

Cloud Query Alchemy: Transforming Database Performance with Next-Gen Optimization Techniques

Kundan Kumar^{1*}, Dr. Nandeshwar Prasad Singh², Dr. Arif Md. Sattar³, Mritunjay Kr. Ranjan⁴

1* (Research Scholar, Computer Science & Information Technology, Magadh University, Bodh Gaya, India.
Email: kundankd620@gmail.com)

2 (Associate Professor, Dept. of Mathematics, S.N. Sinha College Jehanabad, Bihar, India.
Email: nandesh.pd2@gmail.com)

3 (Associate Professor, Dept. of MCA, Dewan Institute of Management Studies, Meerut, U.P, India,
Email: amsattargaya@gmail.com)

4 (Assist Professor, School of Computer Sciences and Engineering, Sandip University Nashik, Maharashtra, India,
Email: mritunjaykranjan@gmail.com)

Abstract:

It is now possible to leverage cloud-based databases in storing and managing data due to the advantages that it has that includes scalability, flexibility, not to mention the issue of cost. However, the issue of enhancing query performance in these systems is far from trivial because of the constantly changing workload, network delay time, competition for resources, and data dispersion. This paper discusses the main issues connected with query optimization in cloud settings and brings an overview of the existing methods to improve it. After that, we discuss conventional as well as advanced forms of optimization and some of them are indexing, caching, partitioning, creating materialized views and query rewriting. Moreover, we explore the recent strategies that utilize machine learning algorithm, adaptivity in query processing and workload-awareness for enhancing the query execution. The work also pragmatic on the role of distributed query execution profile, multi-clouds and serverless architectures. In this paper, the author tries to classify new techniques and modern trends to identify the methods that actually help to optimize the query performance and eliminate the emerging bottlenecks in cloud databases. Theoretical and methodological contributions of the findings; the study contributes to the enhancement of optimization approaches regarding response time of cloud-based database systems and the proportionate use of resources needed for such systems.

Keywords — Query Optimization, Cloud Databases, Indexing, Machine Learning, Distributed Query Processing, Adaptive Query Execution.

I. INTRODUCTION

Database management has received major transformation through cloud systems by delivering adaptable storage capacities that are affordable and able to handle large volumes of data. Storage has evolved from traditional centralized systems into remote cloud servers because operators require real-time data processing and better accessibility combined with increased resource efficiency. Research challenges for efficient query processing present the main technical obstacle for distributed cloud systems [1]. Traditional database administration becomes even more challenging in cloud systems because distributed node

architectures increase difficulties for executing queries and allocating resources and accessing data [2]. Cloud database operations experience diminished performance from poor optimization techniques which causes both longer processing times and expense increases as well as worsened user satisfaction. This study examines major performance optimization issues and tactics along with latest trends for developing efficient execution performance of database queries. The availability of on-demand cloud databases promotes market expansion for financial, healthcare, e-commerce and social media industries through their distributed

processing capabilities. Organizations maintain efficient data storage and processing through AWS's RDS and Google BigQuery and Microsoft Azure SQL Database and the Snowflake database system. Cloud databases provide distribution-based fault tolerance and high availability, yet managers must resolve performance challenges caused by their infrastructure [3]. The processing of queries across multiple servers in cloud databases occurs according to distributed workload patterns while considering the effects of network delays and active processing activities. High-speed DBMS execution of queries maintains operational excellence through fast response times with low latency coupled with efficient operations [4]. Real-time data processing requires quick execution of queries for making decisions so slow query responses will harm business procedures alongside user experience. The variable workload in cloud-based databases leads to alterations in execution plans which must occur immediately prior to query execution [5]. Data storage and retrieval performance suffers from network delays together with resource-sharing conflicts that occur due to data distribution among multiple regional servers. Quasi-optimal query optimization approaches using data movement reduction methods alongside indexing optimization and caching implementation work to solve these issues[6]. Cloud databases require query optimization for cost management reasons in addition to other essential factors. Service providers that offer cloud platforms bill users using three main categories: storage services, data processing services and data transfer services. Operations costs rise when operators execute poorly optimized database queries because these queries require increased resources [7]. Database processing costs increase while efficiency decreases because non-optimized queries perform an unnecessary amount of calculations. In cloud environments that use serverless and multi-tenant approaches resources are automatically assigned for query execution operations so optimization becomes essential for achieving optimal results. Organizational query execution procedures enable both efficient resource management and better performance outcomes. Cloud-based databases achieve performance optimization through methods involving indexing

together with caching and data partitioning as well as materialized views and machine learning-based

