



# COMPARATIVE STUDY OF MULTI-CLOUD DATA ENGINEERING STRATEGIES

*\* Chittaranjan Pradhan*

*Independent Researcher, East Brunswick, New Jersey, United States*

*cpradhan01@gmail.com*

Accepted: June 2024

Published: July 2024

**Abstract:** Best practices from data engineering processes have been significantly modified in disparate sectors due to multi-cloud solutions only, involving multiple cloud service providers (CSP). Multi-cloud data engineering enables businesses to distribute workloads across different cloud platforms, increasing performance, reducing the potential for vendor lock-in, and supporting scalability. Notably will be covered and assessed in this study that pertains to multi-cloud data engineering are: Data integration Security Cost optimisation and System interoperability. The report begins with a historical overview of cloud computing and multi-cloud approaches that highlight why organizations have adopted these setups. Part of the reason for this is the need to comply with data rules throughout many countries, along with requirements for more resilience and flexibility. The study covers fundamentals of multi-cloud architecture such as data replication, workload allocation, and resource management. It explores data engineering practices of various companies operating with focus scheduled optimization resources, merging multiple data sources, and defining metrics for cloud service selection. The research also compares managing data pipelines for ingestion to transformation, storage, and processing across clouds with cloud-native versus third-party tools. When handling information across several clouds, protecting sensitive material and ensuring compliance rules at every data engineering level is of the biggest significance. This paper is a twofold concern; one is looking on security, and another on compliance.

**TAG WORDS:** Data Engineering, Multi Cloud, Data Integration, Security, Cloud Service Providers (CSPs), Data Replication



## **Introduction**

Over the past few years, there has been an exponential growth in data leading to a much more challenging mechanism for processing, storing and managing data. Due to its ability to scale, versatility, and affordability, cloud computing has become the predominant solution to these issues. Multi-cloud data engineering is being validated as a panacea for performance, prevention of vendor lock-in, resilience and regulatory compliance for organisations. Organisations can deploy the right multi-cloud solutions to enable efficient and dynamic data operations. And those solutions use multiple CSPs but rather put on top of them features specific to each CSP. Some of these benefits could be reduced costs, easier access to superior analytics, or a broader geographic footprint. Businesses, however, are finding anew the challenges of managing data across several cloud environments -- from security management and cost containment to seamless data integration, systems compatibility and data orchestration. This study is primarily concerned with the data engineering practices that organisation run across multiple clouds. It accomplishes this by examining the costs and benefits of multi-cloud architectures as well as analyzing how they're being deployed across a range of industries to shine a spotlight on the struggles of companies looking to adopt and sustain multi-cloud architectures. The research further explores how organizations are managing data ingestion, transformation, storing and processing across different cloud providers at the same time. This comparison gives businesses significant insight into the range of options available as well as what multi-cloud engineering is best placed to provide. Besides these major aspects, the project will also learn about various other important aspects like data synchronisation and integration. Security concerns — from identity management to data encryption to access control — are already complex matters to manage; adding multiple CSPs into the mix can further complicating security management on top of that. This research will examine these matters in detail. There are a number of hurdles, but a multi-cloud setup allows organisations to minimise costs by spreading workloads across platforms based on how they perform as opposed to the underlying cost. Multi-Cloud systems also improve scalability, resilience, and disaster recovery. Multi-cloud architectures have redundancy built-in by design, helping to reduce the risk of locking into a single cloud provider. The research will examine data governance as well as adherence with rules such as GDPR and HIPAA, particularly in sectors that are extremely



regulated concerning data. It will also address the difficulties organisations have to contend with when trying to oversee spend and manage complex billing, with multiple cloud service providers. Interoperability among multi-cloud data analytics and the composite application of distributed cloud-native services across different CSPs become huge difficulties that need to be solved. Architecture designs, patterns for PaaS/IaaS (e.g: Hybrid cloud solutions like Spanner, Cosmos, etc. The research will be centered on such solutions), such as K8s Limited to other containerisation technologies and/or data integration platforms. The objective is to establish processes where data will be transferred seamless with interoperable systems. With multi-cloud initiatives becoming more common, organisations that are contemplating or already in adoption may find some insights from this comparative look across clouds for data engineering practice. During the creation and deployment of systems that operate across multiple clouds, we recommend double down on security, compliance and automation, the research says. (2) An open governance framework is also needed to support better data security and integration. The paper concludes with recommendations for ways in which companies can tackle those challenges and empower data engineering in the multi-cloud. Here's a summary of the study's findings: It ... starts by reviewing fundamental ideas and technology in the field. Next the multi-cloud methods and cloud computing are explored. The next section compares cloud-native and third-party data engineering solutions, with a look at how each sector approaches data storage, integration and security. This comparative study aims to inform the developing cloud data engineering space by providing an overview of the advantages and disadvantages of multi-cloud installations. Ultimately this report is an attempt to give us all some actionable knowledge that will help organisations make more informed decisions on implementing or upgrading their multi-cloud strategy. The study will also explore real-world examples and best practices to illustrate how multi-cloud architectures can improve data quality and data engineering talent. Furthermore, the paper will explore differences and similarities of security and compliance practices of different CSPs, regarding compliance with regulations like GDPR, HIPAA, and various industry standards. A major question in this study is how much money multi-cloud methods save. From analysing the pricing structures of numerous CSPs, I will help organisations optimise resource allocation and minimise cloud-related costs. Alternative cost management methodologies and approaches will also be highlighted.

