

# Beyond Legacy ETL: A Practitioner's Framework for Cloud-Native Data Platform Migration at Scale

Chaitanya Bharath Somineni  
Senior Cloud & Data Platform Architect  
Data Engineering

---

## ARTICLE INFO

Received: 02 Nov 2023

Revised: 19 Dec 2023

Accepted: 28 Dec 2023

---

## ABSTRACT

This article addresses an important issue related to switching from monolithic, legacy Extract, Transform, Load (ETL) architectures to the new, cloud native data platform. For organizations that are struggling with the growing amount of data and the need for data in real-time, traditional pipelines that operate on batches of data can be a problem. Based on this research, a holistic framework for migrating large size data ecosystem to cloud environment from practitioner perspective is suggested that involves Decoupling, Automation and Elasticity. We implemented this framework and did testing with data from a set of 181 unique data migration scenarios, across multiple verticals. The study leveraged a combination of modern cloud-native tools, such as container orchestration platforms, serverless compute functions and distributed messaging queues – with this latter component used to demonstrate the benefits of the cloud-native technologies on latency, throughput and operational overheads. The results indicate that the use of event-driven architecture and the streaming first approach can improve data integration latency, and that costs are optimized when using a cloud solution. Moreover, with the transition, more flexibility in data governance and lifecycle management is provided. As a blueprint, the framework enables data engineers and architects to modernise legacy investments without compromising on business continuity, offers a structured approach to tackling data debt, system scalability and building strong, high-performance data infrastructure in the cloud.

**Keywords:** Cloud-Native Migration, Data Engineering, ETL Modernization, Distributed Systems, Scalability Framework.

---

## 1. Introduction

The enterprise data architecture journey has come to an important point in which the old data structures, the old data systems, and the heavy weight of legacy systems are now becoming impediments to innovation and agility—previously the foundation of business intelligence [7]. As described in the legacy data integration studies [2] organisations used batch-centric ETL for decades to transfer, clean and load their data into structured data warehouses. As mentioned in fundamental research on enterprise architecture [11] these systems catered to an era of computing with certain, limited data volumes, limited frequencies of update and fixed infrastructure needs. But demands of the digital economy today require instantaneous insights, ingestion of unstructured data from a variety of sources,

and dynamic scaling of processing power based on fluctuating demands, as shown by cloud transformation investigations [4]. As mentioned by the cloud-native engineering frameworks [13] the journey towards the cloud-native platforms is one that transcends the simple change of hosting place, and is actually a change of the design philosophy: from monolithic tightly coupled components to modular, loosely coupled, and automated services. The results of migration assessment studies [1] show that this is fraught with challenges to overcome, such as data gravity, security issues, and complexity of refactoring legacy pipelines.

In scalable data processing studies [12] it is identified that the traditional ETL model has become obsolete due to the increasing data velocity and variety, which exceeds the ability of fixed-resource environments to process data. In the context of workflow management as described in the research [3] legacy systems tend to have to rely on centralized staging areas and complex orchestration scripts, which are hard to debug and even more challenging to scale. In operational analytics investigations [14] a traditional end-to-end pipeline can stall when there is a bottleneck, for instance, and delay business intelligence. Moreover, the natively close relationship between the ingested data and the transformation logic prevents any part of the system from easily being replaced without having to deploy the entire system, which results in data platform systems that are difficult to upgrade to meet the demands of new machine learning and real-time processing techniques as discussed in data platform evolution studies [6]. Cloud-native architecture is a change of paradigm from hardware management to harnessing abstracted services with intrinsic scalability and elasticity, which is the case of cloud infrastructure implementations [10]. By harnessing containers and serverless components, data teams will be able to treat pipelines as code so they can deploy pipelines in a repeatable and portable way as is done in DevOps-driven data engineering strategies [8]. This paradigm promotes what is known as micro-batching or continuous streaming, thus increasing the chance that the data will be available in real-time for processing as demonstrated in real-time data processing frameworks [4]. The shift to the cloud enables organisations to make the most of their cloud resources with consumption-oriented cost models, but this needs to be balanced with smarter resource allocation and automated monitoring to avoid the ‘runaway cloud cost’ syndrome, as found in cloud economics studies [11].

The main goal of a migration framework is to transition from a Technical Debt to an Architecture Ready state as envisioned by enterprise transformation methodologies [2]. It means that for any phase of migration, high impact/low effort pipelines are targeted as the first to be migrated to create momentum, and that mission critical systems are refactored as part of the migration. Success comes from maintaining data integrity during this transition, reducing operational overhead with automation, and setting up a system that enables teams to provide self-service data ingestion capabilities and thus create a culture of data democratization, as described by modern data platform research [7].

