

# OPTIMIZED HIGH-VOLUME TRANSACTIONS IN AZURE COSMOS DB WITH PARTITIONING TECHNIQUES

Saatwik Gilakattula\*

\*Purdue university, Indianapolis

E mail ID- [saathwik@gmail.com](mailto:saathwik@gmail.com)

## Abstract

The fast-growing surge in cloud-native and globally distributed applications has heightened the need to have database systems that can support high-volume transactional workloads with low latency, scalability according to high workload, and with high reliability. Azure Cosmos DB has become a leading NoSQL database system that is capable of fulfilling these needs by its globally distributed design, elastic provisioning of throughput, and scalable consistency models. The data partitioning used in the central concept of scalability and performance allows the movement of horizontal scaling and efficient utilization of resources. But the best transactional performance in Cosmos DB is very sensitive to the best partitioning strategies and the work load sensitive data models. The given review paper gives an in-depth analysis of optimized high-volume transactions in Azure Cosmos DB with special focus on partitioning techniques. It discusses the architecture of Cosmos DB, the developed methods of partitioning in distributed NoSQL systems, and the effects of partition key selection, workload fit, and transactional boundaries on throughput, latency, and cost efficiency. Practical optimization techniques, challenges that may occur, including cross-partition transactions and partition key rigidity, and trade-offs between consistency, availability, and performance are also discussed in the paper. Moreover, the review has also identified the new research directions such as AI-based adaptive partitioning, event-based architectures, and self-tuning database mechanisms. This paper seeks to offer practical knowledge to both researchers and practitioners in designing and optimizing transactional workloads in the Azure Cosmos DB scalable and cost-effective in both bulk and scale by synthesizing academic research and system-level studies.

**Keywords:** Azure Cosmos DB; High-volume transactions; NoSQL databases; Data partitioning; Distributed systems.

## 1. Introduction

The proliferation of cloud-based globally distributed applications has greatly escalated the need of database systems that have the capability of supporting high volume transactional workloads with low latency and high availability. The modern applications of e-commerce services, financial services, Internet of Things (IoT), and real-time analytics have resulted in huge amounts of simultaneous reads and writes, frequently with geographically distributed users. The conventional relational database management systems (RDBMS) that depends most on vertical scaling and remains on inflexible schema fail to address these demands as they face scalability bottlenecks and high operational costs [1]. NoSQL databases have become a formidable solution to overcome these limitations, with newly introduced databases being able to scale horizontally with any type of schema, and with high availability. Scalability in these systems is mainly brought out by data partitioning also known as sharding where data is shared among more than two nodes to balance the load and enhance

throughput. Yet, the level of success of NoSQL systems within the transactional setting largely relies on the way the partitioning is organized and executed. Thoughtless partitioning techniques may result in hot partitions, poor usage of resources, increased latency, and high costs of operation [2].

Azure Cosmos DB, a globally replicated, multi-model NoSQL database service offered by Microsoft, can be used to run mission-critical applications on a planetary scale. It has capabilities that include automatic indexing, consistent levels, elastic throughput provisioning and transparent global replication. The fundamental feature of the Cosmos DB scalability is its partitioned architecture that separates the data into logical partitions that are remapped to physical partitions of distributed infrastructure. It is a design that helps Cosmos DB to be scalable yet predictable in its performance with high-throughput transactional workloads [3]. Although it has a sophisticated structure, it does not automatically support the best performance of Azure Cosmos DB. The choice of a suitable partition key, knowledge of what transactions can be performed in partitions and matching data models with access patterns in the workloads are all design choices that have a direct relationship with throughput, latency and cost efficiency. Static partitioning strategies can be unsuitable as workloads change, and informed and workload-sensitive partitioning strategies [4] are needed.

The current review paper is dedicated to high-volume transactions, which are optimized in Azure Cosmos DB, and the classification of partitioning techniques is paid specific attention. It consolidates the currently available studies on NoSQL strategies of partitioning, mechanisms of partitioning in the architecture of Cosmos DB and the effects of these mechanisms on the performance of transactions, as well as best practices, challenges, and future research. Combining scholarly knowledge and system-wide research, the paper will equip the researcher and practitioner with all the details on how effective partitioning would greatly reduce transactional scalability in Azure Cosmos DB.