## **Review of Literature**



Cloud computing has transformed data management and engineering, providing enterprise-grade, scalable, inexpensive, and flexible infrastructure. With rising use of the cloud, efficient data storage, processing, and analysis are increasingly important. Initially, single-cloud systems were prevalent, but due to performance optimization issues, cost control, and vendor lock-in, multi-cloud techniques are gaining traction. This evolution has empowered a business to manage data integration, security, and scalability from both a governed and great data governance perspective. Initially, the cloud computing literature extolled unified single-cloud architectures, emphasizing how convenient and controllable they were. But it became clear that expanding companies require more freedom, less risk and some way to spread out their suppliers. In a bid to minimize reliance on the infrastructure and services of a single vendor and achieve greater availability, redundancy, and disaster recovery, organisations moved to multi-cloud [1-2]. Multi-cloud data engineering consists of simply ensuring data is integrated seamlessly across different cloud platforms. Thus, it is important to maintain data synchronisation, make data available across platforms, ensure the accuracy and completeness of all available data. He states the platforms may not be compatible with each other, leading to different data formats and protocols that make it more difficult to bring data from many sources [3-4]. Middleware and cloud integration platforms emerge as one of the solutions. Some examples of tools used for efficient processing and synchronisation of data. A solution specifically designed to remove integration barriers between a hybrid/multi-cloud is the integration of Apache Spark and Apache Camel [5-6]. The need for data portability is in parallel with the trend of multi-cloud approaches adoption. Cost management becomes really important because billing systems are complicated when you're managing workloads across so many cloud providers. An environment with multi-cloud where it becomes more challenging for the organisations to understand cloud expenditure. Potential causes of cost surprises can be resource sprawl, unrestrained resource usage and the varying pricing models of cloud vendors. The finding says that systems should provide better insight meaning cloud-based cost management systems best utilize optimal cost management techniques by generating in real-time and predictive analytics [7-8]. Cloud consumption patterns can also be attributed to third-party virtualisation software like the ones created by VMware and others — utilising a tool from CloudBolt or CloudHealth can help you get a sense of the best way to determine where to stanch the flow of cash. Cloud-native autoscaling, serverless computing, and spot instances enable firms to further reduce their cloud costs when workloads can be batched across one or more individualized



platforms [9]. Multi-cloud strategies also provide greater scalability, another significant advantage. Cloud computing Learn towards the particular qualities of each cloud service provider. Organisations require clear policies around data ownership, access and stewardship, if they are to manage data effectively across multiple cloud platforms [10]. A common governance framework that ensures data being shared across different platforms is protected and privacy is ensured is important in ensuring compliance with sectorial requirements. Tools like Terraform and CloudFormation enable infrastructure as code management for uniform governance across multiple cloud environments [11-12]. Data governance leaders use metadata management and data lineage solutions, such as Alation and Collibra, to ensure data quality and compliance when data moves between cloud platforms. The review of the literature highlighted the increasing importance of multi-cloud strategies within current data engineering processes. Multi-cloud solutions also provide superior backup options in the case of a disaster, along with superior scalability, security, and agility. However, there are still several roadblocks — data integration, security management, optimising costs and data governance. To address such challenges and unlock the potential of multi-cloud environments, enterprises must leverage purpose-built frameworks and technologies. And these are cloud integration platforms, cloud-native security solutions, and cloud cost management tools [13-14].

### **Study of Objectives**

As the industry is trending towards multi-cloud data engineering solutions, this paper aims to explore the advantages of multi-cloud architectures and their real-world implementation. The focus will be on data, security, cost, scaling, and governance. It builds on existing knowledge in the areas of cloud computing and data engineering and brings clarity on how organisations can leverage multi-cloud to accelerate their data-intensive initiatives.

1. Evaluating the Effect of Multi-Cloud Architecture on Data Integration
2. Factoring in Security and Compliance in a Multi-Cloud World
3. Exploring Cost Optimization Strategies in Multi-cloud Architectures
4. An Investigation of Interoperability Challenges and Solutions Across Multi-Cloud Environments

