

DATA ANALYSIS & VISUALIZATION USING AWS

A PROJECT WORK

*Submitted in partial fulfilment of
requirements for the award of the degree of*

**Bachelor of Technology
in
INFORMATION TECHNOLOGY
By**

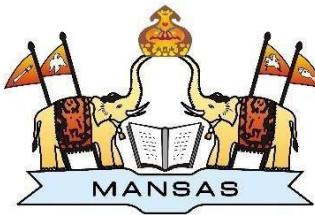
**B.Srija
(21331A1214)**

**N.Lahari
(22335A1206)**

**K.Prabhath Naidu
(21331A1255)**

Under the esteemed guidance of

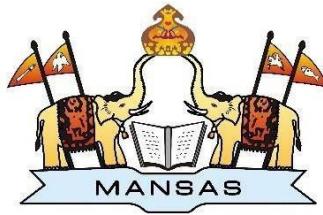
**Dr Anjanadevi B
Associate Professor
Department of Information Technology
M.V.G.R. College of Engineering(A)**



**Department of Information Technology
MAHARAJ VIJAYARAM GAJAPATHI RAJ COLLEGE OF ENGINEERING(A)
(Affiliated to Jawaharlal Nehru Technological University, Gurajada Vizianagaram)
VIZIANAGARAM
2024 - 25**

MAHARAJ VIJAYARAM GAJAPATHI RAJ COLLEGE OF ENGINEERING (A)
VIZIANAGARAM

BONAFIDE CERTIFICATE



Certified that this is a bonafide record of project work entitled “Data analysis & visualization using AWS” ,done by B.Srija (21331A1214), N.Lahari (22335A1206), K.Prabhath Naidu (21331A1255) in partial fulfillment for the award of the degree of “**Bachelor of Technology** “ in Information Technology, M.V.G.R. College of Engineering, Vizianagaram , year 2024 – 25.

Dr Anjanadevi B

Associate professor,

Supervisor,

Department of IT,

MVGR College of Engineering(A),

Vizianagaram.

Dr. P. Srinivasa Rao

Professor,

Head of the Department,

Department of IT,

MVGR College of Engineering(A),

Vizianagaram.

External Examiner

DECLARATION

I hereby declare that the work done on the dissertation entitled "**DATA ANALYSIS & VISUALIZATION USING AWS**" has been carried out by us and submitted in partial fulfillment for the award of credits in **BACHELOR OF TECHNOLOGY** in **INFORMATION TECHNOLOGY** of MVGR College of Engineering(Autonomous) affiliated to the Jawaharlal Nehru Technological University (Vizianagaram).

The various contents incorporated in the dissertation have not been submitted for the award of any other degree of any other institution or university.

B.Srija (21331A1214)
N.Lahari (22335A1206)
K.Prabhath Naidu(21331A1255)

ACKNOWLEDGEMENT

As a note of acknowledgement, with great solemnity and sincerity, we offer our profuse thanks to Dr Anjanadevi B, Associate Professor, Dept. of IT, for guiding us all through our project work, giving a right direction and shape to our learning by extending her expertise and experience in the field of education. Really, we are indebted to her for her excellent and enlightened guidance.

Let us lay it on the line we should thank Dr Anjanadevi B, who remained all through our project work as a great source of inspiration for her resourceful counsel and for inducing in us her empirical knowledge that pivoted the goal and destination of our project.

We consider it our privilege to express our deepest gratitude to Dr. P Srinivasa Rao, Professor and Head of the Department for his valuable suggestions and constant motivation that greatly helped the project work to get successfully completed.

We also thank Prof. P S Sitharama Raju, Director and Dr. R Ramesh, Principal, for extending his utmost support and cooperation in providing all the provisions for the successful completion of the project.

We sincerely thank all the members of the staff in the Department of Information Technology for their sustained help in our pursuits.

We thank all those who contributed directly or indirectly in successfully carrying out this work.

B.Srija (21331A1214)

N.Lahari (22335A1206)

K.Prabhath Naidu(21331A1255)

DATA ANALYSIS & VISUALIZATION USING AWS

Team Name

A6

Objective

This project aims to enhance data analysis and visualization using AWS services like Amazon S3, AWS Glue, Amazon Athena, and Amazon QuickSight. It focuses on efficient data storage, transformation, querying, and visualization to enable faster insights and better decision-making.

Team

Department: Information Technology

Technology

Name: B.Srija

Email: srijabehara5

@gmail.com

Mobile:8309971252



Name: N.Lahari

Email: niddanalahari

@gmail.com

Mobile:6300558941



Name: K.Prabhath Naidu

Email: prabhathkarri0629

@gmail.com

Mobile:7386055602



Select domain(s) where your solution can be implemented

- Cloud Computing
- Data Analytics
- Business Intelligence

Describe how your solution is going to meet the program outcomes

Our project utilizes AWS services for efficient data processing and visualization. By leveraging AWS tools, it ensures cost effective, automated ETL workflows and interactive dashboards. This project aligns with program outcomes by providing hands-on experience with cloud technologies and improving data analysis capabilities.

Describe the engineering solution

This solution leverages AWS services to efficiently manage data analysis and visualization. Amazon S3 ensures secure and scalable data storage, while AWS Glue automates ETL processes, handling data cleaning and transformation. Amazon Athena enables serverless querying, optimizing data retrieval for analysis. Amazon QuickSight provides interactive dashboards, offering insights for decision-making. AWS IAM enforces strict access control, ensuring secure authorization and authentication.

Project Guide

Name : Dr.B.Anjanadevi

Designation: Associate professor

Email: anjanadevi@mvrce.edu.in



End Users / Stake holders of solution

- Data analysts
- Business intelligence teams
- Cloud architects & IT Teams



ABSTRACT

Data-driven decision-making is crucial in today's business landscape, and cloud-based analytics platforms provide scalable and efficient solutions for handling data. This project, "Data Analysis and Visualization using AWS, focuses on analyzing the Fortune 1000 Companies dataset (2024) to extract valuable business insights. The dataset consists of financial and operational metrics of the top 1000 U.S. companies, including revenue, profit, assets, market value, sector classification, and leadership attributes. To process this data, we utilize Amazon S3 for storage, AWS Glue for data transformation, Amazon Athena for serverless querying, and Amazon Quick Sight for interactive dashboard visualizations. Various analytical queries were executed in Athena to uncover trends, correlations, and key performance indicators across different industries. There by visualized using Quick Sight dashboards, providing intuitive insights into company rankings, profitability trends, and sectoral distributions. This project demonstrates the efficiency of AWS based data analytics in handling data without infrastructure management. The approach ensures scalability, cost-effectiveness, and real-time insights, making it a good solution for enterprise-level analytics and business intelligence.

CONTENTS

	Page No
List of Figures	1
1. Introduction	2
1.1. Cloud Computing	2
1.2. Amazon Web Services	2
1.3. Problem definition	3
1.4. Objective	3
1.5. Existing System	4
1.6. Proposed System	4
2. Literature Survey	5
3. System Study and Analysis	8
3.1 User Requirements	8
3.1.1 Functional Requirements	8
3.1.2 Non Functional Requirements	8
3.2 Data Flow Diagram	9
3.3 Hardware and Software Configurations	10
3.3.1 Hardware Requirements	10
3.3.2 Software Requirements	10
4. System Design	11
4.1 Architecture	11
4.1.1 Data ingestion	11
4.1.2 Data storage and organization	13
4.1.3 Data Transformation and ETL Process	13
4.1.4 Data Analysis and Query Processing	14
4.1.5 Data Visualization and Reporting	15
4.1.6 AWS Services Integration and Security	16
5. System Testing	17
5.1 Data Ingestion Testing	17
5.2 ETL Process Testing	19
5.3 Query Processing Testing	22

5.4 Visualization Testing	27
5.5 Results & Performance Comparison	30
5.5.1 Performance Metrics	30
5.5.2 Comparison with Traditional Tools	30
6. Conclusion & Future Work	31
6.1 Conclusion	31
6.2 Future work	31
7. References	32

LIST OF FIGURES

- Figure 3.2 :** The data flow diagram of the project
- Figure 4.1 :** The architecture of the system
- Figure 4.1.1 :** The raw dataset in CSV format.
- Figure 5.1.1 :** The data is uploaded into S3.
- Figure 5.1.2 :** Crawler is created in AWS Glue
- Figure 5.1.3 :** The schema and the tables are stored in glue data catalog
- Figure 5.2.1 :** The ETL Job created in glue to transform the data
- Figure 5.2.2 :** The raw data is transformed into parquet format by ETL Job
- Figure 5.2.3 :** The parquet data schema in the glue data catalog
- Figure 5.2.4 :** The parquet data schema in the glue data catalog
- Figure 5.3.1 :** The query analysis performed in Athena
- Figure 5.3.2 :** The query results that are generated in Athena
- Figure 5.3.3 :** The query analysis performed in Athena
- Figure 5.3.4 :** The query results that are generated in Athena
- Figure 5.3.5 :** The query analysis performed in Athena
- Figure 5.3.6 :** The query results that are generated in Athena
- Figure 5.3.7 :** The query analysis performed in Athena
- Figure 5.3.8 :** The query results that are generated in Athena
- Figure 5.4.1 :** The analysis of KPI's in Quicksight
- Figure 5.4.2 :** The analysis of profits by company and assets, market capital by sector in Quicksight
- Figure 5.4.3 :** The analysis of assets, profits percent change and revenue percent change by company in Quicksight
- Figure 5.4.4 :** The analysis about number of employees by company, industry and change in rank of company in Quicksight
- Figure 5.4.5 :** The analysis of market capital and revenue by company, industry in Quicksight

1. INTRODUCTION

1.1 Cloud Computing?

Cloud computing is a technology that provides on-demand computing resources over the internet, eliminating the need for physical infrastructure. It allows organizations to access storage, processing power, databases, and networking services remotely, reducing operational costs and improving efficiency. Key features of cloud computing include scalability, cost-effectiveness, accessibility, security & reliability, and automation & integration for efficient resource management.

