

Original Article

Serverless ETL: Leveraging AWS Glue and PySpark for Efficient Data Processing

Dharanidhar Vuppu¹, Mounica Achanta²

¹Sr Data Engineer, SurveyMonkey, Texas, United States of America.

²Independent Research at IEE, Texas, United States of America.

¹Corresponding Author : dharanidhar@ieee.org

Received: 03 June 2025

Revised: 26 June 2025

Accepted: 18 July 2025

Published: 29 July 2025

Abstract - In today's cloud-native data landscape, data engineers are expected to build ETL pipelines that can scale effortlessly, remain easy to maintain, and stay within budget. With data volumes growing rapidly and business needs constantly evolving, traditional ETL setups—typically run on provisioned clusters—can become a bottleneck. They often bring challenges like over-provisioned resources, ongoing infrastructure upkeep, and complicated scaling mechanisms. This paper explores a serverless approach using AWS Glue and PySpark, aimed at simplifying ETL development while cutting down significantly on operational complexity.

We share a hands-on implementation of a serverless ETL setup that takes advantage of AWS Glue's built-in orchestration, Spark-based distributed processing, and tight integration with the AWS Data Catalog for managing schemas. This approach simplifies the process of ingesting and transforming data from sources like S3 and RDS, cuts down on setup time, and scales effortlessly without the need for manual tuning.

Through a real-world case study, we benchmark AWS Glue's performance, scalability, and cost-efficiency against traditional Spark clusters hosted on EC2. The results show tangible benefits in terms of time-to-value, fault tolerance, and operational simplicity, particularly for mid-sized batch processing workloads. The paper concludes with practical considerations, limitations, and lessons learned from adopting serverless ETL, offering guidance for data engineers looking to modernize their pipelines using fully managed, cloud-native solutions.

Keywords - Serverless ETL, PySpark, Big Data Processing, AWS Glue, Data Engineering.

1. Introduction

Modern data engineering has evolved beyond traditional ETL models, which were once tightly coupled to fixed infrastructure and required constant resource tuning and maintenance. As organizations accelerate their migration to the cloud and adopt real-time analytics, the demand for flexible, scalable, and low-maintenance data pipelines has grown significantly. In this landscape, serverless computing has emerged as a powerful paradigm, enabling engineers to focus on transformation logic and data modeling rather than provisioning and managing compute resources. (Pogiatzis & Samakovitis, 2020)

ETL pipelines, which remain at the core of most data engineering workflows, have seen a major shift with the advent of serverless frameworks like AWS Glue. AWS Glue abstracts away infrastructure concerns by providing a fully managed, Spark-based environment for running distributed

data processing jobs. It integrates tightly with the AWS ecosystem, allowing seamless access to object storage (S3), relational databases, and schema registries through the AWS Glue Data Catalog. Combined with PySpark, a widely adopted interface for Spark programming in Python, Glue presents an attractive solution for teams looking to build pipelines without the overhead of managing clusters.

While the value proposition of serverless ETL is clear—reduced operational complexity, autoscaling, pay-as-you-go pricing—its practical adoption raises questions about performance, cost-efficiency, debugging capabilities, and suitability for different workload patterns. These concerns are particularly relevant for data engineering teams working with semi-structured or evolving data formats, managing large-scale batch jobs, or integrating across multiple data systems.



This paper examines the architectural and operational benefits of using AWS Glue in combination with PySpark to implement robust serverless ETL pipelines. Drawing from hands-on experience, we walk through the design choices, transformation strategies, and optimization techniques that make Glue a viable alternative to traditional ETL setups. We also evaluate its strengths and limitations through a case study, benchmarking it against an EC2-hosted Spark deployment.

We aim to give data engineers practical, hands-on insights into adopting serverless ETL—showing where AWS Glue shines, its limitations, and how it fits into a modern data stack. By focusing on real-world use cases and the trade-offs engineers face in practice, we hope to share knowledge that's directly applicable to today's data workflows. Plale, B., & Kouper, I. (2017)

2. Evolution of ETL Architectures

As data pipelines become more complex and data volumes continue to grow, the ETL ecosystem has evolved rapidly. Innovations in cloud-native and serverless technologies have reshaped how teams think about data processing. Much of this area's recent work has centred on simplifying pipeline development, improving performance, cutting costs, and reducing maintenance effort. This section covers the foundational ideas and technologies that laid the groundwork for serverless ETL and looks at how AWS Glue fits into this evolving space. (Bussa & Hegde, 2024)

