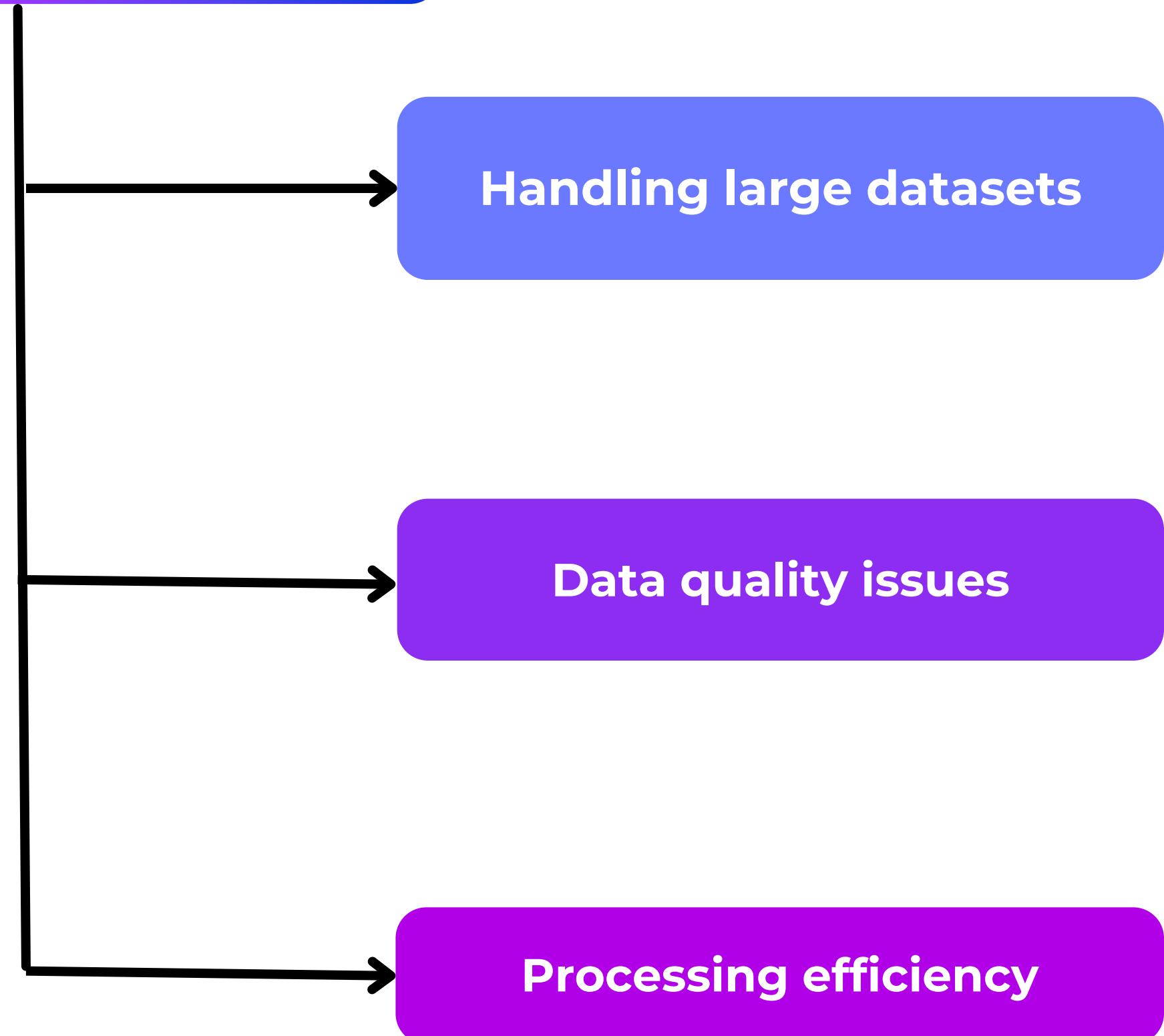


Problem



The project aims to develop a scalable, efficient pipeline to process large CSV files with data quality issues

Pain Points Identified



Significance of the Project



- The project is significant because it addresses a common issue faced by data engineers and analysts—how to handle and process large datasets efficiently.
- By automating data preprocessing and making it scalable, the solution can save significant time and effort, enabling better data-driven decision-making.

Original Objectives

01

Transform data for downstream analysis.

02

Handle data quality issues like missing values, duplicates, and formatting errors.

03

Develop a scalable pipeline for cleaning and preprocessing large CSV files.

