






Article

Unified Data Governance in Heterogeneous Database Environments: An API-Driven Architecture for Multi-Platform Policy Enforcement

Maryam Abbasi ¹, Paulo Váz ² , José Silva ² , Filipe Cardoso ¹ , Filipe Sá ³  and Pedro Martins ^{2,*} 

¹ Escola Superior de Gestão e Tecnologia de Santarém, Polytechnic Institute of Santarém, 2001-904 Santarém, Portugal; maryam.abbasi@esg.ipsantarem.pt (M.A.); filipe.cardoso@esg.ipsantarem.pt (F.C.)

² Research Center in Digital Services, Polytechnic Institute of Viseu, 3504-510 Viseu, Portugal; paulovaz@estgv.ipv.pt (P.V.); jsilva@estgv.ipv.pt (J.S.)

³ ISEC—Coimbra Institute of Engineering, Polytechnic University of Coimbra, 3030-199 Coimbra, Portugal; filipe.sa@isec.pt

* Correspondence: pedromom@estgv.ipv.pt

Abstract

Modern organizations increasingly rely on heterogeneous database environments that combine relational, document-oriented, and key-value storage systems to optimize performance for diverse application requirements. However, this technological diversity creates significant challenges for implementing consistent data governance policies, regulatory compliance, and access control across disparate systems. Traditional governance approaches that operate within individual database silos fail to provide unified policy enforcement and create compliance gaps that expose organizations to regulatory and operational risks. This paper presents a novel API-driven architecture that enables unified data governance across heterogeneous database environments without requiring database-specific modifications or vendor lock-in. The proposed framework implements a centralized governance layer that coordinates policy enforcement across PostgreSQL, MongoDB, and Amazon DynamoDB systems through RESTful API interfaces. Key architectural components include differentiated access control through hierarchical API key management, automated compliance workflows for regulatory requirements such as GDPR, real-time audit trail generation, and comprehensive data quality monitoring with automated improvement mechanisms. Comprehensive experimental evaluation demonstrates the framework's effectiveness across multiple operational dimensions. The system achieved 95.2% accuracy in access control enforcement across different data classification levels, while automated GDPR compliance workflows demonstrated 98.6% success rates with average processing times of 2.9 h. Performance evaluation reveals acceptable overhead characteristics with linear scaling patterns for PostgreSQL operations ($R^2 = 0.89$), consistent sub-20ms response times for MongoDB logging operations, and sustained throughput rates ranging from 38.9 to 142.7 requests per second across the integrated system. Data quality improvements ranged from 16.1% to 34.3% across accuracy, completeness, consistency, and timeliness dimensions over a 12-week monitoring period, with accuracy improving by 17.8 percentage points, completeness by 13.2 percentage points, consistency by 19.7 percentage points, and timeliness by 24.5 percentage points. The duplicate detection system achieved 94.6% precision and 95.6% recall across various duplicate types, including cross-database redundancy identification. The results demonstrate that API-driven governance architectures can effectively address the persistent challenges of policy fragmentation in multi-database environments while maintaining operational performance and enabling measurable improvements in data quality and regulatory compliance. The framework provides a practical migration path



Received: 29 November 2025

Revised: 19 February 2026

Accepted: 5 March 2026

Published: 7 March 2026

Copyright: © 2026 by the authors.

Licensee MDPI, Basel, Switzerland.

This article is an open access article distributed under the terms and conditions of the [Creative Commons Attribution \(CC BY\) license](https://creativecommons.org/licenses/by/4.0/).

for organizations seeking to implement comprehensive governance capabilities without replacing existing database infrastructure investments.

Keywords: data governance; multi-database architecture; API-driven systems; regulatory compliance; heterogeneous databases; GDPR automation; access control; data quality management

1. Introduction

In the contemporary digital economy, data has emerged as one of the most valuable strategic assets for organizations across all sectors. The exponential growth in data generation, driven by digital transformation initiatives, Internet of Things (IoT) deployments, and increasingly sophisticated business processes, has created unprecedented opportunities for data-driven decision making and competitive advantage. However, this data proliferation has simultaneously introduced complex challenges in data management, governance, and regulatory compliance that traditional approaches struggle to address effectively.