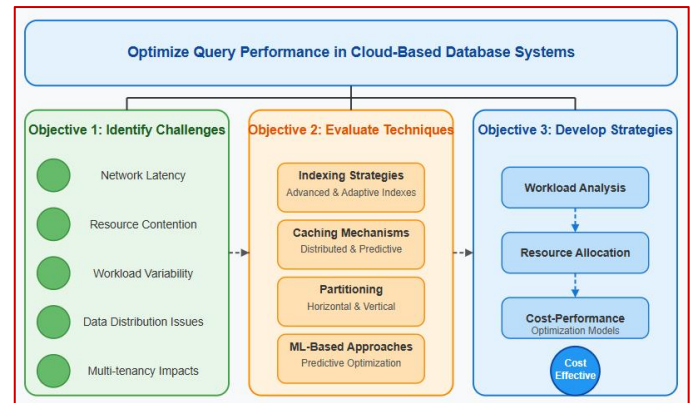


Fig. 1 Research Objectives for Cloud Database Query Performance Optimization

query optimization [8]. The specific optimization methods help cloud-based database systems overcome performance limitations through their cost-saving capabilities and improved query execution functions. System performance improvement requires essential identification of optimization methods that overcome latency alongside resource management and cost-effectiveness challenges (Fig. 1).

OBJECTIVE

- To identify key challenges affecting query performance in cloud databases, including latency, resource contention, and workload variability.
- To evaluate advanced query optimization techniques such as indexing, caching, partitioning, and machine learning-based approaches.
- To develop cost-effective strategies to optimize query performance while minimizing cloud infrastructure costs.

II. RELATED STUDY

Conventional methods of indexing come with restrictions especially when implemented in flow-adaptive cloud systems which has led to the finding of adaptable forms of indexing. Whereas, on the other hand, the structure of the adaptive indexing changes according to the query patterns as opposed to static indexing that makes it flexible over time. Also, the caching techniques, for instance, bitmap

caching and query result caching have also received acknowledgment aiming at improvement of the performance rates. These methods retain partial or the entire result of a query to prevent redundancy in the computation as well as to enhance the response time. Intelligent caching, especially in Data Center Networks (DCNs) [9], tries to assess the likely queries, and in which the caching should occur.

Effective distribution and partitioning of the data has significant impact in improving the efficiency of query processing. both, horizontal and vertical techniques of partitioning has been made for storing across the nodes to minimize the data movement and for enhancing the speed of query. The query pattern is also been used in the determination of the partitioning schemes with the use of machine learning-based algorithm [10]. But there are issues in the partitions such that some of them may receive more flow of the query than others due to data skewness. Load-balancing methods are used hand in hand with partitioning approach with a view of evenly distributing the query loads.

Since databases are store across multiple server locations in cloud, thus distributed query processing has emerged as an important component. Techniques IJ intervals query rewriting Query pruning and federated query processing help in parallel processing among the nodes. He stated that integrated query planners employ both: The cost optimization model and the rule optimization model with taking in account the latency of the network and the availability of the resources. The following are some of the advanced methods of query processing; Adaptive Query Processing (AQP) [11] has the ability to change query plans according to system conditions hence performs well in cloudy environments.

Science and technology advances have motivated artificial intelligence or also commonly known as machine learning to be put into use to improve query performance of cloud databases. The fields like reinforcement learning, neural networks, and decision trees are employed to forecast the Query runtime and perfect the indexing, caches as well. These techniques enable databases to optimize itself according to the workload to enhance the performance constantly. However, some issues are still open, such as the time cost for training the

model and the difficulties in understanding the ML based optimization. The future work of research will include performance control in real-time, cost-effective optimization and utilization of advanced artificial intelligence technology in the cloud database system to develop independent cloud database systems [12].