2.1. Traditional ETL with Spark and Hadoop Ecosystems

Apache Spark has long been a cornerstone for distributed data processing, valued for its speed and flexibility. Its in-memory execution model and support for both batch and streaming workloads made it a top choice for many ETL use cases. Traditionally, Spark was deployed on self-managed clusters using tools like Hadoop YARN, Apache Mesos, or Kubernetes. While these setups are powerful, they often come with a high operational burden—teams need to provision infrastructure, manage scaling, monitor performance, and handle failures, all of which can be time-consuming and resource-intensive. Traditionally, data engineers have deployed Spark on self-managed clusters using resource managers like Hadoop YARN, Apache Mesos, or Kubernetes. While this approach is powerful, it often comes with challenges such as:

- Manual resource tuning
- Cluster provisioning and autoscaling logic
- High DevOps overhead

Research such as Zaharia et al. (2016) emphasized Spark's performance advantages but acknowledged that operational complexity limits accessibility to smaller teams or projects with limited DevOps support.

2.2. Emergence of Cloud-Managed ETL Services

Cloud providers have attempted to abstract cluster management through managed services like:

- Amazon EMR (Elastic MapReduce): Offers Spark with simplified cluster provisioning but still requires node configuration and scaling strategies.
- Google Cloud Dataflow: Focuses on stream and batch processing with autoscaling, but introduces a learning curve with Apache Beam.
- Azure Data Factory: Provides a GUI-based orchestration layer, with some native data transformation capabilities, though lacking the flexibility of Spark for complex logic.

These platforms represent significant steps forward, but still fall short of a fully serverless experience, particularly when Spark-based transformations are involved. (Kriushanth, Arockiam, & Mirobi, 2013)

2.3. AWS Glue and Serverless ETL

AWS Glue was introduced to bridge this gap—offering a Spark-based environment that: (Sudhakar, 2018)

- Runs on demand with no cluster setup
- Integrates natively with S3, RDS, Redshift, and DynamoDB
- Provides schema management via the Glue Data Catalog
- Supports PySpark and dynamic frames for schema evolution and nested data handling

Recent whitepapers and AWS case studies have demonstrated Glue's effectiveness for handling large-scale batch jobs, especially when combined with features like job bookmarking, partition pruning, and DPU (Data Processing Unit) auto-allocation. However, Glue has also been critiqued for:

- Longer cold start times
- Limited job customization compared to EC2-based Spark
- Debugging complexity due to reliance on CloudWatch logs and lack of an interactive UI

2.4. Academic and Industry Insights into Serverless ETL

Academic literature on serverless ETL remains relatively sparse, though industry-driven evaluations are growing. For example:

- Netflix and Expedia have discussed moving to serverless pipelines using Glue and Lambda to reduce infrastructure friction.
- Research by Zhang et al. (2020) on "Serverless Dataflow Systems" evaluates performance trade-offs in serverless orchestration models, highlighting latency, cold starts, and cost granularity.

- Papers on "Function-as-a-Service" paradigms and "Data Engineering at Scale" provide foundational context for understanding the evolution of ETL architectures.

2.5. Positioning this Work

This paper builds on existing work by:

- Taking a hands-on engineering perspective to validate Glue's capabilities with real-world data flows.
- Benchmarking performance, cost, and operability against traditional Spark clusters.
- Providing a practical framework for other data engineering teams evaluating the move to serverless Spark ETL.

Rather than proposing a theoretical model, our goal is to add practitioner-level depth to the serverless ETL conversation—grounded in practical usage, architectural trade-offs, and lessons learned.

3. System Architecture

The architecture behind this serverless ETL pipeline is built with scalability, maintainability, and minimal operational overhead in mind—hallmarks of an efficient, modern data engineering setup. At its core, the pipeline uses AWS Glue to manage and run distributed data transformations written in PySpark, all within a fully serverless environment.

The design is modular and event-driven, and it integrates seamlessly with other AWS services for storage, schema management, and monitoring. (Mehmood & Anees, 2022)

3.1. Data Ingestion and Storage

The pipeline ingests raw data from a variety of sources:

- Amazon S3 serves as the primary storage for clickstream logs, CSV and Parquet files, and data exports from third-party systems.
- Amazon RDS and other JDBC-compatible sources are used to capture structured, transactional data.