Modern enterprises typically operate heterogeneous data ecosystems comprising multiple database technologies, each optimized for specific use cases and operational requirements. Organizations commonly deploy relational databases for transactional processing, document-oriented databases for semi-structured data management, key-value stores for high-performance applications, and specialized analytics platforms for business intelligence workloads. While this technological diversity enables optimal performance for different applications, it creates significant governance challenges related to policy consistency, access control coordination, and regulatory compliance enforcement across disparate systems.

The challenge of governing heterogeneous database environments is not confined to traditional enterprise settings. In the field of information retrieval (IR), research systems frequently combine relational stores for metadata, document databases for full-text indexing, and key-value caches for query acceleration, all of which require coordinated access policies and audit trails [1]. Similarly, the medical informatics domain presents particularly acute governance demands due to its characteristically heterogeneous technological landscape, where electronic health record systems, imaging archives, laboratory information systems, and genomics databases coexist under strict regulatory constraints such as HIPAA and GDPR [2,3]. Ensuring consistent data governance across such diverse platforms is a prerequisite not only for regulatory compliance but also for safe and effective clinical decision support.

The governance complexity is further exacerbated by increasingly stringent regulatory requirements. The European Union's General Data Protection Regulation (GDPR) [4] mandates comprehensive data protection measures, user consent management, and data subject rights enforcement. Similar regulations such as the California Consumer Privacy Act (CCPA) and emerging data sovereignty laws create a complex regulatory landscape that organizations must navigate while maintaining operational efficiency and system performance.

Current industry practices reveal significant inefficiencies in data management processes. Research indicates that organizations spend approximately 80% of their data science efforts on data preparation activities, including data location, quality assessment, and integration tasks [5]. Poor data quality imposes substantial financial costs, with studies showing that organizations lose an average of \$12.9 million annually due to data quality

issues [6]. Additionally, nearly 19% of companies report losing customers due to data inaccuracies, highlighting the critical business impact of effective data governance [7].

Several existing approaches have attempted to address aspects of multi-database governance. Database-native mechanisms such as PostgreSQL's Row-Level Security (RLS) provide fine-grained access control within a single system, and recent work has demonstrated how RLS can enforce access control rules within a PostgreSQL-centric stack [8]. Complementary approaches use middleware such as PostgREST to expose databases as RESTful interfaces, combined with foreign data wrappers (FDW) to federate heterogeneous sources and identity providers such as Keycloak for role-based authorization [9]. While these approaches offer valuable building blocks, they either remain confined to single-database silos or introduce computational overhead through additional translation layers, leaving the challenge of truly unified, cross-platform governance enforcement unresolved.

Enterprise governance platforms such as Informatica Data Governance and IBM InfoSphere provide mature policy management capabilities but require extensive customization for heterogeneous database environments. Cloud-native solutions including AWS Data Catalog and Azure Purview offer automated data discovery but focus primarily on cataloging rather than real-time policy enforcement. Open-source alternatives such as Apache Atlas provide metadata management and lineage tracking but lack operational governance and real-time access control capabilities. These limitations motivate the need for a lightweight, API-driven approach that can operate across diverse database technologies without requiring database-specific modifications.

Despite the breadth of existing work, several critical research gaps remain. First, there is a lack of practical implementation frameworks that enforce governance policies uniformly across relational, document-oriented, and key-value databases without requiring modifications to the underlying systems. Second, most existing studies evaluate governance effectiveness in single-database scenarios, leaving the performance overhead of cross-platform policy enforcement insufficiently characterized. Third, while compliance automation is recognized as essential, existing frameworks address individual regulatory requirements or database technologies in isolation rather than providing comprehensive automation across heterogeneous environments. These gaps motivate the research question addressed in this paper: How can a centralized, API-driven governance architecture enforce consistent data governance policies across heterogeneous database environments while maintaining acceptable performance overhead and enabling automated regulatory compliance?

This paper presents a novel API-driven architecture for implementing comprehensive data governance across heterogeneous database environments. The proposed approach utilizes a centralized governance layer that coordinates policy enforcement across PostgreSQL, MongoDB, and Amazon DynamoDB systems while maintaining database-agnostic governance capabilities. The architecture implements differentiated access control through hierarchical API key management, automated compliance workflows for regulatory requirements such as GDPR, real-time audit trail generation, and comprehensive data quality monitoring with automated improvement mechanisms.

