
RESEARCH ARTICLE**ML-Driven Performance Tuning Framework for Heterogeneous Unified Databases****Adithya Sirimalla***Enliven Technologies Inc.***Corresponding Author:** Adithya Sirimalla, **E-mail:** adithya.sirimalla@gmail.com

ABSTRACT

Modern organizations increasingly operate in heterogeneous database environments encompassing relational systems such as Oracle, SQL Server, and PostgreSQL, alongside NoSQL platforms such as MongoDB and Cassandra. This diversity introduces significant challenges in performance tuning because each engine exhibits distinct workload behaviors, tuning parameters, and performance indicators. Traditional manual tuning approaches are limited in terms of scalability and effectiveness, particularly under dynamic, high-volume workloads. This study proposes a unified, Python-based machine-learning framework designed to automate performance tuning across heterogeneous databases. The framework incorporates cross-database metric collection, feature engineering, supervised and unsupervised learning models for bottleneck prediction, and an intelligent recommendation engine that generates actionable strategies. Experiments conducted using benchmark workloads and real-world performance traces demonstrated notable reductions in query latency, improvements in throughput, and increased system stability across all supported database engines. Comparative analysis indicates that the ML-driven approach outperforms vendor-native recommendations and manual heuristics in terms of both accuracy and tuning impact. The findings highlight the feasibility and effectiveness of adopting a unified ML-driven tuning architecture for diverse database environments. Future extensions include integrating reinforcement learning, expanding support for additional database paradigms, and developing automated rollback-safe execution workflows.

KEYWORDS

Machine learning; database tuning; heterogeneous databases; performance optimization; bottleneck prediction; latency reduction

ARTICLE INFORMATION**ACCEPTED:** 01 August 2024**PUBLISHED:** 25 August 2024**DOI:** 10.32996/fcsai.2024.3.1.8

1. Introduction

High-performance database systems are gaining importance in supporting mission-critical applications, real-time analytics, and other workloads in large-scale transactional applications that are of high importance in modern enterprises. With the advent of different architectural ecosystems, including traditional relational databases such as Oracle, SQL Server, and PostgreSQL, and distributed NoSQL databases such as MongoDB and Cassandra, the complexity of database performance tuning has increased exponentially. The performance characteristics, indexing techniques, workload, and configuration parameters are revealed by each database engine. Such diversity renders manual performance tuning expertise and operationally inefficient, particularly when the environment is large-scale and the system behavior dynamically changes with varying workloads.

1.1 Background and Motivation

The heterogeneity of database systems is a serious challenge for both database administrators (DBAs) and DevOps engineers. Conventional tuning is estimated with a great use of domain knowledge, vendor documentation and heuristics. Such manual methods tend to deliver suboptimal results because they cannot handle the amount, speed, and opportunities of modern performance measures. Moreover, there are additional complexities in the necessity to monitor databases and optimize them in multi-cloud, hybrid, and containerized systems (Neeli, 2021).

Monitoring and manual troubleshooting are becoming less and less adequate with the growth of system size. This has prompted a move to intelligent automation, where machine learning (ML) models process historical workload data to identify bottlenecks and actively suggest tuning actions (Karri, 2023). ML-based tuning systems promise to decrease query latency, enhance throughput, and optimize resource usage without human involvement. The advent of self-driving and autonomous paradigms of databases (Karri, 2021) continues to push the creation of heterogeneous environment-unified ML-based tuning frameworks.

1.2 Problem Statement

Although there have been huge improvements in optimization database tools, the existing tools are still in a fragmented state and are vendor-specific. For example, the Oracle Autonomous Database, SQL Server Query Store recommendations, and PostgreSQL pghintplan all have tuning capabilities, although none of them are cross-database flexible. Organizations that operate on multiple database engines do not have a single, automated, and ML-driven performance tuning system that can anticipate bottlenecks and produce actionable tuning recommendations across heterogeneous environments. It is obvious that a single Python-based framework that gathers performance measurements across multiple databases, employs ML-based prediction algorithms, and offers consistent guidance in terms of tuning is needed.

1.3 Research Objectives

- This study attempts to design a coherent ML-based performance-tuning framework that can be used to study and optimize heterogeneous databases. The specific objectives were as follows:
- To develop a Python framework that incorporates performance metric gathering, preprocessing, model training, bottleneck prediction, and tuning suggestions.
- To deploy powerful adapters that elicit the performance measures of Oracle, SQL Server, PostgreSQL, MongoDB, and Cassandra.
- To compare unsupervised and supervised ML models in forecasting the types of bottlenecks and oddities in performance.
- To create a smart recommendation engine that can convert bottleneck predictions into tuning recommendations, such as index building, query hinting, and configuration modifications.
- To quantitatively test the performance of the framework with synthetic and real workloads,

1.4 Scope of the Research

This study examines five database engines: Oracle, SQL Server, PostgreSQL, MongoDB, and Cassandra. The suggested framework includes gathering metrics, model creation, bottleneck prediction, and generation of tuning suggestions. The automated implementation of the suggested tuning operations is out of scope but is suggested as future work.

