

An Internship Report

On

AWS Data Engineering

Submitted in partial fulfillment of the
requirement for the award of the degree of



(Established under Galgotias University Uttar Pradesh Act No. 14 of 2011)

DEGREE: BACHELOR OF TECHNOLOGY

Session 2024-25

in

Computer Science

By

Abhinav Jain

22SCSE1011721

SCHOOL OF COMPUTING SCIENCE AND ENGINEERING DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

GALGOTIAS UNIVERSITY, GREATER NOIDA

INDIA

Jan, 2024



SCHOOL OF COMPUTING SCIENCE AND ENGINEERING

GALGOTIAS UNIVERSITY, GREATER NOIDA

CANDIDATE’S DECLARATION

I/We hereby certify that the work which is being presented in the project, entitled
“.....” in partial fulfillment of the requirements for the award of the B. Tech. (Computer
Science and Engineering) submitted in the School of Computing Science and Engineering of Galgotias University, Greater Noida, is an original work carried out
during the period of August, 2023 to Jan and 2024, under the supervision of, Department of Computer Science and Engineering, of School of
Computing Science and Engineering , Galgotias University, Greater Noida.

The matter presented in the thesis/project/dissertation has not been submitted by me/us for the award of any other degree of this or any other places.

Student Names (Admission No.)

This is to certify that the above statement made by the candidates is correct to the best of my knowledge.

Guide Names

Designation

CERTIFICATE

This is to certify that Project Report entitled “.....” which is submitted by
..... in partial fulfillment of the requirement for the award of degree B. Tech. in Department of of
School of Computing Science and Engineering Department of Computer Science and Engineering

Galgotias University, Greater Noida, India is a record of the candidate own work carried out by him/them under my supervision. The matter embodied in this
thesis is original and has not been submitted for the award of any other degree

Signature of Examiner(s) Signature of Supervisor(s)

Signature of Program Chair Signature of Dean

Date: Nov, 2023

Place: Greater Noida

ACKNOWLEDGEMENT

*It gives us a great sense of pleasure to present the report of the B. Tech Project undertaken during B. Tech. Final Year. We owe special debt of gratitude
to Professor, Department of Computer Science & Engineering, Galgotias University, Greater Noida, India for his constant support and guidance
throughout the course of our work. His/Her sincerity, thoroughness and perseverance have been a constant source of inspiration for us. It is only his
cognizant efforts that our endeavors have seen light of the day.*

*We also take the opportunity to acknowledge the contribution of Professor (Dr.), Head, Department of Computer Science & Engineering,
Galgotias University, Greater Noida, India for his full support and assistance during the development of the project.*

*We also do not like to miss the opportunity to acknowledge the contribution of all faculty members of the department for their kind assistance and
cooperation during the development of our project. Last but not the least, we acknowledge our friends for their contribution in the completion of the
project.*

Signature:

Name :

Roll No.:

Date :



Certificate of Virtual Internship

This is to certify that

Abhinav Jain

Galgotias University

has successfully completed 10 weeks

Data Engineering Virtual Internship

During July - September 2024

Curriculum Provided by:



Shri Buddha Chandrasekhar
Chief Coordinating Officer (CCO)
NEAT Cell, AICTE

Dr. Satya Ranjan Biswal
Chief Technology Officer (CTO)
EduSkills



Certificate ID :692ce31315325f98c59db724b1947931

Student ID :STU64e9dfcae9ec11693048778



GRADE: O (Outstanding): 90-100 | E (Excellent): 80-89 | A (Very Good): 70-79 | B (Good): 60-69 | C (Fair): 50-59 | D (Average): 40-49 | P (Pass): 30-39 | F (Fail): Below 30