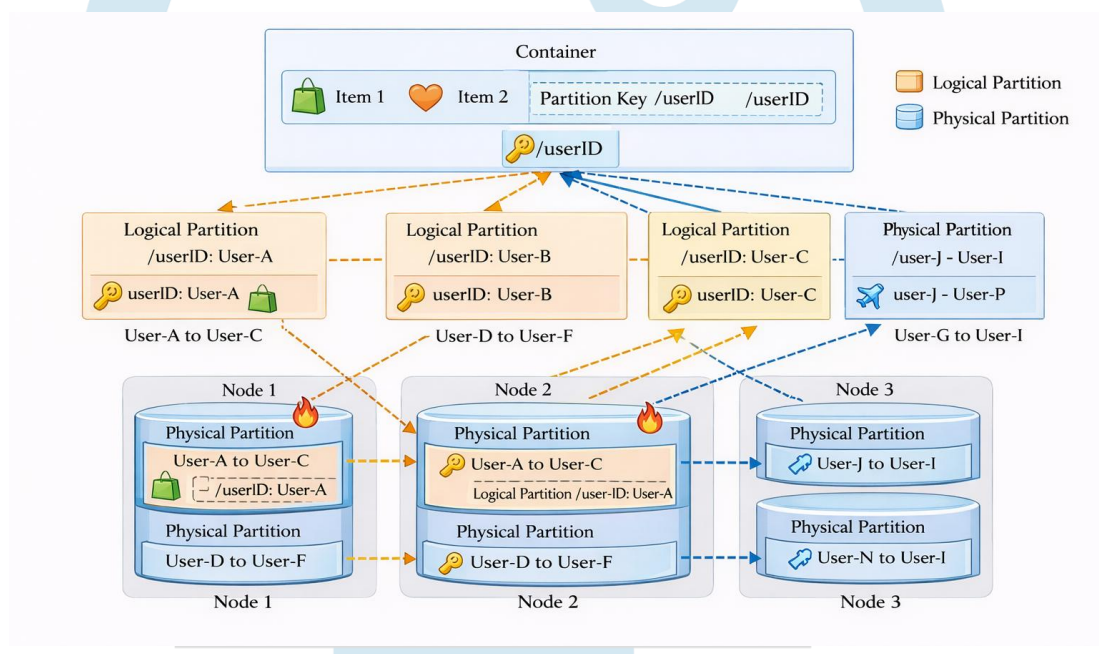
## **2. Azure cosmos DB architecture for high-volume transactions**

Azure Cosmos DB is a cloud-native fully managed NoSQL database that is intended to be used in a high throughput and low latency workload of transactional workload on a global scale. The basic idea behind its architecture is the horizontal scalability, which is implemented by automatic data distribution and replication in different regions. Compared to the old-fashioned databases where vertical scaling was the order of the day, Cosmos DB divides data and workloads into smaller autonomous units that the system can scale down and up accordingly depending on increasing or decreasing transaction volume [5].

One of the fundamental design principles in Cosmos DB is the difference between logical partitions and physical partitions. The application determines logical partitions using a selected key, which is a partition key, and these partitions include items with identical partition key values. These logical partitions are then overlaid to physical partitions which are the actual storage and compute units that the Cosmos DB service has to manage. There are predetermined limits to the storage capacity and throughput of each physical partition and the system automatically adds or divides physical partitions in response to the increase in workload. This structure allows Cosmos DB to accommodate millions of operations in a second and with performance characteristics that are predictable [6].

The Cosmos DB measures throughput provisioning and measurement in Request Units per second (RU/s), which are an abstraction of the cost of CPU, memory and I/O resources necessary to execute database transactions. In case of high-volume transaction systems, RU assignment should be well-adjusted with the data distribution among the partitions. Logical partitions can be evenly distributed so that RU consumption can be uniformly distributed across physical partitions and throughput can be maximized. On the other hand, the skewed data distributions may cause underutilization of resources in one partition and throttling in others,

adversely affecting the transaction latency and system efficiency. The other important architectural characteristic that has impact on the transactional performance is the availability of Cosmos DB to support different consistency models that include strong consistency to eventual consistency. The models enable developers to trade off between latency, availability and correctness depending on application needs. Greater levels of consistency offer more stringent transactional guarantees at the cost of more coordination cost over across replicas, whereas weaker consistency models can offer better write throughput and less latency. When the level of consistency is being selected, it becomes very sensitive to the partitioning design in high-volume transactions situations, due to the ability of cross-partition coordination to substantially impact performance of the system as a whole [7]. In general, the architecture of Azure Cosmos DB offers a platform that supports scalable transaction processing. Nonetheless, the usefulness of this architecture in practice relies a lot on the nature with which partitioning schemes, throughput provisioning and consistency settings are integrated. It is important to understand these architectural elements before analyzing particular partitioning methods and optimization techniques with high volume of transactional workloads (Figure 1).