All incoming data is staged in S3 buckets that are partitioned by event date. This approach helps optimize performance for downstream processing and querying.

3.2. Glue Data Catalog

The AWS Glue Data Catalog serves as the central hub for managing schemas. It supports schema discovery, validation, and versioning. Tables are either automatically registered using Glue crawlers or defined manually through job configurations. This ensures that all downstream jobs interact with consistent, well-documented metadata, reducing the chances of schema-related errors and improving pipeline reliability.

3.3. ETL Jobs (PySpark on Glue)

The heart of the architecture is the Glue job, which:

- Reads from the catalog using dynamic frames (useful for semi-structured or evolving schemas).
- Applies transformation logic in PySpark (e.g., joins, filters, column derivations).
- Writes the output to target destinations (e.g., S3, Snowflake, Redshift) in optimized formats such as Parquet or Delta.

Jobs are configured with appropriate DPUs (Data Processing Units) and leverage features like:

- Job bookmarking for incremental loads
- Partition pushdown for performance optimization
- Retry policies and CloudWatch logging for observability

3.4. Orchestration and Scheduling

Orchestration is handled via:

- AWS Glue Workflows for DAG-like execution with dependencies.
- Optionally, Amazon MWAA (Airflow) or Step Functions for cross-system orchestration.

3.5. Monitoring and Alerts

Operational telemetry is captured using the following:

- AWS CloudWatch Logs and Metrics
- Glue job metrics dashboard for DPU usage and job durations
- Custom alerts for job failures or SLA breaches

This architecture enables data engineers to build robust pipelines without worrying about infrastructure provisioning, while still retaining the power and flexibility of Apache Spark.

4. Implementation and Methodology

The decision to deploy a serverless ETL pipeline with AWS Glue and PySpark was driven by the need for simplified orchestration, scalable compute, and flexible transformation logic. This section explores the core design principles, the technology stack used, and the practical steps taken to build and optimize the ETL workflow efficiently. (Warneke & Kao, 2009)

4.1. Job Configuration and Environment Setup

AWS Glue jobs were built using version 3.0 or higher, which supports Spark 3.x and Python 3. Writing the jobs in PySpark allowed the team to take advantage of Spark's distributed processing while also benefiting from Python's clean syntax and rich ecosystem. (Batmaci, 2022)

Here are some key configuration choices that helped tailor performance to the workload:

- Worker Type: G.1X was used for standard workloads, while G.2X was chosen for jobs that required more memory.
- Number of DPUs: Adjusted depending on the size and complexity of the data to balance cost and performance.
- Job Bookmarks: Enabled to support incremental processing, ensuring only new or updated data was picked up in each run.

Jobs were version-controlled via Git and deployed using AWS Glue's script editor or through automated CI/CD pipelines using AWS CLI and boto3.

4.2. Data Ingestion Strategy

Data ingestion was source-dependent:

- For batch files on Amazon S3, the job used `glueContext.create_dynamic_frame.from_catalog()` to load data registered in the Glue Data Catalog.
- JDBC connections were established using connection objects defined within Glue for RDS and external systems. Credentials were managed via AWS Secrets Manager.

The ingestion logic included schema validation, null-handling, and deduplication steps using dynamic frame and Data Frame operations.

4.3. Transformation Logic in PySpark

The transformation layer was written in PySpark, incorporating:

- Joins between multiple datasets (e.g., user logs with metadata)
- Derivation of new columns based on business logic
- Conditional filtering, grouping, and aggregations
- Schema normalization (flattening nested fields, resolving schema drift)

Where schema evolution was expected, dynamic frames were used for flexibility. In other cases, data was cast into static Spark Data Frames to leverage stricter typing and better performance.

4.4. Output Handling and Storage

Transformed data was written back to Amazon S3 in columnar formats (Parquet or ORC), partitioned by date and other relevant keys. Data was optionally pushed to Snowflake or Amazon Redshift for analytical querying or BI reporting using external connectors.

To avoid redundant writes, job outputs were idempotent and often included:

- Overwrite modes for daily partitions
- Upserts are simulated using Spark's window functions and deduplication before writing

4.5. Optimization Techniques

To improve runtime and reduce cost:

- Partition pruning was applied during reads to limit scan size
- Predicate pushdown and `.repartition()` calls were used judiciously
- Glue job metrics were monitored regularly to tune DPU usage
- Logging was instrumented using `print()` statements and CloudWatch custom metrics

Jobs were profiled with sample datasets before scaling up to full runs, ensuring transformations were logically correct and performant.