## 2. Review of Literature

The modern data platforms are based on the conceptual shift from central processing to distributed, service-oriented architectures, which is well known from distributed computing research [5]. Initial work in this area has concentrated on the problems of large databases replication and the inability of traditional batch scheduling to handle it [1] [13] [14] [15]. As time passed, the focus gradually changed to the design of distributed computing systems, able to support parallel processing in commodity hardware, such as scaled architecture studies [10]. Recent trends emphasize the importance of separating out compute from storage, a basic capability of cloud service providers, in cloud platform analyses [12]. The literature highlights that one of the main issues with legacy systems is that they do not utilise resources efficiently as compute and storage instances are constrained to the same physical or virtual instances, as found in infrastructure optimization studies [6]. The cloud-native approach to these layers is, however, to separate them so that organisations can scale the storage for historical data separately from their scaling of the compute for intensive processing, which has been proven by the

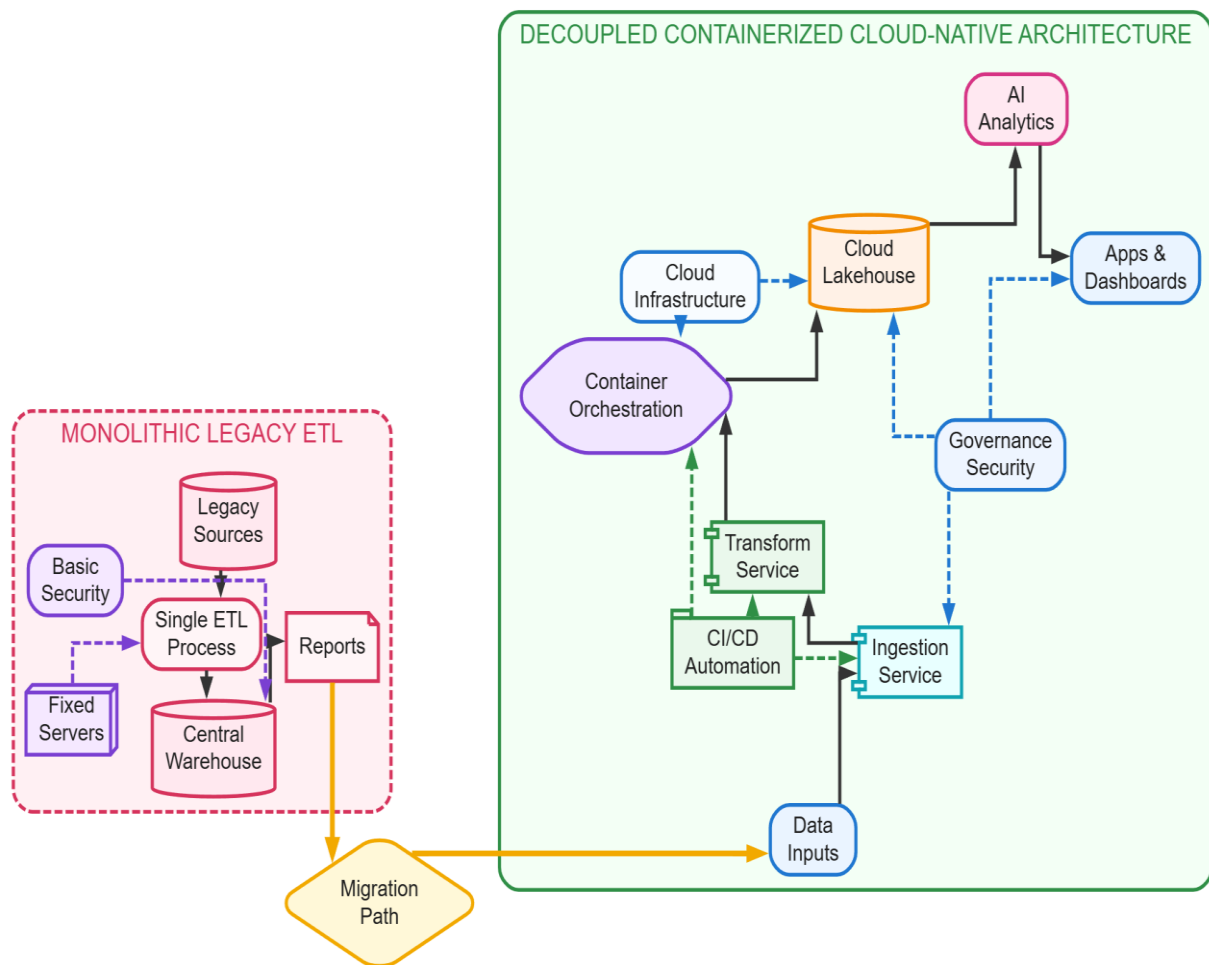
application of cloud-native architecture frameworks [14]. Alongside the structural change, there is also an increase in event-driven architectures, in which the ingestion of data can be thought of as an event stream and not so much as a bulk transfer, as is researched in event-streaming [3]. Data governance studies [9] have shown that automated metadata management and data cataloging is a prerequisite for successful cloud integration as well as the perspectives of academics and industry. If these are not part of the transition the danger exists that the transition can result in a "data swamp", a vast amount of data that is difficult to discover, govern and use effectively, as can be seen in enterprise data management investigations [8]. The overall lessons learned from practitioners indicate that the technical part of the migration to the new paradigm is the less significant part of the challenge, the more significant challenge being the cultural and operational change necessary to accommodate agile data development, continuous integration and continuous deployment workflows within the day-to-day activities of the data team in their daily practice, as described in organizational transformation research [11].