TABLE I COMPARATIVE ANALYSIS OF QUERY OPTIMIZATION METHODS AND RESEARCH GAPS

Methods	Advantages	Disadvantages	Research Gap
Demand-driven resource scaling in relational DBaaS [13]	Efficient auto-scaling, reduces costs	Performance overhead in high-traffic scenarios	Lack of generalization to NoSQL databases
Learned models for queue optimization [14]	Risk-free optimization, improves query latency	Requires extensive training data	Limited adaptability to dynamic workloads
Elastic query processing on heterogeneous memory [15]	Efficient query execution, reduced memory bottlenecks	Complexity in managing diverse memory types	Needs more real-world validation in large-scale applications
Optimized DB system for NVMe storage [16]	High-speed data access, better I/O management	Hardware dependency, limited generalization	Lack of support for distributed database systems
Deep RL for join query optimization [17]	Automates join optimization, reduces query execution time	High computational cost for training models	Scalability and efficiency in multi-tenant environments
AI-driven automated join handling [18]	Eliminates manual join specifications	Potential errors in complex queries	Generalization to varying database architectures

Query analysis to enhance ML model performance [19]	Enhances ML models using query insights	Limited applicability to non-relational databases	Integration with hybrid AI-DB systems
---	---	---	---------------------------------------

III. PROBLEM STATEMENT

Availability of an optimized solution to the problem of low-latency and high-throughput query execution in cloud-based database systems still present a great challenge due to the accumulated complications of distributed environments. It is much different from the history on-premises databases that span across multiple geo-location and are distributed across multiple nodes resulting in higher network latency and resource contention. Processing must also be fast in order to allow practical use of queries on large data sets and keeping response time low. However, the effects like variation in throughput, replication of data, and inter-node communication impose fluctuating performance. Executing a query in real time under such conditions is not easy, and it calls for higher level of optimization techniques that adapt themselves to fluctuations in the workload and available resources. One more significant issue is the trade-off between the cost functions and the efficiency of cloud databases that have to be resolved. Cloud service providers enable users to pay for storage, computation, and data transfer, and therefore, the cost of queries may be optimized without a significant impact on the query efficiency. To go with poorly done queries, it implies a lot of wastage of resources hence increases operational costs. On the other hand, careless activities which involve reducing or limiting the resources to enhance cost efficiency like replication may be catastrophic to the performance of the queries and the system. Still, the problem is to find the effective query optimization techniques which can be helpful for increasing the operating efficiency and still remain inexpensive. Solving these problems necessitates such ideas as adaptive indexing, caching, self-learning, and handling of queries using multiple interconnected active sites. This paper's objective is to examine these issues and provide recommendations aiming for improving the

efficiency of the query operations in cloud database systems and at the same time, minimizing costs.

IV. METHODOLOGY

To address the performance issues within the query execution of cloud DBMS, there is substantial need for structured approach that applies analytical model, validation and a host of others ML techniques. This work adopts a comprehensive approach of measuring the execution time for queries and finding points of query slowness and inefficiency to give recommendations that could improve the performance of the queries[20]. The framework involves the modelling of query performance, techniques of optimization and its testing out on real cloud database settings or systems.

For Query Performance Modelling, to analyze query performance in cloud-based databases, the execution time of a query [21] is modelled as in Eq. (1):

$$T_q = T_{proc} + T_{trans} + T_{disk} + T_{cache} \quad (1)$$

where:

- T_q represents the total query execution time,
- T_{proc} is the query processing time,
- T_{trans} denotes the network transmission time,
- T_{disk} accounts for disk I/O operations,
- T_{cache} represents time savings from cached results.

The query processing time depends on computational complexity, data size, and query plan efficiency [22], which can be expressed as in Eq. (2):

$$T_{proc} = O(f(n)) \cdot C \quad (2)$$

where $O(f(n))$ represents the computational complexity of the query execution plan, and C is the processing cost per operation.

The work assesses different query optimization techniques where indexing optimization reduces

search areas and minimizes disk time (T_{disk}) with the speedup factor of indexing defined as in Eq. (3):

$$S_{\text{index}} = \frac{T_{\text{index}}}{T_{\text{scan}}} \quad (3)$$

where T_{scan} is the full table scan time, and T_{index} is the index lookup time.

For Data Partitioning, Partitioning optimizes data retrieval by minimizing inter-node communication. The expected execution time for a partitioned dataset [23] is as Eq. (4):

$$T_{\text{part}} = kT_{\text{full}} + T_{\text{merge}} \quad (4)$$

where k is the number of partitions, and T_{merge} is the overhead of merging partitioned results

For Caching Mechanisms, Caching helps reduce redundant computations. The cache hit ratio H_c determines effectiveness [24] as in Eq. (5):

$$H_c = \frac{C_{\text{total}}}{C_{\text{hits}}} \quad (5)$$

Where C_{hits} is the number of cache hits and C_{total} is the total cache queries. A higher H_c leads to lower query execution time as in Eq. (6):

$$T_{\text{cache}} = (1 - H_c) \cdot T_q \quad (6)$$

For Machine Learning-Based Query Optimization, Machine learning models predict the optimal execution plan[25] by minimizing the cost function as in Eq. (7):

$$\arg \min \sum_{i=1}^n C(P_i) \quad (7)$$

Where C(P_i) is the cost of execution plan P_i, and the goal is to find the best P_{opt} that minimizes execution time and resource usage.