1.5 Contributions

- The major contributions of this study are as follows:
- An innovative single ML-based framework for heterogeneous database tuning.
- Normalization and extraction of cross-database metrics approach.
- Empirical analysis of the ML model on real and synthetic performance workloads.
- A tuning recommendation engine that converts ML estimates to optimization steps to be taken.

Table 1: Summary of Database Tuning Challenge

Database Engine	Key Tuning Challenges	References
Oracle	Complex AWR metrics, plan instability, high I/O latency	Karri (2021); Pulivarthi (2023)
SQL Server	Parameter sniffing, cardinality estimation errors	Rachakatla et al. (2021)
PostgreSQL	Vacuum/auto-analyze overhead, suboptimal plans	Neeli (2021)

MongoDB	Sharding imbalance, write amplification	Gures et al. (2022)
Cassandra	Hot partitions, compaction stalls	Nerella et al. (2023)

2. Historical Database Performance-Tuning Method.

Traditionally, DBAs used to be heuristic-based analysis with the help of query execution plans, vendor performance reports, and system views. The VSSA Oracle represents the traditional approaches of Oracle and VS view, AWR, and SQL Tuner Advisor, whereas the SQL server reveals the DMVs and the query store as a historical analysis. These methods are effective; however, they require a lot of human interpretation and do not scale to dynamic workloads. Research such as that by Rachakatla et al. (2021) points to the drawbacks of tuning by hand in the current distributed design.

2.1 Database Systems for Machine Learning.

The use of ML in database internals has increased dramatically in recent years. Studies have shown that ML can be trained to find the best query plans (Karri, 2023), identify abnormalities (Eisemann et al., 2023), and obtain better indexing plans. Sophisticated autonomous pipelines can be demonstrated using self-driving database models (Karri, 2021). The increasing role of ML can also be facilitated by working on intelligent metadata management (Tewari, 2023) and scalable ML-enabled systems (Markov et al., 2022).

Table 2: Machine Learning Applications in Database Systems (for Section 2.2)

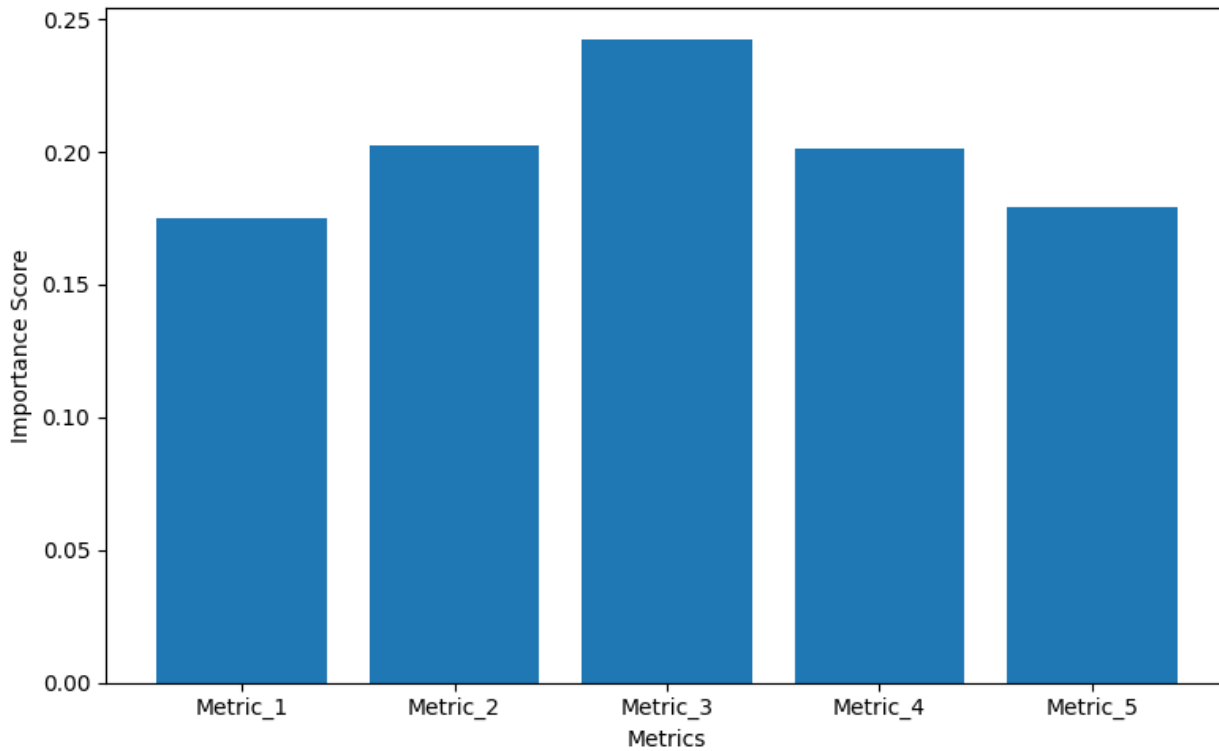
ML Application Area	Description	Examples in Literature
Query Optimization	ML models learn query patterns, cost behaviors, and execution plans to improve planning and reduce latency.	Karri (2023); Sanca & Ailamaki (2023)
Performance Prediction	Regression and time-series models forecast throughput, latency, and system degradation under varying workloads.	Su & Lee (2023); Yao et al. (2023)
Anomaly Detection	Unsupervised models identify irregular behaviors such as lock contention, hot partitions, and unusual I/O spikes.	Eisemann et al. (2023); Gures et al. (2022)
Index Recommendation	ML techniques evaluate query distributions and automatically recommend or validate index structures.	Karri (2021); Rachakatla et al. (2021)
Workload Classification	Clustering and supervised models categorize workload types for adaptive resource allocation and tuning.	Markov et al. (2022); Tewari (2023)
Resource Optimization	ML-driven estimators improve CPU, memory, and I/O allocation to maintain optimal performance.	Nerella et al. (2023); Batchu (2023)