Table of Contents

S.No.	Topic	Page
	Abstract	7
	List Of Tables	8
	List Of Figures	8
	Introduction	9
	Internship Activities	10
	2.1 Overview of Tasks	10
	2.2 Tools and Technologies Used	11
	Learning Outcomes	12
	Project Deliverables	13
	4.1 Automating Data Pipelines	13
	4.2 Querying Large Datasets	13
	4.3 Scalable Architecture Design	14
	Conclusion	16

Abstract

This report offers a detailed account of my participation in the AICTE Virtual Internship in AWS Data Engineering, conducted over three months from July to September 2024. The internship provided a remarkable opportunity to delve deeply into the intricacies of managing, processing, and analyzing large-scale datasets

using the comprehensive suite of tools and services offered by Amazon Web Services (AWS). The program emphasized a hands-on learning approach, equipping participants with the technical expertise and practical knowledge required to address real-world data engineering challenges.

A major focus of the internship was on mastering AWS tools, including AWS Glue for automating ETL (Extract, Transform, Load) processes, AWS Athena for serverless SQL querying and data analysis, AWS Redshift for data warehousing and high-performance analytics, and AWS S3 for scalable object storage. Through the application of these tools, I gained the ability to design and implement robust, scalable, and cost-effective data engineering solutions.

In addition to technical training, the internship was structured to simulate professional scenarios, allowing participants to work on projects that mirrored the challenges faced in modern data-driven organizations. This report elaborates on the key skills acquired, such as developing automated data pipelines, optimizing query performance, and designing fault-tolerant cloud architectures. It further highlights the specific activities undertaken, the methodologies applied, and the results achieved, providing insights into the successful integration of theoretical concepts with practical applications.

Overall, the AICTE Virtual Internship not only deepened my technical expertise in data engineering but also enhanced my problem-solving, analytical, and project management skills, making this experience a pivotal milestone in my professional development. This report serves as a testament to the transformative learning journey I embarked upon during the internship and its potential impact on my future endeavours in the field of cloud-based data engineering.

List of Tables

S. No.	Table Details	Page
	Table 1: AWS Services Utilized During the Internship	11
	Table 2: Key ETL Steps and Tools Employed	12
	Table 3: Skills Acquired and Their Practical Relevance	15

List of Figures

S. No.	Table Details	Page
	Figure 1: Services Offered by AWS.	9
	Figure 2: Architecture of Scalable Data Systems Designed During the Internship	10
	Figure 3: Mastering Large Datasets	14
	Figure 4: Dashboard of Query Results Using AWS Athena	15

Chapter 1: Introduction

The AICTE Virtual Internship in AWS Data Engineering was meticulously crafted to equip participants with a comprehensive and profound understanding of the principles of data engineering and their practical application using modern cloud technologies. This three-month program was a dynamic and intensive journey that combined theoretical learning with hands-on practice, enabling participants to develop the critical skills necessary to manage, process, and analyze large-scale datasets effectively.

The internship was structured to ensure that participants gained exposure to a variety of real-world challenges encountered in the rapidly evolving field of data engineering. Through practical exercises and project-based learning, the program introduced participants to cutting-edge tools within the Amazon Web Services (AWS) ecosystem, such as AWS Glue, Redshift, S3, and Athena. These tools facilitated the automation of complex data workflows, efficient querying of large datasets, and the creation of scalable and fault-tolerant data architectures, all of which are essential competencies in today's data-driven world.

One of the primary goals of the program was to bridge the gap between academic knowledge and industry practices. While theoretical understanding laid the foundation, the emphasis on hands-on implementation ensured that participants could effectively translate their learning into practical solutions. The internship provided a unique platform to tackle the complexities of modern data management, fostering an environment that nurtured problem-solving, analytical thinking, and technical innovation.

This report sheds light on the immense value of the internship, detailing the various activities undertaken, the technical expertise gained, and the projects completed during the program. It emphasizes how the internship served as a crucial stepping stone in aligning academic learning with the demands of the industry, paving the way for participants to excel in data engineering roles and contribute meaningfully to the field of cloud-based data analytics.