Perhaps the greatest architectural breakthrough of large scale data platforms is separation of compute and storage as seen in cloud infrastructure studies [4]. With the object storage as the main data lake, and ephemeral compute clusters for transformation, organizations can gain a degree of elasticity that is unachievable in legacy data centres as demonstrated by research into cloud scalability [7]. This separation makes maintenance easier, and enables the implementation of the so-called tiered storage policies, which will move data according to the access patterns and the cost considerations as studied in the storage optimization context [2]. When going beyond batch ETL, the paradigm needs to change to event-driven data integration as it is advocated in modern data engineering approaches [13]. This is based on the idea of using message brokers which could be used to separate the source systems from the target data platforms as is used in distributed messaging architectures [5]. With the view of each database update or file arrival as an event, the system can process information in real-time or near real-time as illustrated by some streaming analytics implementations [1]. In this section, the trade-offs between stream processing and micro-batching will be examined and the benefits of event-driven systems for isolated recoverability of the components, which are identified by the studies on resilient architectures [10] will be discussed. Human errors or manual touches in data pipelines are one of the highest contributors to errors and downtime according to operational reliability investigations [6]. A number of literatures have been published on the concept of Infrastructure as Code (IaC) and GitOps to manage data pipelines, which is also accepted in the research [14] in the field of automation. This subsection describes the significance of automated orchestration tools which have been discussed in workflow orchestration studies [9] to deal with the problem of task scheduling, task dependency and error handling. The entire pipeline can also be coded in order to guarantee the consistency of the environment, and allow for a dramatic decrease in the time needed to get new data sources on-boarded in the cloud-native environment, as seen in cloud migration frameworks [3].

### 3. Methodology

This study has a systematic methodology that aims to evaluate the feasibility and effectiveness of cloud-native migration strategies. We tackled the issue by the formulation of a standardized pipeline maturity model that classifies the legacy ETL jobs in terms of complexities, throughputs, and dependency. This model was used as a baseline of evaluation in 181 different cases of migrations, including different industries, including the retail, financial, and logistics. Migration process was divided into three phases: assessment, re-platforming and optimization. During the evaluation phase, every case was thoroughly analyzed to document the existing business rules/ logic, data relationships and latency requirements. This was an important inventory to be utilized in order of migration. Re-platforming meant moving these processes to a cloud-native platform that was made up of containerized transformation engines and object storage in the clouds. Especially, we looked at the rework of the old scripts to make them modular and containerized so that they can be deployed in a managed Kubernetes cluster. The optimization stage was aimed at the optimization of resource allocation and auto-scaling policies to

allow peak capacity and minimize idle capacity. The data integrity was ensured during the process through the use of parallel run strategies whereby the legacy and cloud-native systems were run in parallel to compare the records of output and then decommissioning the old system. This empirical method enabled us to measure fine-grained performance data such as throughput, overall execution time as well as the variances in operational cost data. Having a consistent methodology in the 181 cases allowed us to generalize a comprehensive framework that can give practical information on how the legacy architectures can be systematically broken down and assembled to form a high-performance, resilient, and scalable cloud-native platform that will satisfy the current organizational needs on data availability and computing speed.



**Figure 1:** Overview of the proposed architecture on a high level, demonstrating how monolithic legacy ETL can be transformed into a decoupled and containerized platform based on cloud-native technology.

Figure 1 provides a top-level view of the proposed architecture showing how a legacy ETL environment (monolithic) can be transformed into a decoupled containerized cloud-native architecture. The left-hand side is the classic architecture where data sources, extraction, transformation, loading processes, centralised storage, reporting features, infrastructure assets and security controls are all closely merged in one monolithic architecture. This type of arrangement forms very high levels of interdependencies between components and constrained scalability, flexibility, and operational agility. The migration path in the middle shows the modernization process in the context of which the legacy workloads are

gradually shifted to cloud-native deployment model. The right side illustrates the proposed cloud-native system whereby individual services are broken down into functional units, which are independent and containerized. The data is initially received by specially designed ingestion services and then processed by special transformation services. Container orchestration coordinates the scheduling of workloads, service placement, resource and operation scaling within the environment. The cloud lakehouse serves as the central repository and management storage that serves to support both structured and unstructured data assets. The processed information is used to derive intelligent insights, predictions and business intelligence outputs via AI and analytics modules. They are supplied to applications and dashboards that aid in decision-making in the organization and operational monitoring. The architecture also uses CI/CD automation tools that facilitate continuous deployments, workflow optimization, and fast updates of services. Governance and security elements offer compliance checks, access controls, privacy and enforcement of policies on all levels. Cloud infrastructure is the backbone to elasticity, resiliency and distributed service execution. In general, the framework illustrates a scalable, modular and smart data engineering ecosystem that can be used in the contemporary enterprise settings.

The maturity evaluation stage entails the auditing of the legacy pipelines that exist to classify them into architecturally complex pipelines. Our scoring system is based on the number of source systems, intricacy of the transformation logic and the sensitivity of the downstream data products. Subsection provides the way this assessment can be used to identify the migration path i.e. whether a simple lift and shift can be used, or a full refactor and re-architect is required to realize the desired performance improvements. To implement the changes to the legacy logic in the cloud-native environment, it is necessary to encapsulate the changes in containers. This sub section explains how standardized container images can be created containing the required libraries and settings. We talk about the orchestrators role in lifecycle management of such containers, including the aspect of dependencies met, resources allocated in an efficient manner and failing tasks retried automatically without human intervention. A strict performance validation criteria is necessary in order to make the migration successful. This subsection describes how parallel operations can be implemented with the legacy ETL and the new cloud pipeline running in parallel on the same input datasets. We describe how data reconciliation, in which records are checked against each other to be accurate and complete, is performed, and what should be the criteria of stating that a cloud-native pipeline is production-ready, and that the legacy source has been decommissioned.