Updated Objectives

01

02

03

04

Data Transformation: Enhance the pipeline to support more complex data transformations beyond just cleaning

Real-Time Processing : Explore the potential of adding real-time processing capabilities if the data comes from live systems

Output Formats: Enhance the pipeline to output cleaned data in different formats like JSON, SQL, and other structured formats besides just CSV.

Parallel Processing: Add support for parallel processing to handle large datasets faster.

Methodology & Approach

Methodologies:

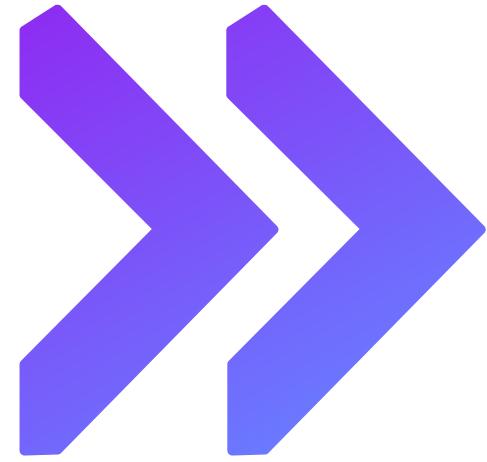
Reading and Cleaning Data

Handling Missing Values

Data Transformation

Duplicates Removal

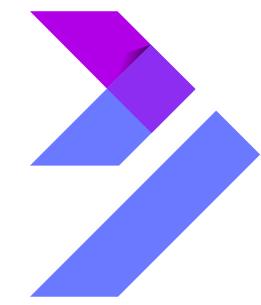
Data Collection & Preprocessing



Data Sources



The data comes in the form of raw CSV files sourced from various internal databases or system logs, potentially containing millions of rows and multiple columns.