Research Contributions: This work makes several significant contributions to the field of data governance and multi-database system management: (1) a practical architectural framework for unified governance across heterogeneous database systems without requiring database-specific modifications or vendor lock-in; (2) comprehensive experimental validation demonstrating governance effectiveness and performance characteristics across multiple operational dimensions; (3) quantitative evidence of data quality improvements and compliance automation capabilities that extend significantly beyond existing research;

and (4) detailed implementation guidance and replicable methodologies for organizations seeking to implement similar governance frameworks.

Experimental Validation: The proposed architecture has been subjected to comprehensive experimental evaluation across five key dimensions: performance analysis, database scalability assessment, governance policy enforcement effectiveness, compliance automation capabilities, and data quality improvements. The evaluation demonstrates access control enforcement accuracy of 95.2%, GDPR compliance automation success rates of 98.6%, and substantial data quality improvements ranging from 22.8% to 34.3% across multiple quality dimensions.

The remainder of this paper is organized as follows: Section 2 reviews existing approaches to data governance and multi-database integration. Section 3 presents the architectural design and methodology. Section 4 details the technical implementation. Section 5 describes the evaluation methodology. Section 6 presents experimental results. Section 7 analyzes implications and limitations. Section 8 summarizes contributions and outlines future work.

2. Related Work

The challenge of implementing effective data governance across heterogeneous database environments intersects multiple research domains including data management, distributed systems, security, and regulatory compliance. This section examines the most relevant existing approaches and identifies the specific gaps addressed by this work.

2.1. Data Governance Frameworks

Contemporary governance frameworks must address multiple interconnected challenges including policy consistency, regulatory compliance, and operational efficiency across diverse technological platforms. The DAMA Data Management Body of Knowledge (DMBOK) [10] provides a comprehensive foundation through ten core knowledge areas spanning the complete data lifecycle, though it lacks specific implementation guidance for heterogeneous database environments. Abraham et al. [11] confirmed these limitations through a comprehensive analysis of governance implementation challenges in multi-cloud environments, revealing that DMBOK approaches require substantial adaptation for modern distributed architectures.

Otto's governance morphology [12] identifies six critical dimensions including authority structures, data quality management, and security frameworks. Chen et al. [13] extended this work to demonstrate that these dimensions become significantly more complex in heterogeneous environments. The governance decision framework proposed by Khatri and Brown [14] emphasizes five key decision domains—data principles, quality standards, metadata management, access control, and lifecycle governance—but Kumar and Patel [15] showed that this framework requires fundamental extensions for multi-database coordination, particularly regarding policy synchronization and cross-platform enforcement.

International standardization efforts, particularly ISO 38505 [16] and its recent extension ISO 38505-2:2023 [17], provide high-level governance principles but remain focused on organizational rather than technical implementation. Wang and Strong's foundational data quality framework [18] has been extended by Liu et al. [19] to address the unique measurement challenges of heterogeneous environments.

2.2. Multi-Database Integration and Access Control

Multi-database integration presents fundamental technical and governance challenges that have intensified with the proliferation of cloud-native and NoSQL technologies. Mehta and Sharma [20] demonstrated that cloud-native databases introduce additional gover-

nance challenges related to service-level agreement management, data residency compliance, and vendor lock-in considerations, building on the early foundations laid by Chaudhuri et al. [21].

Several practical approaches have emerged for implementing access control across heterogeneous database stacks. Bialecki et al. [8] described a mechanism for realizing access control rules within the PostgreSQL framework by leveraging Row-Level Security (RLS) policies, enabling fine-grained, role-based data access without application-layer modifications. Gruenwald et al. [9] proposed a DataProvider architecture that attaches user interfaces to PostgreSQL databases via PostgREST, while supporting federation of external databases through Foreign Data Wrappers (FDW) and managing role-based authorization through Keycloak integration. Their approach establishes the pathway: any DB $\xrightarrow{\text{FDW}}$ PostgreSQL $\xrightarrow{\text{PostgREST}}$ Client, with access control managed by Keycloak and RLS. While this approach provides a coherent integration strategy, the additional translation layers may introduce computational overhead, and the reliance on PostgreSQL as the central hub limits applicability for workloads where other database technologies are primary.