4.6. Resilience and Monitoring

Error handling was implemented through:

- Try-catch blocks around transformation logic
- Job parameters for rerun flexibility
- Retry policies configured at the Glue job level

Glue job logs were streamed to CloudWatch, and alerts were set up for failures or SLA breaches. For production pipelines, job state and run metadata were stored in audit logs to track lineage and detect anomalies. This methodology reflects a balance between leveraging Glue's serverless advantages and applying Spark-native optimization practices. It enables repeatable, reliable data transformations without the burden of managing infrastructure, allowing data engineers to focus on business logic and data quality. (Singh P., 2021)

5. Case Study / Experiment

To evaluate the effectiveness of a serverless ETL approach in real-world data engineering scenarios, we developed and benchmarked a generic pipeline using AWS Glue and PySpark. The objective was to assess performance, scalability, and operational simplicity when handling large-scale batch data processing tasks without provisioning or managing infrastructure.

5.1. Objective

- Build a resilient, serverless ETL pipeline capable of ingesting structured and semi-structured data from cloud-based sources.
- Apply common transformation logic, including joins, aggregations, and type casting.
- Deliver clean, partitioned output ready for downstream analytics consumption.

5.2. Input and Output Characteristics

- Input Data:
 - Simulated daily batch of log and transaction data (~300–500 million records)

- Stored in Amazon S3 in CSV and JSON formats
- **Output Data:**
 - Transformed and normalized
 - Written back to S3 in partitioned Parquet format for efficient querying

5.3. Implementation Details

- **Ingestion**
 - Glue DynamicFrames are used to read schema-flexible data from the Glue Catalog
 - JDBC connectors are employed to pull structured reference data from cloud-based relational databases
- **Transformation Logic:**
 - Conversion from dynamic frames to Spark Data Frames for performance-critical operations
 - Applied transformations such as:
 - Complex joins across datasets
 - Type casting and schema enforcement
 - Row-level filtering, aggregations, and calculated columns
- **Job Configuration:**
 - Glue 3.0 environment with Spark 3.x and Python 3.x
 - Worker Type: G.1X for baseline testing; G.2X for larger loads
 - DPU allocation: 10–15 for moderate-volume test runs
 - Job bookmarks are enabled to support incremental data loads

5.4. Orchestration and Scheduling

- Glue Workflows are used for chaining multiple ETL jobs
- Triggered daily on a fixed schedule
- Optional integration tested with Amazon EventBridge and Airflow (MWAA) for custom orchestration logic

5.5. Observations and Insights

- **Performance:**
 - Daily batch processing (~400M records) completed within 10–12 minutes on average
 - Autoscaling and memory optimization are handled transparently by Glue
- **Scalability:**
 - Easily scaled up by adjusting DPU allocation without architectural changes
- **Operational Simplicity:**
 - No cluster setup or tuning required

- Schema evolution is handled gracefully using Glue's dynamic frames

- **Cost Management:**

- The pay-per-second billing model made it cost-effective for intermittent workloads

- **Monitoring**

- CloudWatch provided job-level visibility, though deeper Spark-level diagnostics were more limited than in traditional Spark setups

5.6. EMR Comparison Snapshot

To contextualize the benefits and limitations of AWS Glue, a high-level comparison was made against a typical self-managed Spark setup on Amazon EMR. The comparison reflects commonly observed characteristics across similar batch ETL use cases:

Factor	AWS Glue (Serverless Spark)	Amazon EMR (Self-managed Spark)
Setup & Provisioning	Fully managed, no infrastructure	Manual cluster provisioning is required
Autoscaling	Native and automatic	Configurable, but requires tuning
Cold Start Latency	Moderate (30–60 sec)	High (5–10 min for cluster spin-up)
Cost Model	Pay-per-second for job runtime	Pay-per-instance, even when idle
Performance Tuning	Abstracted from the user	Full control over Spark configs
Debugging & Logs	CloudWatch + basic logs	Spark UI, SSH access, full logs
Operational Overhead	Minimal	High (maintenance, patching, scaling)
Use Case Fit	Best for batch, periodic jobs	Better for persistent, fine-tuned jobs

5.7. Key Takeaways

- For batch ETL workloads that prioritize rapid development and minimal infrastructure management, AWS Glue with PySpark offers a highly efficient solution.
- The managed nature of Glue helps reduce engineering overhead while retaining the flexibility and power of Spark for data transformation.
- While less suitable for streaming or ultra-low latency use cases, it excels in scheduled data prep tasks where schema flexibility and operational agility are essential.