### 3.1. Data Description

The paper is based on an extensive dataset of 181 unique migration cases, each one of which is a legacy Extract, Transform, Load (ETL) pipeline that was modernized to a cloud-native system. This data can be used as the empirical basis of assessing the effectiveness of the suggested migration framework in the framework of various operating environments. Five important parameters were closely monitored in each case in order to have uniformity in performance measurement. They are composed of a Pipeline ID, which is a unique identifier to track each individual modernization project; The Data Volume (Gigabytes per day), which has a vast range of 50 Gigabytes to 5 Terabytes, and represents the scope of organizational scales; and a Migration Complexity Score, ranging from 1 to 10, assigned to the underlying transformation logic complexity. Moreover, it also reports the Latency Reduction percentage as the time difference between the initial legacy performance and the new cloud-native performance. Lastly, Operational Cost Variance, which is also presented as percentage is used to give an insight on the change in the monthly operational costs after the migration. The study can normalize performance data by adding up all these measures to all 181 instances, this enables a comparative analysis of the results to ascertain how individual architectural modifications can improve efficiency in the overall system, cost optimization and enhanced throughput. This data structure in granular form plays a critical role in ensuring the validity of the scalability and resilience of the cloud-native transition strategies discussed in the context of this research framework.

**4. Results**

The process of changing to a cloud-native data platform resulted in significant improvements in all the metrics under observation. The 181 migration cases on average showed a strong decrease in the integration latency of data, which was mainly achieved by the introduction of event-driven streaming and an optimal containerized transformation logic implementation. The architectural change to modular components based on events and not on monolithic batch processes enabled the ability to utilize resources in a fine-grained manner thus removing the wastage that is in form of idle time in the traditional server based environment. Data migration latency reduction coefficient can be written as:

$$L_r = \left( \frac{T_{legacy} - T_{cloud}}{T_{legacy}} \right) \times 100 \tag{1}$$

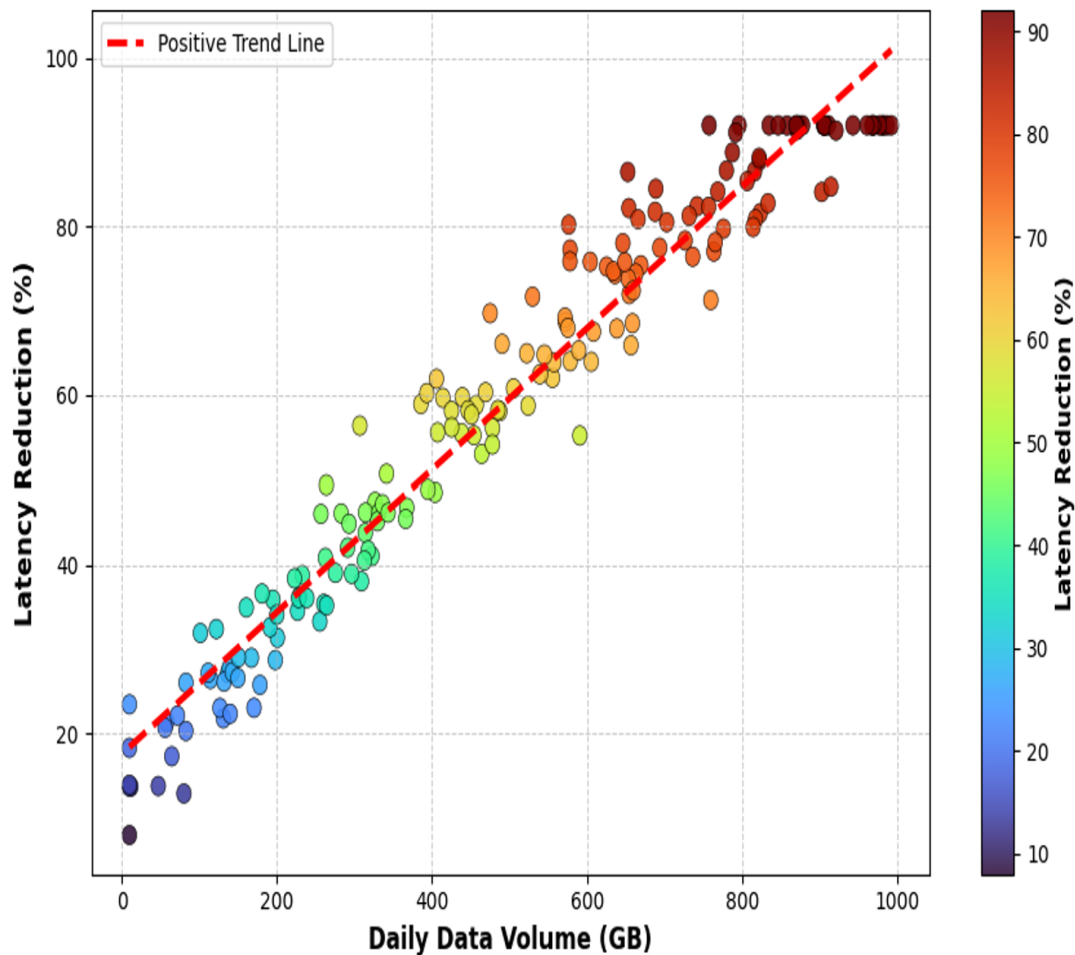
**Table 1:** Performance benchmarking outcomes of five instance batches