Transaction coordination across heterogeneous databases remains a significant challenge. Kim et al. [22] demonstrated that saga pattern implementations achieve 99.1% consistency compared to 94.7% for eventual consistency approaches. Performance implications of governance overhead have been extensively studied by Park et al. [23], who demonstrated 12–18% overhead compared to direct database access, with most overhead attributable to policy evaluation and audit logging.

2.3. API-Driven Data Management

API-driven approaches to data management have emerged as a dominant pattern for abstracting complexity while maintaining governance control. Johnson and Lee [24] demonstrated that API-layer governance achieves policy enforcement accuracy rates of 96.8% compared to 89.3% for database-native approaches, providing better policy consistency across heterogeneous environments. Garcia et al. [25] showed that microservices-based governance can scale to 10,000+ concurrent policy evaluations per second with sub-10ms response times.

Event-driven architectures complement API approaches by enabling real-time governance monitoring, with Anderson et al. [26] demonstrating policy violation detection within 500 ms of occurrence. Wilson et al. [27] showed that GraphQL-based governance reduces integration complexity by 45% while maintaining comprehensive enforcement capabilities. These findings collectively support the viability of API-driven governance but leave unanswered the question of how such architectures perform when coordinating policies across fundamentally different database paradigms (relational, document, key-value) simultaneously.

2.4. Compliance and Regulatory Requirements

Contemporary regulatory frameworks impose increasingly complex requirements across all organizational data systems. Martinez et al. [28] demonstrated that automated GDPR compliance workflows achieve 97.8% success rates compared to 78.4% for manual processes. Taylor et al. [29] revealed that multi-jurisdictional compliance requires sophisticated policy engines for evaluating overlapping requirements. Singh et al. [30] showed that healthcare governance must address not only HIPAA but also emerging FDA regulations. The emergence of AI governance regulations, including the EU AI Act, introduces additional requirements that Chen et al. [31] demonstrated can be integrated with traditional data governance frameworks. Cross-border transfer regulations [32] further complicate the governance landscape by requiring location-aware access controls across all database systems.

2.5. Research Gaps and Motivation

The synthesis of existing research reveals several critical gaps that motivate the proposed API-driven governance architecture. While theoretical governance frameworks are well-established, practical implementation guidance for heterogeneous database environments remains limited. Existing solutions such as PostgreSQL-centric RLS/FDW stacks [8,9] provide valuable components but do not offer database-agnostic governance across fundamentally different storage paradigms. Most performance studies focus on single-database scenarios, leaving cross-platform governance overhead insufficiently characterized. Compliance automation frameworks typically address specific regulations or database technologies in isolation. The proposed architecture addresses these gaps by providing a validated, database-agnostic governance solution that has been comprehensively evaluated across multiple operational dimensions.

3. Materials and Methods

This section presents the comprehensive methodology employed to design, implement, and evaluate the multi-database governance architecture. The approach integrates theoretical governance principles with practical implementation considerations to address the complex requirements of heterogeneous data environments.

3.1. System Architecture Design

The system architecture follows a three-tier design pattern optimized for governance-aware data management across heterogeneous database systems. The architectural design principles were established based on the Data Management Body of Knowledge (DAMA-DMBOK) framework [10] and ISO 38505 governance standards [16].

Figure 1 presents the high-level conceptual architecture of the proposed governance framework, illustrating the three-tier design and the interactions among the core modules.