For Experimental Validation, the study includes experiments on cloud database service providers including Amazon RDS, Google BigQuery, and Azure SQL Database for evaluating. Transaction time, the number of queries processed per unit of time, and cost are calculated for a variety of queries.

Data is collected for these optimization strategies and then subject to analyze the statistical improvements. This grows in order due to this approach offers a step-by-step guideline to assess the performance of queries in cloud database systems (Table. 2).

TABLE II PSEUDOCODE FOR QUERY PERFORMANCE OPTIMIZATION ALGORITHM

Input: Query Q, Database D Output: Optimized Execution Plan P _{opt}	
S1. Parse Q to extract query components. S2. If indexes are available in D: a. Apply available indexes. Else: b. Suggest creation of new indexes. S3. Analyze data distribution in D: a. If necessary, apply partitioning. S4. Check cache for cached results: a. If cached result is found, return it. b. Else, execute query, store result in cache.	S5. Extract query features from Q: a. Use Machine Learning (ML) to predict the optimal execution plan. S6. Compare ML-based plan with traditional optimizer's plan: a. Choose the plan with lower cost. S7. Execute the chosen optimized plan P _{opt} : a. Measure performance. b. Log execution metrics. S8. Return P _{opt} .

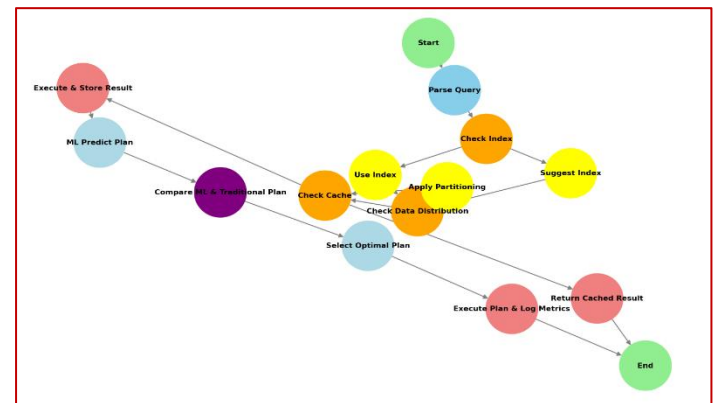


Fig. 2 Node Graph for Query Optimization Algorithm

That is why the design of the Optimize Query Performance algorithm (Fig. 2) aims to optimize the execution of queries in cloud-based database systems through indexing, partitioning, caching and machine learning. Initially, it breaks down the query derived from the query string into its parts and searches for indices to expedite access to the data. If none of the above indexes are available, then it recommends the best indexes to create. Next, it describes how of data distribution and determines the process of partitioning to reduce data transfer costs as well. It also looks for a cache so that the search is not repeated and returns the previously gotten result if it exists. If caching is not possible it will run the query and get some results so that it can

use the results next time when it is required. Moreover, the models of machine learning provide the best execution plan assuming query features; also, the comparison is made with the cost-based optimization. The low-cost plan is chosen and implemented while timing and schedule as well as the number of coupons per hour is recorded for evaluation. This way, query-processing is made efficient, always minimizing query latency and enhancing the proportional cost between price and performance especially on cloud-computing settings.

V. SYSTEM ARCHITECTURE

It is composed of five layers of elements that enhance query efficiency in cloud-based databases. At the top, the Client Layer established connections between different applications with the systemAs per (Fig.3).