**Figure 1: Azure cosmos DB logical and physical partition architecture.**

This diagram describes the process of partitioning items in a container with the help of partitions keys in logical partitions that are subsequently partitioned to physical partitions and nodes.

### 3. NoSQL distributed databases partitioning methods

Partitioning is a basic tool that allows distributed NoSQL databases to scale horizontally and provides reasonable performance under high transaction loads. Partitioning improves fault tolerance, contention, and parallelism by dividing large dataset into smaller and manageable parts and allocating these parts to several nodes. When dealing with high-volume transactional setting, the partitioning technique impacts directly throughput, latency, and rate of distributed transactions [8]. The hash-based partitioning technique is one of the most widely used partitioning techniques whereby a hash function is used to define the location of data using a partition key. This is a good load balancing method with uniform distribution of data and is applicable in the workloads with random access patterns. Nevertheless, hash-based partitioning usually does not take into account the relationship between the data and as a result, more cross-partition operations may be done in case the transaction needs related data, which is in a different partition. Distributed transactions create overhead in coordination and may lead to worse performance in a write-intensive system [9]. The range-based

partitioning is another widely used method whereby data is partitioned based on the ordered key ranges. It works especially well with workloads that deal with range queries or time-series data, since it maintains data locality. However the range based partitioning is vulnerable to data and workload skewness particularly where access pattern becomes skewed or varies with time. Hotspots may be formed when some critical ranges get a disproportionately large amount of traffic causing uneven use of resources and limited transaction throughput.

In an effort to overcome these constraints of the static partitioning methods workload aware and graph partitioning methods have been suggested. The techniques examine patterns of access and relationships between data to co-locate commonly used or related data in the same partition. Specifically, graph-based partitioning represents data items as nodes and the transactional relationships as edges and strives to decrease cross-partition edges. Workload-aware partitioning can greatly enhance transactional NoSQL systems response time and throughput by decreasing the quantity of transactions sent out [10]. Despite the high performance advantages of the advanced partitioning techniques in both experimental and academic environments, the use of this technique in production systems is difficult. Unpredictable workloads, changing data structure, and operational complexity make it hard to keep changing partitioning schemes without service interruption. This has prompted most cloud-native NoSQL databases, such as Azure Cosmos DB, to use less complex application-level abstractions around partitioning alongside physical data distribution in the internal implementation. The knowledge of these general partitioning methods offers important background to the exploration of how Cosmos DB offers and optimizes partitioning of high-volume transactional loads in practice (Table 1).

**Table 1. Comparison of partitioning techniques for high-volume transactional workloads**

Partitioning Technique	Partitioning Basis	Key Advantages	Major Limitations	Suitability for High-Volume Transactions	Relevance to Azure Cosmos DB
Hash-Based Partitioning	Hash function applied to partition key	Uniform data distribution; good load balancing; simple implementation	Ignores data relationships; increases cross-partition transactions	High throughput for random access workloads	Commonly used through partition key hashing for balanced RU utilization
Range-Based Partitioning	Ordered key ranges (e.g., time, ID intervals)	Preserves data locality; efficient range queries	Prone to hotspots and workload skew	Moderate, dependent on evenly distributed access patterns	Applicable for time-series containers but requires careful key design