1.2 Amazon Web Services(AWS)

What is AWS?

Amazon Web Services (AWS) is a comprehensive cloud computing platform that provides a wide range of on-demand cloud services, including storage, computing power, databases, networking, and machine learning tools. Launched by Amazon in 2006, AWS has become the leading cloud service provider, offering scalable, secure, and cost-effective solutions for businesses of all sizes.

Why Use AWS?

AWS is widely used due to its flexibility, global presence, and extensive suite of services. Some key benefits include:

1. **Scalability** – AWS allows automatic scaling of resources based on demand, ensuring optimal performance.
2. **Cost-Effectiveness** – With a pay-as-you-go model, users only pay for what they use, reducing operational expenses.
3. **Security** – AWS provides built-in security measures, including encryption, firewalls, identity access management (IAM), and compliance certifications.

4. **Reliability** – AWS offers 99.99% uptime with global data centers, ensuring minimal downtime.
5. **Flexibility & Integration** – AWS services can be seamlessly integrated with various tools, databases, and analytics platforms.
6. **Fully Managed Services** – AWS automates many processes like data storage, computing, and analytics, reducing manual effort.

1.3 Problem Definition

As organizations generate massive amounts of data, traditional tools like SQL, Spark, Excel, Power BI, and Tableau struggle to keep up due to scalability limitations, slow processing speeds, and complex infrastructure requirements. SQL and Excel are effective for small datasets but become inefficient with large-scale data, leading to slow query execution. Spark, while powerful, requires complex configurations and high computational resources. Power BI and Tableau, though strong in visualization, lack built-in scalability and automated data integration, making real-time analysis difficult. To overcome these challenges, this project adopts a cloud-based, serverless approach using AWS services—Amazon S3 for scalable storage, AWS Glue for automated ETL, Amazon Athena for fast SQL-based querying, and Amazon QuickSight for dynamic visualization. This solution ensures faster data processing, real-time insights, reduced infrastructure costs, and seamless scalability, enabling organizations to make data-driven decisions efficiently.

1.4 Objective

The objective of this project is to develop a cloud-based, scalable, and cost-effective data analytics and visualization system using AWS services. By leveraging AWS tools, the project aims to improve data storage, processing, querying, and visualization while overcoming the limitations of traditional methods. The key objectives are:

- Enhance data processing efficiency

- Improve scalability and performance
- Enable fast and cost-effective querying
- Provide real-time, interactive visualization
- Minimize operational costs & infrastructure maintenance

The goal is to create a fully cloud-based analytics system that simplifies data management, speeds up insights, and provides a seamless visualization experience for businesses across various industries.

1.5 Existing System

Before the adoption of cloud-based solutions like AWS, organizations relied on SQL databases, Spark, Excel, Power BI, and Tableau for data processing and visualization. These tools struggle with scalability, real-time processing, and high infrastructure costs, making them inefficient for large-scale datasets. SQL databases slow down with complex queries, Spark requires expensive computational resources, and Excel lacks automation. Power BI and Tableau offer strong visualization but lack real-time data integration. These limitations lead to slow data retrieval, high costs, and inefficiencies, highlighting the need for a scalable, automated, and cost-effective solution.

1.6 Proposed System

To overcome the limitations of traditional tools, this project adopts a cloud-based, serverless architecture using AWS for efficient data processing and visualization. Amazon S3 provides scalable storage, eliminating infrastructure management. AWS Glue automates ETL (Extract, Transform, Load) processes, reducing manual effort. Amazon Athena enables fast, serverless SQL-based querying, improving real-time data retrieval. Amazon QuickSight offers interactive visualizations, enhancing decision-making with real-time insights. This fully managed AWS solution ensures scalability, automation, cost efficiency, and real-time analytics, enabling organizations to process large datasets effortlessly while reducing operational overhead.

2. LITERATURE REVIEW

A) Big Data Analytics on AWS Cloud

Mishra and Kumar (2021), in their study published in the *International Journal of Engineering Research & Technology (IJERT)*, explored the evolving landscape of big data analytics in AWS cloud environments. The paper provides a comprehensive overview of AWS services such as Amazon S3, Glue, Athena, and QuickSight, demonstrating their role in enabling efficient data storage, querying, and visualization. The study highlights how cloud-based analytics platforms offer scalability, cost-effectiveness, and ease of deployment compared to traditional data processing systems. Furthermore, the authors discuss the four dimensions of Big Data (Volume, Velocity, Variety, and Veracity) and the challenges associated with managing large-scale datasets in cloud environments.

B) Serverless Data Processing and Visualization using AWS

Caliwag and Caliwag (2021), in their conference paper on *ResearchGate*, presented an AWS-powered data visualization framework leveraging DynamoDB and Lambda functions. The proposed architecture integrates AWS services for real-time data ingestion, transformation, and visualization using QuickSight. The study emphasizes the benefits of serverless computing, including reduced infrastructure overhead and cost efficiency, making it a compelling choice for data-driven applications.

C) Sentiment Analysis using AWS NLP Services

Satyanarayana et al. (2020), in their *IEEE International Conference Paper*, investigated the application of AWS Comprehend for sentiment analysis of voice-based data. The study outlines an automated workflow where speech data is transcribed using AWS Transcribe, processed via AWS Glue and Athena, and visualized in QuickSight. The research underscores the efficacy of AWS-based NLP tools in extracting meaningful insights from unstructured data, particularly in business intelligence applications.

D) Cloud-Based Weather Data Analytics

Wavhale et al. (2020), in their publication in the *International Research Journal of Engineering and Technology (IRJET)*, proposed a real-time weather forecasting and analytics system utilizing AWS Lambda, S3, and API Gateway. The study integrates various meteorological APIs with AWS cloud services to collect, store, process, and analyze weather patterns. The findings highlight how serverless computing and cloud storage improve efficiency in large-scale data processing and decision-making applications.

E) Amazon Athena: A Serverless Data Querying Approach

Kulkarni (2023), in the *International Journal of Computer Trends and Technology (IJCTT)*, examined Amazon Athena's serverless architecture and its impact on big data querying and analytics. The study explores Athena's integration with S3 and Glue, highlighting its ability to perform cost-effective, high-speed SQL-based queries on large datasets. Additionally, the paper discusses troubleshooting techniques for optimizing query performance and cost management in serverless environments.

F) Digital Architecture for Cloud-Based Microgrid Monitoring

Patsidis et al. (2023), in their article published in *Energies (MDPI)*, explored a cloud-hosted data analytics framework for monitoring microgrid operations. The study utilizes AWS Lambda, DynamoDB, and QuickSight to analyze time-series data from energy grids. By leveraging real-time analytics and predictive modeling, the research demonstrates significant cost reductions and efficiency improvements in energy management.

G) AWS-Powered Data Warehousing for Structured Data Processing

Gupta et al. (2023), in their *IEEE International Conference Paper*, proposed an approach for ingesting and visualizing CSV files using AWS services. The research highlights the effectiveness of AWS Glue for ETL (Extract, Transform, Load), Redshift for structured data storage, and QuickSight for visualization. The study emphasizes scalability, automation, and cost reduction in cloud-based data warehousing.

H) AI-Driven Big Data Analytics on AWS

Kaniganti (2018), in the *International Journal of Science and Research (IJSR)*, presented an AI-integrated big data analytics framework utilizing AWS Kinesis, DynamoDB, and SageMaker. The study showcases the potential of machine learning and AI models for predictive analytics, demonstrating their application in finance, manufacturing, and retail industries. The research concludes that AWS-powered AI frameworks enable real-time decision-making and improved business intelligence.

I) Serverless Architecture for Big Data Analytics

Rahman and Hasan (2019), in their *IEEE Conference Paper*, examined the serverless computing paradigm for big data analytics. The study compared traditional and cloud-based architectures, highlighting how AWS Lambda, S3, Glue, and Athena streamline data processing while eliminating infrastructure management complexities. The paper concludes that serverless models significantly enhance efficiency, cost savings, and scalability in data-intensive applications.

3. SYSTEM STUDY AND ANALYSIS

3.1 User Requirements

3.1.1 Functional Requirements

1. **Scalable Data Storage** – The system should handle large datasets efficiently using Amazon S3.
2. **Automated Data Processing** – Data should be cleaned, transformed, and loaded using AWS Glue.
3. **Fast Query Execution** – The system should enable quick SQL-based queries using Amazon Athena.
4. **Interactive Data Visualization** – The processed data should be presented in a clear and dynamic manner using Amazon QuickSight.
5. **Serverless Architecture** – The system should operate without requiring dedicated infrastructure, reducing maintenance costs.

3.1.2 Non-Functional Requirements

1. **Security** – Data should be stored securely with IAM roles, encryption, and access controls.
2. **Cost Optimization** – The system should use pay-as-you-go AWS services to reduce costs.
3. **Backup & Disaster Recovery** – AWS S3 versioning and backup policies should ensure data recovery in case of failures.
4. **Reliability** – The system should ensure high availability and fault tolerance with AWS-managed services.
5. **Performance Optimization** – Athena queries should be optimized using partitioning and compression to reduce query execution time.
6. **Ease of Use** – The system should be user-friendly, allowing analysts to interact with dashboards without complex coding.