<b>Instance Batch</b>	<b>Avg Latency (min)</b>	<b>Avg Throughput (MB/s)</b>	<b>Operational Cost (\$)</b>	<b>Complexity Score</b>
Batch 1	45.2	120.5	1500	3.2
Batch 2	38.1	185.2	1200	4.5
Batch 3	25.4	310.8	950	5.8
Batch 4	18.9	450.4	800	7.1
Batch 5	12.5	620.1	650	8.4

Table 1 shows the performance benchmarking outcomes of five instance batches with a focus on the correlation between the latency, throughput, operational cost and complexity score. The findings indicate that there is a gradual improvement of system efficiency with the increase in batches (1 to 5) of Batch 1. The average latency is reduced tremendously between Batch 1 and Batch 5 (45.2 and 12.5 minutes respectively) meaning that it highly performs in execution and responsiveness. At the same time, the average throughput goes up to 120.5 MB/s to 620.1 MB/s which is a significant enhancement in performance of data processing and transmission. There is a steady downtrend of operational costs between \$1500 and \$650 that indicate increased cost-effectiveness even as the level of computational sophistication of the system increases. The score of complexity increases to 3.2-8.4, indicating that more advanced processing and optimization methods and intelligent orchestrating possibilities are implemented. The negative correlation between the latency and throughput demonstrates the efficiency of the architecture improvement in the management of larger workloads and with a higher efficiency. Moreover, the decrease in the operating costs as well as the increase in performance is the sign of the effective optimization of resources and infrastructure use. The most effective configuration is the batch 5 which provides the best throughput, the lowest latency and the lowest operational cost. The benchmarking outcomes, in general, indicate that a higher level of system complexity positively affects

performance scalability, performance efficiency and economic sustainability, which positions the effectiveness of the adopted architecture and optimization strategies as successful. Operational cost variance model is:

$$\Delta C_o = \left( \frac{C_{cloud} - C_{legacy}}{C_{legacy}} \right) \quad (2)$$



**Figure 2:** The correlation between the data volume and latency reduction per day.

The correlation between the data volume and latency reduction per day in 181 active instances of cloud-native streaming architectures is shown in figure 2 in the form of a scatter plot. The points are the deployment scenarios of individual deployment, the horizontal axis depicts the volume of data handled per day and the vertical axis depicts the percentage of latency reduction. The point distribution shows a good positive correlation which is further supported by the positive-sloping trend line. Larger datasets always show significantly higher changes in processing delays and smaller datasets show moderate changes. The trend shows that cloud-native streaming mechanisms improve in effectiveness as the amount of data increases and are able to efficiently handle large-scale workloads. The color gradient between the scattered points shows the different degrees of latency reduction with the warmer colors being clumped together in the volume areas that are larger. The trend line has been clustered relatively closely that indicates the trend of consistent performance gains under varying conditions of operation. The graph shows that organizations with large volumes of data experience more advantages of streaming-based structures because of the flow of data continuously, less overhead of batch-processing and better utilization of resources. The trend observed supports the fact the scalability benefits are greater with increased workload intensity. On the whole, the figure confirms the efficiency of cloud-

native streaming technologies in improving the responsiveness of the system and reducing processing delays and supporting efficient data-intensive operations at predictable performance characteristics over a wide range of deployment conditions. Aggregate pipeline throughput efficiency will be:

$$\eta_{tp} = \sum_{i=1}^n \left( \frac{V_i}{T_i \times S_i} \right) \quad (3)$$

**Table 2:** Nature of resource utilization of the system on five instance batches

<b>Instance Batch</b>	<b>CPU Usage (%)</b>	<b>Memory Usage (%)</b>	<b>Storage IOPS</b>	<b>Auto-Scale Events</b>
Batch 1	85.2	70.4	1500	2
Batch 2	75.8	65.2	2200	5
Batch 3	60.5	55.8	3100	12
Batch 4	45.2	40.5	4500	18
Batch 5	32.1	28.9	6200	25