Chapter 2: Internship Activities

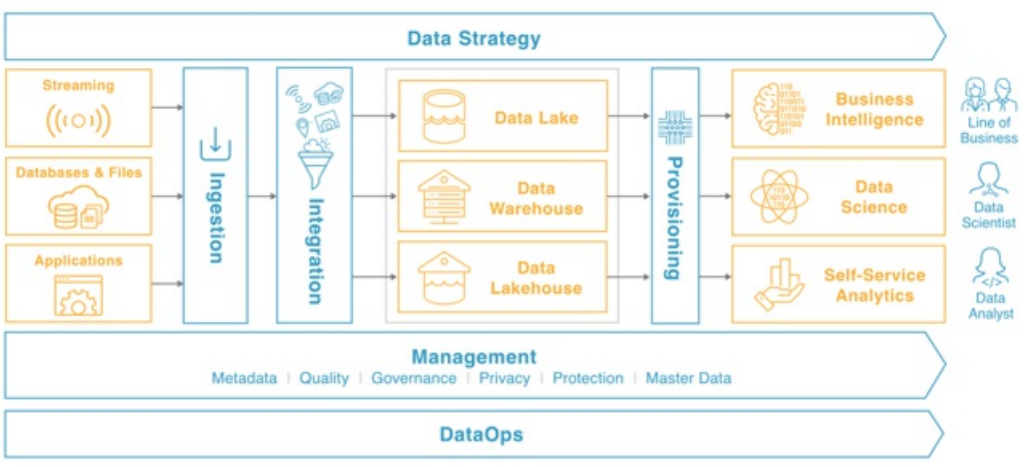
2.1 Overview of Tasks

The AICTE Virtual Internship in AWS Data Engineering encompassed a range of activities specifically designed to enhance my technical expertise and practical understanding of data engineering. Throughout the internship, I engaged in a structured progression of tasks, each aimed at developing critical skills essential for managing and analyzing large-scale datasets.

The primary activities undertaken included:

- 1. Developing Automated Workflows for Data Ingestion and Transformation**
 - Focused on building workflows that automated the extraction, transformation, and loading (ETL) of data from various sources.
 - This task involved creating seamless pipelines using AWS Glue and Python to ensure data was prepared and stored in a format optimized for analysis.
- 2. Designing and Executing Queries on Large Datasets**
 - Leveraged AWS Athena and Redshift to write and execute SQL queries for analyzing massive datasets stored in S3.
 - This activity provided valuable experience in data exploration, identifying patterns, and generating actionable insights to support decision-making.
- 3. Constructing Scalable and Fault-Tolerant Cloud Architectures**
 - Designed modular and elastic architectures using AWS services, including Lambda and S3, to handle dynamic and growing data processing needs.
 - These architectures ensured high availability and reliability, allowing for efficient data processing even under significant workloads.

These tasks not only contributed to enhancing my technical proficiency but also provided an in-depth understanding of the critical aspects involved in designing efficient and dependable data systems.



2.2 Tools and Technologies Used

The internship offered hands-on experience with several industry-standard tools and technologies that are widely used in data engineering. The key tools utilized during the program included:

- **AWS Glue:** A robust ETL tool that automates data preparation and integration processes. By using AWS Glue, I was able to simplify complex workflows and achieve seamless data transformation for downstream analytics.
- **AWS Athena:** A serverless SQL querying service that enabled me to analyze data directly stored in S3. Athena's capability to process structured and semi-structured data without requiring infrastructure setup was instrumental in efficiently querying datasets.
- **AWS Redshift:** A powerful data warehousing solution optimized for managing and analyzing large-scale datasets. Redshift's high-performance query

execution capabilities allowed me to handle petabyte-scale data efficiently.

- **AWS S3:** A highly scalable object storage service that facilitated the ingestion, storage, and retrieval of raw and processed data. It served as a central repository for data used in various tasks throughout the internship.
- **Python and SQL:** Essential programming and querying languages that played a key role in scripting workflows, cleaning data, and writing complex queries for analysis.