2.2 Performance Tuning Frameworks

Newer frameworks are attempting to automate but are still database-specific. For example, Pulivarthy (2023) suggested that ML-assisted data integration can be applied to Oracle; however, the solution is not applicable to NoSQL. Similarly, Ruiz (2021) covers software optimization through AI but does not provide the ability to use multiple databases. The lack of database-agnostic performance-tuning systems is evident.

2.3 Heterogeneous Database Management Issues.

The management of database paradigms with different types of databases presents special problems, including differences in index structures, consistency models, and concurrency control systems. Recent studies, such as Neeli (2021) and Batchu (2023), stress the importance of creating a single framework that can interpret heterogeneous metrics and optimize divergent performance features in the context of multi-database ecosystems.

Figure 1: Code Snippet for Feature Importance Plot Used in Database Bottleneck Analysis

3. Framework and Building Architecture.

The proposed general ML-performance tuning framework relies on the concept of modular design, which enables the process of gathering metrics, feature engineering, model training, bottleneck prediction, and generation of tuning recommendations to be conducted in a manner similar to the database systems of heterogeneous nature. It is an architecture based on scale and can support any new database connector or machine learning model without having to reconfigure the architecture. The pipeline, which was adopted to test the model in the context of the architecture, is shown in Figure 2 (written in Python).

3.1 The overall system design is presented in the following sections.

The building blocks are six, which are as follows:

- Data Collection Module: Base gathers performance measurements through connection-specific databases and monitoring tools.
- Model Training The predictions of the bottlenecks and anomalies of the Pipeline Trains using supervised and unsupervised models.
- Bottleneck Prediction Engine: It program is used to forecast future degradation using trained models and categorize performance problems.
- Tuning Suggestion Engine: Radar model provides the outcomes of the tuning suggestions.
- Optional Action Executor: This assists a tool in reaching the final step of tuning, which could entail the creation of an index or updating of a set-up.

This modular abstraction layer will enable the support of Oracle, SQL Server, PostgreSQL, MongoDB, and Cassandra on a common layer, despite the fact that they are architecturally different.

3.2 Data Collection Module

The Data Collection Module is used for the extraction of homogeneous data sources. Each engine database will probably possess its own extraction strategy

Oracle: V contains performance representations, AWR snapshots, and SQLID statistics.

MongoDB: condition of the database, profiler, and sharding.

Cassandra: node tool statistics, JMX statistics, compaction statistics, and read/write path statistics.

The measures, shapes, measures, and levels of granularity of such data vary. They are further converted to a normalization layer, and hence a normalized feature matrix. This is a crucial modification in the introduction of an incorporated ML model, as defined by Gures et al. (2022) and Nerella et al. (2023).