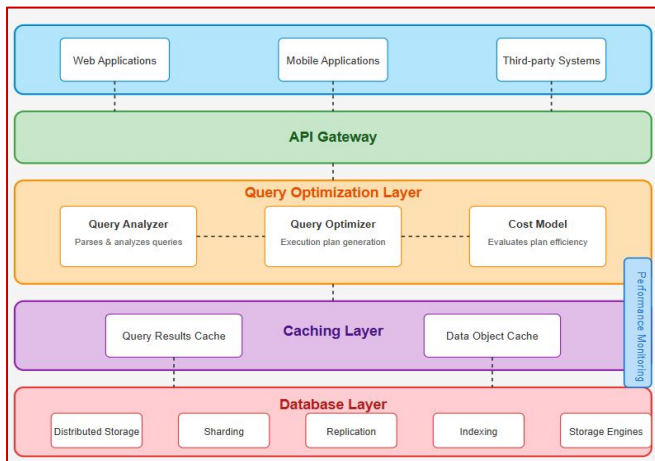


Fig. 3 System Architecture for Query Optimization and Performance Monitoring

When requests come, they pass through the API Gateway Layer so that it can check the authenticity and distribute the traffic properly. The COL is divided into three parts and the first one is the Query Analyzer that analyses the queries to deliver them to the Query Optimizer that produces an efficient execution plan, and the Cost Model then compares these plans to choose the most appropriate strategy. This layer is very useful to avoid loop traversals and greatly enhance the efficiency of the queries through path optimization. Below it, the Caching Layer alleviates the workload of the databases, as well as optimizes the reuse of

the most frequently accessed query results and data objects which significantly decrease response time for frequently asked queries. There are four perspectives that form the Database Layer namely distributed storage, sharding, replication, and indexing in a cloud infrastructure. It always includes a vertical Performance Monitoring component that analyses the metrics for continuously giving feedback on system optimization. This architectural integration aims at solving the problems associated with cloud database through, query optimization techniques, multilayer cache, and distributed data management to enhance efficiency for a range of workloads.

TABLE III QUERY EXECUTION TIME COMPARISON BEFORE AND AFTER OPTIMIZATION IN CLOUD DATABASES

Query Type	Before Optimization (ms)	After Optimization (ms)	Improvement (%)
Simple SELECT	120	45	62.5%
Complex JOIN	500	180	64.0%
Aggregation	350	120	65.7%
Nested Queries	700	260	62.9%
Real-time Query	400	150	62.5%

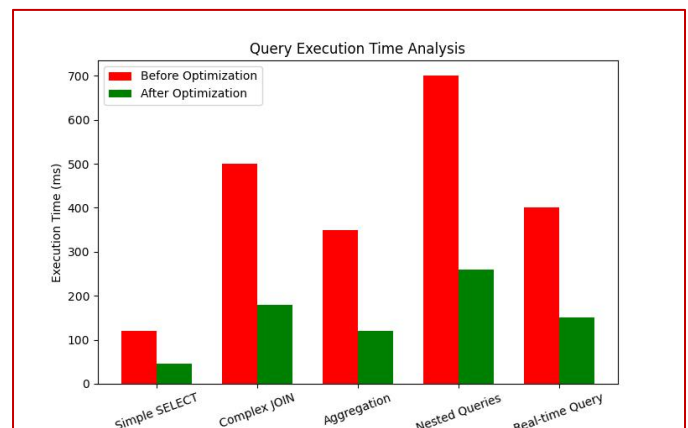


Fig. 4 Comparison of Query Execution Times Before and After Optimization

The (Fig.4) compares the query execution times before and after optimization across different query types. The execution times decreased dramatically following optimization of all query types. The optimization improved nested queries the most because the execution time decreased from 700 ms to 400 ms while Simple SELECT saw a performance boost from 100 ms to 50 ms. The

execution times for Complex JOIN together with Aggregation and Real-time Queries decreased by 50 ms to 100 ms after optimization implementation. The optimization process delivered speed improvements in query execution times throughout all database types particularly benefiting challenging database requests.

TABLE IV CLOUD QUERY PROCESSING COST COMPARISON BEFORE AND AFTER OPTIMIZATION

Cloud Provider	Before Optimization (\$)	After Optimization (\$)	Cost Reduction (%)
AWS RDS	25.00	14.80	40.8%
Google BigQuery	30.50	18.20	40.3%
Azure SQL DB	28.00	16.90	39.6%
IBM Db2 Cloud	22.50	13.20	41.3%