3.2 Data Flow Diagram

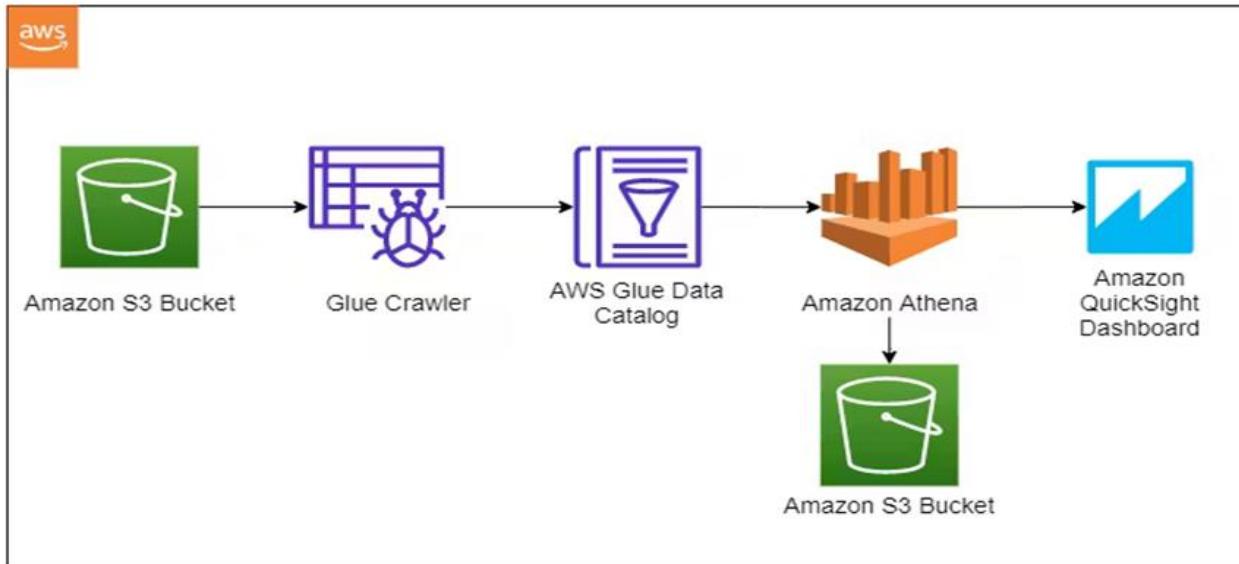


Fig.3.2 Data Flow Diagram of the Project (Source:[10]).

The data flow diagram represents the systematic processing of data from ingestion to visualization using AWS services. The workflow begins with Amazon S3, where raw data (CSV, JSON, or other formats) is stored. To enhance query performance, the data is converted into Parquet format, a columnar storage format that optimizes storage and speeds up querying in Amazon Athena.

An AWS Glue Crawler scans the raw and Parquet-formatted data in S3 and updates the AWS Glue Data Catalog, which maintains metadata such as table structure and schema details. This enables Amazon Athena to efficiently query the data using standard SQL. Athena processes user queries and retrieves relevant datasets, which can optionally be stored back in S3 for further processing.

The final step involves Amazon QuickSight, which connects to Athena to generate interactive dashboards and reports, providing insights into the processed data. This pipeline ensures scalability, cost efficiency, and near real-time analytics, leveraging serverless AWS services to streamline big data processing and visualization.

3.3 Hardware and Software Configurations

To ensure optimal performance, scalability, and cost-effectiveness, our system leverages AWS cloud infrastructure. Below are the key hardware and software configurations required for implementing the project.

3.3.1 Hardware Requirements

To support seamless execution and data handling, the following hardware specifications are required:

- **RAM:** Minimum 8 GB (recommended 16 GB for high-performance processing).
- **Processor:** Multi-core processor (Intel i5/i7 or AMD Ryzen equivalent).
- **Storage:** Minimum 500 GB SSD/HDD for local backup and processing.
- **Internet:** High-speed and stable internet connection for accessing AWS services.

3.3.2 Software Requirements

Cloud Services (AWS Stack)

The project is implemented using AWS cloud services, providing serverless analytics and data visualization. The required AWS services include:

- **Amazon S3** – Cloud storage for raw and processed datasets.
- **AWS Glue** – ETL (Extract, Transform, Load) service for data integration.
- **Amazon Athena** – Serverless SQL-based query engine for data analysis.
- **Amazon QuickSight** – Data visualization and business intelligence tool.

System Requirements

- **Operating System:** Windows 10/11, macOS, or Linux.
- **Web Browser:** Latest version of Google Chrome, Mozilla Firefox, or Microsoft Edge for accessing the AWS Management Console.

4. SYSTEM DESIGN

4.1 Architecture

The architecture of our AWS-based solution is designed to handle the Fortune 1000 companies dataset, which includes structured data about companies, their financials, and other attributes. The system leverages Amazon S3 for storage, AWS Glue for ETL (Extract, Transform, Load), Amazon Athena for querying, and Amazon QuickSight for visualization.

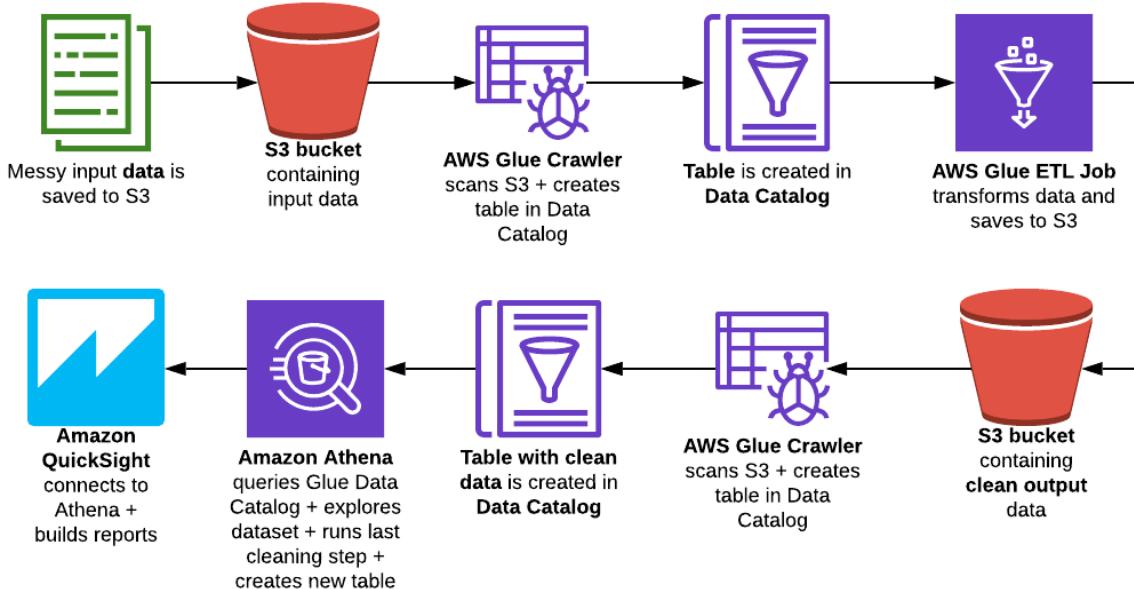


Fig.4.1 The architecture of the system (Source:[11]).

4.1.1 Data Ingestion

Description:

Data ingestion is the process of collecting and importing data into the system. In our aws project, the Fortune 1000 companies dataset is ingested into Amazon S3(simple storage service).

It serves as the primary data storage for our project and is the first step in the data pipelining of our project.

The dataset is initially uploaded in raw format (e.g., CSV) to an S3 bucket.

Tools and Services:

- Amazon S3:** Used as the primary storage for raw data. The dataset is uploaded to an S3 bucket.
- AWS Glue Crawler:** Automatically scans the raw data in S3 to infer the schema and store it in the AWS Glue Data Catalog.

Process:

1. The raw dataset of fortune companies (CSV format) is uploaded to the S3 bucket.

Rank	Company	Ticker	Sector	Industry	Profitable	Founder	Female	Growth	In Change	Gained	In Dropped	Newcome	Global500	Worlds	M Best	Com Number	c Market	Caj Revenues	RevenueP	Profits_M	ProfitsPer Assets	M	CE
1	Walmart	WMT	Retailing	General M	yes	no	no	0 no	no	no	yes	yes	yes	yes	2100000	484852.8	648125	6	15511	32.8	252399	C	
2	Amazon	AMZN	Retailing	Internet S	yes	no	no	0 no	no	no	yes	yes	yes	no	1525000	1873676	574785	11.8	30425	527854	Ar		
3	Apple	AAPL	Technolog	Computer	yes	no	no	1 yes	no	no	yes	yes	yes	no	161000	2647974	383285	-2.8	96995	-2.8	352582	Ti	
4	UnitedHealth UNH		Health Car	Health Car	yes	no	no	yes	1 yes	no	no	yes	yes	no	440000	456080.8	371622	14.6	22381	11.2	273720	An	
5	Berkshire H BRKA		Financials	Insurance	yes	no	no	yes	2 yes	no	no	yes	yes	no	396500	908919.7	364482	20.7	96223	1069978	W		
6	CVS Health CVS		Health Car	Health Car	yes	no	yes	no	0 no	no	no	yes	yes	no	259500	100373.9	357776	10.9	8344	101.1	249728	Ka	
7	Exxon Mo EXM		Energy	Petroleum	yes	no	no	-4 no	yes	no	yes	no	no	no	61500	461222.2	344582	-16.7	36010	-35.4	376317	De	
8	Alphabet GOOGL		Technolog	Internet S	yes	no	no	0 no	no	no	yes	yes	yes	no	182502	1884633	307394	8.7	73795	23	402392	Su	
9	McKesson MCK		Health Car	Wholesale	yes	no	no	no	0 no	no	no	yes	no	no	48000	70546.5	276711	4.8	3560	219.6	62320	Br	
10	Concra COR		Health Car	Wholesale	yes	no	no	yes	1 yes	no	no	yes	no	no	44000	48472.1	262173.4	9.9	1745.3	2.7	625587	St	
11	Costco WI COST		Retailing	General M	yes	no	yes	1 yes	no	no	yes	yes	yes	no	316000	324924.4	242290	6.8	6292	7.7	68994	Ro	
12	JPMorgan JPM		Financials	Commerce	yes	no	yes	11 yes	no	no	yes	yes	yes	no	309926	576938.4	239425	54.7	49552	31.5	3875393	Ja	
13	Microsoft MSFT		Technolog	Computer	yes	no	no	0 no	no	no	yes	yes	no	no	221000	3126133	211915	6.9	72361	-0.5	411976	Sa	
14	Cardinal H CAH		Health Car	Wholesale	yes	no	yes	0 no	no	no	yes	no	no	no	47520	27217.8	205012	13	261	43417	Ja		
15	Chevron CVX		Energy	Petroleum	yes	no	yes	-5 no	yes	no	yes	no	no	no	45600	292956.6	200949	-18.4	21369	-39.7	261632	Mi	
16	Cigna CI		Health Car	Health Car	yes	no	yes	-1 no	yes	no	yes	yes	no	no	71413	103017.9	195265	8.2	5164	-22.6	152761	Da	
17	Ford Moto F		Motor Veh	Motor Veh	yes	no	yes	2 yes	no	no	yes	no	no	no	177000	53017.8	176191	11.5	4347	273310	Ja		
18	Bank of A BAC		Financials	Commerce	yes	no	no	14 yes	no	no	yes	yes	yes	yes	212985	299213	171912	49.4	26515	-3.7	3180151	Br	
19	General M GM		Motor Veh	Motor Veh	yes	no	yes	no	2 yes	no	no	yes	no	no	163000	52353.5	171842	9.6	10127	1.9	273064	Mi	
20	Elevance I-ELV		Health Car	Health Car	yes	no	yes	2 yes	no	no	yes	yes	yes	yes	104900	120619.6	171340	9.4	5987	-0.6	108928	Ge	
21	Citigroup C		Financials	Commerce	yes	no	yes	15 yes	no	no	yes	no	no	no	237925	121122.2	156820	55.1	9228	-37.8	2411834	Ja	
22	Centene CNC		Health Car	Health Car	yes	no	yes	no	3 yes	no	no	yes	yes	yes	67700	41979.4	153999	6.5	2702	124.8	84641	Se	
23	23 Home Dep HD		Retailing	Specialty F	yes	no	no	-3 no	yes	no	yes	yes	yes	no	463100	380153.7	152669	-3	15143	-11.5	76530	Ed	
24	Marathon MPC		Energy	Petroleum	yes	no	yes	-8 no	yes	no	yes	no	no	no	18200	72607.7	150307	-16.5	9681	-33.3	85987	Mi	
25	Kroger KR		Food & Dr	Food & Dry	yes	no	no	-1 no	yes	no	yes	no	no	no	414000	41100.7	150039	1.2	2164	-3.6	50505	W	
26	Phillips 66 PSX		Energy	Petroleum	yes	no	yes	-9 no	yes	no	yes	no	no	no	14000	69880.8	149890	-14.7	7015	-36.4	75501	Mi	
27	27 Englehardt E		Financials	Finance	yes	no	yes	-4 no	yes	no	yes	no	no	no	6400	1610.2	144340	-15.3	13409	-34.7	1235437	Or	

Fig.4.1.1 The raw dataset in CSV Format (Source:[12]).

2. An AWS Glue Crawler is triggered to scan the raw data and infer the schema.
3. The schema is stored in the AWS Glue Data Catalog, making the data ready for transformation and querying.

4.1.2 Data Storage and Organization

Description:

The data is stored in Amazon S3 and organized into different layers to facilitate efficient processing and analysis. The storage architecture follows a 4-layer data lake model:

1. **Raw Layer:** Contains the original, unprocessed data (e.g., CSV files).
2. **Curated Layer:** Stores cleaned and enriched data after ETL processing.
3. **Processed Layer:** Contains data transformed into an optimized format (e.g., Parquet) for efficient querying.
4. **Provisioning Layer:** Stores final datasets ready for visualization and reporting.

Tools and Services:

- **Amazon S3:** Used for storing data in all layers (Raw, Curated, Processed, and Provisioning).
- **AWS Glue Data Catalog:** Stores metadata about the datasets, including schema information.

Data Organizing:

- The raw data is stored in the S3 bucket.
- After ETL job processing, the data is transformed into Parquet format and stored in the S3 bucket.

4.1.3 Data Transformation and ETL Process

Description:

The ETL process is responsible for transforming raw data into a format suitable for analysis. In our project, AWS Glue is used to perform ETL operations, including data cleaning, enrichment, and conversion to Parquet format.

Workflow:

1. **Extract:** The raw data is extracted from the fortune1000-raw-data S3 bucket.
2. **Transform:** AWS Glue performs the following transformations:
 - Cleansing: Removing null values and duplicates.
 - Enrichment: Changing data types for better format and analysis.
 - Conversion: Converting the data from CSV to Parquet format for efficient storage and querying.
3. **Load:** The transformed data is loaded into the S3 bucket.

Tools and Services:

- **AWS Glue:** Used for ETL jobs, including data cleaning, enrichment, and conversion to Parquet.
- **AWS Glue Data Catalog:** Stores metadata about the transformed datasets.

4.1.4 Data Analysis and Query Processing

Description:

Once the data is transformed and stored in Parquet format, it is ready for analysis. Amazon Athena is used to run SQL queries on the processed data, enabling fast and efficient data analysis.

Process:

1. The processed data in Parquet format is queried using Amazon Athena.
2. SQL queries are executed to analyze key metrics such as revenue, profit, and market capitalization by sector and industry.
3. Query results are stored temporarily in S3 for further visualization.

Tools and Services:

- **Amazon Athena:** Used for running SQL queries on the processed data.
- **Amazon S3:** Stores query results temporarily.

Query Optimization:

- The use of **Parquet format** and **data partitioning** significantly reduces query execution time and cost.

4.1.5 Data Visualization and Reporting

Description:

The final step in the architecture is data visualization and reporting. Amazon QuickSight is used to create interactive dashboards and reports based on the analyzed data.

Process:

1. Query results from Amazon Athena are imported into Amazon QuickSight.
2. Dashboards are created to visualize key metrics such as:
 - Revenue by Industry
 - Profit by Sector
 - Market Capitalization Trends
3. Users can interact with the dashboards to explore data and gain insights.

Tools and Services:

- **Amazon QuickSight:** Used for creating interactive dashboards and reports.
- **Amazon Athena:** Provides the query results for visualization.

UI/UX:

- The dashboards are designed to be user-friendly, with interactive charts, filters, and drill-down capabilities.
- Users can customize views based on their analysis needs.

4.1.6 AWS Services Integration and Security

Security Measures:

- **AWS IAM (Identity and Access Management):** Used to define roles and permissions for accessing AWS resources. For example:
 - A specific IAM role is created for the AWS Glue ETL job to access S3 buckets.
 - Another IAM role is created for Amazon QuickSight to query data from Athena.
- **Encryption:** Data stored in S3 is encrypted using AWS Key Management Service (KMS) to ensure security.
- **Access Control:** S3 bucket policies are configured to restrict access to authorized users only.

Service Integration:

- Amazon S3 serves as the central storage for raw and processed data.
- AWS Glue integrates with S3 and Athena to perform ETL and store metadata in the Glue Data Catalog.
- Amazon Athena queries the processed data stored in S3 and provides results to QuickSight.
- Amazon QuickSight integrates with Athena to visualize query results in interactive dashboards.

5. SYSTEM TESTING

5.1 Data Ingestion Testing

The data ingestion process ensures that raw data is successfully uploaded to Amazon S3, maintaining integrity and accessibility for further processing. The Fortune 1000 dataset, stored in CSV format, was uploaded to an Amazon S3 bucket, and its structure and completeness were verified. AWS Glue Crawler was executed to infer the schema, ensuring proper integration with downstream processing tools.

Key Considerations:

- The dataset must be ingested without errors, missing values, or corruption.
- The schema must be accurately inferred to maintain consistency across data processing stages.

Testing Approach:

1. Uploaded the raw dataset to Amazon S3 and verified object metadata, including file format and structure.
2. Executed AWS Glue Crawler to scan the dataset and generate schema information in the AWS Glue Data Catalog.
3. Cross-checked the schema structure with the expected attributes to ensure accuracy.

Observations:

- The dataset was successfully ingested, and all records were intact.
- The Glue Crawler accurately identified column names, data types, and relationships without inconsistencies.
- No data loss or corruption was observed during the ingestion process.

Result: The ingestion process was completed successfully, ensuring data readiness for transformation and analysis.

Amazon S3

General purpose buckets

Directory buckets

Table buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

fortune_reports/

Objects (1)

Name	Type	Last modified	Size	Storage class
csv_reports/	Folder	-	-	-

Fig.5.1.1 The data is uploaded into S3 bucket.

