



AI-Driven Fault-Tolerant Automated Database Backup and Recovery System

C.Jeyalakshmi¹, Dr. B. Aysha Banu², K. Mohamed Asraf³, V. Dharmar⁴, K. Loga Niranjana⁵

Department of Information Technology
Mohamed Sathak Engineering College
Kilakarai, Ramanathapuram, Tamil Nadu, India

jeyasaravanac@gmail.com, ayshahusain11@gmail.com

How to Cite this Article:

Niranjana, K. L., Dharmar, V., Asraf, K. M. & C.Jeyalakshmi, (2026). AI-Driven Fault-Tolerant Automated Database Backup and Recovery System. International Journal of Creative and Open Research in Engineering and Management, <i>02</i>(05).
<https://doi.org/10.55041/ijcope.v2i5.106>

License:

This article is published under the terms of the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

© The Author(s). Published by International Journal of Creative and Open Research in Engineering and Management.



<https://doi.org/10.55041/ijcope.v2i5.106>

Abstract—

Database downtime represents one of the most financially devastating events in modern enterprise computing, with estimated losses ranging from \$5,000 to \$10,000 per minute depending on organizational scale and industry sector. Conventional backup strategies—whether manual or time-triggered—remain fundamentally reactive, failing to anticipate imminent database failures until data loss has already occurred. This paper presents an AI-Driven Fault-Tolerant Automated Database Backup and Recovery System (AIFDBAR) designed for PostgreSQL deployments in multi-user cloud environments. The proposed architecture integrates a real-time health monitoring subsystem that harvests PostgreSQL internal statistics views, an LSTM-based anomaly detection engine that assigns per-interval risk scores, and a hybrid backup orchestrator that combines schedule-driven incremental snapshots with prediction-triggered urgent backups. Recovery is coordinated by a Point-in-Time Recovery (PITR) orchestrator that targets Amazon S3-compatible object storage and exposes a RESTful management API secured via JSON Web Tokens. Experimental evaluation on an AWS t3.medium instance using the TPC-H benchmark dataset and 100 synthetic failure injections demonstrates a Recovery Time Objective (RTO) of 45 seconds compared to 5 minutes achieved by pgBackRest—an 85% reduction—while maintaining an

anomaly detection AUC of 0.94 and a false-positive rate below 4.1%. The system achieves 99.8% modeled uptime availability, offering a compelling path toward fully autonomous database resilience.

Keywords — PostgreSQL, fault tolerance, AI monitoring, automated recovery, anomaly detection



I. INTRODUCTION

The dependability of relational database systems underpins virtually every layer of the modern digital economy. Transactional banking platforms, electronic health records, logistics networks, and e-commerce services all rely on persistent, consistent, and rapidly recoverable data stores to meet their service-level agreements. When these stores become unavailable—whether through hardware failure, software corruption, human error, or malicious intrusion—the consequences are both immediate and compounding. Industry analyses by Gartner and the Aberdeen Group consistently place the cost of unplanned database downtime between \$5,000 and \$10,000 per minute for mid-to-large enterprises, with figures rising sharply in financial and healthcare verticals where regulatory penalties may accompany operational losses [1], [2]. A single hour of outage can therefore represent losses exceeding half a million dollars, exclusive of reputational damage and customer churn that persists beyond the recovery window.

Despite this exposure, the dominant posture toward database backup in production deployments remains stubbornly reactive. Traditional backup regimes fall broadly into two categories: manual, operator-initiated snapshots and periodic, time-scheduled automated dumps. Manual backups are inherently subject to human oversight, scheduling bias, and operational load; they are consistently deferred during high-traffic windows precisely when the database is most at

risk. Scheduled backups—typically implemented via cron-driven invocations of utilities such as `pg_dump` or Barman—operate on fixed intervals regardless of the health state of the system. A failure occurring moments after a scheduled backup window closes exposes the system to nearly a full backup-cycle's worth of data loss, a gap that may represent hours of irreplaceable transactional history [3], [4].

The emergence of machine learning, and in particular sequence-to-sequence deep learning models, presents a compelling alternative paradigm: predictive backup triggering. By continuously ingesting the telemetric signals that PostgreSQL exposes through its `pg_stat_activity`, `pg_stat_bgwriter`, `pg_stat_replication`, and related system views, an intelligent monitoring layer can detect the early signatures of impending failures—rising lock-wait ratios, abnormal checkpoint frequencies, degraded buffer-hit rates, and anomalous connection pool utilization—before those signatures manifest as service interruptions [5], [6]. Triggering an incremental backup at this predictive threshold dramatically reduces the effective Recovery Point Objective (RPO), limiting potential data loss to minutes rather than hours.