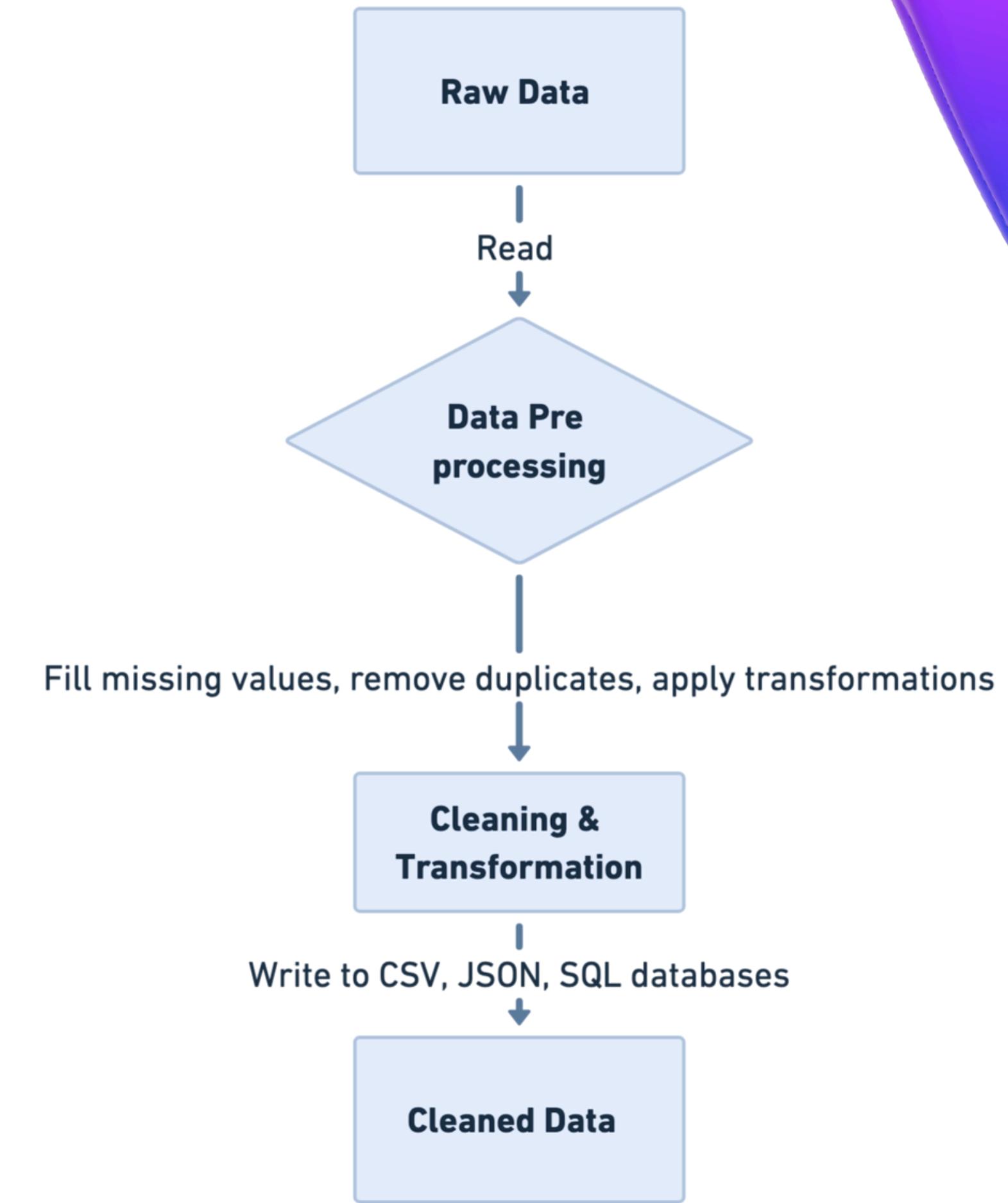
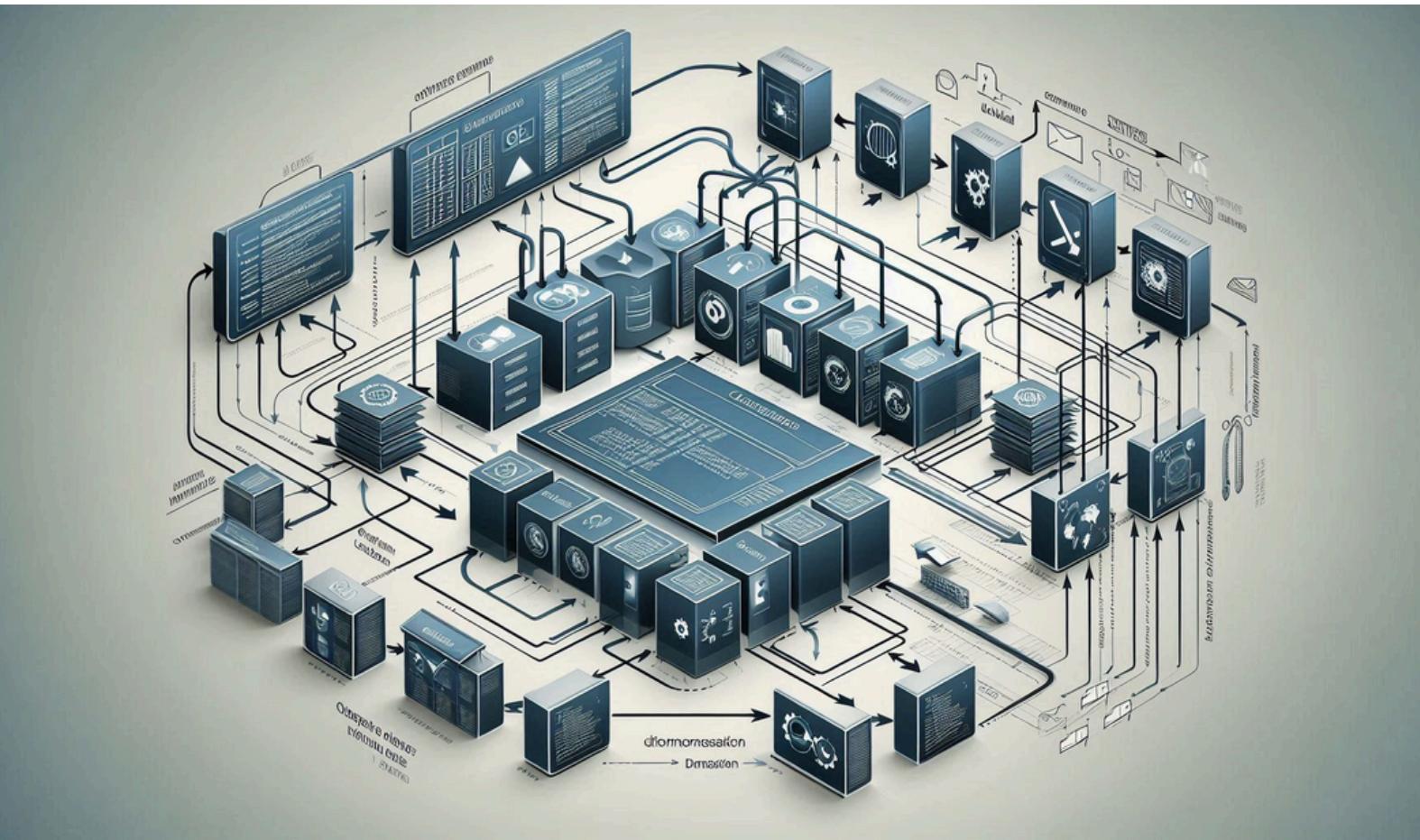
Data Cleaning Techniques