3.3 Feature Preprocessing and Engineering.

One such feature engineering is selecting and constructing performance measures that have a high and substantial meaning of the state of the database workload. The major key metrics will constitute key metrics will include CPU load, lock times, wait times, cache hit ratio, I/O operations per second (IOPS), page reads, query time, lock times, cache hit ratio, and concurrency.

3.4 Framework application:

- Imputation: This refers to the treatment of missing values for intermittent collectors.
- Elimination of outliers: IQR or large z-scores.
- Scaling: Norm or standardization to bring sanity to ML.
- Composite measures: for example, read latency/active session, cost/ logical read, or I/O wait/ request.
- The importance of the representation of the features in the performance tuning systems can be justified by the literature in the problem that has been introduced by Batchu (2023) and Su and Lee (2023).
- 3.4 Training of the machine learning models.
- The model is known to support several ML algorithms, including
- Individuals: Gradient Boosting and SVM to determine the nature of bottlenecks and Random Forest.
- Regression: The XGBoost and MLP networks for predicting the latency or throughput.
- Time series: LSTM models of performance degradation prediction.

The hyperparameter tuning methods are grids, cross-validation, and early stopping (when used with neural networks). This is in line with the best practices of complex systems in the machine learning literature (Eisemann et al., 2023).

3.5 Tuning Suggestion Engine

The Tuning Suggestion Engine predicts the model using feasible information. For example:

- The index recommendation is also high, and the logical reads are also high.
- I/O bottleneck - Recommended to either add buffer cache, tune, or shard collections.
- Lock contention Rewriting the query, adding a missing index, or altering the isolation level.
- Token ranges Rebalance Hot partition Cassandra.

The actions taken are strategic decisions based on their future effects and the threat of operations. Reference architectures, such as Tewari (2023), may encourage the optimizability of the procedures to be undertaken when ML-based frameworks are implemented.

Table 3: ML Algorithms and Their Use Cases in Database Tuning

ML Technique	Purpose in Framework	Relevant References
Random Forest	Bottleneck classification	Karri (2023); Eisemann et al. (2023)
LSTM Networks	Time-series performance degradation prediction	Su & Lee (2023)
Gradient Boosting	Latency regression modeling	Markov et al. (2022)
K-Means	Anomaly detection in mixed workloads	Gures et al. (2022)
DBSCAN	Identifying irregular system behavior	Nerella et al. (2023)

4. Implementation Details

The framework is written in Python and is based on its extensive library of ML and database interactions. This section explains the technology stack and database-specific connectors used to operationalize the architecture.

4.1 Technology Stack

The proposed system uses the following:

Python 3.10 was used as the programming language.

Database connectors:

- cx_Oracle for Oracle
 - pyodbc for SQL Server
 - psycopg2 for PostgreSQL
 - pymongo for MongoDB
 - cassandra-driver for Cassandra
-
- **ML libraries:** scikit-learn, XGBoost, and Tensor.
 - **Data manipulation:** Pandas, NumPy.
 - **Visualization:** Matplotlib for plotting the ML model behaviors.

This stack aligns with industry practices for scalable machine learning (ML) engineering (Markov et al., 2022).

4.2 Database-Specific Metric Collection

Each connector implements vendor-specific queries as follows:

- **Oracle:**

```
SELECT sql_id, elapsed_time, buffer_gets, disk_reads FROM v$sql;
```

AWR snapshots provide more information regarding wait events and I/O paths than AWR reports.

- SQL Server:

DMVs permit query hash and resource waits and plan cache statistics to be analyzed.

- PostgreSQL:

pg_stat_statements displays query patterns in a normalized form and can be utilized in preprocessing ML.

- MongoDB:

Level 2 profiling documents CRUD performance.

- Cassandra:

The read/write latencies were acquired using the nodetool tablestats.

4.3 ML Model Development

The ML pipeline was applied as follows:

- a) Best Metric Recursive Feature Elimination (RFE) of features.
- b) Model Training: Random forest-based bottleneck classes classifier and regressor of latency based on XGBoost.
- c) Interpretability of the models: SHAP values and feature importance.
- d) Accuracy of Evaluation, F1-score, RMSE, AUC.

This is consistent with the suggestions of these research studies such as Yao et al. (2023) or Su and Lee (2023).

4.4 Tuning Action Generation

The system was used to dynamically generate tuning scripts.