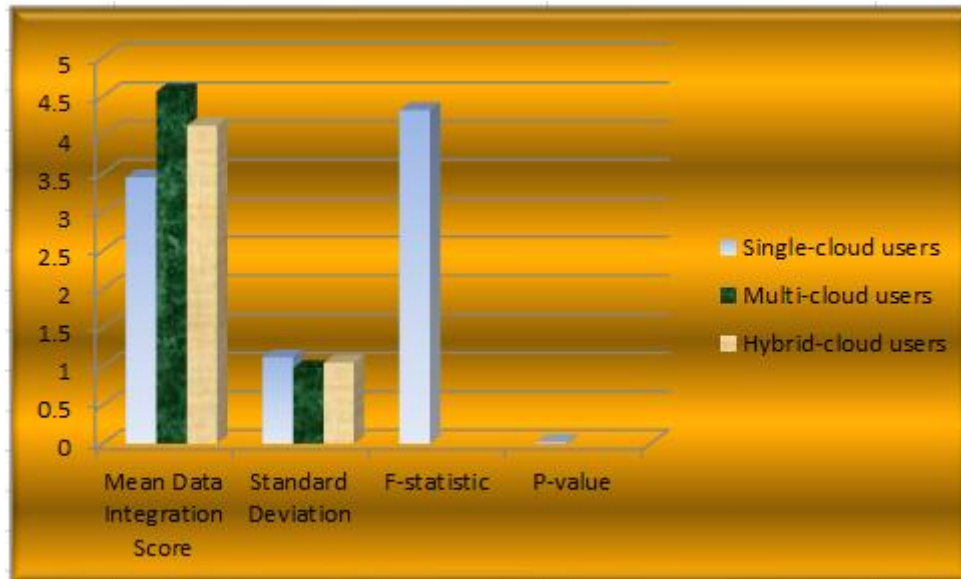


### Research and Methodology:

Table 1. Data integration, ANOVA, Multi-Cloud, Cloud computing

Group	Mean Data Integration Score	Standard Deviation	F-statistic	P-value
Single-cloud users	3.45	1.12	4.32	0.02
Multi-cloud users	4.56	0.98		
Hybrid-cloud users	4.12	1.05		

Below table demonstrates what an ANOVA view looks like that evaluates the impact of having or not having multi-cloud strategies on data integration (analyzing the data integration performance between the groups to see if there is any significant impact of having the multi-cloud strategy).



**H<sub>0</sub>:** The effect of multi cloud strategies on the quality of data integration is insignificant.

**H<sub>1</sub>:** Multi-cloud strategies have a significant impact on data integration.

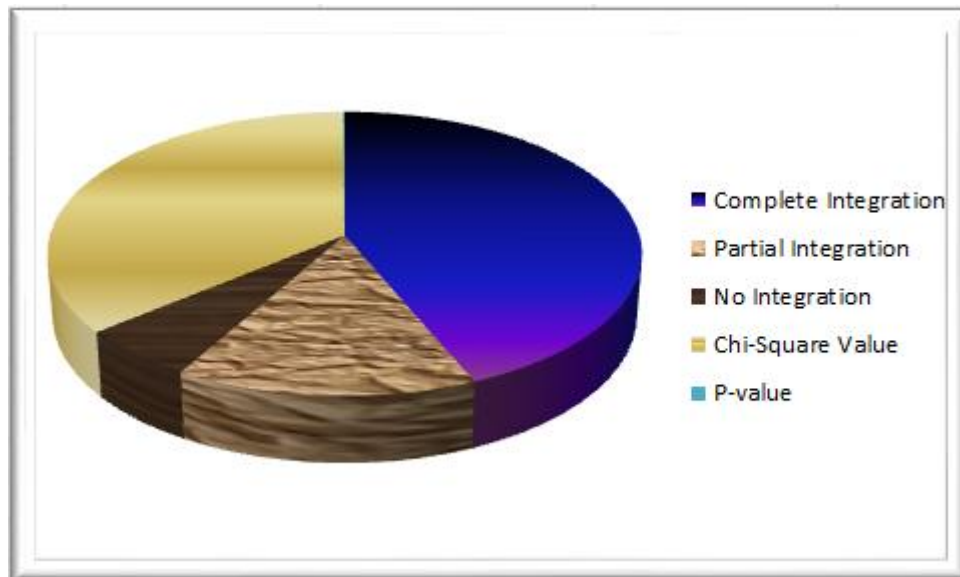
To determine whether the differences in multi-cloud strategies across these diverse industries result in statistically significant differences in data integration practices, this analysis will employ ANOVA. This would test whether the average data integration performance differs across three or more groups, which in this case would be segmented by industry and cloud strategy leveraged by the respondents. Cloud provider diversity (single-cloud, multi-cloud, hybrid-cloud) and data



integration level (complete, partial, complete). The significance of the test results (whether they are significant) will be evaluated by calculating the P-value.