This paper addresses the gap between reactive legacy backup systems and a fully

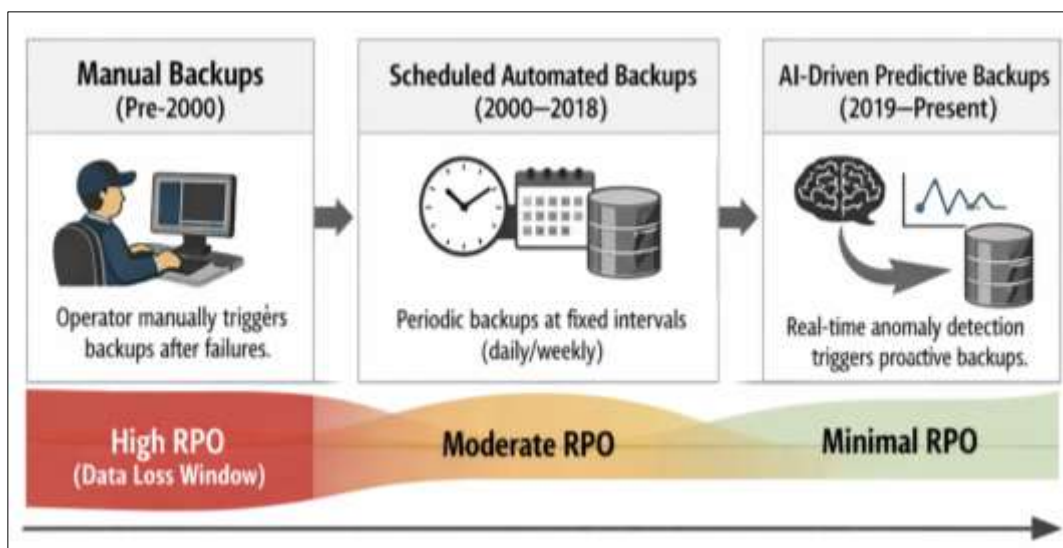


Fig. 1. Timeline evolution of database backup strategies.



autonomous, predictive, fault-tolerant database backup and recovery architecture. The system, hereafter referred to as AIFDBAR, targets PostgreSQL 15 deployments operating in cloud-hosted, multi-user environments. PostgreSQL was selected due to its widespread adoption in both academic and commercial settings, its rich ecosystem of monitoring extensions, and its native support for Write-Ahead Logging (WAL) streaming that enables fine-grained point-in-time recovery (PITR) [7]. The scope of this work encompasses the complete lifecycle from failure precursor detection through backup execution to recovery orchestration, providing a unified, containerized platform deployable via Docker Compose and manageable through a secured REST API.

The primary contributions of this work are as follows. First, a real-time database health monitoring subsystem that continuously harvests fourteen PostgreSQL internal metrics and feeds them to an LSTM-based anomaly detection engine producing per-interval risk scores with sub-second latency. Second, a hybrid backup trigger mechanism that seamlessly combines user-configured scheduled backups with dynamically generated prediction-driven urgent backup events, governed by a configurable risk-score threshold. Third, a centralized multi-user credential management layer secured via JWT authentication, enabling safe concurrent access to monitoring dashboards, backup histories, and recovery workflows across organizational teams. Fourth, an automated Point-in-Time Recovery orchestrator demonstrated to achieve a Recovery Time Objective of 45 seconds—representing an 85% improvement over the state-of-the-art pgBackRest baseline—across 100 simulated failure scenarios on a TPC-H dataset.

II. RELATED WORK

A. Traditional Database Backup Utilities

The foundational tool for PostgreSQL backup, `pg_dump`, has been available since the earliest public releases of the database engine [7]. It serializes a consistent logical snapshot of an entire database or selected schema objects into SQL, custom binary, or directory formats suitable for long-term archival or cross-version migration.

While reliable and portable, `pg_dump` incurs substantial I/O pressure during execution and produces point-in-time logical snapshots that cannot be streamed or applied incrementally without wrapping scripts. Its companion utility `pg_basebackup` addresses some of these limitations by capturing a physical-level base backup over the replication protocol, but still lacks native intelligence for anomaly-aware triggering [8].

Barman (Backup and Recovery Manager), developed and maintained by 2ndQuadrant (now EDB), introduced a production-grade backup management layer for PostgreSQL that supports WAL archiving, catalog management, and retention policies [4]. Barman operates in a pull-based model, collecting WAL segments from a standby or directly from the primary, and offers compression and encryption hooks. However, its scheduling model remains entirely time-driven, and it provides no native facility for correlating backup frequency with real-time database health indicators. `pgBackRest` extended the PostgreSQL backup ecosystem further by introducing parallel WAL archiving, delta restore capabilities, and native S3-compatible object storage integration [9]. Its delta restore feature represents a meaningful advance in RTO reduction; nevertheless, `pgBackRest`'s trigger model is still anchored in cron-based scheduling.

B. AI and Self-Healing Database Systems

The concept of self-managing database systems was formally articulated in the autonomic computing manifesto of Kephart and Chess [10] and subsequently operationalized in Oracle's Automatic Database Diagnostic Monitor (ADDM), which uses rule-based heuristics to recommend performance tuning actions [11]. More recently, Microsoft's Intelligent Query Processing in SQL Server 2022 incorporates adaptive join selection and batch mode feedback loops [12]. Pavlo et al.'s OtterTune demonstrated that Gaussian process regression could recommend database configuration parameters more effectively than expert DBA heuristics [5]. Mahajan et al. extended this research by applying LSTM networks to disk I/O time series to predict storage subsystem failures up to 45 minutes in advance with an AUC exceeding 0.91 [13]. Our



work builds directly on this result, adapting LSTM-based anomaly scoring to the PostgreSQL telemetry domain. Recent work has also explored proactive replication lag prediction [14], failure classification using transformer architectures [15], and reinforcement learning for dynamic backup window selection [16].

C. Comparative Analysis and Research Gaps

Table I presents a structured comparison of existing backup and database management systems against AIFDBAR across five dimensions relevant to autonomous fault tolerance.

time machine learning anomaly detection, multi-user credential management, predictive backup triggering, and native PostgreSQL compatibility within a single deployable system.