AWS Glue

Crawlers

Crawlers (2) Info

Name	State	Schedule	Last run	Last run times...	Log	Table changes fr...
fortune_reports_c...	Ready	(Succeeded)	December 12, 202...	View log	1 created	
parquet_fortune...	Ready	(Succeeded)	December 12, 202...	View log	-	

Fig.5.1.2 Crawler is created in AWS Glue

The screenshot shows the AWS Glue Data Catalog Tables page. On the left, there's a navigation sidebar with sections like AWS Glue, Data Catalog, and Data Integration and ETL. The main area is titled 'Tables' and contains a table with two rows. The columns are Name, Database, Location, Classification, Deprecated, View data, Data quality, and Column stats... . The first row has 'csv_reports' as the name, 'fortune_reports' as the database, 's3://aws-glue-demo' as the location, 'CSV' as the classification, and '-' as the deprecated status. The second row has 'parquet_parquet_re' as the name, 'fortune_reports' as the database, 's3://aws-glue-demo' as the location, 'Parquet' as the classification, and '-' as the deprecated status. There are buttons for 'Add tables using crawler' and 'Add table' at the top right of the table view.

Name	Database	Location	Classification	Deprecated	View data	Data quality	Column stats...
csv_reports	fortune_reports	s3://aws-glue-demo	CSV	-	Table data	View data quality	View statistics
parquet_parquet_re	fortune_reports	s3://aws-glue-demo	Parquet	-	Table data	View data quality	View statistics

Fig.5.1.3 The crawler infers the schema and the tables are stored in glue data catalog

5.2 ETL Process Testing

The ETL (Extract, Transform, Load) pipeline is responsible for structuring and optimizing data for analysis. AWS Glue was used to process the dataset, converting it into a format that enhances performance in subsequent query and visualization stages.

Key Considerations:

- Data transformation must preserve integrity and structure.
- The ETL job should enhance data accessibility while optimizing storage efficiency.

Testing Approach:

1. Executed AWS Glue ETL job to clean, transform, and convert the dataset into an optimized format.
2. Stored the transformed dataset in Amazon S3 and updated metadata in the AWS Glue Data Catalog.
3. Validated the record count and schema consistency between raw and transformed datasets.

Observations:

- The ETL job processed the dataset successfully without introducing inconsistencies.
- The transformed dataset was structured correctly, ensuring compatibility with Amazon Athena for querying. The optimized format improved data retrieval without data loss.

Result: The ETL process was executed successfully, ensuring an optimized dataset for analysis.

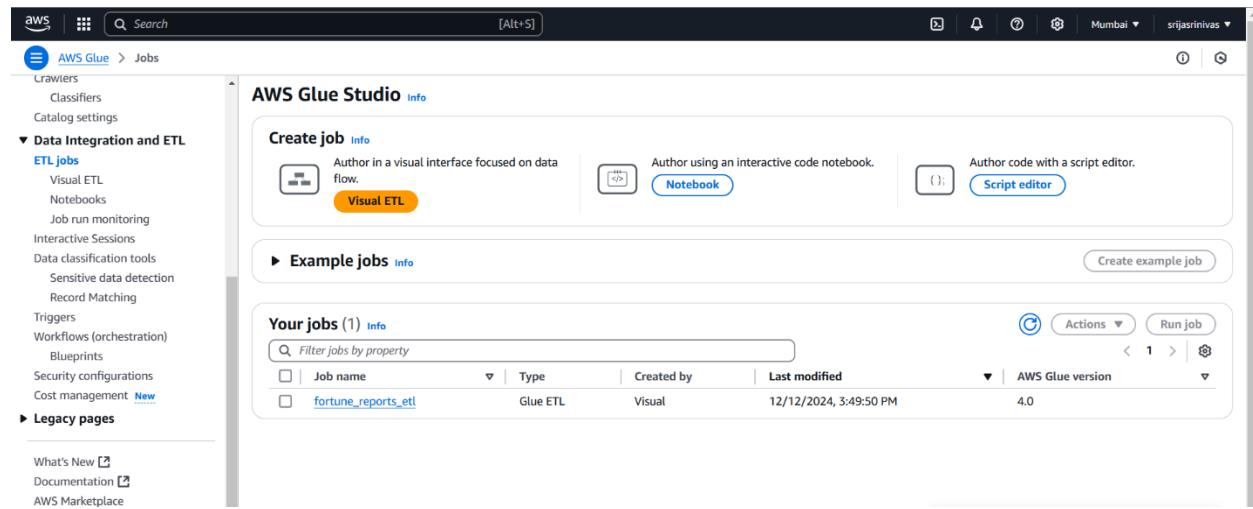


Fig.5.2.1 The ETL Job is created in glue to transform the data

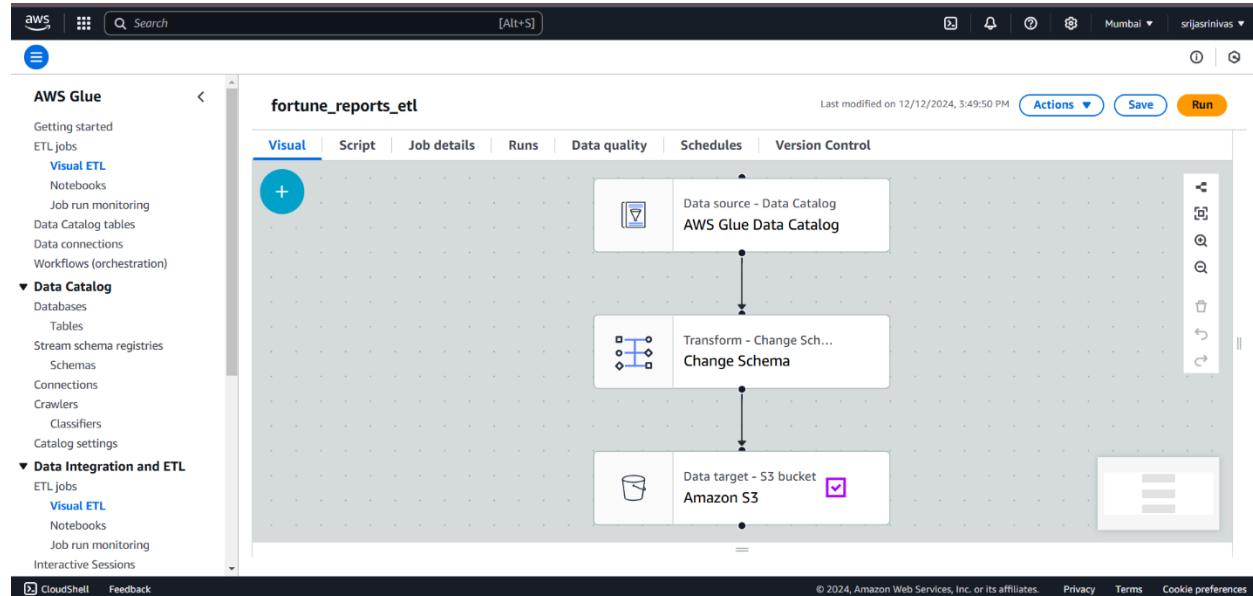


Fig.5.2.2 The raw data is transformed into parquet format by ETL Job

The screenshot shows the AWS Glue Data Catalog interface. On the left, there's a navigation sidebar with links for AWS Glue, Data Catalog tables, Data Catalog, Data Integration and ETL, and Legacy pages. The main area is titled 'Schema (32)' and displays a table of 32 columns. The columns are numbered 1 to 17, and then 18 to 32. The columns include 'rank', 'company', 'ticker', 'sector', 'industry', 'profitable', 'founder_is_ceo', 'femaleceo', 'growth_in_jobs', 'change_in_rank', 'gained_in_rank', 'dropped_in_rank', 'newcomer_to_the_fortune500', 'global500', 'worlds_most_admired_comp...', 'best_companies_to_work_for', and 'number_of_employees'. The data types for these columns are bigint, string, string, string, string, string, string, string, string, double, string, string, string, string, string, string, string, string, string, double, and bigint respectively. There are also columns 18 through 32 which are mostly empty or have '-' in them.

#	Column name	Data type	Partition key	Comment
1	rank	bigint	-	-
2	company	string	-	-
3	ticker	string	-	-
4	sector	string	-	-
5	industry	string	-	-
6	profitable	string	-	-
7	founder_is_ceo	string	-	-
8	femaleceo	string	-	-
9	growth_in_jobs	string	-	-
10	change_in_rank	double	-	-
11	gained_in_rank	string	-	-
12	dropped_in_rank	string	-	-
13	newcomer_to_the_fortune500	string	-	-
14	global500	string	-	-
15	worlds_most_admired_comp...	string	-	-
16	best_companies_to_work_for	string	-	-
17	number_of_employees	bigint	-	-
18				
19				
20				
21				
22				
23				
24				
25				
26				
27				
28				
29				
30				
31				
32				

Fig.5.2.3 The parquet data schema in the glue data catalog

This screenshot shows the AWS Glue Data Catalog interface, similar to Fig.5.2.3 but with a different set of columns. The schema table has 32 columns, numbered 17 to 32. The columns include 'number_of_employees', 'marketcap_march28_m', 'revenues_m', 'revenuepercentchange', 'profits_m', 'profitspercentchange', 'assets_m', 'ceo', 'country', 'headquarterscity', 'headquartersstate', 'website', 'companytype', 'footnote', 'marketcap_updated_m', and 'updated'. The data types for these columns are bigint, double, double, double, double, double, double, string, string, string, string, string, string, string, string, string, double, and string respectively.