This experiment validates serverless ETL as a practical and scalable option for modern data engineering pipelines, particularly when time-to-delivery, maintenance effort, and cost control are key considerations.

6. Results and Discussion

The evaluation of AWS Glue with PySpark for serverless ETL highlighted several strengths and a few limitations when applied to large-scale batch processing in a real-world data engineering context. Lee, D. (2020)

6.1. Key Observations

6.1.1. Execution Time

- Daily batch jobs (processing ~400M records) consistently completed in 9–12 minutes.
- Performance scaled linearly with data volume when DPU allocation was increased appropriately.

6.1.2. Scalability and Flexibility

- Jobs scaled seamlessly by adjusting DPUs or switching worker types—no infrastructure changes were needed.
- Schema changes in upstream data were handled gracefully using dynamic frames, reducing breakage risk.

6.1.3. Cost Efficiency

- Pay-per-second pricing led to lower costs for infrequent or bursty workloads compared to EMR.
- No charges incurred when jobs were idle, unlike provisioned clusters.

6.1.4. Operational Overhead

- No cluster setup, patching, or autoscaling configuration required.
- Integration with CloudWatch provided sufficient monitoring for most use cases.

6.1.5. Limitations

- Cold start latency (~30–60 seconds) introduced slight delays but was acceptable for batch workloads.
- Debugging was less interactive than traditional Spark environments—limited visibility into execution plans or memory usage.

6.2. Overall Assessment

AWS Glue turned out to be a solid, low-maintenance choice for running batch ETL jobs—especially in cases where fast development, cost efficiency, and tight integration with the AWS ecosystem are bigger priorities than fine-tuning Spark internals.

For teams that want the power of Spark without the hassle of managing infrastructure, Glue strikes a good balance.

7. Challenges and Lessons Learned

While AWS Glue simplifies many aspects of ETL development, working with it in a production-like environment surfaced several important challenges and practical lessons for data engineering teams considering serverless ETL.

7.1. Challenges Encountered

7.1.1. Cold Start Latency

- Initial job startup time ranged from 30 to 60 seconds, adding overhead to short-duration tasks.
- Not ideal for low-latency or on-demand, user-triggered pipelines.

7.1.2. Limited Debugging Visibility

- CloudWatch logs were useful for basic error tracking but lacked the granularity of Spark UI or cluster-level metrics.
- Diagnosing memory pressure, skewed joins, or shuffles required additional instrumentation in code.

7.1.3. Schema Evolution Management

- Dynamic frames are flexible but can silently accommodate upstream changes, making it easy to miss critical schema shifts.
- Requires explicit validation logic or version control on schema definitions.

7.1.4. Job Configuration Nuances

- Optimal DPU settings, worker types, and partitioning strategies were not always intuitive.
- Performance tuning often requires trial and error due to limited documentation on edge cases.

7.1.5. Resource Quotas and Limits

- Default Glue limits (e.g., max DPUs per account, concurrent job runs) require adjustment via AWS support for scaling use cases.

7.2. Lessons Learned

- Build transformation logic to be idempotent and partition-aware from day one to avoid reprocessing issues.
- Use Data Frames over DynamicFrames for better performance and more predictable schema handling where possible.
- Custom logging and metric tracking should be added to supplement CloudWatch, especially for critical production jobs.
- Maintain a versioned schema registry or enforce data contracts to catch breaking changes early.
- Treat Glue as part of a larger architecture—use tools like Airflow (MWAA), Step Functions, or EventBridge for orchestration and alerting.

Overall, while Glue reduces operational overhead significantly, teams still need to apply solid engineering practices around monitoring, schema governance, and performance tuning to ensure reliable and maintainable pipelines.

8. Future Trend

While AWS Glue with PySpark has proven effective for batch ETL at scale, several areas remain open for further exploration to enhance its utility and address current limitations.