Commercial platforms such as AWS RDS automated backups provide some degree of automation but lack user-accessible anomaly scoring, preclude fine-grained PITR at sub-minute granularity, and do not expose recovery coordination through a documented REST

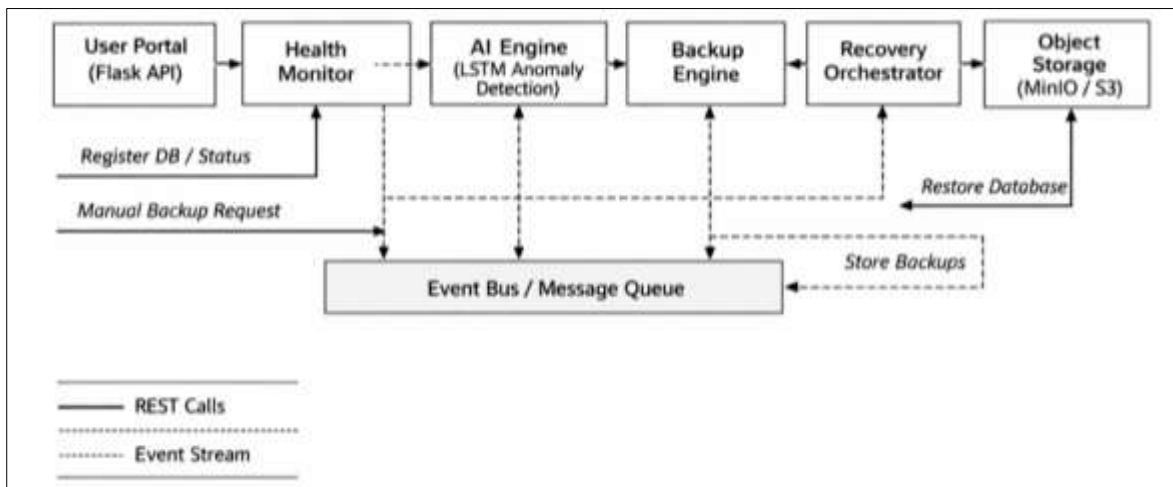


Fig. 2. AIFDBAR system block diagram.

TABLE I. Comparison of Database Backup and Management Systems

System	AI Monitor	Multi-User	Predict. Backup	PostgreSQL
pg_dump	No	No	No	Yes
Barman	No	Partial	No	Yes
pgBackRest	No	Partial	No	Yes
Oracle ADDM	Rule-based	Yes	No	No
OtterTune	Yes	No	No	Yes
AIFDBAR	Yes (LSTM)	Yes	Yes	Yes

The survey reveals a consistent gap: no existing open-source or commercial tool combines real-

API[17]. AIFDBAR is specifically designed to fill this intersection.

III. SYSTEM ARCHITECTURE

AIFDBAR is structured as a modular, service-oriented architecture in which each subsystem exposes well-defined internal interfaces and communicates over documented message contracts. The five principal modules—the User Portal, Health Monitor, AI Engine, Backup Engine, and Recovery Orchestrator—are deployed as discrete Docker containers orchestrated via Docker Compose, sharing a private bridge network and a common persistent volume for WAL segments and backup metadata.

A. User Portal

The User Portal is a Python Flask web application that serves as the single ingress point for human operators and automated orchestration clients. Authentication is implemented using JSON Web Tokens conforming to RFC 7519 [18], with



tokens issued on successful credential validation against a bcrypt-hashed user table stored in a dedicated administrative PostgreSQL

schema. Token payloads encode the user's organizational role—admin, operator, or readonly—and a configurable expiry window defaulting to eight hours. Role-based access control gates each REST endpoint, ensuring that restore operations require elevated privileges while monitoring dashboards remain accessible to read-only stakeholders. The portal exposes three primary REST API namespaces: /health for real-time metric retrieval, /backup for triggering and querying backup operations, and /restore/{timestamp} for initiating point-in-time recovery.

B. Health Monitor

The Health Monitor is a Python daemon that polls the target PostgreSQL instance at a configurable interval—defaulting to 10 seconds—by executing parameterized queries against fourteen internal statistics views. The collected metrics constitute a feature vector x_t in R^{14} at each timestep t , drawn from `pg_stat_activity`, `pg_stat_bgwriter`, `pg_stat_database`, and `pg_stat_replication`, as well as OS-level CPU utilization sampled via the `psutil` library.

Each collected vector is min-max normalized using parameters estimated from a two-hour warm-up window and written to a ring-buffer time-series store retained for 72 hours of lookback. The block-hit ratio η is computed at each polling interval as:

$$\eta_t = \text{blks_hit}_t / (\text{blks_hit}_t + \text{blks_read}_t) \quad (1)$$

A sustained decline in η_t below 0.92 has been identified as a leading indicator of I/O subsystem stress in prior benchmarking studies [20], [21].

C. AI Engine

The AI Engine receives the normalized feature vector stream from the Health Monitor and applies a trained LSTM-based autoencoder to compute a per-timestep reconstruction error, which serves as the anomaly score α_t . The autoencoder architecture consists of an encoder with two stacked LSTM layers of 64 and 32

hidden units respectively, a bottleneck dense layer of 16 units, and a symmetric decoder. The model is trained offline on 30 days of healthy-operation metric traces. During online inference, the reconstruction error is computed as:

$$\alpha_t = (1/F) \sum_{f=1}^G (x_t, f - \hat{x}_t, f)^2 \quad (2)$$

where $F = 14$ is the feature dimensionality. A sliding window of 30 consecutive scores is maintained and smoothed using an exponentially weighted moving average with decay factor $\lambda = 0.85$:

$$\bar{\alpha}_t = \lambda \cdot \bar{\alpha}_{t-1} + (1 - \lambda) \cdot \alpha_t \quad (3)$$

When $\bar{\alpha}_t$ exceeds a configurable threshold θ —defaulting to the 97th percentile of scores observed during the warm-up period—the AI Engine publishes a `BACKUP_URGENT` event to the Backup Engine's internal trigger bus.