17	number_of_employees	bigint	-	-
18	marketcap_march28_m	double	-	-
19	revenues_m	double	-	-
20	revenuepercentchange	double	-	-
21	profits_m	double	-	-
22	profitspercentchange	double	-	-
23	assets_m	double	-	-
24	ceo	string	-	-
25	country	string	-	-
26	headquarterscity	string	-	-
27	headquartersstate	string	-	-
28	website	string	-	-
29	companytype	string	-	-
30	footnote	string	-	-
31	marketcap_updated_m	double	-	-
32	updated	string	-	-

Fig.5.2.4 The parquet data schema in the glue data catalog

5.3 Query Processing Testing

Query processing is a critical aspect of data analysis, ensuring that information can be retrieved accurately and efficiently. Amazon Athena was used to execute SQL-based queries on the transformed dataset to extract meaningful insights.

Key Considerations:

- Queries must return correct results without missing or incorrect values.
- Query performance should be optimized for large datasets to ensure efficiency.

Testing Approach:

1. Executed multiple analytical queries to assess revenue, profitability, sector distribution, and other key metrics.
2. Compared query results against expected values to validate accuracy.
3. Assessed query response times and retrieval performance after data optimization.

Observations:

- All queries returned accurate and consistent results, matching expected dataset attributes.
- The optimized dataset structure improved query execution performance, enhancing retrieval speed.
- Queries involving aggregations and filters executed efficiently without processing delays.

Result: The query processing phase demonstrated efficiency in retrieving insights, supporting seamless data exploration.

The screenshot shows the AWS Athena Query editor interface. On the left, there's a sidebar titled 'Data' with sections for 'Data source' (set to 'AwsDataCatalog') and 'Database' (set to 'fortune_reports'). Below these are 'Tables and views' and a list of two tables: 'csv_reports' and 'parquet_parquet_reports'. The main area is titled 'Query 7' and contains the following SQL code:

```

1 SELECT company, sector, number_of_employees
2 FROM "fortune_reports"."parquet_parquet_reports"
3 WHERE number_of_employees > 100000
4 ORDER BY number_of_employees DESC;

```

Below the code, it says 'SQL Ln 4, Col 35'. At the bottom of the editor are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. To the right, there's a checkbox for 'Reuse query results up to 60 minutes ago'.

Fig.5.3.1 This is the query analysis performed in Athena

The screenshot shows the 'Query results' tab in the Athena Query editor. It displays the results of the previously run query. The results are listed in a table with columns: '#', 'company', 'sector', and 'number_of_employees'. There are 140 rows of data. The top of the table includes a search bar and buttons for 'Copy' and 'Download results CSV'.

#	company	sector	number_of_employees
1	Walmart	Retailing	2100000
2	Walmart	Retailing	2100000
3	Amazon	Retailing	1525000
4	Home Depot	Retailing	463100
5	Home Depot	Retailing	463100
6	FedEx	Transportation	446400
7	FedEx	Transportation	446400
8	UnitedHealth Group	Health Care	440000
9	Concentrix	Technology	440000

Fig.5.3.2 This is the query results that are generated in Athena

The screenshot shows the AWS Athena Query editor interface. On the left, there's a sidebar titled 'Data' with sections for 'Data source' (set to 'AwsDataCatalog'), 'Catalog' (set to 'None'), and 'Database' (set to 'fortune_reports'). Below this is a 'Tables and views' section with a 'Create' button and a dropdown menu. Under 'Tables and views', there are two entries: 'csv_reports' and 'parquet_parquet_reports'. The main area is titled 'Query 1' and contains the following SQL code:

```

1 SELECT rank, company, sector, industry, marketcap_march28_m
2 FROM parquet_parquet_reports
3 ORDER BY marketcap_march28_m DESC
4 LIMIT 10;

```

Below the code, it says 'SQL Ln 5, Col 1'. At the bottom of the editor are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. To the right, there's a note about 'Reuse query results up to 60 minutes ago'. The footer includes links for 'CloudShell', 'Feedback', and copyright information.

Fig.5.3.3 This is the query analysis performed in Athena

The screenshot shows the results of the query execution. The top bar indicates the query is 'Completed' with a green status bar. It also shows 'Time in queue: 65 ms', 'Run time: 545 ms', and 'Data scanned: 103.16 KB'. There are buttons for 'Copy' and 'Download results CSV'. The results table is titled 'Results (10)' and has columns: '#', 'rank', 'company', 'sector', 'industry', and 'marketcap_march28_m'. The data rows are:

#	rank	company	sector	industry	marketcap_march28_m
1	13	Microsoft	Technology	Computer Software	3126133.1
2	13	Microsoft	Technology	Computer Software	3126133.1
3	3	Apple	Technology	Computers, Office Equipment	2647973.8
4	3	Apple	Technology	Computers, Office Equipment	2647973.8
5	65	Nvidia	Technology	Semiconductors and Other Electronic Components	2258900.0
6	65	Nvidia	Technology	Semiconductors and Other Electronic Components	2258900.0
7	8	Alphabet	Technology	Internet Services and Retailing	1884633.0
8	8	Alphabet	Technology	Internet Services and Retailing	1884633.0
9	2	Amazon	Retailing	Internet Services and Retailing	1873675.8
10	30	Meta Platforms	Technology	Internet Services and Retailing	12357940.1

The footer includes links for 'CloudShell', 'Feedback', and copyright information.

Fig.5.3.4 This is the query results that are generated in Athena

The screenshot shows the AWS Athena Query editor interface. On the left, there's a sidebar titled 'Data' with dropdowns for 'Data source' (AwsDataCatalog), 'Catalog' (None), and 'Database' (fortune_reports). Below this is a 'Tables and views' section with a 'Create' button and a table listing 'Tables (2)': csv_reports and parquet_parquet_reports. The main area is titled 'Query 1' and contains the following SQL code:

```

1 SELECT rank, company, sector, profits_m
2 FROM parquet_parquet_reports
3 ORDER BY profits_m DESC
4 LIMIT 10;

```

Below the code, the status bar shows 'SQL Ln 1, Col 1'. At the bottom of the editor are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. To the right, there's a note about reusing query results up to 60 minutes ago.

Fig.5.3.5 This is the query analysis performed in Athena

The screenshot shows the results of the query execution. The top part of the interface is identical to Fig.5.3.5. The results are displayed in a table titled 'Results (10)'. The table has columns: #, rank, company, sector, and profits_m. The data is as follows:

#	rank	company	sector	profits_m
1	3	Apple	Technology	96995.0
2	3	Apple	Technology	96995.0
3	5	Berkshire Hathaway	Financials	96223.0
4	8	Alphabet	Technology	73795.0
5	8	Alphabet	Technology	73795.0
6	13	Microsoft	Technology	72361.0
7	13	Microsoft	Technology	72361.0
8	12	JPMorgan Chase	Financials	49552.0
9	12	JPMorgan Chase	Financials	49552.0
10	30	Meta Platforms	Technology	39098.0

Fig.5.3.6 This is the query results that are generated in Athena

The screenshot shows the Amazon Athena Query editor interface. On the left, there's a sidebar titled 'Data' with sections for 'Data source' (set to 'AwsDataCatalog'), 'Catalog' (set to 'None'), and 'Database' (set to 'fortune_reports'). Below this is a 'Tables and views' section with a 'Create' button and a dropdown menu. A table listing 'Tables (2)' is shown, with 'csv_reports' and 'parquet_parquet_reports' selected. The main area is titled 'Query 1' and contains the following SQL code:

```

1 SELECT rank, company, sector, industry, change_in_rank
2 FROM parquet_parquet_reports
3 WHERE change_in_rank > 0
4 ORDER BY change_in_rank DESC;

```

Below the code, it says 'SQL Ln 5, Col 1'. At the bottom of the editor are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. To the right, there's a checkbox for 'Reuse query results up to 60 minutes ago'. The footer includes links for 'CloudShell', 'Feedback', and copyright information.

Fig.5.3.7 This is the query analysis performed in Athena

The screenshot shows the Amazon Athena Query results page. At the top, it says 'Views (0)'. Below that are tabs for 'Query results' (which is selected) and 'Query stats'. The 'Query results' tab shows a green bar indicating the query is 'Completed' with a time of 65 ms, a run time of 619 ms, and data scanned of 105.16 KB. It also has 'Copy' and 'Download results CSV' buttons. The results table is titled 'Results (1,010)' and has columns: '#', 'rank', 'company', 'sector', 'industry', and 'change_in_rank'. The data is as follows:

#	rank	company	sector	industry	change_in_rank
1	182	First Citizens BancShares	Financials	Commercial Banks	420.0
2	182	First Citizens BancShares	Financials	Commercial Banks	420.0
3	387	Las Vegas Sands	Hotels, Restaurants & Leisure	Hotels, Casinos, Resorts	326.0
4	387	Las Vegas Sands	Hotels, Restaurants & Leisure	Hotels, Casinos, Resorts	326.0
5	188	KKR	Financials	Securities	308.0
6	418	FM Global	Financials	Insurance: Property and Casualty (Stock)	290.0
7	473	Interactive Brokers Group	Financials	Securities	263.0
8	473	Interactive Brokers Group	Financials	Securities	263.0
9	681	Agilon Health	Health Care	Health Care: Pharmacy and Other Services	260.0

Fig.5.3.8 This is the query results that are generated in Athena

5.4 Visualization Testing:

Data visualization is essential for interpreting analytical results effectively. Amazon QuickSight was used to generate interactive dashboards, incorporating various chart types to represent different aspects of the dataset.