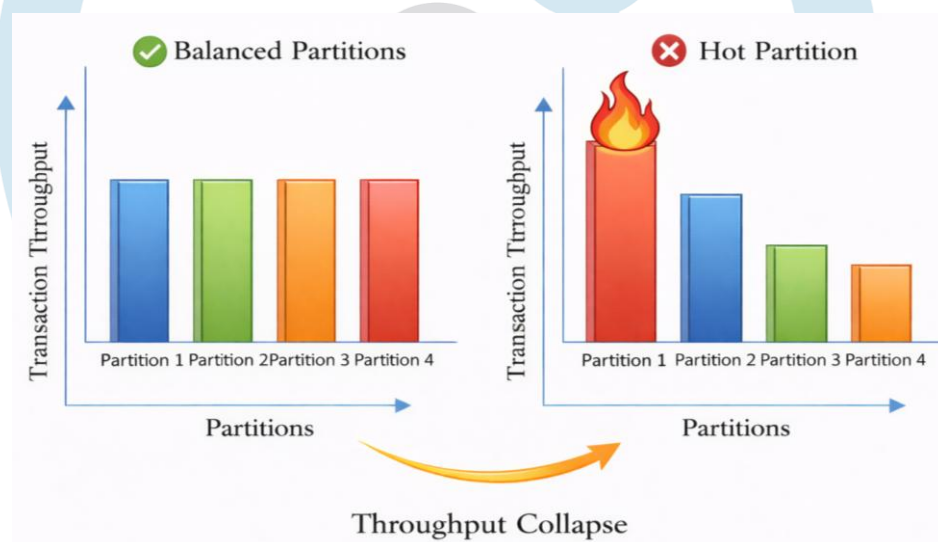
Round-Robin Partitioning	Sequential distribution across nodes	Simple distribution; avoids initial hotspots	Poor data locality; high coordination cost	Low suitability due to frequent cross-partition operations	Not natively supported due to inefficiency for transactional workloads
Workload-Aware Partitioning	Access patterns and transaction frequency	Reduces distributed transactions; improves latency	High implementation complexity; dynamic overhead	High, especially for write-intensive workloads	Conceptually aligns with Cosmos DB best practices for partition key selection
Graph-Based Partitioning	Data relationships modeled as graphs	Co-locates related data; minimizes cross-partition edges	Computationally expensive; limited production adoption	Very high for relationship-heavy transactions	Useful for data modeling decisions though not directly implemented

#### 4. Optimizing high-volume transactions using azure cosmos DB partitioning

Maximizing the high-volume transactional workloads in the Azure Cosmos DB cannot be achieved through the mere activation of the partitioning feature, but rather requires a fine balance between the design of the partition key, access patterns of the workloads, and the transactional delimits. Although Cosmos DB hides much of the underlying infrastructure management, the application-level design choices continue to be the major factor in system performance and cost efficiency. The platform provides a scaling advantage missing in write intensive and latency sensitive applications where suboptimal partitioning is destructive of the scaling advantage [11]. One key aspect of optimization that should be critical is the choice of partition key. A partition key must have high cardinality and uniform distribution of data as well as support the most common transactional access patterns. Cosmos DB is able to provide atomicity and isolation when transactional operations are bounded by a single logical partition without introducing the overhead of distributed coordination. On the other hand, ill-selected keys tend to result in cross-partition transactions and this is more costly and adds further latency. These inefficiencies are amplified as the volumes of transactions grow hence the direct influence on throughput and response times[12].

Workload-aware data modeling can also be another key optimization strategy. Applications with a high volume of use do tend to have predictable access patterns (e.g. user-centric work or time-based work). By matching data models with these patterns, related items are positioned in the same logical partition, minimizing the number of times cross-partition queries and writes occur. Experiments of real-time and collaborative systems implemented on Cosmos DB show that the use of the functionality of change feed processing and event-driven architectures can further enhance the efficiency of transactions with a carefully

designed partitioning strategy. Unlike in the past where the workloads remain constant, it may be the case that the current partitioning strategies prove ineffective as workloads vary. New studies also point to the promise of adaptive and smart optimization methods such as the analysis of access logs and transactional activity by AI. These strategies seek to prescribe or program partitions which evolve with changes in volume and characteristics of workloads. Although these techniques remain mostly experimental, they are an indication of how future optimization frameworks may operate in which Cosmos DB deployments can dynamically scale to maintain large volume transaction demands in the long term without human intervention [13]. To conclude, high-volume operations in Azure Cosmos DB require a comprehensive solution combining 1) partition key design, 2) workload-sensitive modeling, and 3) workload-sensitive optimization concepts. Although the platform offers a strong architectural support on scalability, ultimately the platform offers informed design decisions based on the system capabilities and application-level needs in order to achieve a consistent and cost-effective transactional performance (Figure 2).



**Figure 2: Impact of choice of partition key on transaction.**

The idea of this conceptual graph is that balanced partitions, which leave throughput stable, are compared to skewed partitions, which creates a hot partition. It points out the tendency of throughput stalling because of lopses in partition key distribution, which increases the system overall performance.