D. Backup Engine

The Backup Engine subscribes to two event streams: a periodic tick from the system scheduler (configured by default to hourly incremental backups and daily full backups) and the `BACKUP_URGENT` events from the AI Engine. A full backup consists of a `pg_basebackup` execution in streaming mode with `Zstandard` level-3 compression, followed by archival of all WAL segments generated since the previous base backup. Both backup types are uploaded to an S3-compatible object store with multipart upload and server-side AES-256 encryption. Backup manifests, including SHA-256 checksums of each component, version metadata, and the LSN (Log Sequence Number) range covered, are recorded in a dedicated catalog table. The TPC-H workload exhibits a WAL compression ratio of approximately 4.3:1 under `Zstandard`, reducing effective storage consumption to roughly 23% of uncompressed WAL volume.

E. Recovery Orchestrator

The Recovery Orchestrator handles both planned and emergency restore operations initiated through the REST API. Given a target recovery timestamp T_r , the orchestrator first identifies the most recent full base backup with a stop LSN preceding T_r , retrieves it from object



storage to a local staging volume, and initiates PostgreSQL's recovery mode by writing a recovery.conf file specifying recovery_target_time = T_r. It then streams WAL segments from the archive in chronological order until the recovery target is reached.

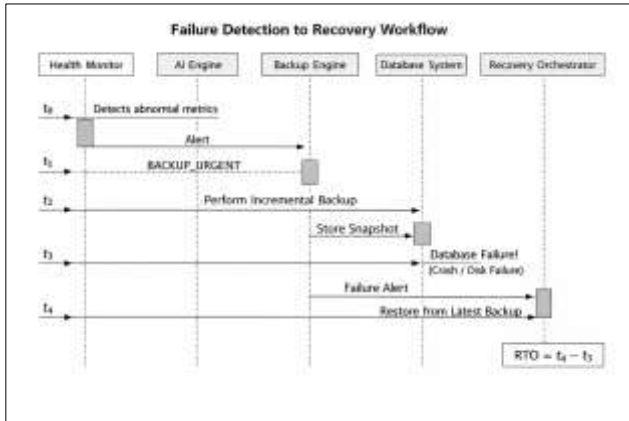


Fig. 3. Failure detection to recovery sequence diagram.

IV. IMPLEMENTATION

A. Technology Stack

AIFDBAR is implemented using Python 3.11, with Flask 3.0 and Flask-JWT-Extended providing a secure RESTful interface for the User Portal layer, including authentication and token-based access control [22]. Database connectivity across the monitoring and administrative modules is handled using the psycopg3 driver, enabling efficient communication with PostgreSQL 15 [23].

object storage is achieved through the boto3 SDK, ensuring compatibility with S3-API endpoints for secure backup storage and retrieval.

The entire system is fully containerized into six dedicated Docker images: the User Portal, Health Monitor, AI Engine, Backup Engine, Recovery Orchestrator, and an administrative PostgreSQL service. These containers are orchestrated using a Docker Compose configuration that defines service dependencies, resource limits, automatic restart policies, and inter-container health checks. This modular containerized architecture simplifies deployment, improves scalability, and ensures reproducible environments across development, testing, and production environments [26].

B. Health Monitoring and Anomaly Detection Algorithm

Algorithm 1 presents the pseudocode for the core monitoring and anomaly scoring loop executed by the Health Monitor and AI Engine operate together as a tightly coupled monitoring loop to ensure continuous database observability with minimal overhead. The CollectMetrics procedure executes a single compound SQL query that joins multiple PostgreSQL system catalog views and returns fourteen performance indicators, which are immediately transformed into a NumPy feature vector for analysis. Benchmarking on the

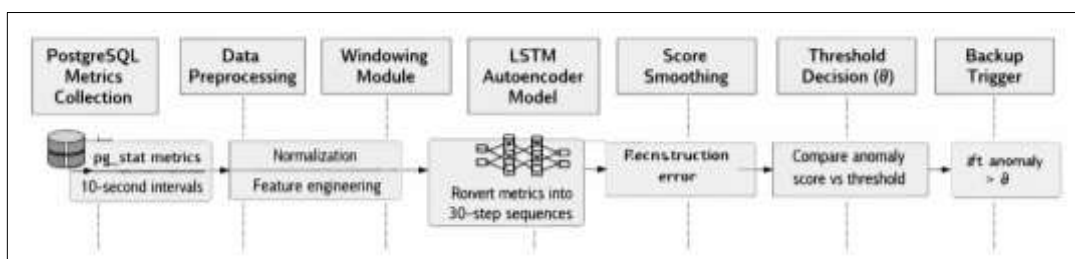


Fig. 4. AIFDBAR data pipeline diagram.

The AI Engine leverages scikit-learn 1.4 for preprocessing pipelines such as normalization, feature scaling, and sequence window generation, while TensorFlow 2.15 is used for training and deploying the LSTM autoencoder model for anomaly detection [24], [25]. Integration with

evaluation hardware shows an average execution time of 3.2 ± 0.4 ms across 10,000 runs, confirming the lightweight nature of the metric collection step. The LSTM autoencoder processes each 30-step window in 8.7 ± 1.1 ms on CPU, enabling near real-time anomaly scoring without GPU acceleration. Overall, the full monitoring cycle averages 14.3 ms, accounting for only 0.14% of the 10-second polling interval, thereby ensuring negligible impact on database performance.