Key Considerations:

- Visualizations should accurately reflect data insights without discrepancies.
- The dashboard should support interactive exploration, enabling users to filter and analyze data dynamically.

Testing Approach:

1. Configured Amazon QuickSight to connect with Athena and fetch query results.
2. Created multiple visualizations, including bar charts, pie charts, donut charts, KPIs, area charts, and line charts, to represent different financial and operational metrics.
3. Validated the correctness of visualizations by cross-referencing displayed values with query results.

Observations:

- All visual elements accurately represented the dataset, ensuring alignment with query results.
- The dashboard provided interactive exploration capabilities, allowing users to filter and analyze data seamlessly.
- Performance remained stable, ensuring smooth visualization updates without lag.

Result: The visualization phase successfully translated analytical insights into meaningful graphical representations, enhancing data interpretation and decision-making.

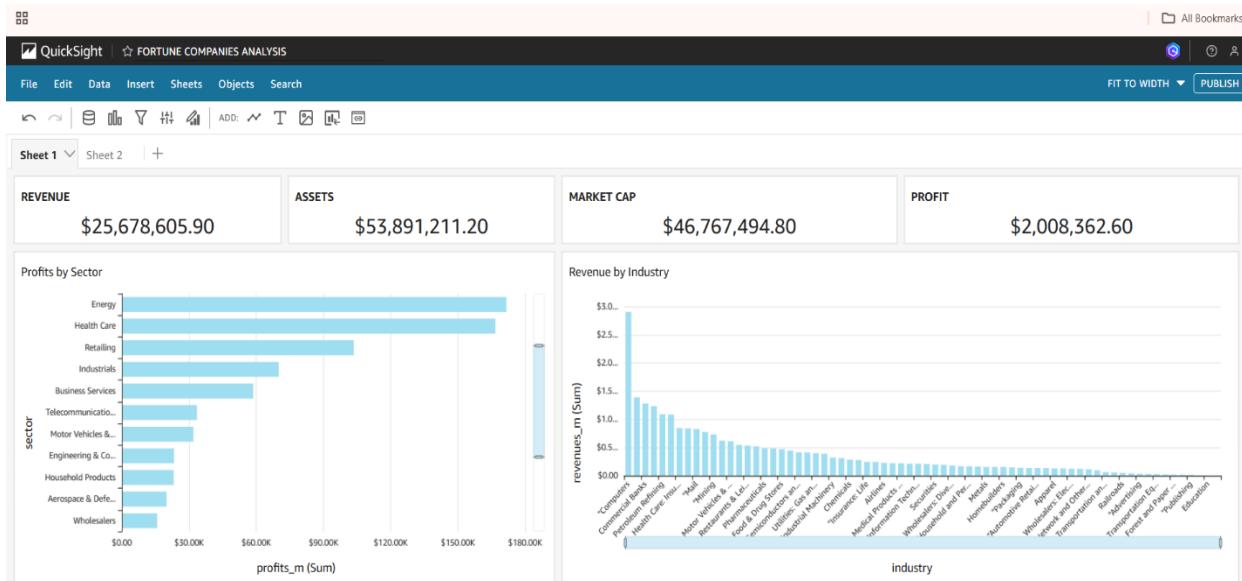


Fig.5.4.1 The data visualization in Quicksight which represents the analysis of key performance indicators related to companies using bar charts.

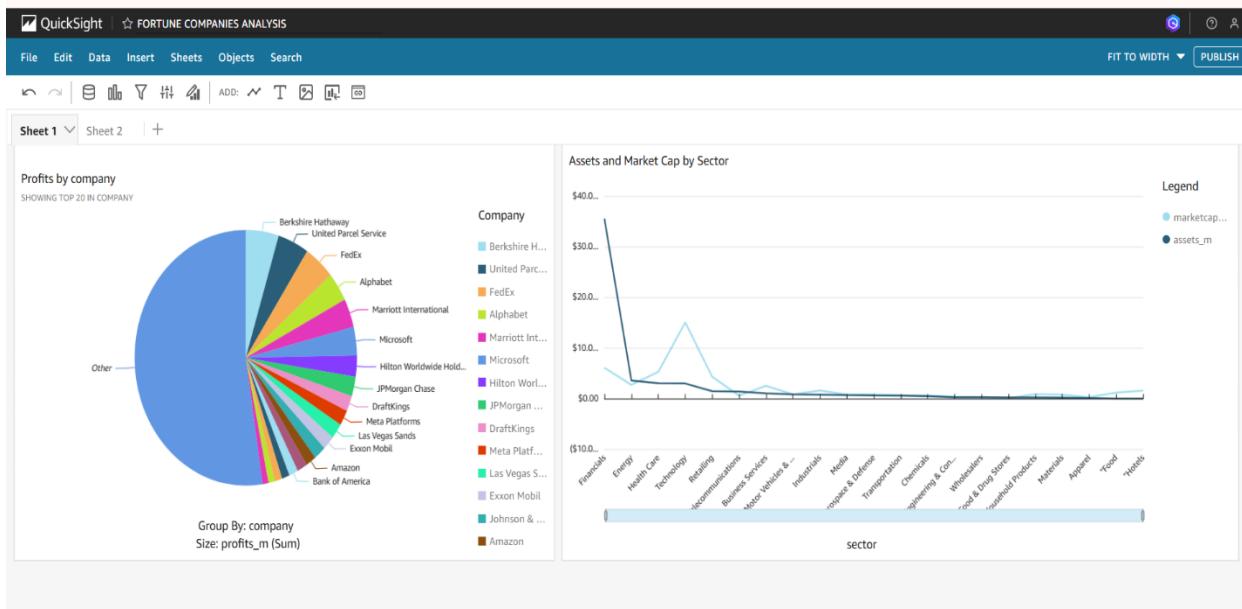


Fig.5.4.2 This represents the analysis of profits by company and assets, market capital by sector performed in Quicksight using pie chart and stacked area chart.

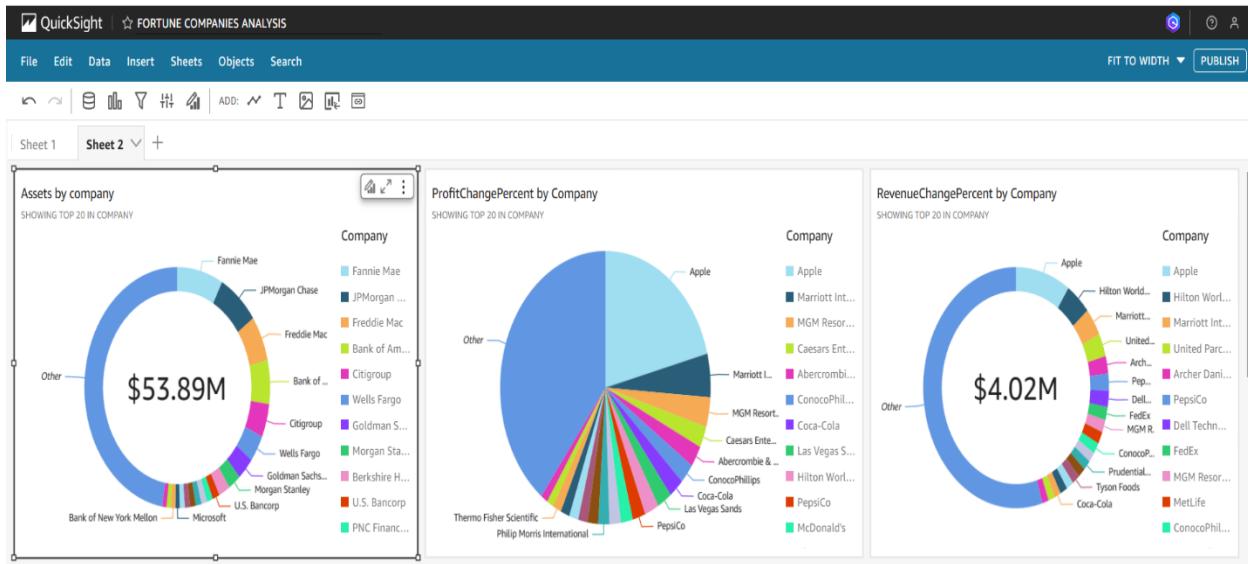


Fig.5.4.3 This represents the analysis of assets, profits percent change and revenue percent change by company using pie chart and donut chart in Quicksight