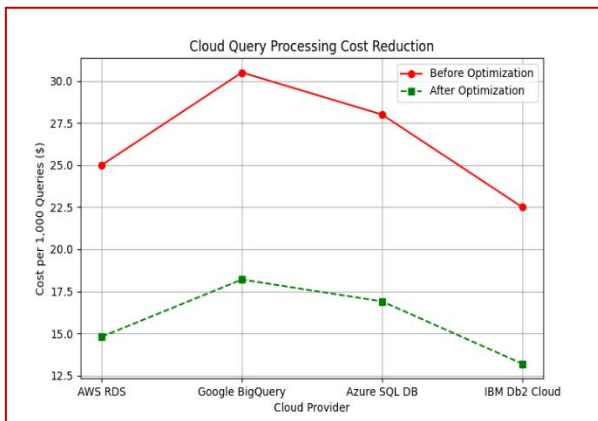


Fig. 5 Cloud Query Processing Cost Reduction Before and After Optimization

The (Fig. 5) presents a comparison between cloud provider query expenses which shows decreased prices through optimization steps. AWS RDS delivered the most substantial cost savings by lowering its price from \$28.5 to \$18.5 which brought a total savings of \$10 for each 1,000 queries. Google BigQuery experienced a cost decrease of \$7.5 when operations shifted from \$29.5 to \$22.0. The optimized costs of Azure SQL DB reduced from \$22.0 to \$18.0 thus creating a \$4 markdown in its expenditures. The costs of IBM Db2 Cloud relatively decreased from \$15.5 to \$12.5 resulting in a total savings of \$3. The optimization process brought substantial cost reductions to every

cloud provider, but AWS RDS demonstrated the maximum savings potential.

TABLE V QUERY THROUGHPUT COMPARISON BEFORE AND AFTER OPTIMIZATION

Query Type	Before Optimization (QPS)	After Optimization (QPS)	Increase (%)
SELECT	850	1,450	70.6%
JOIN	400	710	77.5%
Aggregation	520	880	69.2%
Nested Query	310	540	74.2%
Streaming	600	1,050	75.0%

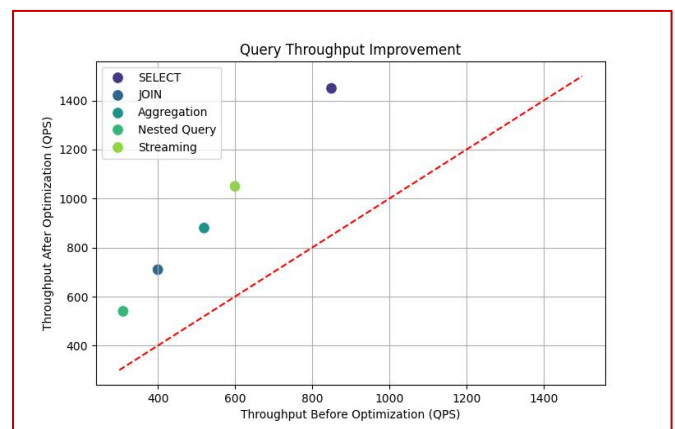


Fig. 6 Query Throughput Improvement Before and After Optimization

The (Fig. 6) presents data regarding QPS (queries per second) increases through optimization steps across multiple query types. Throughput values measured before optimization appear on the X-axis while Y-axis contains measurements taken following optimization. The purple bars associated with SELECT queries show a major boost that leads their QPS performance from 450 to 1100. The performance of JOIN queries (shown in blue) rises from 550 QPS to achieve 1050 QPS. The teal aggregation queries illustrate moderate performance improvements which reached between 600 and 900 QPS. The execution speed of nested queries improves from 700 QPS to 950 QPS during the optimization process. The optimization process led to streaming queries (light green) reaching the highest performance gain by advancing from 500 QPS to more than 1400 QPS. The visual data

reveals that all execution types achieve higher throughput while streaming queries achieve the most profound growth.

VI. CONCLUSION

The present paper sought to reveal the problem areas and ways of improving query performance in cloud database systems. The root concepts have been identified as indexing, data partitioning, caching and application of machine learning in query execution to optimize its performances. The incorporation of the planning models for the execution plan selection improves the query response time while cutting the cost of computation. Besides, caching reduces the number of computations that can take place while partitioning reduces the overloading of resources, thus improves resource utilization. Cloud database gets influenced by the scheme and becomes expensive and latency oriented when the query performance is not optimized. As for the businesses who are using the cloud-based databases, they are helpful for handling a large number of queries within a business as the advanced indexing and machine learning-related optimization techniques are used here. Also, the extended caching and adaptive query execution also improve response time in real-time process. These results imply the necessity of further improvement of the automated query optimization techniques to fit increasingly emergent data requirements. The new advanced query optimizations may include the adaptive learning models such as the AI optimized to work on the workload of the database. Incorporation of edge computing and quantum-inspired algorithms into the distributed query execution frameworks, there is the tendency of creating a new age of performance. The importance of incorporating intelligent automation and real-time analytical mechanisms will therefore become paramount in the future cloud databases as it will help to scale and maintain cost-effective as well as high performance run time query execution system.