- Handling Missing Values
 - Removing Duplicates
 - Type Casting



System Architecture

Flow Diagram



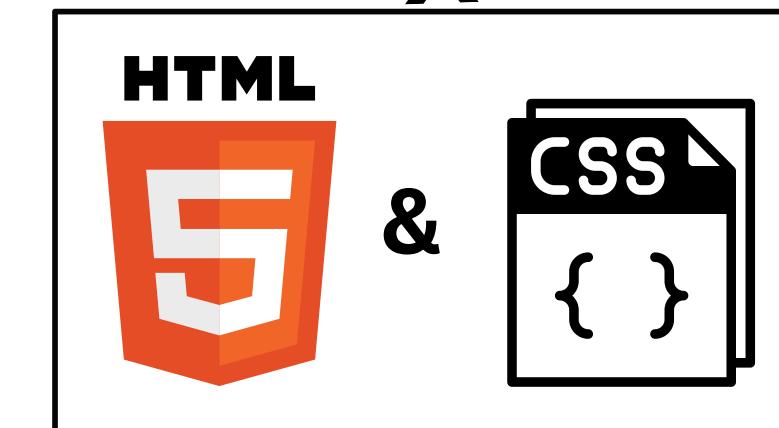
Technology Stack

Language:

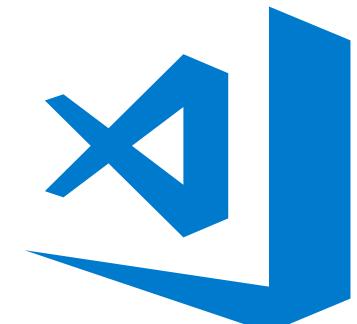


Backend
programming

Frontend
programming



Tools:



Visual Studio Code



Google Sheets

Challenges & Solutions

Challenges:

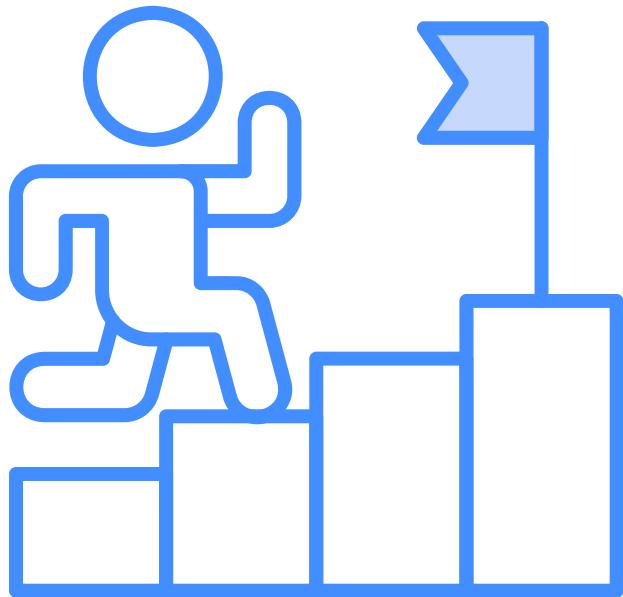


Data Cleaning Challenges

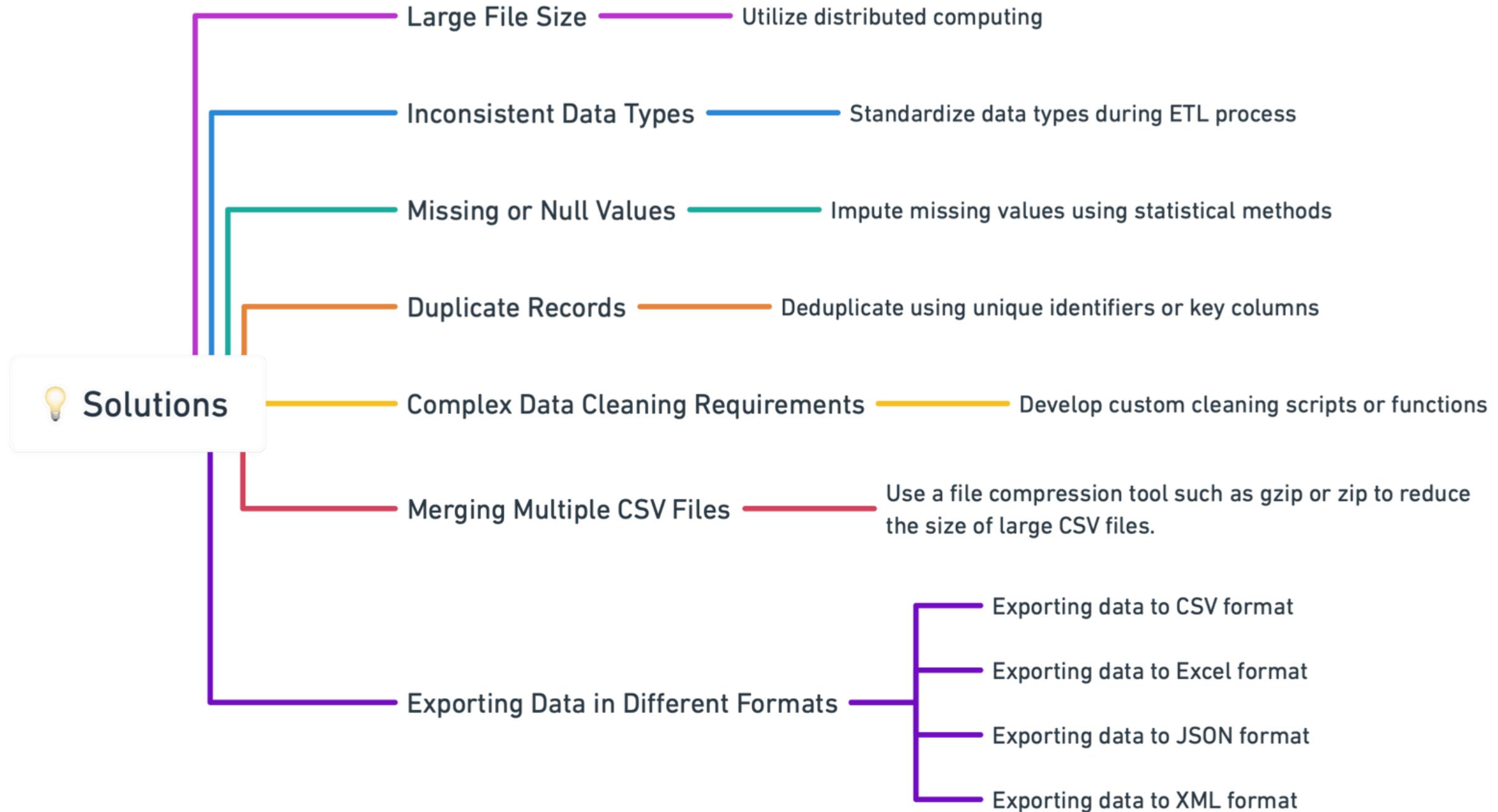
- Large File Size
- Inconsistent Data Types
- Missing or Null Values
- Duplicate Records
- Complex Data Cleaning Requirements
- Merging Multiple CSV Files
- Exporting Data in Different Formats



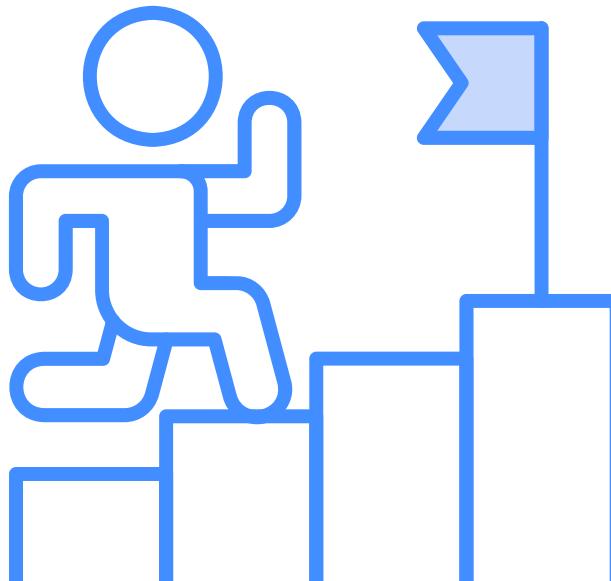
Google Sheets



Solutions:



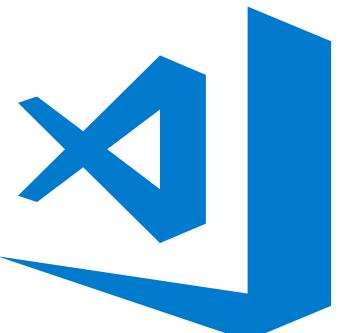
Google Sheets



**Project
Code
100%**

index.html

```
< index.html U >
templates > < index.html > html
1  <!DOCTYPE html> Contains emphasized items
2  <html lang="en">
3  <head>
4      <meta charset="UTF-8">
5      <meta name="viewport" content="width=device-width, initial-scale=1.0">
6      <title>Upload CSV</title>
7      <style>
8          body {
9              font-family: Arial, sans-serif;
10             margin: 0;
11             padding: 0;
12             background: #f4f4f4;
13             color: #333;
14         }
15         .container {
16             max-width: 800px;
17             margin: 50px auto;
18             background: #fff;
19             padding: 20px;
20             box-shadow: 0px 4px 8px #rgba(0, 0, 0, 0.1);
21             border-radius: 8px;
22         }
23         h1 {
24             text-align: center;
25             color: #007BFF;
26         }
27         form {
28             display: flex;
29             flex-direction: column;
30             align-items: center;
31         }
32         label {
33             font-size: 18px;
34             margin-bottom: 10px;
35         }
36         input[type="file"] {
37             margin-bottom: 20px;
38         }
39         .upload-button {
40             background: #007BFF;
41             color: white;
42             border: none;
43             padding: 10px 20px;
44             font-size: 16px;
45             border-radius: 4px;
46             cursor: pointer;
47         }
48         .upload-button:hover {
49             background: #0056b3;
50         }
51     </style>
52 </head>
53 <body>
54     <div class="container">
55         <h1>Upload Your CSV File</h1>
56         <form action="/clean" method="POST" enctype="multipart/form-data">
57             <label for="file">Select CSV File:</label>
58             <input type="file" id="file" name="file" accept=".csv" required>
59             <button type="submit" class="upload-button">Upload and Clean</button>
60         </form>
61     </div>
62 </body>
63 </html>
```



Visual Studio Code

result.html [1] ×

result.html



Visual Studio Code

```
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4      <meta charset="UTF-8">
5      <meta name="viewport" content="width=device-width, initial-scale=1.0">
6      <title>Cleaned Data</title>
7      <style>
8          body {
9              font-family: Arial, sans-serif;
10             margin: 0;
11             padding: 0;
12             background: #f4f4f4;
13             color: #333;
14         }
15         .container {
16             max-width: 1000px;
17             margin: 30px auto;
18             background: #fff;
19             padding: 20px;
20             box-shadow: 0px 4px 8px #007BFF;
21             border-radius: 8px;
22         }
23         h1, h2 {
24             text-align: center;
25             color: #007BFF;
26         }
27         .table-container {
28             margin: 20px auto;
29             overflow-x: auto;
30         }
31         .styled-table {
32             border-collapse: collapse;
33             margin: 20px auto;
34             font-size: 16px;
35             width: 100%;
36             border: 1px solid #ddd;
37         }
38         .styled-table th, .styled-table td {
39             text-align: left;
40             padding: 10px;
41         }
42         .styled-table th {
43             background: #007BFF;
```

```
        color: white;
    }
.styled-table tr:nth-child(even) {
    background: #f4f4f4;
}
.styled-table tr:hover {
    background: #ddd;
}
.chart-container {
    margin: 20px auto;
    text-align: center;
}
.button-group {
    text-align: center;
    margin-top: 20px;
}
.button-group a {
    display: inline-block;
    text-decoration: none;
    margin: 0 10px;
    background: #007BFF;
    color: white;
    padding: 10px 20px;
    border-radius: 4px;
    font-size: 16px;
}
.button-group a:hover {
    background: #0056b3;
}
.search-bar {
    text-align: center;
    margin-bottom: 20px;
}
.search-bar input {
    padding: 10px;
    width: 50%;
    font-size: 16px;
    border: 1px solid #ccc;
}
```

```
1  templates > <> result.html > ...
2      <html lang="en">
3          <head>
4              <style>
5                  .search-bar input {
6                      border: 1px solid #ccc;
7                      border-radius: 4px;
8                  }
9              </style>
10         </head>
11     <body>
12         <div class="container">
13             <h1>Cleaned Data Preview</h1>
14
15             <h2>Preview of Cleaned Data</h2>
16             <div class="search-bar">
17                 <input type="text" id="search" placeholder="Search in"
18             </div>
19             <div class="table-container">
20                 {{ cleaned_preview|safe }}
21             </div>
22
23             <h2>Summary Statistics</h2>
24             <div class="table-container">
25                 {{ stats_preview|safe }}
26             </div>
27
28             <h2>Data Visualization</h2>
29             <div class="chart-container">
30                 {{ chart_html|safe }}
31             </div>
32
33             <div class="button-group">
34                 <a href="{{ download_link }}>Download Cleaned Data</a>
35                 <a href="/">Back to Upload</a>
36             </div>
37
38         </div>
39         <script>
40             function filterTable() {
41                 let input = document.getElementById("search");
42                 let filter = input.value.toUpperCase();
43             }
44         </script>
45     </body>
46 </html>
```