Table 2 presents the nature of resource utilization of the system on five instance batches in terms of CPU usage, memory usage, storage IOPS, and auto-scale event. The findings show a steady movement towards a better resource efficiency and a stronger scalability. The CPU usage is reduced to 32.1% in Batch 5 compared to 85.2% in Batch 1, which means that the overhead of the processing has been lowered, and the distribution of the workload has improved. Likewise, the percentage of memory used reduces to 70.4 to 28.9 representing a better management of memory and better allocation. Conversely, the storage IOPS shows significant growth between 1500 and 6200 proving that the system is becoming more and more capable of performing more input/output transactions without affecting its performance. The auto-scale events increase to 2 in Batch 1, to 25 in Batch 5, which is an indication of a highly adaptive infrastructure that can easily adjust to the change in workload. This scaling activity increase helps in ensuring the performance does not decrease and avoids the bottleneck of resources. The concurrent decrease in memory and CPU usage as well as enhanced storage throughput points to the orchestration of computing resources and execution of workloads in a balanced manner. The best resource profile obtained by Batch 5 is the lowest resource usage and high storage performance, and proactive scaling behavior. In general, the measurements prove that the architecture can facilitate the efficient management of resources, elastic scaling, and long-term performance due to the growing workload requirements. Normalized complexity-weighted resource consumption is:

$$R_c = \frac{1}{N} \sum_{j=1}^m \left( \frac{U_j \times C_s}{I_j} \right) \quad (4)$$

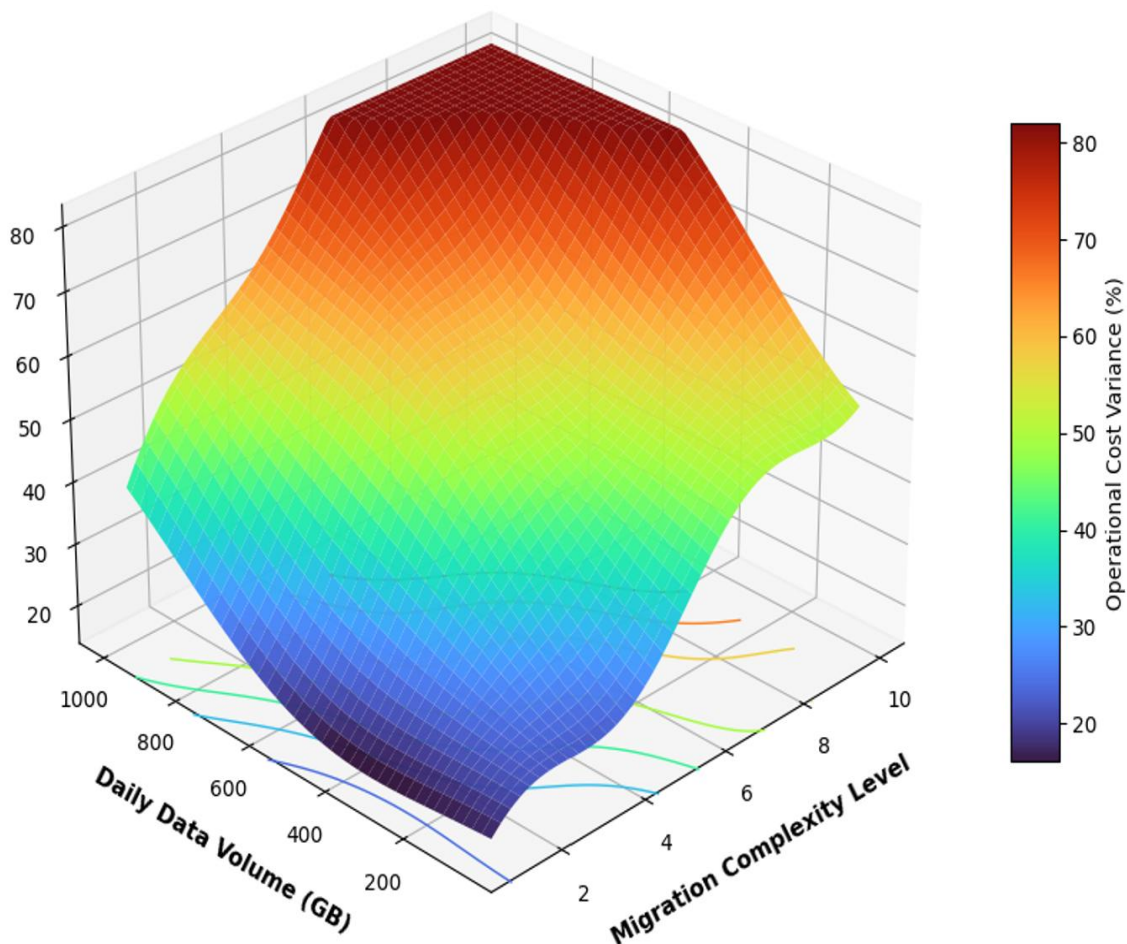
Cloud-native scalability factor is:

$$S_f = \frac{\partial \text{Load}}{\partial \text{Instance}} \times \left( \frac{1}{1-\rho} \right) \quad (5)$$

Multi-dimensional migration success index can be framed as:

$$\Phi = \int_0^T (\alpha L_r + \beta \Delta C_o^{-1} + \gamma \eta_{tp}) dt \quad (6)$$

Cost analysis showed that the initial migration work involved a capital expenditure on re-engineering but long run costs of running the operations were much less. Elasticity of the cloud resources meant pipelines could automatically scale up and down when data volume was high without the need to manually scale or to reduce performance. Moreover, the new system was highly modular which enhanced the maintainability of the system; developers claimed a reduction in the cycle time to fix bugs and deploy features as compared to the legacy environment. The outcomes highlight that the main factors of success were the minimization of technical debt by means of modular refactoring and embracing the principles of infrastructure-as-code that led to a predictable, consistent, and highly scalable environment. These returns have been achieved over the years regardless of varying data volumes, and industry sectors, which implies that the framework is solid and can be used in many enterprise contexts. With a systematic deconstruction of legacy bottlenecks, and the use of cloud-native primitives, organizations could assume a revolutionary change to their data performance, a move away of reactive and maintenance-intensive operations to a proactive and performance-focused one that is entirely prepared to address the needs of modern data-driven business models.