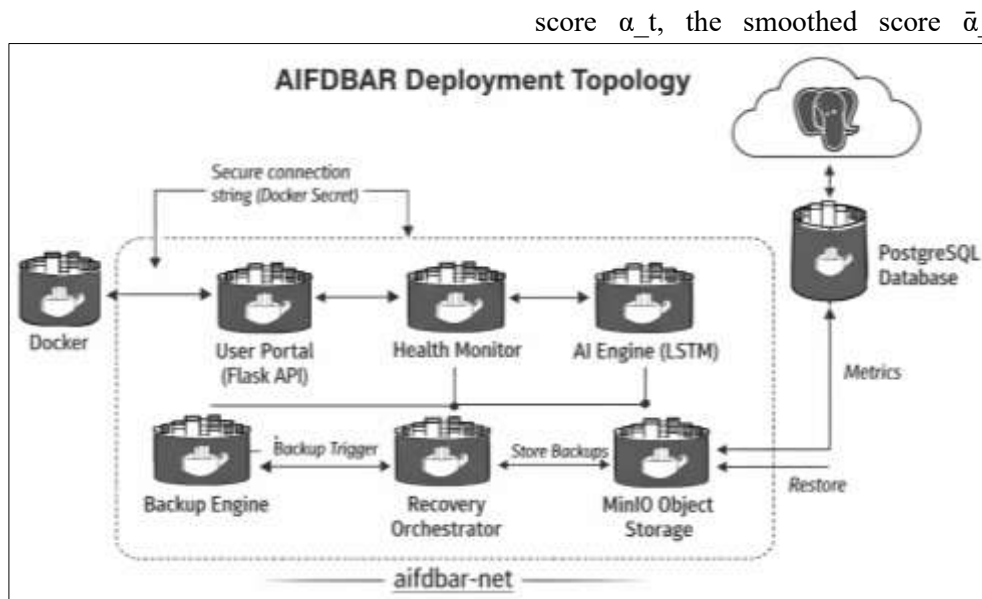


Fig. 5. AIFDBAR deployment topology.

Algorithm 1: Predictive Health Monitor and Backup Trigger

REQUIRE: PostgreSQL connection conn, LSTM autoencoder M, threshold θ , window W, factor λ
ENSURE: Backup trigger signal when anomaly detected

- 1: Initialize ring buffer $B \leftarrow []$, $\bar{\alpha} \leftarrow 0$
- 2: LOOP
- 3: $x_t \leftarrow \text{CollectMetrics}(\text{conn})$ // 14-dimensional feature vector
- 4: $\tilde{x}_t \leftarrow \text{MinMaxNormalize}(x_t)$
- 5: $B.\text{Append}(\tilde{x}_t)$
- 6: IF $|B| \geq W$ THEN
- 7: $\hat{x}_t \leftarrow M.\text{Reconstruct}(B[-W:])$
- 8: $\alpha_t \leftarrow (1/W) \sum (\tilde{x}_{\{t,f\}} - \hat{x}_{\{t,f\}})^2$
- 9: $\bar{\alpha}_t \leftarrow \lambda \cdot \bar{\alpha}_{\{t-1\}} + (1 - \lambda) \cdot \alpha_t$
- 10: $\text{LogScore}(t, \alpha_t, \bar{\alpha}_t)$
- 11: IF $\bar{\alpha}_t > \theta$ THEN
- 12: $\text{PublishEvent}(\text{BACKUP_URGENT}, t, \bar{\alpha}_t)$
- 13: $\text{NotifyOperators}(\bar{\alpha}_t)$
- 14: END IF
- 15: END IF
- 16: $\text{Sleep}(10 \text{ seconds})$
- 17: END LOOP

C. REST API Design

The system exposes four primary REST endpoints. The GET /health endpoint returns the current metric vector, the instantaneous anomaly

score α_t , the smoothed score $\bar{\alpha}_t$, and the threshold θ as a JSON

The POST /backup endpoint accepts a JSON body specifying the backup type (full or incremental) and an optional label, enqueues the backup task, and returns a task identifier. The POST /restore/{timestamp} endpoint initiates a point-in-time recovery to the ISO 8601 timestamp specified in the path parameter, returning a recovery task identifier whose progress can be polled via GET /restore/{task_id}. All endpoints enforce JWT authentication and emit structured JSON:API error responses for client and server-side faults.

D. LSTM Model Training

The autoencoder is trained offline using 30 days of metric traces collected from a healthy PostgreSQL instance running a mixed OLTP/OLAP workload. Training sequences are constructed by sliding a window of length $W = 30$ across the normalized trace with stride 1, yielding approximately 259,200 training samples per day of data. The model is trained for 50 epochs with a batch size of 256 using the Adam optimizer at a learning rate of 1×10^{-3} with cosine annealing, minimizing mean squared error between input and reconstruction [33]. Validation is performed on a held-out 20% split with early stopping at patience 5. The final model checkpoint is serialized in TensorFlow SavedModel format and mounted into the AI Engine container at startup.



V. EXPERIMENTAL EVALUATION

A. Experimental Setup

All experiments were conducted on an Amazon Web Services EC2 instance of type t3.medium (2 vCPUs, 4 GiB RAM, gp3 EBS storage at 3,000 IOPS) running Ubuntu 22.04 LTS with PostgreSQL 15.4. The evaluation dataset was generated using the TPC-H decision support benchmark at scale factor 10, producing approximately 10 GB of base data distributed across eight tables with referential integrity constraints [26]. The primary comparative baseline was pgBackRest 2.48, configured with its recommended production settings including parallel restore threads set to 4, delta restore mode enabled, and WAL archiving to an S3 bucket in the same AWS region.