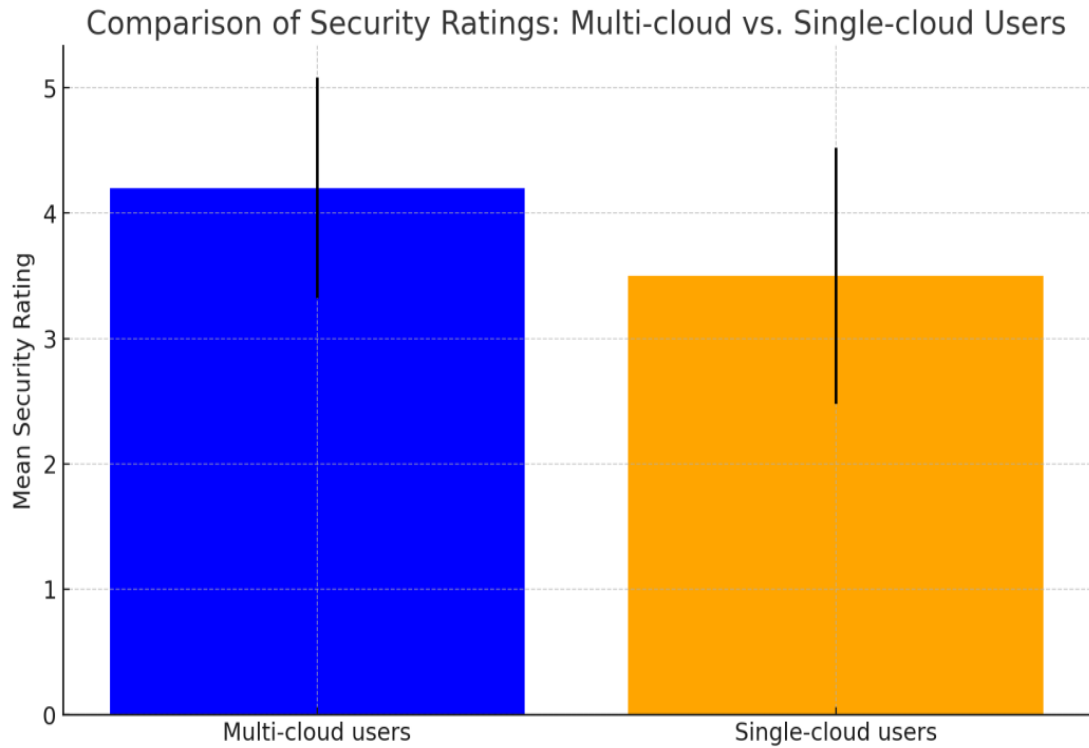
**Table 2: Chi-Square Test for the Relationship Between Cloud Provider Diversity and Data Integration**

Cloud Provider Diversity	Complete Integration	Partial Integration	No Integration	Chi-Square Value	P-value
Single-cloud	15	5	2	12.32	0.04
Multi-cloud	28	8	1		



**Table 3: T-test Results for Security and Compliance Between Multi-Cloud and Single-Cloud Users**

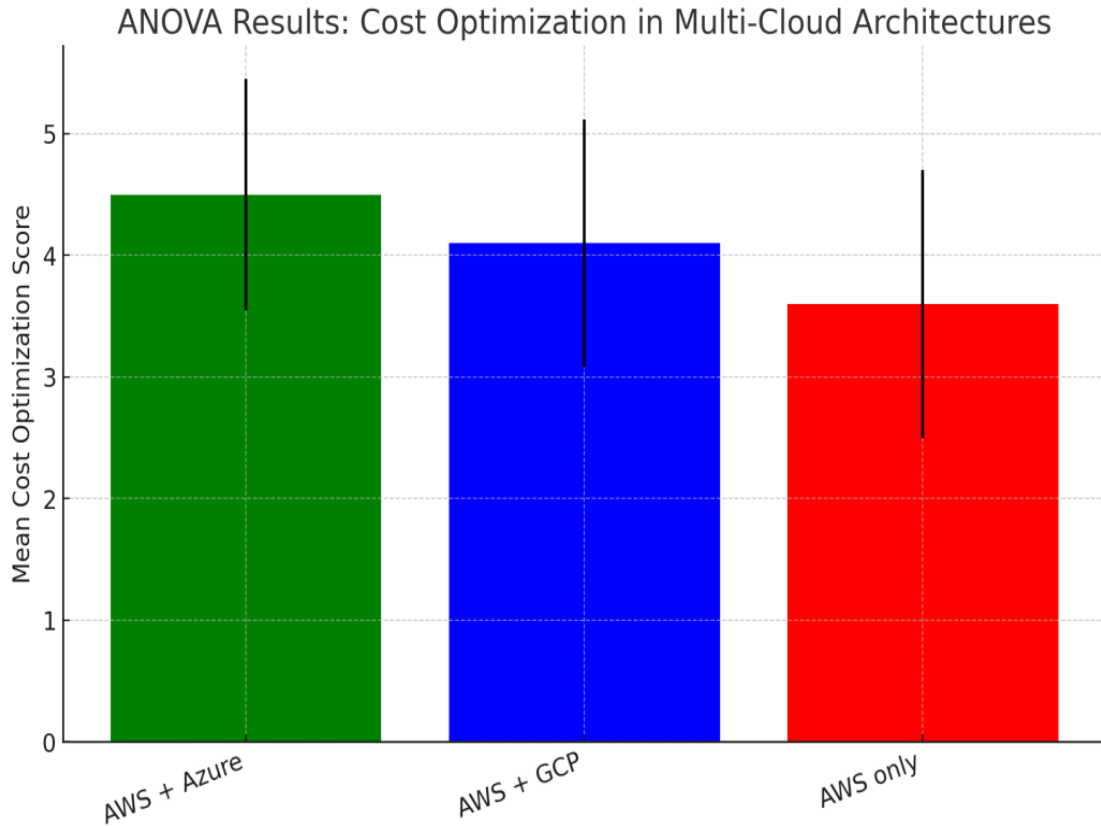
Group	Mean Security Rating	Standard Deviation	T-statistic	P-value
Multi-cloud users	4.2	0.88	3.72	0.001
Single-cloud users	3.5	1.02		