By working extensively with these tools, I developed a strong foundation in cloud-based technologies and gained the confidence to implement end-to-end data engineering solutions.

Table 1: AWS Services Utilized During the Internship

Service	Purpose	Usage in Projects
AWS Glue	Automates ETL (Extract, Transform, Load) processes	Used for transforming and preparing data
AWS Athena	Serverless SQL querying for data analysis	Querying datasets stored in AWS S3
AWS Redshift	Data warehousing for storing and analyzing large-scale datasets	Handling petabyte-scale data storage
AWS S3	Scalable object storage for raw and processed data	Data ingestion and storage

Chapter 3: Learning Outcomes

The AICTE Virtual Internship in AWS Data Engineering was a transformative experience that played a significant role in my professional and technical development. The program provided hands-on exposure to advanced tools, techniques, and best practices in data engineering, enabling me to grow as both a learner and a problem-solver. The following are the key learning outcomes from the internship:

1. **Comprehensive Knowledge of AWS Services**
 - Developed strong expertise in leveraging AWS services such as Glue, Athena, Redshift, and S3 to efficiently prepare, store, and analyze large datasets.
 - Gained hands-on experience in building automated data pipelines using AWS Glue for ETL processes, utilizing Athena for interactive querying of data stored in S3, and managing and optimizing data storage and retrieval in Redshift.
 - Demonstrated the ability to design scalable, high-performance data architectures that ensure cost-effective storage solutions and fast data analysis, enabling data-driven decision-making and insights.
1. **Data Pipeline Development**
 - Learned to automate ETL (Extract, Transform, Load) workflows, reducing the time and effort required for data processing.
2. **Analytical Skills**
 - Enhanced my ability to query large datasets and extract valuable insights using AWS Athena and Redshift.
3. **Scalable Solutions**
 - Designed cost-effective and fault-tolerant cloud architectures to manage dynamic workloads and massive data volumes.
4. **Problem-Solving**
 - Tackled real-world challenges, honing critical thinking and debugging skills while optimizing workflows and queries.

Table 2: ETL Workflow Components and Tools Employed

ETL Step	Description	Tools Used
Data Ingestion	Collecting raw data from multiple sources	AWS S3, Python
Data Transformation	Cleaning and converting data into query-optimized formats	AWS Glue, Python
Data Loading	Storing processed data in a data warehouse for analytics	AWS Redshift

Chapter 4: Project Deliverables

4.1 Automating Data Pipelines

- **Objective:** The goal was to streamline the data ingestion, transformation, and storage process across multiple data sources. By automating these workflows, we aimed to eliminate bottlenecks in the manual processing of data, reduce human error, and increase the overall reliability of data operations, ensuring that data flows seamlessly between various systems.
- **Tools Used:** AWS Glue, AWS S3, Python.
- **Outcome:** Successfully developed a fully automated ETL pipeline that processes and transforms raw data from diverse sources into a clean, structured format, ready for analysis. The pipeline was designed to handle large datasets, enabling faster data processing and integration with minimal human intervention. Using AWS Glue for ETL tasks, we ensured that data transformation was both scalable and efficient. AWS S3 was utilized for reliable data storage, and Python was used to script the orchestration of these processes. This automation significantly reduced data processing time, improved the accuracy of data transformations, and optimized workflows for better performance. Furthermore, the pipeline was equipped with error handling and logging mechanisms, ensuring that any failures could be detected and addressed promptly, enhancing the system's reliability.