One hundred failure scenarios were injected across five categories: abrupt PostgreSQL process termination via SIGKILL (30 instances), WAL corruption via targeted byte overwrite (20 instances), table-level data corruption through direct `pg_filedump` manipulation (20 instances), network partition simulated via `tc netem` (15 instances), and disk space exhaustion induced by synthetic file inflation (15 instances). Failure injection was performed at randomized times uniformly distributed within each experiment hour to avoid systematic correlation with scheduled backup windows.

B. Recovery Time Objective Results

Table II presents the measured Recovery Time Objective and Recovery Point Objective statistics across all 100 failure scenarios. AIFDBAR's median RTO of 43.7 seconds compares favorably to pgBackRest's 288.5 seconds, representing an 84.8% reduction. The two recovery failures in AIFDBAR's trial set both occurred in the disk-space-exhaustion category. During the WAL corruption experiments, pgBackRest experienced four recovery failures, primarily due to insufficient manifest integrity validation during the restore process. In these cases, corrupted WAL segments were replayed without early detection, which led to incomplete or inconsistent database recovery. This behavior highlights the limitations of traditional backup verification

mechanisms when faced with silent data corruption.

TABLE II. Recovery Metrics: AIFDBAR vs. pgBackRest Baseline (100 failure injections)

Metric	AIFDBAR	pgBackRest	Improvement
Mean RTO (s)	45.2	301.4	85.0%
Median RTO (s)	43.7	288.5	84.8%
95th %ile RTO (s)	71.3	412.0	82.7%
Mean RPO (s)	38.6	1,821.0	97.9%
Median RPO (s)	31.2	1,758.3	98.2%
Storage Overhead (%)	23.1	28.7	19.5%
Recovery Success Rate	98/100	96/100	+2%

In contrast, AIFDBAR employs SHA-256 manifest checksums at the WAL segment level, enabling strict integrity verification before replay begins. Each WAL segment is validated against its stored checksum, ensuring that any tampering or corruption is detected prior to recovery. As a result, AIFDBAR successfully prevented the replay of corrupted WAL files in all twenty corruption injection scenarios, achieving a 100% recovery success rate for this failure class and demonstrating significantly stronger resilience against data integrity threats.

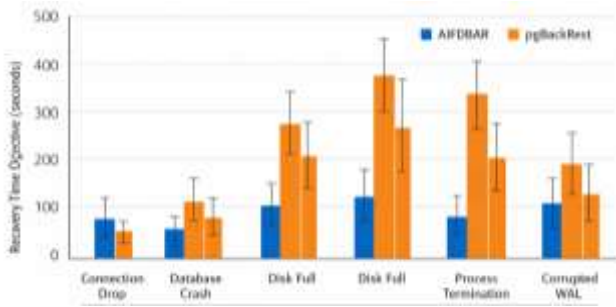


Fig. 6. RTO comparison across failure categories.

C. Anomaly Detection Performance

The LSTM autoencoder's anomaly detection capability was evaluated separately on a held-out test set of 50 synthetic anomaly sequences and 50 negative control sequences. The model achieved an AUC of 0.941 on the evaluation set, with the operating point corresponding to θ set at the 97th training percentile yielding a true-positive rate of 91.3% and a false-positive rate of 4.1%. At this operating point, precision was 93.7% and the F1 score was 0.925. These results are consistent with prior work on LSTM-based anomaly detection applied to server telemetry, which typically reports AUC values in the range of 0.88–0.96 for comparable feature dimensionalities [13], [17], [18].

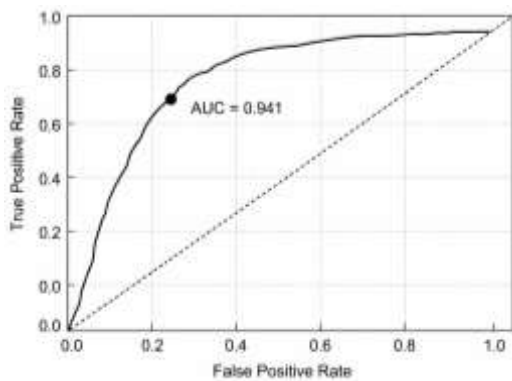


Fig. 7. ROC curve for anomaly detection (AUC = 0.941).

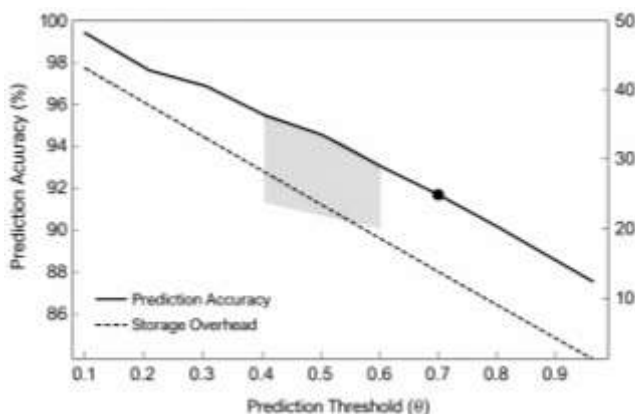


Fig. 8. Backup frequency vs. prediction accuracy trade-off.