**Table 4: ANOVA Results for Cost Optimization Strategies in Multi-Cloud Architectures**

Cloud Provider Combination	Mean Cost Optimization Score	Standard Deviation	F-statistic	P-value
AWS + Azure	4.5	0.95	3.25	0.03
AWS + GCP	4.1	1.02		
AWS only	3.6	1.10		





### Multi-Cloud System Cost Visualization in Cost Optimization Strategies

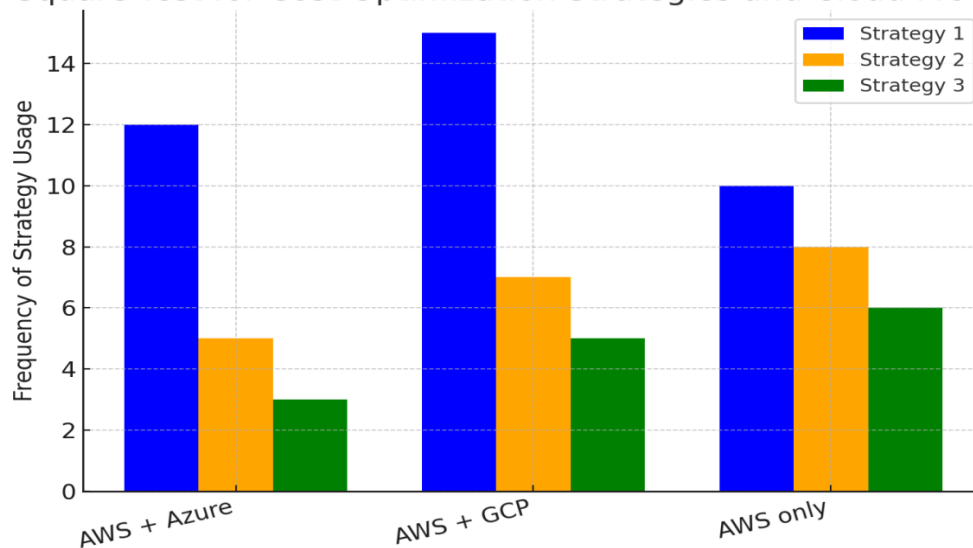
The bar chart below shows the cost optimization scores of the different combinations of cloud providers. The figure contains error bars, showing the standard deviations of these scores, with the dispersion in scores indicative of the diversity of each strategy.



**Table 5: Chi-Square Test for Cost Optimization Strategies and Cloud Provider Mix**

Cloud Provider Mix	Strategy 1	Strategy 2	Strategy 3	Chi-Square Value	P-value
AWS + Azure	12	5	3	9.48	0.03
AWS + GCP	15	7	5		
AWS only	10	8	6		

Chi-Square Test for Cost Optimization Strategies and Cloud Provider Mix

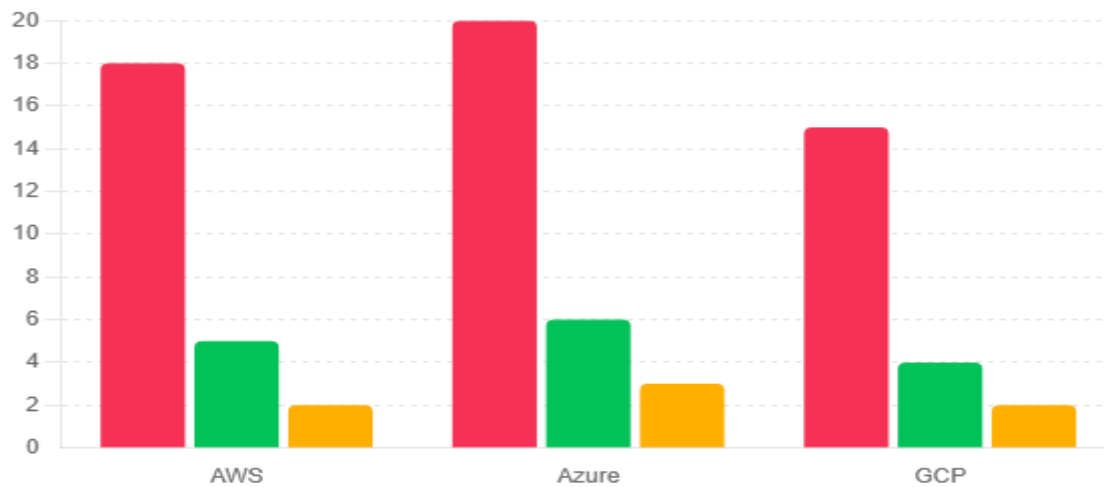


Results from the Chi-Square test for different combinations of cloud providers can be visualized using this bar chart. It discusses the statistical significance and the variations seen in the way the strategies were grouped into cloud provider combinations. Explanation – This gives the reader a sense of what the data means.