main.py

```
main.py > clean_data
1  from flask import Flask, request, render_template, send_file
2  import pandas as pd
3  import io
4  import matplotlib
5  matplotlib.use('Agg') # Use a non-GUI backend
6  import matplotlib.pyplot as plt
7  import base64
8  import os
9
10 app = Flask(__name__)
11
12 # Ensure the downloads directory exists
13 os.makedirs('downloads', exist_ok=True)
14
15 @app.route('/')
16 def index():
17     return render_template('index.html')
18
19 @app.route('/clean', methods=['POST'])
20 def clean_data():
21     try:
22         # Retrieve uploaded file
23         file = request.files['file']
24         if not file:
25             return "No file uploaded", 400
26
27         try:
28             # Read the CSV file with UTF-8 encoding using the python engine
29             df = pd.read_csv(file, on_bad_lines='skip', engine='python')
30         except UnicodeDecodeError:
31             # If UTF-8 fails, try with ISO-8859-1 encoding
32             file.seek(0) # Reset file pointer to the beginning
33             df = pd.read_csv(file, encoding='ISO-8859-1', on_bad_lines='skip', engine='python')
34
35         # Data Cleaning
36         df.fillna(value='N/A', inplace=True) # Fill missing values
37         df.drop_duplicates(inplace=True) # Remove duplicates
38
39         # Generate cleaned data preview (first 10 rows)
40         cleaned_preview = df.head[100000].to_html(classes='styled-table', index=False)
```

```
41
42     # Summary Statistics
43     numeric_summary = df.describe().transpose()
44     numeric_summary.reset_index(inplace=True)
45     if numeric_summary.shape[1] == 8:
46         numeric_summary.columns = ['Column', 'Count', 'Mean', 'Std', 'Min', '25%', '50%', '75%', 'Max']
47     else:
48         numeric_summary.columns = ['Column'] + [f'Stat_{i}' for i in range(1, numeric_summary.shape[1])]
49
50     stats_preview = numeric_summary.to_html(classes='styled-table', index=False)
51
52     # Data Visualization (example: column distribution)
53     fig, ax = plt.subplots(figsize=(8, 6))
54     if not df.select_dtypes(include=['number']).empty:
55         df.select_dtypes(include=['number']).mean().plot(kind='bar', ax=ax, color='skyblue')
56         ax.set_title('Mean of Numeric Columns', fontsize=14)
57         ax.set_ylabel('Mean Value', fontsize=12)
58         plt.tight_layout()
59
60     # Save plot to a string buffer
61     buf = io.BytesIO()
62     plt.savefig(buf, format='png')
63     buf.seek(0)
64     chart_html = f''
65     buf.close()
66
67     # Save plot to a file in the downloads directory
68     graph_path = os.path.join('downloads', 'mean_numeric_columns.png')
69     plt.savefig(graph_path)
70     plt.close(fig)
71     else:
72         chart_html = "<p>No numeric columns available for visualization.</p>"
73
74     # Save cleaned data for download
75     cleaned_data_path = os.path.join('downloads', 'cleaned_data.csv')
76     df.to_csv(cleaned_data_path, index=False)
77
78     # Provide a downloadable link
79     download_link = "/download"
80     download_graph_link = "/download-graph"
```

OUTPUT

INDEX PAGE

Upload Your CSV File

Select CSV File:

no file selected

AFTER UPLOADING THE FILE

Upload Your CSV File

Select CSV File:

 IOT_Assignme...r_range.csv

Cleaned Data Preview

Preview of Cleaned Data

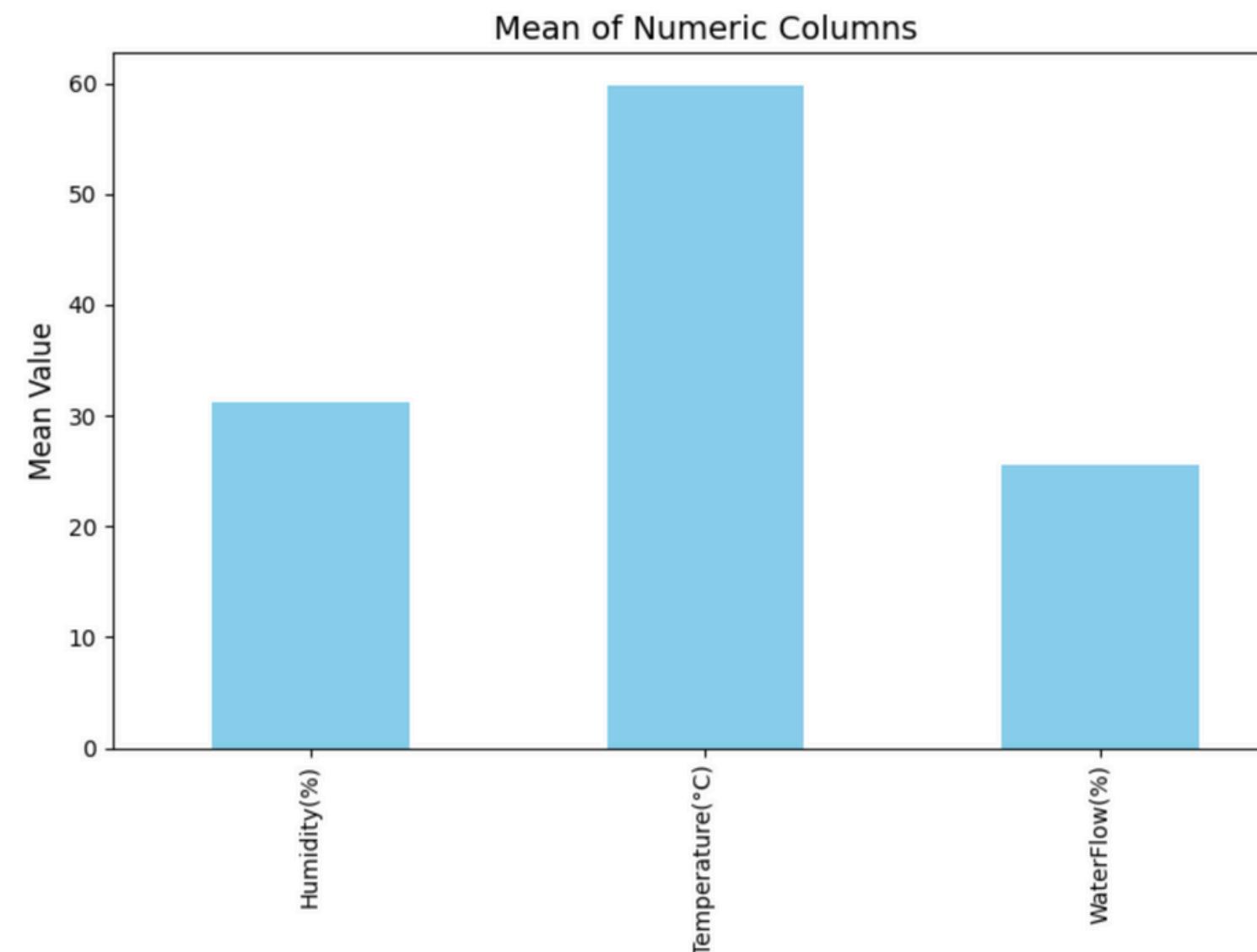
Search in table...

Humidity(%)	Temperature(°C)	WaterFlow(%)
25.045045	76.856857	67.432607
8.880881	31.771772	0.000000
24.956957	75.735736	66.105904
40.328328	88.308308	0.000000
38.742743	57.157157	0.000000
45.613614	24.164164	0.000000
33.237237	43.383383	0.000000
23.415415	81.421421	73.999716
10.026026	96.316316	100.000000
15.399399	58.838839	50.670417

Summary Statistics

Column	Count	Mean	Std	Min	25%	50%	75%	Max
Humidity(%)	200.0	31.137437	12.386344	8.132132	21.015015	30.616617	41.528529	51.867868
Temperature(°C)	200.0	59.747347	23.112767	20.000000	40.860861	58.638639	78.998999	100.000000
WaterFlow(%)	200.0	25.496364	34.953661	0.000000	0.000000	0.000000	61.884829	100.000000

Data Visualization



<THANK YOU>