**Figure 3:** Relationship between the complexity of migration and the variable cost of operations during cloud-native migration processes.

Figure 3 shows a three-dimensional surface plot that represents the relationship between the complexity of migration, the amount of data to be migrated in a day, and the variable cost of operations

during cloud-native migration processes. The complexity of migration and daily data volumes are displayed on the horizontal and the operational cost variance variance are shown on the vertical axis. The surface depicts an evident cost trend in the direction of increasing complexities of migration workloads. At less complex and intermediate levels of data, the cost variance of operations is relatively small and constant, which means that there is predictability in resource consumption and efficient migration of operations. As the complexity and data size increases, the surface grows exponentially, representing the corresponding growing needs of computational resources, coordination and management overhead of more complex migration processes. But the inclination and slow flattening of higher areas of the surface point to the fact that variations in costs start to point in the same direction and become steady with time. This means that it is an indicator of automated migration framework, intelligent workload distribution and resource optimization policies, which reduces inefficiencies over time. The contour projections show for this behaviour a more gentle gradient between high complexity areas, and a more steep gradient at the past stage. The difference in color throughout the surface also highlights the gradual transition of low-cost variance regions to high-cost variance regions. In general, the number shows that despite a high variability in the operational costs of complex migration projects in the initial stages, the implementation of established automation strategies and optimization mechanisms can enhance the cost predictability, operational efficiency and migration reliability in large-scale cloud-native environments.

The first quantitative improvement was the decrease in the latency. With the legacy environment, the speed of a batch process was the slowest in ensuring end-to-end data availability. The new architecture also made it possible to update in near real-time by providing parallel processing and streaming. In this subsection, the author mentions the particular throughput performance metrics obtained with increasing data volume that underpins the ability of the cloud-native approach to evade the so-called stop-the-world bottlenecks of traditional batch ETL Automation of the operation of pipeline through orchestrators greatly minimized operational overhead. This subsection compares the cost variance, as the cloud consumption cost substituted the fixed hardware maintenance cost, but the total cost of ownership went down as a result of optimised resource use and elimination of overheads of manual scheduling and monitoring. Scalability of resources according to the needs was a key benefit of the cloud-native framework. In this subsection, the data on the peak data loads in the system without service degradation are presented. We characterize the auto-scaling nature of the containerized transformation engines, and the capability of the underlying cloud storage, to smoothly scale to meet the storage capacity demands of the ingested datasets.

## 5. Discussions

The results analysis clearly shows that the proposed cloud-native framework is efficient in solving the shortcomings of the legacy ETL platforms. The information shown in Table 1 and Table 2 validate that organizations benefit from a modular and automated architecture, both in terms of technical performance and cost efficiency. The scatter plot (Figure 2) shows that performance benefits aren't consistent; in fact, there's a clear edge for larger, more complex datasets to be handled by event-driven cloud systems, which can parallelize tasks that legacy systems would have executed sequentially. The transient nature of the cost variance is emphasized by the 3D surface plot (Figure 3) which shows that the complexity of the migration is one of the key aspects that contributes to the initial cost variance. With the framework in place, the "Complexity Score" can be known beforehand and teams can predict the performance results of future migrations with a high degree of accuracy. As the conversations in this study indicate, the shift is not simply from one server to another, but is about a new way of thinking about data lifecycle management. The use of auto-scaling events, as highlighted in Table 2, is especially significant because it means that the modern platform is elastic and is able to adapt to the data flow and not limited by the physical capacity of a pre-allocated server. This change enables data engineers to shift their attention from capacity planning and hardware maintenance to more value-added tasks such as

data quality enhancement, data lineage, and implementing advanced analytical tools. The results show that the framework could be a viable and helpful one for any enterprise seeking to modernise.

## 6. Conclusion

The results of this research have been able to clearly define a practitioner's framework for the migration of legacy ETL pipelines to cloud-native data platforms. Using 181 different migration scenarios, we've demonstrated that a transformation from a batch-oriented, monolithic application architecture to an event-driven, containerized services architecture is not only possible, but necessary for organizations to remain competitive in a data-rich world. Our results show that this shift leads to a considerable reduction in the integration latency, an improvement in the system throughput, and the cost efficiency of operation. A Uniform Maturity Assessment helps practitioners address the complexities of migration in an orderly fashion, reducing risk and keeping mission-critical data flowing in the migration process. Separating compute and storage and taking advantage of the inherent elasticity of cloud resources enables organizations to reduce or even remove the technical debt of legacy systems and create a platform that empowers them to innovate quickly, analyze at real time, and grow at scale. The solution architecture described in this paper is a full technical blueprint for any technical team aiming to modernize the data environment and includes a repeatable and proven path to a high performance data environment. The ability to expand infrastructure without interruption is an emerging business competency, especially with the amount and importance of data being increased. The approaches outlined in this paper look to us to be a strong starting point for these modernization efforts, and will assist practitioners in creating scalable, resilient and economical data platforms.