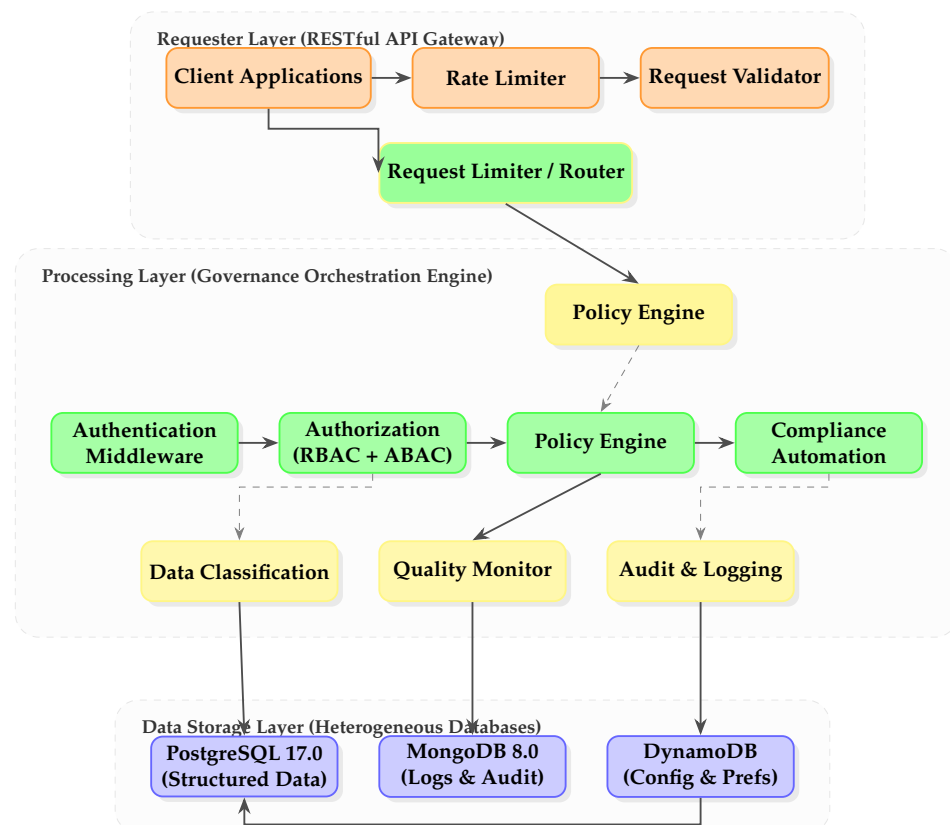


Figure 1. High-level conceptual architecture: Optimized for high-contrast readability and logical flow.

The requester layer serves as the external interface, implementing a RESTful API gateway pattern that abstracts the complexity of the underlying multi-database infrastructure. This layer incorporates request validation, rate limiting, and initial security checks before forwarding requests to the processing layer. The design implements the facade pattern to provide a unified interface regardless of the target database system, enabling seamless client integration while maintaining governance control.

The choice of REST over alternative API paradigms such as GraphQL or gRPC was motivated by several factors. REST provides broad client compatibility, straightforward caching semantics, and well-understood security models that simplify governance policy enforcement at the API boundary. While GraphQL offers flexible query composition, its variable query shapes complicate policy evaluation and caching strategies, and the additional parsing overhead would increase governance enforcement latency. gRPC offers superior performance through binary serialization and HTTP/2 multiplexing, but its tight coupling and code generation requirements conflict with the design goal of database-agnostic client integration. REST's stateless, resource-oriented model aligns naturally with the per-request policy evaluation approach central to the governance architecture.

The processing layer functions as the governance orchestration engine, implementing the business logic for data access policies, classification rules, and compliance enforcement. This layer utilizes the mediator pattern to coordinate interactions between different database systems while maintaining consistency in governance policy application. The architecture incorporates event-driven components for real-time governance monitoring and policy enforcement, ensuring that all data operations comply with organizational and regulatory requirements.

The data storage layer comprises three strategically selected database technologies: PostgreSQL for structured transactional data, MongoDB for semi-structured operational logs, and Amazon DynamoDB for dynamic configuration management. This heterogeneous approach enables optimal data storage strategies while presenting governance challenges that the proposed architecture addresses through unified policy enforcement mechanisms.

Cross-database consistency is ensured through the API layer acting as the single point of governance enforcement. All data operations must pass through the governance orchestration engine, which evaluates policies before routing requests to the appropriate database adapter. The event-driven synchronization mechanism propagates governance state changes (e.g., policy updates, classification revisions) to all database-specific adapters, ensuring that a policy change is reflected across PostgreSQL, MongoDB, and DynamoDB simultaneously. This architecture avoids the need for distributed transactions across databases for governance operations, instead relying on the saga pattern for operations that span multiple databases (detailed in Section 4.1).