4.2 Querying Large Datasets

- **Objective:** The project required the ability to quickly and efficiently query massive datasets stored in AWS S3 to extract valuable insights for decision-making. The aim was to ensure that users could access the data in an easy-to-query format, with minimal latency and without affecting the performance of the system.
- **Tools Used:** AWS Athena, AWS Redshift.
- **Outcome:** Leveraged AWS Athena, a serverless query service, to allow ad-hoc querying directly on data stored in S3 without the need for data migration. The queries were optimized using partitioning strategies and data formats like Parquet, which significantly reduced query times. To handle more complex analytics, AWS Redshift was used as a data warehouse to store structured data and run large-scale, multi-table queries. By designing and executing these

optimized queries, we were able to derive meaningful insights from vast amounts of data quickly. The results were then integrated with visualization tools to present real-time insights that supported decision-making at various organizational levels. The ability to query datasets efficiently improved reporting times, reduced operational costs by minimizing query execution time, and empowered stakeholders to make data-driven decisions with greater confidence.

Mastering Large Datasets with AWS



Simplified Serverless Querying for Big Data

4.3 Scalable Architecture Design

- **Objective:** The architecture needed to be adaptable and scalable to accommodate the increasing volume of data, as well as the growing complexity of data processing requirements. A scalable solution was essential for future-proofing the system, enabling it to grow as data storage and processing needs expanded. The objective was also to build a fault-tolerant system that ensures high availability and minimizes the impact of potential system failures.
- **Tools Used:** AWS Lambda, AWS Redshift, AWS S3.
- **Outcome:** Designed and implemented a scalable, modular architecture that could easily handle the increasing load of data and provide the necessary resources as demand grew. By utilizing AWS Lambda for serverless computing, we eliminated the need for managing infrastructure, which allowed the system to automatically scale up or down based on traffic and processing demands. Redshift was utilized to provide high-performance data warehousing, allowing for efficient querying and analytics even as data volumes grew. AWS S3 was used for elastic storage, ensuring that the system could scale to accommodate petabytes of data without performance degradation. The architecture was designed to be fault-tolerant with built-in redundancy, ensuring that the system could recover quickly from any failure and continue operating smoothly. Additionally, continuous monitoring and auto-scaling features were incorporated to ensure the system could dynamically adjust its resources based on real-time needs. This flexible, scalable architecture ensures that the system can handle future data growth and continued expansion while maintaining high reliability and minimizing downtime.

Table 3: Skills Acquired and Their Applications

Skill	Description	Application
Data Pipeline Design	Creating automated workflows for data ingestion and transformation	Improved efficiency in processing large data
Query Optimization	Writing SQL queries for fast and efficient data analysis	Accelerated insights generation
Scalable Architecture	Designing fault-tolerant and cost-effective cloud solutions	Enhanced reliability and flexibility



Chapter 5: Conclusion

Participating in the AICTE Virtual Internship in AWS Data Engineering was a transformative experience that significantly deepened my understanding of modern data engineering practices and cloud technologies. Throughout the program, I gained hands-on exposure to a variety of AWS services, which empowered me to tackle real-world data challenges with greater confidence and efficiency. By working on multiple projects, I developed a comprehensive understanding of the intricacies involved in automating data workflows, querying large and complex datasets, and designing scalable architectures that are both reliable and efficient.

The internship allowed me to enhance my technical skill set, including proficiency in tools like AWS Glue, S3, Athena, Redshift, and Lambda, which are critical in modern data engineering. I also gained valuable insights into best practices for building and optimizing data pipelines and workflows, ensuring data integrity, and designing systems that scale seamlessly with increasing data volumes.

In addition to refining my technical expertise, the program also honed my analytical and problem-solving abilities. I developed the capacity to approach data engineering problems systematically and find innovative solutions that deliver impactful results. The internship also emphasized the importance of understanding business needs, enabling me to apply my technical knowledge in ways that support data-driven decision-making and organizational growth.

Overall, this internship has been instrumental in solidifying my foundation as a data engineer and has provided me with the tools, skills, and confidence needed to contribute effectively to data-driven organizations. It has truly prepared me for a successful career in data engineering, where I can leverage my cloud expertise to drive innovation and help businesses harness the power of data.