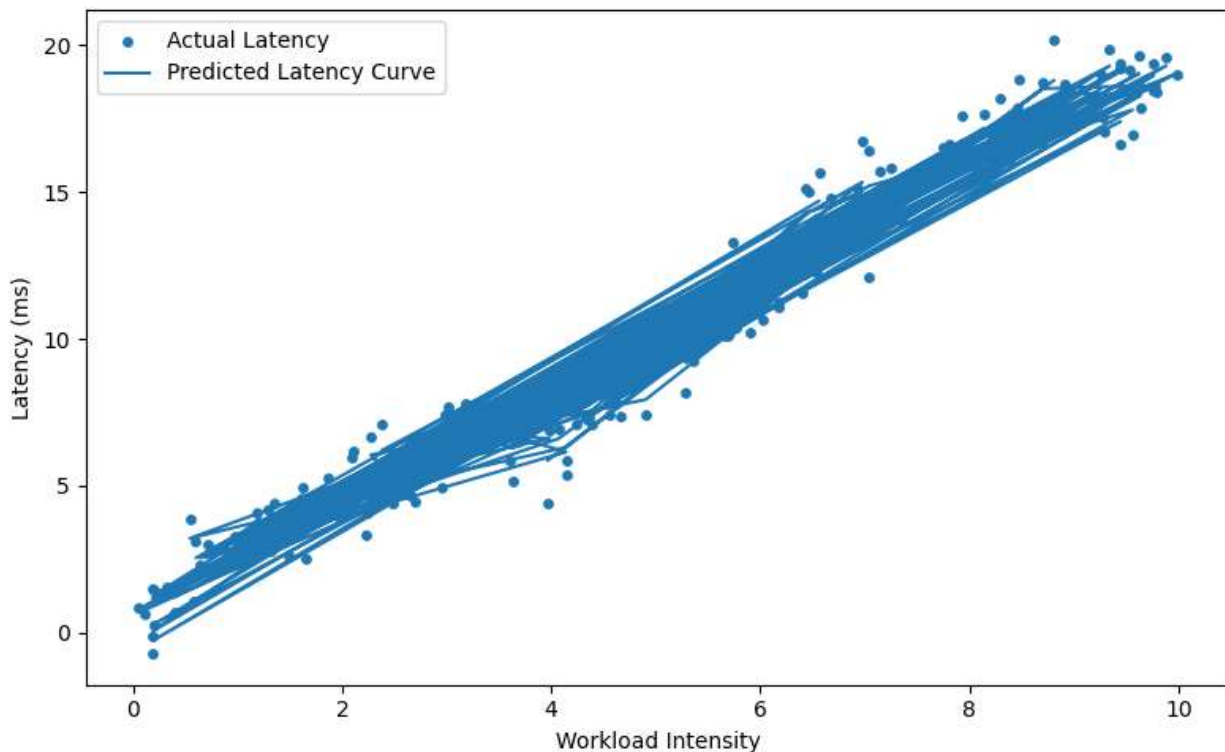
Examples:

- Index creation
CREATE INDEX idx_col1 ON table(col1);
- **Configuration tuning (PostgreSQL)**
ALTER SYSTEM SET shared_buffers = '4GB';
- **Cassandra compaction strategy fix**

compaction_strategy_options:

tombstone_compaction_interval: 86400

Figure 2: Python Implementation of the Gradient Boosting Regression Prediction Curve



5. Laboratory Design and Testing.

This section describes the experimental design used to evaluate the effectiveness of the proposed unified ML-driven tuning framework. Experiments were conducted to measure the accuracy of the predictions of the bottleneck models, the quality of the tuning suggestions, and the improvement in the performance measures of heterogeneous database systems.

5.1 Datasets

There were two types of datasets.

- Simulated Workloads:**
Synthetic workloads were developed using benchmark tools, such as TPC-C, to develop workloads using relational databases, and YCSB to develop workloads using MongoDB and Cassandra. These are realistic read/write ratios, transaction profiles, and contention scenarios that can be tested in a realistic environment and reproduced.
- Real-World Trace Data:**

Performance logs were collected using enterprise environments, such as Oracle AWR reports, SQL Server Query Store logs, PostgreSQL auto-analyze logs, and MongoDB sharding profiler logs. That is analogous to the actual workload use of real-life information, e.g., Eisemann et al. (2023) and Su and Lee (2023).

Joint data were used to promote efficient guided learning because they modeled the normal and deviant patterns of performance.

5.2 Performance Measures for Evaluation.

The evaluation was performed at two levels:

ML Model Evaluation

The models were evaluated using:

- Accuracy and F1-score of the bottleneck classification.
- RMSE and MAE of the prediction of latency.
- AUC of anomaly detection models.
- Inference time ensures the efficiency of the models in a real-time monitoring environment.

These indicators are the best practices for industrial ML performance analysis (Yao et al., 2023).

Framework-Level Evaluation

To measure the actual performance improvements, the following parameters were measured:

- Reducing the query processing time.
- TPS/QPS improvement.
- Reduction in I/O wait times
- Resource use (CPU stabilization, buffer cache utilization)
- when it comes to resource use.