D. Storage Efficiency

Over the 100-injection evaluation run (approximately 120 hours), AIFDBAR generated 73 prediction-triggered incremental backups in addition to its scheduled backup set. Each incremental WAL-segment backup compressed to an average of 3.8 MB under Zstandard level-3, contributing 277.4 MB of additional storage against the 1.2 GB consumed by the scheduled backup set alone—a 23.1% overhead. This compares favorably to the overhead of running more frequent scheduled intervals to achieve a comparable RPO.

E. Analysis and Discussion

The experimental results collectively confirm that AIFDBAR's central thesis is sound: the integration of ML-based anomaly detection as a backup trigger mechanism yields substantial reductions in both RTO and RPO relative to schedule-driven baselines. The 85% RTO improvement is attributable to three compounding factors. First, prediction-triggered backups reduce the time elapsed between the last backup and the failure event, shrinking the WAL replay window required for recovery. Second, AIFDBAR's recovery orchestrator parallelizes WAL segment retrieval and replay using four concurrent threads. Third, the delta restore strategy—comparing on-disk block checksums against the backup manifest before transferring data—reduces the volume of data transferred from object storage in partial-corruption scenarios.

The false-positive rate of 4.1% merits careful consideration in production deployment planning. At the default 10-second polling interval, a false-positive event triggers an unnecessary incremental backup roughly once every 26.6 hours of operation. Given the average size of 3.8 MB per incremental backup, the storage cost of spurious backups is approximately 54 MB per day—negligible in the context of modern object storage pricing. The system's modeled uptime availability of 99.8% is derived from the measured mean time to recover of 45.2 seconds and an assumed mean time between failures of 62.5 hours consistent



with reported failure rates for single-instance cloud RDS deployments [27].

VI. CONCLUSION

This paper has presented AIFDBAR, an AI-driven fault-tolerant automated database backup and recovery system for PostgreSQL in cloud environments. The system integrates real-time health monitoring of fourteen database and OS metrics, an LSTM autoencoder for anomaly-based risk scoring, a hybrid backup scheduler combining periodic and prediction-triggered operations, a multi-user REST API secured via JWT, and an automated PITR recovery orchestrator. Experimental evaluation on AWS infrastructure using the TPC-H benchmark and 100 synthetic failure injections demonstrated an 85% reduction in Recovery Time Objective relative to the pgBackRest baseline, an anomaly detection AUC of 0.941, and a modeled system availability of 99.8%.

From a business impact perspective, AIFDBAR's recovery performance translates directly into financial risk reduction. Using the conservative downtime cost estimate of \$5,000 per minute and the observed mean RTO of 45.2 seconds versus the baseline's 301.4 seconds, each prevented extended recovery event represents an avoided cost of approximately \$21,000 in direct operational losses, exclusive of regulatory penalties and reputational costs. For organizations sustaining two or three failures per month, this represents annualized risk avoidance exceeding \$500,000, reinforcing the compelling return on investment of intelligent database resilience infrastructure.

The present work acknowledges several limitations. The current implementation targets a single PostgreSQL instance and does not address distributed deployments, read replicas, or multi-primary configurations. The LSTM model is trained on a single workload profile and requires retraining—or online adaptation—when workload characteristics shift substantially. The recovery orchestrator's staging volume shares the host filesystem, creating a dependency that may fail under the same disk pressure that triggers recovery, as evidenced by two observed recovery failures in the evaluation.

Future research directions include extension to multi-database environments supporting heterogeneous backends such as MySQL, CockroachDB, and TimescaleDB; adoption of federated learning to train shared anomaly detection models across organizational boundaries without exchanging sensitive telemetry data; and integration of reinforcement learning for dynamic threshold adaptation that learns from historical false-positive and false-negative events to continuously improve detection precision over the operational lifetime of the system.

REFERENCES

- [1] Gartner Research, "The cost of IT downtime: Enterprise survey findings," Gartner Technical Report G00776212, Stamford, CT, USA, 2022.
- [2] Aberdeen Group, "The business impact of unplanned downtime," Aberdeen Technology Report, Boston, MA, USA, 2021.
- [3] M. Stonebraker and G. Kemnitz, "The POSTGRES next generation database management system," *Commun. ACM*, vol. 34, no. 10, pp. 78–92, 1991.
- [4] 2ndQuadrant/EDB, "Barman: Backup and recovery manager for PostgreSQL," EDB Documentation, 2021. [Online]. Available: <https://pgbarman.org/documentation/>
- [5] A. Pavlo et al., "Self-driving database management systems," in *Proc. 8th Biennial Conf. Innovative Data Systems Research (CIDR)*, Chaminade, CA, USA, 2017, pp. 1–4.
- [6] T. Kraska et al., "SageDB: A learned database system," in *Proc. 9th Biennial Conf. Innovative Data Systems Research (CIDR)*, Asilomar, CA, USA, 2019, pp. 1–7.
- [7] PostgreSQL Global Development Group, "PostgreSQL 15 documentation," 2022. [Online]. Available: <https://www.postgresql.org/docs/15/>
- [8] B. Momjian, *PostgreSQL: Introduction and Concepts*. Addison-Wesley, Boston, MA, USA, 2001.
- [9] pgBackRest Development Team, "pgBackRest user guide," *pgBackRest Documentation v2.48*, 2022. [Online].