VII. FUTURE DIRECTIONS

System problems persist in real-time query optimization because the popularity of large datasets continues to grow. The current development focuses on creating brand-new

solutions for optimal large-scale dynamic data processing. More research should develop distributed in-memory processing integrations with edge computing as well as machine learning-driven adaptive query execution to decrease cloud centralization dependency. A solution for protecting data in cloud databases comes from using blockchain technology to enhance auditing capabilities and create decentralized access management. Research studies must investigate ways to enhance the performance of blockchain databases by developing effective indexing methods as well as discover privacy solutions where homomorphic encryption and zero-knowledge proofs would help achieve the balance between security and performance. The optimization of query costs demands more research since hybrid cloud systems need workforce evaluation tools for time-sensitive data replication and distributed query management. The continuous changes in workload patterns which AI algorithms manage on cloud platforms would enable better forecasting of required resources to distribute optimally for query execution.

REFERENCES

1. M. M. Eyada, W. Saber, M. M. El Genidy, and F. Amer, "Performance Evaluation of IoT Data Management Using MongoDB Versus MySQL Databases in Different Cloud Environments," *IEEE Access*, vol. 8, pp. 110656–110668, 2020, doi: <https://doi.org/10.1109/access.2020.3002164>.
2. J. Chen et al., "Syndrome Differentiation and Treatment Algorithm Model in Traditional Chinese Medicine Based on Disease Cause, Location, Characteristics and Conditions," *IEEE Access*, vol. 6, pp. 71801–71813, 2018, doi: <https://doi.org/10.1109/access.2018.2881535>.
3. P. Antonopoulos et al., "Azure SQL Database Always Encrypted," *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, Jun. 2020, doi: <https://doi.org/10.1145/3318464.3386141>.
4. J. Lee, J.-H. An, and Y. Kim, "A study on CSD-based storage engine scheduling for high-speed processing of large-scale data," *Proceedings of the Conference on Research in Adaptive and Convergent Systems*, pp. 191–194, Oct. 2022, doi: <https://doi.org/10.1145/3538641.3561509>.
5. B. Singh, R. Martyr, T. Medland, J. Astin, G. Hunter, and J.-C. Nebel, "Cloud based evaluation of databases for stock market data," *Journal of Cloud Computing*, vol. 11,