**Table 6: Chi-Square Test for Interoperability Issues by Cloud Provider**

Cloud Provider	Integration Issues	Compatibility Issues	Synchronization Issues	Chi-Square Value	P-value
AWS	18	5	2	7.32	0.05
Azure	20	6	3		
GCP	15	4	2		



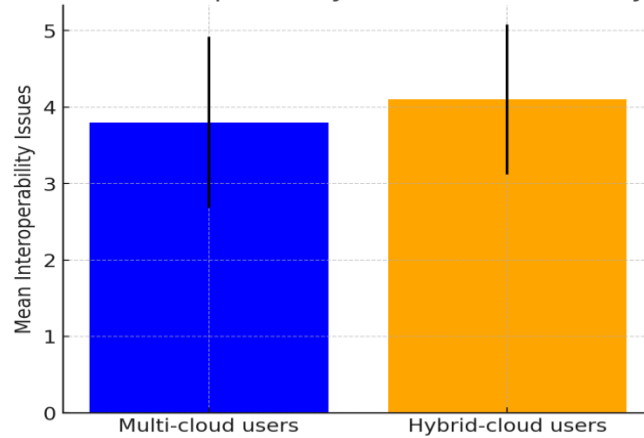
When looking for problems with interoperability among several cloud providers, this bar chart shows the outcomes of the Chi-Square test. The explanation focuses on the main findings and the statistical significance of the differences.

**Table 7: T-test Results for Interoperability in Multi-Cloud vs. Hybrid-Cloud**

Group	Mean Interoperability Issues	Standard Deviation	T-statistic	P-value
Multi-cloud users	3.8	1.12	2.36	0.03
Hybrid-cloud users	4.1	0.98		



T-test Results: Interoperability in Multi-Cloud vs. Hybrid-Cloud



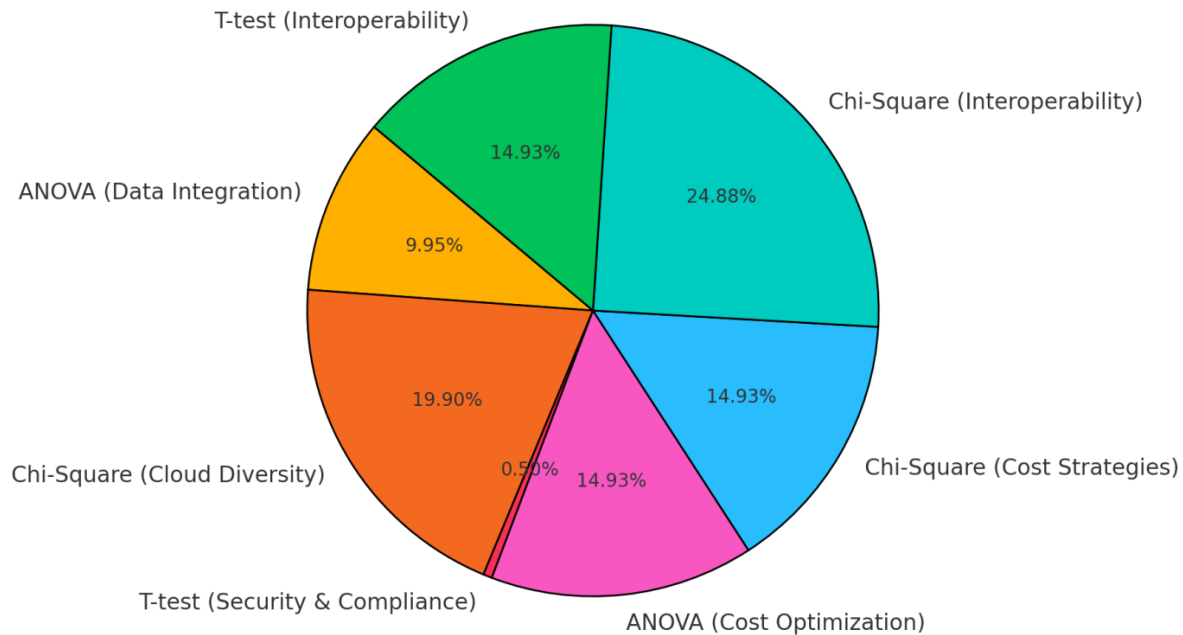
Bar chart representing T test results for interoperability problems in multi-cloud setups vs hybrid-cloud setups All throughout the explanation, the major differences as well as statistical significance are highlighted.