## 5. Issues and constraints of partitioned high-volume transactions in Azure cosmos DB

Although Azure Cosmos DB has strong architecture and scalability capabilities, it has a number of limitations as far as it can support high-volume transactional workloads, especially with regard to partitioning. Although partitioning allows a horizontal scaling, it creates design complexities that may affect performance, consistency, and the operational cost negatively unless handled carefully. Such issues do not exist only in Cosmos DB, but are also inherent to distributed NoSQL systems at scale [14]. The processing of cross-partition transactions is one of the weaknesses. Cosmos DB only ensures that ACID transactions are performed in a single logical partition. In cases where the transactional operations in one partition interact with those in other partitions, developers have to use application level coordination or compensating transaction, making them more complex and less reliable. This weakness is particularly troublesome with workloads having complicated relationships or multi-entity updates, in which the imposition of strict transactional semantics is essential [15].

Partition key immutability is yet another important obstacle. After a container has been created the partition key cannot be modified without a data migration to a new container. Access patterns can change as the application workloads change with time, making the initial partitioning strategy no longer optimal. This rigidity may result in performance loss, higher request unit (RU) usage, and downtime of operations in the process of migration or re-architecture work [16].

In high-volume cases, predictability of the costs also is a problem. The RU-based pricing model of Cosmos DB is transparent but leaves the possibility of inefficient use of RU because of an unequal distribution of data or workload. Hot partitions can be throttled and the rest of the partitions might be underutilised, resulting in higher costs with no proportionate improvement in performance. Continuous monitoring and proactive workload tuning is thus needed to design cost-efficient scalability [17]. Lastly, there is the trade-off between consistency, availability, and latency that has an intrinsic constraint on transactional optimization. Stronger consistency guarantees can lead to a higher latency to applications, as well as the reduced throughput, because of the greater coordination overhead between replicas and regions. Although Cosmos DB also provides several consistency models, choosing the wrong level may negatively affect the advantages of partitioned scalability [18].

## 6. Future research directions

With the ongoing evolution of cloud-native applications, future studies of high-volume transactions optimization in Azure Cosmos DB will shift away with respect to fixed partitioning schemes to more adaptive, intelligent, and application-aware schemes. Though present-day best practice has focused on cautious initial partition key design, dynamic loads and unpredictable access patterns are increasingly requiring systems that can dynamically adapt themselves as changing transactional behavior [19]. The area of AI-based workload analysis integration to aid adaptive partitioning and throughput optimization is one promising area. Machine learning models can be used to detect which hotspots are emerging by examining past access logs, transaction counts, and latency profiles, and suggest changes to partitioning or redistribute throughput. The initial tests of AI-controlled solutions based on Cosmos DB prove that smart data placement and automation can considerably decrease the operational overhead at the cost of preserving high-level transaction performance [20]. The other significant area of research is the hybrid and event-driven architecture, where Cosmos DB is integrated with other supplementary Azure services. Event streams, serverless functions, and change feed-based processing allow systems to take transactional load off, and their responses under heavy loads become more responsive. Next generation transactional systems will probably be based on modes of asynchronous processing, loosely coupled and which minimize the contention on the database level and maintain logical consistency at the application level [21]. Lastly, self-tuning and autonomous database management mechanisms to be used with globally distributed NoSQL systems should be attempted in the future. These are automated repartitioning, workload sensitivity consistency-based adaptation, and request unit (RU/s) predictive scaling. With the growing global and latency-sensitive nature of transactional workloads, these capabilities will be needed to ensure the provision of performance guarantees without human intervention. The innovations in these directions will support further the power of the Azure Cosmos DB as a platform of the next generation high-volume transactional applications.

## 7. Conclusion

The review has discussed how partitioning techniques can be used to optimize high-volume transactional workloads in Azure Cosmos DB. Although Cosmos DB offers a highly scaled and distributed architecture worldwide, the performance of the system in real life is heavily conditioned by the application level architecture choices, especially the choice of partition key and workload responsive data modeling. With proper partitioning, it is possible to have balanced usage of the resources, reduce the number of cross-partition transactions, and have predictable throughput and latency during heavy transactional loads. Other essential issues such as transactional limitations among partitions, the immutability of partition keys, predictable costs, and trade-offs of consistency among distributed NoSQL systems are also pointed out in the analysis. Moreover, new studies indicate that adaptive and AI-based optimization and event-driven architecture has high potential in enhancing transactional efficiency in dynamical workloads. All in all, this paper highlights that a comprehensive view of optimal performance in Azure Cosmos DB needs to integrate the architectural knowledge, accurate partitioning plans, and ongoing workload monitoring to facilitate the successful use of transactional application scale and resiliency.