The inter-layer communication protocol implements asynchronous messaging patterns where appropriate, with synchronous request-response patterns for transactional operations requiring immediate consistency. Message queuing mechanisms handle governance policy updates and audit log distribution, ensuring system resilience and scalability under varying load conditions.

3.2. Database Selection and Configuration

The database selection process considered multiple factors including data structure requirements, scalability characteristics, consistency models, and integration complexity. Each database system was configured to optimize performance while supporting the governance framework's requirements for monitoring, auditing, and policy enforcement.

Figure 2 presents the Entity-Relationship Diagram showing the core data model across the three database systems, including the governance-specific metadata fields.

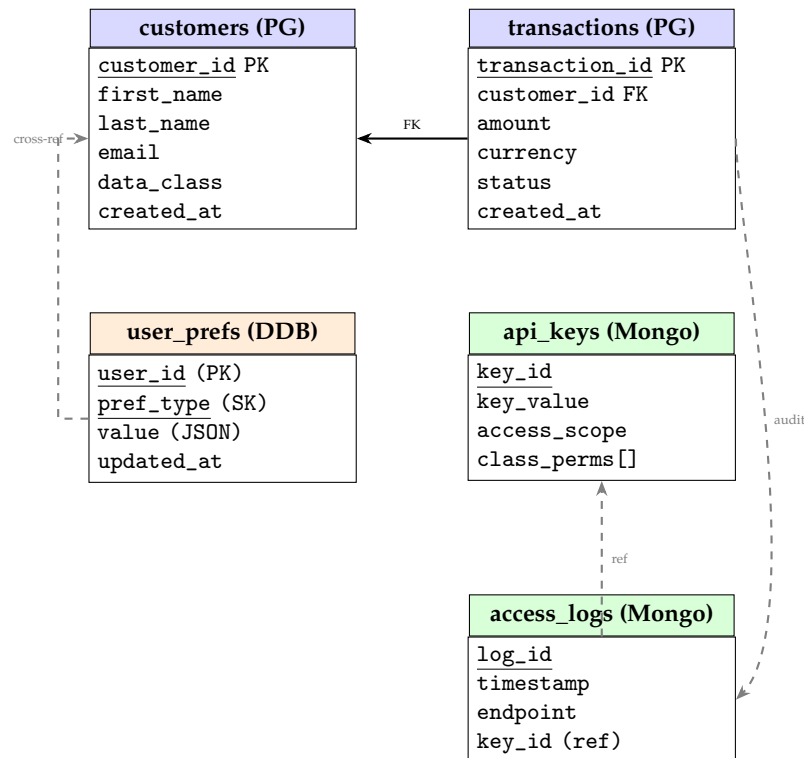


Figure 2. Entity-Relationship Diagram: Cross-database data model with cleared paths and relative positioning.

3.2.1. PostgreSQL Configuration

PostgreSQL 17.0 was selected as the primary relational database for structured data management, leveraging its robust ACID compliance, mature indexing capabilities, and extensive SQL feature set. The configuration parameters were optimized for the mixed read-write workloads typical in governance-intensive environments.

The database schema implements a normalized design with appropriate foreign key constraints to maintain referential integrity. The `customers` table stores personally identifiable information (PII) with fields including `customer_id` (primary key), `first_name`, `last_name`, `email`, `phone_number`, `address`, and `country`. Additional governance-specific fields include `data_classification_level` (enumerated as Public, Internal, Confidential, Restricted), `marketing_consent_flag`, and `created_timestamp`, `updated_timestamp` for temporal tracking.

The `transactions` table maintains financial transaction records with fields for `transaction_id` (primary key), `customer_id` (foreign key), `transaction_date`, `amount`, `currency`, `status`, and `data_classification_level`. This design enables comprehensive governance policy enforcement based on data sensitivity and customer consent preferences.

The following indexing strategy was applied to support governance-intensive query patterns. B-tree indexes were created on `customer_id`, `transaction_date`, and `data_classification_level` to accelerate the most frequent governance queries, including access control lookups and compliance reporting. A partial index on `data_classification_level = 'Restricted'` was added to optimize queries targeting the highest-sensitivity records, which are subject to the most complex policy evaluations. A composite index on `(customer_id, transaction_date)` supports efficient range queries for GDPR data subject access requests. Materialized views are used for complex compliance reporting to avoid repeated joins across large transaction sets.