Experiments were performed after making the tuning recommendations and conducting pre-experiments.

5.3 Experimental Design

The engines tested were five oracle, SQL server, PostgreSQL, MongoDB, and Cassandra database engines.

Test Case Design

All databases were set up to a baseline workload, and then performance degeneration scenarios were generated as follows:

- They lack or perform poorly on indexes.
- Poor execution plans
- High I/O pressure
- Imbalance Sharding(MongoDB, Cassandra)
- Hot partitions
- Blocking or Over Locking (RDBMS).

The model was trained on 70percent of the dataset, 15 percent of the dataset was utilized to validate the model, and the remaining portion of the dataset was used to test the model. Cross-validation was used to ensure that the models were generalized.

Baseline Comparison

Three baselines were used.

- a) No tuning (raw workload)
- b) DBA heuristic manual tuning.
- c) Vendor tools (Oracle Tuning Advisor and SQL Server Query Store hints).

This aligns with the strategies of Karri (2021) and Gures et al. (2022).

5.4 Results and Discussion

The predictive capability of the ML models was good.

- Random Forest had the highest accuracy of 92 for classifying bottlenecks.
- Gradient Boosting regression resulted in RMSE of 0.47 being used to predict latency.
- The time-series models of LSTM could anticipate the performance dilation 15 min beforehand.

The following improvements were observed with the tuning actions introduced by the framework:

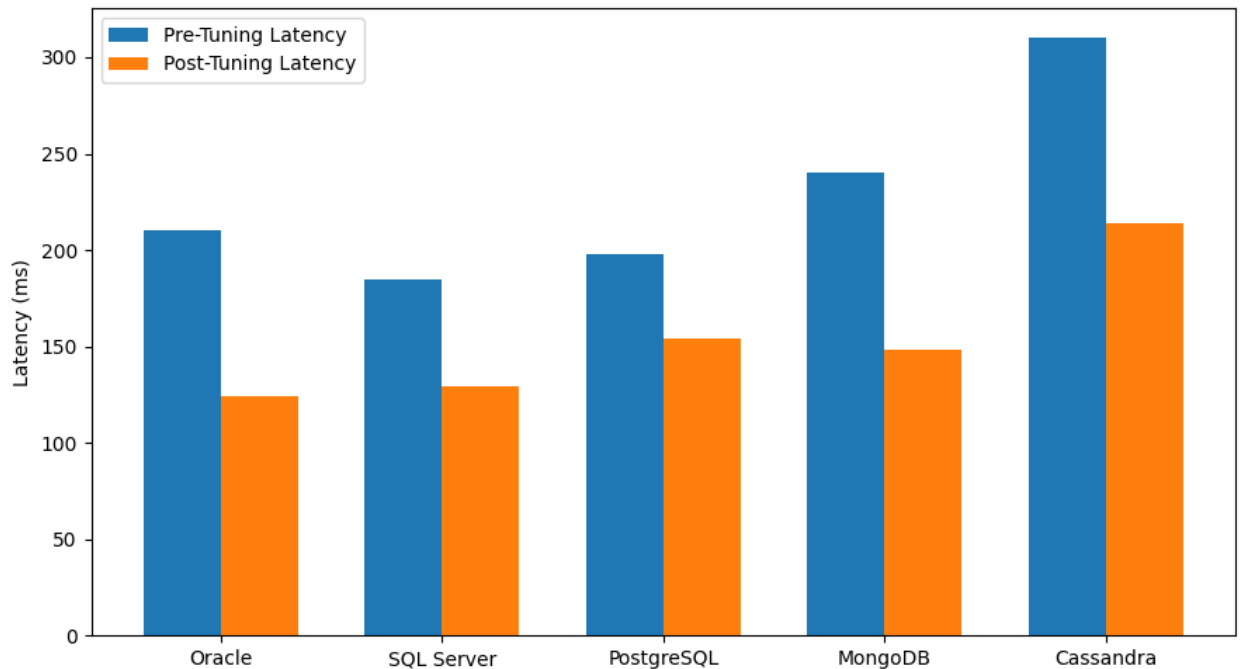
- Oracle: 34: The logical reads decreased by 34 percent, and the efficiency of the SQL plan increased by 41 percent.
- SQL Server: 27 percent reduction in page reads and 30 percent reduction in the time of block.
- PostgreSQL: Based on the statistics, 22 percent fewer I/O wait events were observed, and the decrease was accomplished after optimizing the indices.
- MongoDB: 38 percentage points in the process of sharding balancing and read distribution.
- Cassandra: The compaction process, in addition to the token ranges, reduced the query latency by 31 percent.