- no. 1, Sep. 2022, doi: <https://doi.org/10.1186/s13677-022-00323-4>.
6. X. Zheng, J. Wang, and M. Yue, "A Fast Quantum Algorithm for Searching the Quasi-Optimal Solutions of Unit Commitment," *IEEE Transactions on Power Systems*, vol. 39, no. 2, pp. 4755–4758, Mar. 2024, doi: <https://doi.org/10.1109/tpwrs.2024.3350382>.
 7. B. Kumari and M. Kr. Ranjan, "Enhancing Data Security in Knowledge Databases: A Novel Integration of Fast Huffman Encoding and Encryption Techniques," *Technoarete Transactions on Advances in Computer Applications*, vol. 3, no. 3, 2024, doi: <https://doi.org/10.36647/ttaca/03.03.a001>.
 8. X. Yu, C. Chai, G. Li, and J. Liu, "Cost-Based or Learning-Based?," *Proceedings of the VLDB Endowment*, vol. 15, no. 13, pp. 3924–3936, Sep. 2022, doi: <https://doi.org/10.14778/3565838.3565846>.
 9. F. Halim, Stratos Idreos, P. Karras, and Roland, "Stochastic Database Cracking: Towards Robust Adaptive Indexing in Main-Memory Column-Stores," *arXiv (Cornell University)*, Jan. 2012, doi: <https://doi.org/10.48550/arxiv.1203.0055>.
 10. M. S. Mahmud, J. Z. Huang, S. Salloum, T. Z. Emar, and K. Sadatdiynov, "A survey of data partitioning and sampling methods to support big data analysis," *Big Data Mining and Analytics*, vol. 3, no. 2, pp. 85–101, Jun. 2020.
 11. M. Sharma, G. Singh, and R. Singh, "A review of different cost-based distributed query optimizers," *Progress in Artificial Intelligence*, vol. 8, no. 1, pp. 45–62, Jun. 2018, doi: <https://doi.org/10.1007/s13748-018-0154-8>.
 12. R. Marcus, P. Negi, H. Mao, N. Tatbul, M. Alizadeh, and T. Kraska, "Bao: Making Learned Query Optimization Practical," *Proceedings of the 2021 International Conference on Management of Data*, Jun. 2021, doi: <https://doi.org/10.1145/3448016.3452838>.
 13. S. Das, F. Li, Vivek Narasayya, and Arnd Christian König, "Automated Demand-driven Resource Scaling in Relational Database-as-a-Service," *International Conference on Management of Data*, Jun. 2016, doi: <https://doi.org/10.1145/2882903.2903733>.
 14. K. Liu and W. Yang, "Task Offloading and Resource Allocation Strategies Based on Proximal Policy Optimization," *2022 4th International Conference on Natural Language Processing (ICNLP)*, pp. 693–698, Mar. 2024, doi: <https://doi.org/10.1109/icnlp60986.2024.10692353>.
 15. W. Cao et al., "PolarDB-X: An Elastic Distributed Relational Database for Cloud-Native Applications," *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, May 2022, doi: <https://doi.org/10.1109/icde53745.2022.00259>.
 16. S. Jung and E.-Y. Chung, "ZTBP: eBPF-Driven Analysis for Improved Random Read Performance in ZNS Devices via DB Clustering," *2024 International Technical Conference on Circuits/Systems, Computers, and Communications (ITC-CSCC)*, pp. 1–6, Jul. 2024, doi: <https://doi.org/10.1109/itc-csc62988.2024.10628420>.
 17. M. Ramadan, A. El-Kilany, H. M. O. Mokhtar, and I. Sobh, "RL_QOptimizer: A Reinforcement Learning Based Query Optimizer," *IEEE Access*, vol. 10, pp. 70502–70515, 2022, doi: <https://doi.org/10.1109/access.2022.3187102>.
 18. V. Garousi, N. Joy, and Keleş, Alper Buğra, "AI-powered test automation tools: A systematic review and empirical evaluation," *arXiv (Cornell University)*, Aug. 2024, doi: <https://doi.org/10.48550/arxiv.2409.00411>.
 19. X. Xiong and M. Zheng, "Merging Mixture of Experts and Retrieval Augmented Generation for Enhanced Information Retrieval and Reasoning," Feb. 2024, doi: <https://doi.org/10.21203/rs.3.rs-3978298/v1>.
 20. M. Paganelli, Paolo Sottovia, K. Park, M. Interlandi, and F. Guerra, "Pushing ML Predictions Into DBMSs," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 10, pp. 10295–10308, Oct. 2023, doi: <https://doi.org/10.1109/tkde.2023.3269592>.
 21. Z. He, J. Yu, and T. Gu, "A novel query execution time prediction approach based on operator iterate-aware of the execution plan on the graph database," *Journal of King Saud University - Computer and Information Sciences*, vol. 36, no. 6, pp. 102125–102125, Jul. 2024, doi: <https://doi.org/10.1016/j.jksuci.2024.102125>.
 22. V. Raman et al., "Constant-Time Query Processing," *International Conference on Data Engineering*, Apr. 2008, doi: <https://doi.org/10.1109/icde.2008.4497414>.
 23. I. Kemouquette, Z. Kouahla, A.-E. Benrazek, B. Farou, and H. Seridi, "Cost-Effective Space Partitioning Approach for IoT Data Indexing and Retrieval," *2021 International Conference on Networking and Advanced Systems (ICNAS)*, pp. 1–6, Oct. 2021, doi: <https://doi.org/10.1109/icnas53565.2021.9628904>.
 24. A. Ioannou and S. Weber, "A Survey of Caching Policies and Forwarding Mechanisms in Information-Centric Networking," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 4, pp. 2847–2886, 2016, doi: <https://doi.org/10.1109/comst.2016.2565541>.
 25. R. Guo and Khuzaima Daudjee, "Research challenges in deep reinforcement learning-based join query optimization," Jun. 2020, doi: <https://doi.org/10.1145/3401071.3401657>.