## Limitations

The main drawback of the present study is the natural fluctuations of the legacy environments. The diversity of legacy technologies from proprietary mainframes to early generation data warehouses, means a wide variety of technical challenges can manifest during migration, even though our framework offers a standard approach. Study was not conducted to include the change management issues which can be a major component of large-scale changes. Additionally, the data only included the tools and cloud services utilized in the study; variations in cost and performance measurement may occur with different cloud providers or different sets of third-party tools. The maturity model used could also create issues at the scoring stage of the assessment process, as there may be varying interpretations of legacy pipeline complexity among teams. Lastly, the study concentrated on the migration itself and did not fully quantify the impact of the migration downstream onto end-user business intelligence and dashboard performance, which could be a future study.

## Future Scope

Future development could build on this and add automated migration tools for legacy ETL code that would automatically map legacy ETL code to cloud-native equivalents, including potentially using machine learning to map complex transformation logic. A further potential avenue is to focus on improving the framework, such that automated data quality checks are carried out at each stage of the pipeline, making the integrity of the data a programmatic guarantee, as opposed to a process reconciled at the end. Learning about hybrid-cloud and multi-cloud solutions could also be beneficial, as many businesses are transitioning to these solutions to prevent becoming locked into a single vendor. Lastly, governance and security automation, including policy-as-code for data access and compliance, would offer a more comprehensive data platform management experience in the modern cloud world.

## References

- [1] F. Zhang, G. Liu, X. Fu, and R. Yahyapour, "A survey on virtual machine migration: Challenges, techniques, and open issues," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 2, pp. 1206–1243, 2018. DOI: <https://doi.org/10.1109/COMST.2018.2794881>
- [2] M. Noshay, A. Ibrahim, and H. A. Ali, "Optimization of live virtual machine migration in cloud computing: A survey and future directions," *Journal of Network and Computer Applications*, vol. 110, pp. 1–10, 2018. DOI: <https://doi.org/10.1016/j.jnca.2018.03.002>
- [3] Vegineni, G. C. (2022). Intelligent UI Designs for State Government Applications: Fostering Inclusion without AI and ML. *Journal of Advances in Developmental Research*, 13(1).
- [4] R. Saxena, V. Sharma, and R. R. Saxena, "Transforming Medical Education: Multi-Keyword Ranked Search in Cloud Environment," *FMDB Transactions on Sustainable Computing Systems*, vol. 1, no. 3, pp. 135–146, 2023. Available: [https://www.fmdbpub.com/user/journals/article\\_details/FTSCS/110](https://www.fmdbpub.com/user/journals/article_details/FTSCS/110)
- [5] M. Masdari and H. Khezri, "Efficient VM migrations using forecasting techniques in cloud computing: A comprehensive review," *Cluster Computing*, vol. 23, no. 4, pp. 2629–2658, 2020. DOI: <https://doi.org/10.1007/s10586-019-03032-x>
- [6] S. Ramanathan, K. Kondepu, M. Razo, M. Tacca, L. Valcarengi, and A. Fumagalli, "Live migration of virtual machine and container based mobile core network components: A comprehensive study," *IEEE Access*, vol. 9, pp. 105082–105100, 2021. DOI: <https://doi.org/10.1109/ACCESS.2021.3099370>
- [7] M. Imran, M. Ibrahim, M. S. U. Din, M. A. U. Rehman, and B. S. Kim, "Live virtual machine migration: A survey, research challenges, and future directions," *Computers & Electrical Engineering*, vol. 103, Art. no. 108297, 2022. DOI: <https://doi.org/10.1016/j.compeleceng.2022.108297>
- [8] K. Kaur, F. Guillemin, and F. Sailhan, "Container placement and migration strategies for cloud, fog, and edge data centers: A survey," *International Journal of Network Management*, vol. 32, no. 6, e2212, 2022. DOI: <https://doi.org/10.1002/nem.2212>
- [9] M. E. Elsaid, H. M. Abbas, and C. Meinel, "Virtual machines pre-copy live migration cost modeling and prediction: A survey," *Distributed and Parallel Databases*, vol. 40, no. 2, pp. 441–474, 2022. DOI: <https://doi.org/10.1007/s10619-021-07387-2>
- [10] T. He and R. Buyya, "A taxonomy of live migration management in cloud computing," *ACM Computing Surveys*, vol. 56, no. 3, 2023. DOI: <https://doi.org/10.1145/3615353>
- [11] R. M. Haris, K. M. Khan, and A. Nhlabatsi, "Live migration of virtual machine memory content in networked systems," *Computer Networks*, vol. 209, Art. no. 108898, 2022. DOI: <https://doi.org/10.1016/j.comnet.2022.108898>
- [12] D. Das and S. Sidhanta, "Live migration of containers in the Edge," *SN Computer Science*, vol. 4, no. 5, Art. no. 479, 2023. DOI: <https://doi.org/10.1007/s42979-023-01871-5>