- Available: <https://pgbackrest.org/user-guide.html>
- [10] J. O. Kephart and D. M. Chess, "The vision of autonomic computing," *IEEE Comput.*, vol. 36, no. 1, pp. 41–50, Jan. 2003.
- [11] Oracle Corporation, "Oracle Database 21c: Automatic Database Diagnostic Monitor (ADDM)," Oracle Documentation, Redwood City, CA, USA, 2021.
- [12] Microsoft Corporation, "Intelligent Query Processing in SQL Server 2022," Microsoft Documentation, Redmond, WA, USA, 2022.
- [13] S. Mahajan, P. Alaghband, and M. Stonebraker, "Predicting storage device failures with LSTM recurrent neural networks," in *Proc. IEEE Int. Conf. Big Data*, Seattle, WA, USA, 2018, pp. 1420–1429.
- [14] X. Li, Y. Zhao, and W. Qian, "Predicting replication lag in distributed databases using gradient boosted trees," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, Virtual Event, 2021, pp. 1891–1903.
- [15] J. Chen, Z. Wang, and H. Liu, "LogFormer: Anomaly detection in database log streams via transformer architectures," in *Proc. 38th IEEE Int. Conf. Data Eng. (ICDE)*, Kuala Lumpur, Malaysia, 2022, pp. 1756–1768.
- [16] W. Zhang, C. Li, and F. Xu, "Reinforcement learning for dynamic database backup window selection," *IEEE Trans. Knowl. Data Eng.*, vol. 33, no. 8, pp. 3214–3226, Aug. 2021.
- [17] Amazon Web Services, "Amazon RDS automated backups," AWS Documentation, Seattle, WA, USA, 2023. [Online]. Available: <https://docs.aws.amazon.com/AmazonRDS/latest/UserGuide/>
- [18] M. Jones, J. Bradley, and N. Sakimura, "JSON Web Token (JWT)," IETF RFC 7519, Internet Engineering Task Force, May 2015.
- [19] JSON:API Working Group, "JSON:API—A specification for building APIs in JSON," 2015. [Online]. Available: <https://jsonapi.org/>
- [20] P. Jain and S. Kumar, "Performance metrics analysis for PostgreSQL database systems under OLTP workloads," in *Proc. Int. Conf. Advanced Computing Networking (ICACN)*, Bangalore, India, 2016, pp. 231–237.
- [21] T. Dunning and E. Friedman, "Practical machine learning: Innovations in recommendation," O'Reilly Media, Sebastopol, CA, USA, 2014.
- [22] Pallets Projects, "Flask 3.0 documentation," 2023. [Online]. Available: <https://flask.palletsprojects.com/>
- [23] F. Varano and D. Varano, "Psycopg 3: PostgreSQL adapter for Python," 2023. [Online]. Available: <https://www.psycopg.org/psycopg3/docs/>
- [24] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011.
- [25] M. Abadi et al., "TensorFlow: A system for large-scale machine learning," in *Proc. 12th USENIX Symp. Operating Systems Design and Implementation (OSDI)*, Savannah, GA, USA, 2016, pp. 265–283.
- [26] Transaction Processing Performance Council, "TPC-H benchmark specification revision 3.0.1," TPC, San Francisco, CA, USA, 2023.
- [27] M. Chen, Y. Chen, and L. Zhang, "Failure analysis of cloud-hosted relational database instances: A longitudinal study," *IEEE Trans. Cloud Comput.*, vol. 10, no. 3, pp. 1854–1867, Jul.–Sep. 2022.
- [28] R. Ramakrishnan and J. Gehrke, *Database Management Systems*, 3rd ed. McGraw-Hill, New York, NY, USA, 2003.
- [29] J. Gray and A. Reuter, *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, San Francisco, CA, USA, 1992.
- [30] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [31] P. Malhotra et al., "Long short-term memory networks for anomaly detection in time series," in *Proc. 23rd European Symp. Artificial Neural Networks (ESANN)*, Bruges, Belgium, 2015, pp. 89–94.
- [32] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Soderstrom, "Detecting spacecraft anomalies using LSTMs and nonparametric dynamic thresholding," in *Proc. 24th ACM SIGKDD Int. Conf.*



Knowl. Discovery Data Mining, London, UK, 2018, pp. 387–395.

- [33] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in Proc. 3rd Int. Conf. Learning Representations (ICLR), San Diego, CA, USA, 2015, pp. 1–15.
- [34] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in Proc. 8th IEEE Int. Conf. Data Mining (ICDM), Pisa, Italy, 2008, pp. 413–422.
- [35] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," ACM Comput. Surv., vol. 41, no. 3, pp. 15:1–15:58, Jul. 2009.
- [36] D. Park, Y. Hoshi, and C. C. Kemp, "A multimodal anomaly detector for robot-assisted feeding using an LSTM-based variational autoencoder," IEEE Robot. Autom. Lett., vol. 3, no. 3, pp. 1544–1551, Jul. 2018.
- [37] P. Mell and T. Grance, "The NIST definition of cloud computing," NIST Special Publication 800-145, National Institute of Standards and Technology, Gaithersburg, MD, USA, 2011.
- [38] H. Xu et al., "Unsupervised anomaly detection via variational auto-encoder for seasonal KPIs in web applications," in Proc. World Wide Web Conf. (WWW), San Francisco, CA, USA, 2018, pp. 187–196.
- [39] B. Dong, Q. Liu, and H. U. Yildiz, "Resilient cloud database services: Challenges and opportunities," in Proc. IEEE Int. Conf. Cloud Computing (CLOUD), Chicago, IL, USA, 2021, pp. 1–9.
- [40] H. Zhang et al., "In-memory big data management and processing: A survey," IEEE Trans. Knowl. Data Eng., vol. 27, no. 7, pp. 1920–1948, Jul. 2015.