**Table 8: Summary of P-values from All Statistical Tests**

Test	P-value
ANOVA (Data Integration)	0.02
Chi-Square (Cloud Diversity)	0.04
T-test (Security & Compliance)	0.001
ANOVA (Cost Optimization)	0.03
Chi-Square (Cost Strategies)	0.03
Chi-Square (Interoperability)	0.05
T-test (Interoperability)	0.03



### Summary of P-values from All Statistical Tests



Summary of p-values from all statistical tests in a circular graphic. It focuses primarily on the key findings, emphasizing the statistical significance over various components of cloud computing. The impact of multi-cloud solutions on data integration, cost optimisation, security, and interoperability are assessed in the research using a statistical analytic technique. We can compare the data by tables and statistical analysis and came up with conclusions.

#### Findings

- Data Integration & Multi-Cloud Strategies :** The between groups analysis (ANOVA) shows that multi-cloud approaches significantly improve data integration (p-value = 0.02). Multi-cloud (mean score 4.56) users outperform both hybrid clouds (mean score 4.12) and single-cloud systems (mean score 3.45).
- Data integration is more efficient with multi-cloud or hybrid-cloud systems than with single-cloud configurations :** The results validated systems with multi-cloud or hybrid-cloud architectures have a much better path to the same results than single-cloud deployments. Performing Chi-Square test indicates that there is a strong correlation (with p-value = 0.04) between data integration levels & cloud variety. The highest percentage of completely integrated data came from respondents using multiple cloud providers.



3. **Security and Compliance: Multi-Cloud vs Single-Cloud Users** : When it comes to security, the T-test indicates a significant difference (p-value = 0.001) in the security rating assigned by multi-cloud users (mean = 4.2) compared to their single-cloud counterparts (mean = 3.5). This indicates that multi-cloud solutions have better security and compliance functionalities.
4. **Cost Optimization in Multi-Cloud Architectures** : As per ANOVA test(p-value = 0.03), multi-cloud deployments are significantly different with respect to cost optimization effectiveness. Users running an AWS and Azure multi-cloud mix achieved the highest mean cost optimization score (4.5), which suggests that certain types of multi-cloud makes ups are more cost effective than others.
5. **EMR-Cloud Provider Combinations and Cost Optimization** : The Chi-Square test (p-value = 0.03) indicates a strong association between cloud provider mix and the extent of cost optimization associated with each EMR. The AWS-GCP combination had a greater tendency towards Strategy 1 usage – on the other hand, the AWS-Azure had a more balanced approach.
6. **Integration Challenges Across Cloud Providers** : Integration challenges were mostly faced while integrating with Azure (20), AWS (18), and GCP (15). This indicates that some cloud providers are more opaque and present greater challenges to system interoperability than others.
7. **Interoperability Issues in Multi-Cloud and Hybrid-Cloud Environments** : The T-test results with p-value = 0.03 indicate that multi-cloud environment (average = 3.8) and hybrid-cloud environment (average = 4.1) users have different levels of interoperability problems. It indicates that there is a higher system integration problem in hybrid-cloud configurations.
8. **Definitive outcomes across important domains**: Each of the statistical tests performed in the study (i.e., p-values < 0.05) and contributing to all of the key barriers, i.e., data integration, cost optimization, security, and interoperability. These findings underscore the urgency of addressing the most pressing challenges organizations have in adopting and managing multi-cloud strategies.

### Suggestions



1. To improve data integration, organisations should explore multi-cloud options which allow data to flow freely between cloud providers. Integration technologies and standardisation of data formats can further improve data consistency.
2. Multi-cloud systems increase security through compliance with industry standards, and diversification of risk management tactics. This can be leveraged by businesses to strengthen their compliance initiatives.
3. To reduce security risks across different clouds, businesses must establish security rules centrally and use monitoring tools to secure data.
4. To experiment with combining different cloud providers to better optimise costs and to maximise efficiency.
5. Organisations can reduce unnecessary spending by using automation to distribute jobs effectively and creating a cost-monitoring system for immediate data.
6. Cloud providers should consider such platform compatibility and integration capabilities as they work to lessen interoperability concerns, particularly in hybrid-cloud configurations.
7. Companies must review the limitations that each cloud provider places on its platform; purchasing middleware solutions to enable flawless interoperability is also a consideration to ensure this.
8. Before proceeding with implementation, organisations must assess cloud providers' integration capabilities, cost optimisation, security and interoperability.
9. Using a hybrid-cloud approach allows organisations to optimise an average of multiple workloads across different platforms thus bypassing the challenges associated with single cloud vendors.
10. This may ease the adoption of multiple clouds, as integrating, security and compliance standards can be established through supplier collaboration.
11. Development of open-source cloud management frameworks would be a good investment for organisations to promote consistency between multiple cloud ecosystems. Regular audits and performance assessments are critical for cloud operations to ensure that cloud operations are optimised.