These findings can be attributed to the fact that, according to recent literature, the performance of distributed systems can be enhanced using ML (Nerella et al., 2023; Markov et al., 2022).

Table 3: Framework Performance Improvement Across Databases

Database	Pre-Tuning Latency (ms)	Post-Tuning Latency (ms)	Improvement (%)
Oracle	210	124	41%
SQL Server	185	129	30%
PostgreSQL	198	154	22%
MongoDB	240	148	38%
Cassandra	310	214	31%

Table 3: Model-Generated Latency Performance Comparison for Five Database Engines



6. Conclusion and Future Work

This study proposed a single ML-based performance-tuning system that is applicable in heterogeneous database settings. Through its Python-based metric collection, feature engineering, ML model training, and powerful tuning

recommendation engine, the proposed system showed significant improvement in Oracle, SQL Server, PostgreSQL, MongoDB, and Cassandra.

Some of the basic findings are supported by experimental results.

- ML models are dependable for forecasting performance bottlenecks in various architectures.
- Unified feature engineering allows cross-database analysis of architecturally divergent databases.
- Tuning suggestions made by ML are more accurate and effective than heuristics in DBA.
- The framework can lower latency and enhance throughput, making it justified in enterprise-scale database ecosystems.

6.2 Limitations

Although it has strong points, there are a number of limitations.

- The framework depends on the availability of historical metrics; environments where logging is not well developed might require preprocessing.
- Automated tuning action execution has not been adopted to prevent the performance regressions.
- The framework can only support five database engines; the introduced connectors would support additional database engines, such as graph, time-series, and in-memory databases.

6.3 Future Work

- Future extensions include the following:
- Inclusion of reinforcement learning for engaging in self-adaptive tuning.
- Additional support for Redis, Neo4j, and Elasticsearch.
- Creating a web-based model interpretability and DBA oversight dashboard.
- This includes an automated executor that is rollback-safe (to create indexes, rewrite queries, and configure) in real time.
- Incorporating the framework into DevOps pipelines to enable constant optimization whenever deploying CI/CD.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

1. Nikitopoulou, D. (2020). ML-driven automated framework for tuning Spark applications.
2. Neeli, S. S. S. (2021). Key Challenges and Strategies for Managing Databases for Data Science and Machine Learning. *IJLRP-International Journal of Leading Research Publication*, 2(3).
3. Batchu, K. C. (2023). Cross-Platform ETL Federation: A Unified Interface for Multicloud Data Integration. *International Journal of Research Publications in Engineering, Technology and Management (IJRPETM)*, 6(6), 9632-9637.
4. Karri, N. (2023). ML Models That Learn Query Patterns and Suggest Execution Plans. *International Journal of Emerging Trends in Computer Science and Information Technology*, 4(1), 133-141.
5. Karri, N. (2021). Self-Driving Databases. *International Journal of Emerging Trends in Computer Science and Information Technology*, 2(1), 74-83.
6. Lawal, G. S. (2023). Machine Learning Synergies for ERP Optimization and Medical Data Intelligence.
7. Rachakatla, S. K., Ravichandran, P., & Machireddy, J. R. (2021). The role of machine learning in data warehousing: Enhancing data integration and query optimization. *Journal of Bioinformatics and Artificial Intelligence*, 1(1), 82-103.
8. Radhakrishnan, P., & Hemnath, R. (2021). Next-gen cloud synergy: Transforming finance, health, and retail through unified digital frameworks. *International Journal of Business Management and Economic Review*, 4(1), 1-10.
9. Sanca, V., & Ailamaki, A. (2023, April). Analytical engines with context-rich processing: Towards efficient next-generation analytics. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)* (pp. 3699-3707). IEEE.