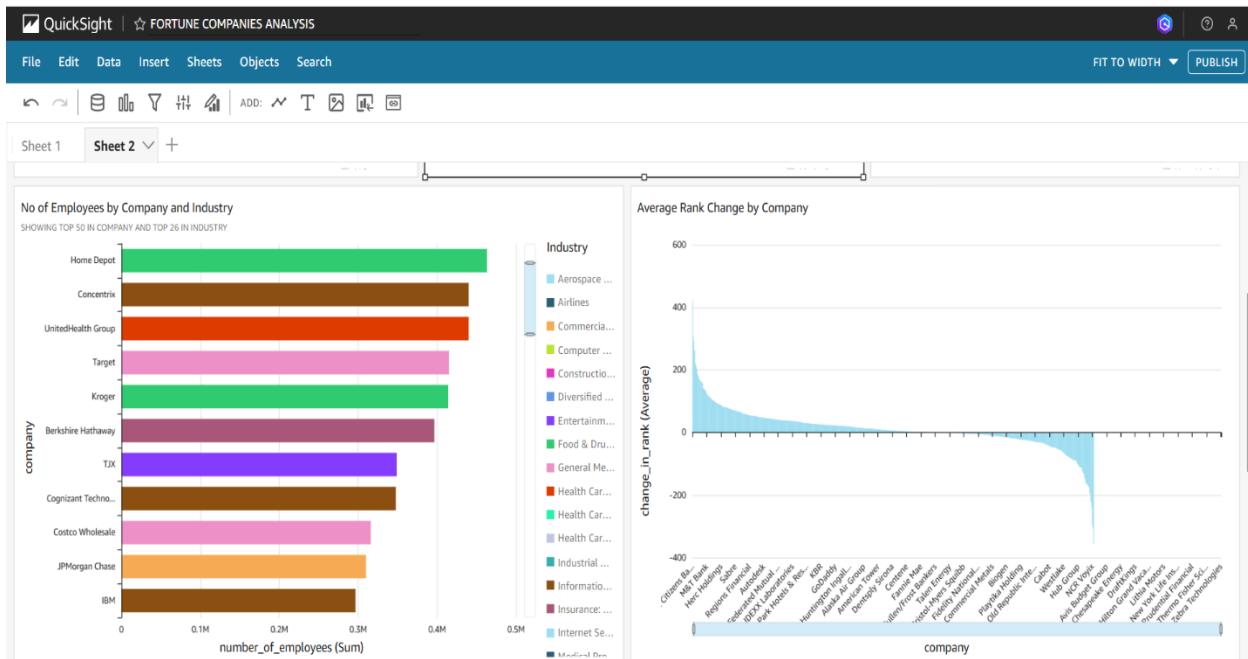


Fig.5.4.4 This represents the analysis about number of employees by company, industry and change in rank of company that are performed using bar graph and area chart in Quicksight

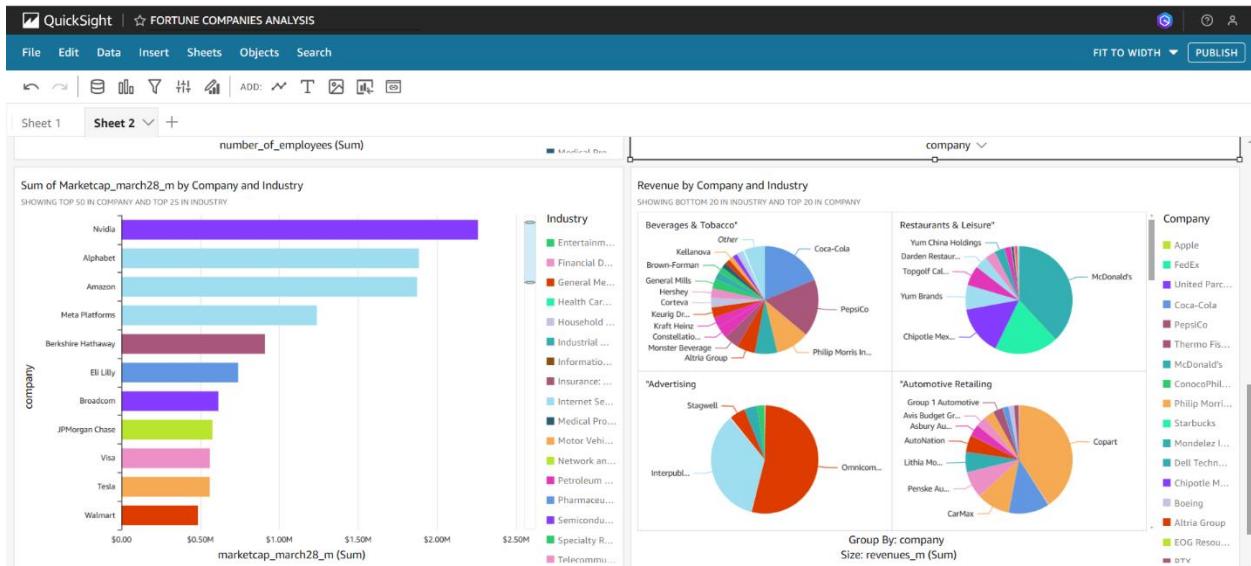


Fig.5.4.5 This represents the analysis of market capital and revenue by company, industry that are performed using bar graph and pie chart in Quicksight

5.5 Results & Performance Comparison

5.5.1 Performance Metrics:

- Data Ingestion:** The AWS Glue Crawler successfully inferred the schema of the Fortune 1000 dataset in under 5 minutes.
- ETL Process:** The transformation of raw CSV data into Parquet format reduced storage size and improved query performance.
- Query Processing:** Queries on Parquet data were faster compared to raw CSV data, with an average query execution time.
- Visualization:** Dashboards in Amazon QuickSight loaded in under 5 seconds, providing real-time insights into the dataset.

5.5.2 Comparison with Traditional Tools:

- Scalability:** AWS tools (S3, Glue, Athena) outperformed traditional tools like SQL and Spark in handling large datasets.
- Cost:** The pay-as-you-go pricing model of AWS reduced infrastructure costs compared to on-premise solutions.
- Ease of Use:** AWS Glue and QuickSight provided a more user-friendly interface.

6. CONCLUSION & FUTURE WORK

6.1 Conclusion

In this project, we explored AWS services such as S3, Glue, Athena, QuickSight. We worked with S3 buckets and objects to store our data, and worked with AWS Glue to create Glue Crawlers and build a Glue Data Catalog. This allowed us to generate tables and define schemas for our dataset. We also configured IAM roles, policies and permissions on who can access the data. We used Athena for querying data and QuickSight for visualizations. This provided us with a solid understanding of foundational AWS services, preparing us to efficiently utilize these tools for scalable data analytics and visualization in aws cloud environment.

6.2 Future Work

- **Real-Time Data Processing:** Integrate Amazon Kinesis for real-time data streaming and analysis.
- **Machine Learning Integration:** Use Amazon SageMaker to build predictive models for forecasting company performance.
- **Advanced Visualizations:** Explore advanced visualization techniques in Amazon QuickSight, such as geospatial analysis and anomaly detection.
- **Multi-Region Deployment:** Extend the solution to support multi-region data storage and processing for global datasets.
- **Cost Optimization:** Implement cost-monitoring tools like AWS Cost Explorer to further optimize resource usage.

7. REFERENCES

- [1] Mishra, A. and Kumar, G. (2021) ‘Big Data Analytics on AWS Cloud Using AWS Athena and QuickSight’, *International Journal of Engineering Research & Technology (IJERT)*, 10(4). Available at: www.ijert.org.
- [2] Caliwag, E.M.F. and Caliwag, A. (2021) ‘AWS Data Visualization using DynamoDB and Lambda’, *Conference Paper*, July 2021. Available at: [ResearchGate](#).
- [3] Satyanarayana, G., Bhuvana, J. and Balamurugan, M. (2020) ‘Sentiment Analysis on Voice Using AWS Comprehend’, *2020 International Conference on Computer Communication and Informatics (ICCCI -2020)*, IEEE. DOI: 10.1109/ICCCI48352.2020.
- [4] Wavhale, V.D., Bira, S., Kumar, V. and Choudhari, V.R. (2020) ‘Weather Data Forecast and Analytics’, *International Research Journal of Engineering and Technology (IRJET)*, 07(08), pp. 651–652.
- [5] Kulkarni, A. (2023) ‘Amazon Athena: Serverless Architecture and Troubleshooting’, *International Journal of Computer Trends and Technology*, 71(5), pp. 57–61. DOI: [10.14445/22312803/IJCTT-V71I5P110](https://doi.org/10.14445/22312803/IJCTT-V71I5P110).
- [6] Patsidis, A., Dyśko, A., Booth, C., Rousis, A.O., Kalliga, P., and Tzelepis, D. (2023) ‘Digital Architecture for Monitoring and Operational Analytics of Multi-Vector Microgrids Utilizing Cloud Computing, Advanced Virtualization Techniques, and Data Analytics Methods’, *Energies*, 16, p. 5908. DOI: [10.3390/en16165908](https://doi.org/10.3390/en16165908).
- [7] Gupta, A., Dhanda, N. and Gupta, K.K. (2023) ‘Ingest and Visualize CSV Files using AWS Platform For Transition from Unstructured to Structured Data’, *11th International Conference on Emerging Trends in Engineering & Technology - Signal and Information Processing (ICETET - SIP)*, IEEE.

- [8] Kaniganti, S.T. (2020) ‘Architecting Real-Time Big Data Analytics: An AWS-Powered Framework Integrating AI and ML for Predictive Insights’, *International Journal of Science and Research (IJSR)*, 9(2). DOI: [10.21275/SR24716230925](https://doi.org/10.21275/SR24716230925).
- [9] Rahman, M.M. and Hasan, M.H. (2019) ‘Serverless Architecture for Big Data Analytics’, *Global Conference for Advancement in Technology (GCAT)*, IEEE.
- [10] Medium (n.d.) "Data Flow Diagram of the Project." Available at Medium: https://miro.medium.com/v2/resize:fit:705/1*Sw_L7u1TKza2CK-gFKTJxQ.png.
- [11] Pochetti, F. (2020) "AWS System Architecture." Available at Francesco Pochetti Blog: <https://francescopochetti.com/wp-content/uploads/2020/03/AWS-2.png>.
- [12] Duval, J.N. (2024) "2024 Fortune 1000 Companies Dataset." Available at Kaggle: <https://www.kaggle.com/datasets/jeannicolasduval/2024-fortune-1000-companies>.