## References

- [1] Mihai, G. (2020). *Comparison between relational and NoSQL databases*. Annals of Dunarea de Jos University of Galati, Fascicle I: Economics and Applied Informatics, 26(1), 45–54.
- [2] Devashish Ghanshyambhai Patel. (2025). *Supply Chain Security in Cloud: Implementing Tamper Resistant Image Life Cycle Management*. International journal of innovative research in technology, 12(1), 530-537.
- [3] Hillar, G., & Yondem, D. (2018). *Guide to NoSQL with Azure Cosmos DB*. Birmingham, UK: Packt Publishing.
- [4] Rowe, J., Horal, M., Sundar, H. S., Arumugam, M., & Köse, B. (2025). *Implementing decentralized per-partition automatic failover in Azure Cosmos DB*. arXiv, abs/2505.14900.
- [5] Upreti, N., Sundaram, K., Sundar, H. S., Boshra, S., Perumalswamy, B., & Simhadri, H. (2025). *Cost-effective, low latency vector search with Azure Cosmos DB*. Proceedings of the VLDB Endowment, 18, 5166–5183.
- [6] Mazurova, O., Andrushchenko, M., & Shirokopetleva, M. (2023). *Research of methods of software implementation of the Cosmos DB API on the .NET platform*. Innovative Technologies and Scientific Solutions for Industries, 24, 118–129.
- [7] Diogo, M., Cabral, B., & Bernardino, J. (2019). *Consistency models of NoSQL databases*. Future Internet, 11(2), 43.
- [8] Krstić, J., & Krstić, S. (2018). *Testing the performance of NoSQL databases via the database benchmark tool*. Vojnotehnički Glasnik, 66(3), 614–639.
- [9] Gidado, A. A., & Ezeife, C. I. (2025). *UniqueNOSD: A novel framework for NoSQL over SQL databases*. Journal of Big Data, 12, Article 130.

- [10] Sahu, A., & Ahirrao, S. (2018). *Graph-based workload driven partitioning system by using MongoDB*. International Journal of Advances in Applied Sciences, 7(1), 29–37.
- [11] Tadi, S. R. C. C. T. (2024). *Interactive document editing and distributed synchronization using Azure Cosmos DB and WebSockets*. International Journal of Advanced Research in Science, Communication and Technology, 4(3), 1–10.
- [12] Devashish, GT. (2025). *Generative AI in DevOps: Enhancing Cloud Workflow Automation*. International journal of innovative research in technology, 12(2), 3732- 42.
- [13] Malakar, A. (2025). *Data-driven migration strategies: Leveraging GenAI for relational to NoSQL cloud database migrations*. International Journal of Engineering and Computer Science, 14(9), 5266–5275.
- [14] Ahmed, J., Karpenko, A., Tarasyuk, O., Gorbenko, A., & Sheikh-Akbari, A. (2023). *Consistency issue and related trade-offs in distributed replicated systems and databases: A review*. Radioelectronic and Computer Systems, 2, 14–32.
- [15] Mahfoud, Z., & Nouali-Taboudjemat, N. (2019). *Consistency in cloud-based database systems*. Informatica, 43(3), 1–14.
- [16] Devashish G.T. (2025). *Automating multi-cloud workflows with packer and terraform*. International Journal of Creative Research Thoughts, 13(6), 987-99.
- [17] Valavandan, R., Gothandapani, B., Gnanavel, A., Ramamurthy, N., & Balakrishnan, M. (2023). *Leveraging Azure platform data services for efficient data analytics in the oil & gas industry*. International Journal of Research Publication and Reviews, 4(6), 45595–45602.
- [18] Devashish G.T. (2025). *Securing Cloud Infrastructure Through Ancestry Tracking in Machine Images*. International Journal on Science and Technology, 16(3), 1-19.
- [19] Reagan, R. (2017). *Web applications on Azure*. Berkeley, CA: Apress.
- [20] Nuriev, M., & Lapteva, M. (2024). *Facilitating efficient energy distribution and storage: The role of data consistency technologies in Azure Cosmos DB*. E3S Web of Conferences, 541, 02003.
- [21] Paz, J. R. G. (2018). *Learning Azure Cosmos DB concepts*. In Microsoft Azure Cosmos DB Revealed (pp. 25–59). Apress.