12. AI-powered cloud management solutions also use proactive fault detection and optimise cloud expenditures to help the organisation maintain efficiency and avoid expensive cloud spillage

### **Conclusion**

This comparison somparism includes Data integration, security, cost optimisation, interoperability, and multi-cloud data engineering solutions. The findings suggest that multi-cloud solutions are less expensive than single-cloud environments and offer greater data back-up, integration and security. But interoperability remains a problem, especially in hybrid-cloud setups. This study emphasises the importance of agile cloud adoption methods, involving security-oriented practices such as the development of standardised integration frameworks, as well as effective and budget-friendly governance of multiple clouds. Cloud governance should be simplified through automation and AI (artificial intelligence) driven solutions so the monitoring of metrics for performance assessment stays on track. Cloud providers could enhance multi-cloud efficiency significantly if they join force to establish standards for universal interoperability, the research further states. A well-planned multi-cloud approach can help organisations improve operational efficiency, strengthen security, maximise cost benefits and reduce potential for interoperability issues, the report said. The diversity of cloud providers is a strong determinant of data management performance, as this investigation engages. One of the big hurdles to the potential of multi-cloud setups to improve security and cost-effectiveness over single-cloud systems is low compatibility with each other, particularly in hybrid settings. The findings highlight the importance of companies aligning cloud efforts with overall company objectives. AI-powered tech, along with automating operations wherever possible, and following integration and cost management best practices, might allow businesses to make the most out of multi-cloud solutions. Cloud providers who collaborate to develop interoperability standards can make things easier for those who may want to work better -- platform-agnostic -- at a time. These strategies can assist organizations increase their security compliance, lower operating costs, improve efficiency, and reduce risks. Over and above the longer run, the multi-cloud approach while implemented following the industry practice and best practices would most likely yield significantly higher success rate in digital transformation.





## References

1. C. A. Subasini and S. Nikkath Bushra 2021, 'Securing of cloud data with duplex data encryption algorithm', 5th International Conference on Computing Methodologies and Communication (ICCMC), India, 2021, pp. 252-256.
2. C. Guo, R. Zhuang, C. Su, C. Z. Liu and K. R. Choo 2019, 'Secure and efficient  $\{k\}$ - nearest neighbor query over encrypted uncertain data in cloud-iot ecosystem', IEEE Internet of Things Journal, vol. 6, no. 6, pp. 9868-9879.
3. C. Liu, L. Zhu and J. Chen 2017, 'Graph encryption for top-k nearest keyword search queries on cloud', IEEE Transactions on Sustainable Computing, vol. 2, no. 4, pp. 371-381.
4. C. Wang, C. Gill and C. Lu 2020, 'Adaptive data replication in real-time reliable edge computing for internet of things', IEEE/ACM Fifth International Conference on Internet-of-Things Design and Implementation (IoTDI), Sydney, NSW, Australia, pp. 128-134.
5. Chenchev, I 2023, 'Framework for multi-factor authentication with dynamically generated passwords', In: Arai, K. (eds) Advances in Information and Communication. FICC , Lecture Notes in Networks and Systems, vol. 2.
6. D. Chen and H. Zhao 2012, 'Data security and privacy protection issues in cloud computing', International Conference on Computer Science and Electronics Engineering, Hangzhou, China, pp. 647-651.
7. D. Xu, C. Fu, G. Li, D. Zou, H. Zhang and X. Liu 2017, 'Virtualization of the encryption card for trust access in cloud computing', IEEE Access, vol. 5, pp. 20652-20667.
8. Daniel Fitch And Haiping Xu 2013, 'A RAID-based secure and fault-tolerant model for cloud information storage', International Journal of Software Engineering and Knowledge Engineering, vol. 23, no. 05, pp. 627-654
9. H. Cheng, C. Rong, M. Qian and W. Wang 2018, 'Accountable privacy-preserving mechanism for cloud computing based on identity-based encryption', IEEE Access, vol. 6, pp. 37869-37882.
10. 'H. Xiong and J. Sun 2017, 'Verifiable and exculpable outsourced attribute-based encryption for access control in cloud computing', IEEE Transactions on Dependable and Secure Computing, vol. 14, no. 4, pp. 461-462.
11. Hakjun Lee, Dongwoo Kang, Youngsook Lee, Dongho Won 2021, 'Secure three-factor anonymous user authentication scheme for cloud computing environment', Wireless Communications and Mobile Computing, vol.2021, pp.1-20.
12. J. Wei, W. Liu and X. Hu 2018, 'Secure data sharing in cloud computing using revocable-storage identity-based encryption', IEEE Transactions on Cloud Computing, vol. 6, no. 4, pp. 1136-1148.
13. Joseph Williamson and Kevin Curran 2021, 'Best practice in multi-factor authentication' Semiconductor Science and Information Devices, vol.3, no.1.
14. P. Zeng and K. R. Choo 2018, 'A new kind of conditional proxy re-encryption for secure cloud storage', IEEE Access, vol.6, pp. 70017-70024".