Areas for Expansion:

- **Real-time and Streaming ETL**
 - Investigate the use of AWS Glue Streaming for near real-time ingestion of event-based data (e.g., clickstreams, logs).
 - Evaluate how it compares with alternatives like Kinesis Data Analytics or Apache Flink in terms of latency and cost.
- **Advanced Orchestration**
 - Integrate Glue jobs more deeply with Amazon MWAA (Airflow) or Step Functions for complex dependency management and conditional execution logic.
 - Enable dynamic job chaining, retries, and notification hooks using external orchestration layers.
- **Observability Enhancements**
 - Implement custom metrics and structured logging for deeper visibility into job internals (e.g., row counts, partition stats, memory usage).
 - Explore the feasibility of proxying Spark UI data to improve debugging and performance insights.
- **Data Quality and Validation Frameworks**
 - Embed validation layers using tools like Deequ or custom PySpark checks to enforce schema conformity, null checks, and distribution rules.
 - Automate alerts for anomalies in data volume or processing time.
- **Integration with ML Pipelines**
 - Extend Glue's output to serve as a clean data layer for SageMaker or other ML workflows, enabling end-to-end feature pipelines.

- Explore automated feature extraction and model scoring using the same serverless stack.

- **Cross-cloud and Multi-region Strategies**

- Evaluate Glue's role in multi-region data replication and hybrid-cloud ETL, especially for organizations operating in regulated environments.

As data architectures become more modular and event-driven, the role of serverless ETL will likely expand beyond batch processing. Future exploration will focus on how Glue can evolve into a central data transformation backbone across streaming, ML, and multi-cloud ecosystems.

9. Conclusion

Using AWS Glue with PySpark shows that serverless ETL is not just hype but a practical and reliable approach that works in real-world data engineering workflows. Glue takes care of the heavy lifting of infrastructure while still giving you the power and flexibility of Spark.

This makes it a great fit for teams looking to speed up pipeline development while minimizing the burden of managing infrastructure.

Key takeaways

- AWS Glue handles large-scale batch ETL workflows well, offering solid performance and scalability across various data volumes.
- Serverless design enables rapid iteration, lower cost of ownership, and minimal platform maintenance—benefits especially valuable to lean data teams.
- While some limitations exist (e.g., cold starts and limited debugging depth), these are manageable within most batch-oriented contexts.
- Best results are achieved when Glue is integrated into a broader architecture that includes robust orchestration, monitoring, and schema validation.

In short, AWS Glue strikes a good balance between ease of use and robust capabilities. It is a solid choice for teams that value agility, want to keep costs in check, and are already working within the AWS ecosystem.

As serverless technology advances, Glue is well-positioned to play an even bigger role in building scalable, event-driven, machine learning-ready data platforms.

References

- [1] Plale, B., & Kouper, I. (2017). The centrality of data: data lifecycle and data pipelines. In *Data analytics for intelligent transportation systems*. Elsevier, 91-111. [[Google Scholar](#)] [[Publisher Link](#)]
- [2] Lee, D. (2020). Data transformation: a focus on the interpretation. *Korean journal*, 503-508. [[Google Scholar](#)] [[Publisher Link](#)]

- [3] Kriushanth, M., Arockiam, L., & Mirobi, G. (2013). Auto scaling in Cloud Computing: an overview. IJARCC, 2278-1021. [[Google Scholar](#)] [[Publisher Link](#)]
- [4] Pogatzi, A., & Samakovitis, G. (2020). An event-driven serverless ETL pipeline on AWS. Applied Sciences, 191. [[Google Scholar](#)] [[Publisher Link](#)]
- [5] Sudhakar, K. (2018). Amazon web services (aws) Glue. International Journal of Management, IT and Engineering, 108-122 [[Google Scholar](#)] [[Publisher Link](#)]
- [6] Singh, P. (2021). Manage data with PySpark. In Machine Learning with PySpark. 15-37. [[Google Scholar](#)] [[Publisher Link](#)]
- [7] Batmaci, G. (2022). Etl Data Pipelines Configurations in Spark. [[Google Scholar](#)] [[Publisher Link](#)]
- [8] Mehmood, E., & Anees, T. (2022). Distributed real-time ETL architecture for unstructured big data, 3419-3445. [[Google Scholar](#)] [[Publisher Link](#)]
- [9] Warneke, D., & Kao, O. (2009). Efficient parallel data processing in the cloud. 1-10. [[Google Scholar](#)] [[Publisher Link](#)]
- [10] Bussa, S., & Hegde, E. (2024). Evolution of Data Engineering in Modern Software Development. Journal of Sustainable Solutions, 116-130. [[Google Scholar](#)] [[Publisher Link](#)]