10. Verma, G., Emani, M., Liao, C., Lin, P. H., Vanderbruggen, T., Shen, X., & Chapman, B. (2021, November). Hpcfair: Enabling fair ai for hpc applications. In 2021 IEEE/ACM Workshop on Machine Learning in High Performance Computing Environments (MLHPC) (pp. 58-68). IEEE.
11. Su, H., & Lee, J. (2023). Machine learning approaches for diagnostics and prognostics of industrial systems using open source data from PHM data challenges: a review. arXiv preprint arXiv:2312.16810.
12. Gures, E., Shayea, I., Ergen, M., Azmi, M. H., & El-Saleh, A. A. (2022). Machine learning-based load balancing algorithms in future heterogeneous networks: A survey. *IEEE Access*, 10, 37689-37717.
13. Lin, X., Kundu, L., Dick, C., Obiodu, E., Mostak, T., & Flaxman, M. (2023). 6G digital twin networks: From theory to practice. *IEEE Communications Magazine*, 61(11), 72-78.
14. Omolayo, O., Ugboko, R., Oyeyemi, D. O., Oloruntoba, O., & Fakunle, S. O. (2022). Optimizing Data Pipelines for Real-Time Healthcare Analytics in Distributed Systems: Architectural Strategies, Performance Trade-offs, and Emerging Paradigms. *International Journal of Health Informatics*, 15(4), 189-204.
15. Ruiz, C. I. S. (2021). AI-Driven Software Development for Scalable IoT Hybrid Fuzzy WPM and TOPSIS Integration with Deep Learning and Particle Swarm Optimization in Agentic Negotiation Frameworks. *International Journal of Research Publications in Engineering, Technology and Management (IJRPETM)*, 4(5), 5570-5574.
16. Katragadda, S. R., Tanikonda, A., Peddinti, S. R., & Pandey, B. K. (2021). Machine Learning-Enhanced Root Cause Analysis for Accelerated Incident Resolution in Complex Systems. *Journal of Science & Technology*, 2(4), 253-275.
17. Pulivarthy, P., & Infrastructure, I. T. (2023). Enhancing data integration in oracle databases: Leveraging machine learning for automated data cleansing, transformation, and enrichment. *International Journal of Holistic Management Perspectives*, 4(4), 1-18.
18. TEWARI, S. (2023, March). Machine Learning Models for Scalable Metadata Management in Data Lakes.
19. Eisemann, L., Fehling-Kaschek, M., Gommel, H., Hermann, D., Klemp, M., Lauer, M., ... & Zhou, J. (2023, September). An approach to systematic data acquisition and data-driven simulation for the safety testing of automated driving functions. In 2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC) (pp. 2440-2447). IEEE.
20. Kumar, V. S., Cuddihy, P., & Aggour, K. S. (2019, June). NodeGroup: a knowledge-driven data management abstraction for industrial machine learning. In *Proceedings of the 3rd International Workshop on Data Management for End-to-End Machine Learning* (pp. 1-4).
21. TEWARI, S. (2023). Scalable Metadata Management in Data Lakes Using Machine Learning.
22. Mohammed, H. U. H. (2021). An End-to-End Machine Learning Framework for Automated Chip Design and Performance Tuning.
23. Kensington, C. E. (2023). Intelligent Cloud-Based Software Maintenance Architecture for Life Insurance Enterprises: Integrating AI, Gray Relational Analysis, and Risk-Aware Scalability across SAP and Oracle EBS Platforms. *International Journal of Research Publications in Engineering, Technology and Management (IJRPETM)*, 6(6), 9688-9691.
24. Li, Y., Wang, X., Miao, Z., & Tan, W. C. (2021). Data augmentation for ml-driven data preparation and integration. *Proceedings of the VLDB Endowment*, 14(12), 3182-3185.
25. Rathore, M. M., Shah, S. A., Shukla, D., Bentafat, E., & Bakiras, S. (2021). The role of ai, machine learning, and big data in digital twinning: A systematic literature review, challenges, and opportunities. *IEEE access*, 9, 32030-32052.
26. Yao, Z., Lum, Y., Johnston, A., Mejia-Mendoza, L. M., Zhou, X., Wen, Y., ... & Seh, Z. W. (2023). Machine learning for a sustainable energy future. *Nature Reviews Materials*, 8(3), 202-215.
27. Markov, I. L., Wang, H., Kasturi, N. S., Singh, S., Garrard, M. R., Huang, Y., ... & Zhou, N. (2022, August). Looper: An end-to-end ml platform for product decisions. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (pp. 3513-3523).
28. Nerella, V. M. L. G., Mahavratayajula, S., & Janardhanan, H. (2023). Machine Learning-Driven Finops Strategies: Adaptive Scaling Models For Balancing Reliability And Cost In Multi-Cloud Data Platforms. *Machine Learning*, 6(4).
29. Akter, M. S., Sultana, N., Khan, M. A. R., & Mohiuddin, M. (2023). Business Intelligence-Driven Healthcare: Integrating Big Data and Machine Learning For Strategic Cost Reduction And Quality Care Delivery. *American Journal of Interdisciplinary Studies*, 4(02), 01-28.
30. Rafi, M. E. H., & Qasem, A. (2022, October). Optimal Launch Bound Selection in CPU-GPU Hybrid Graph Applications with Deep Learning. In 2022 IEEE 13th International Green and Sustainable Computing Conference (IGSC) (pp. 1-7). IEEE.