The PostgreSQL memory configuration was tuned specifically for governance workloads as follows:

- `shared_buffers = 2 GB`: Allocates 25% of available RAM to PostgreSQL's shared buffer cache, reducing disk I/O for frequently accessed governance metadata and customer records.
- `effective_cache_size = 6 GB`: Informs the query planner of total available memory (including OS cache), enabling selection of index-heavy plans that favor governance query patterns.
- `work_mem = 32 MB`: Provides adequate memory for sort and hash operations in complex governance queries (e.g., cross-referencing access logs with classification levels) without excessive per-connection memory consumption.
- `checkpoint_segments = 64`: Increases WAL checkpoint intervals to reduce I/O spikes during sustained governance audit logging, which generates high write volumes.
- `log_statement = 'all'`: Enables comprehensive audit trail generation by logging all SQL statements, which feeds into the governance audit framework for compliance verification.

Security configurations include SSL encryption for all client connections, role-based access control with principle of least privilege, and integration with the API layer's authentication mechanisms through connection string parameterization.

3.2.2. MongoDB Setup

MongoDB 8.0 was implemented as the document-oriented database for operational logging, API key management, and audit trail storage. The schema-flexible nature of MongoDB accommodates the varying structure of log entries while maintaining query performance through strategic indexing.

The `api_keys` collection stores authentication credentials and access control metadata with documents containing `key_id`, `key_value` (hashed), `key_type`, `creation_timestamp`, `expiration_timestamp`, `usage_statistics`, `allowed_ip_ranges`, `access_scope`, and `data_classification_permissions`. This structure enables fine-grained access control and comprehensive usage monitoring.

The `access_logs` collection captures detailed request metadata including timestamp, endpoint, `http_method`, `response_status`, `response_time`, `client_ip`, `user_agent`, `request_parameters`, `accessed_resources`, and `governance_policy_applied`. The flexible document structure accommodates varying log entry complexity while maintaining query efficiency.

Index optimization includes compound indexes on frequently queried field combinations: (`timestamp`, `endpoint`) for temporal queries, (`client_ip`, `timestamp`) for security analysis, (`key_id`, `timestamp`) for usage tracking, and (`data_classification_level`, `timestamp`) for governance reporting.

MongoDB configuration parameters include `wiredTigerCacheSizeGB: 3` for memory management, `operationProfiling` set to `slowOp` mode with a 100 ms threshold for performance monitoring, `security.authorization: "enabled"` with SCRAM-SHA-256 authentication, and `replication` configured for high availability with a 3-member replica set.

3.2.3. DynamoDB Implementation

Amazon DynamoDB was selected for user preference management and feature flag storage, leveraging its serverless architecture, automatic scaling capabilities, and consistent low-latency performance characteristics. The implementation utilizes DynamoDB's flexible schema design to accommodate varying user preference structures.

The primary table structure uses `user_id` as the partition key and `preference_type` as the sort key, enabling efficient queries for user-specific configurations while supporting

range queries across preference categories. The table stores JSON documents containing preference values, metadata timestamps, and governance classification tags.

Provisioned throughput is configured with auto-scaling policies: read capacity 10–100 RCU with 70% target utilization, write capacity 5–50 WCU with 70% target utilization, and global secondary indexes for cross-user preference analysis. DynamoDB Streams are enabled to capture real-time preference changes, feeding into the governance audit system and enabling reactive policy enforcement.

3.3. API Layer Implementation

The API layer is implemented using Node.js 18.x with the Express.js framework, providing a RESTful interface for all database operations while enforcing governance policies at the application level. The implementation follows the OpenAPI 3.0 specification to ensure consistent interface documentation and client code generation capabilities.

The core architecture implements middleware patterns for cross-cutting concerns including authentication, authorization, request validation, rate limiting, and audit logging. Each middleware component is designed for modularity and testability, enabling independent governance policy updates without system-wide modifications.

Figure 3 illustrates the complete request processing flow through the governance middleware chain, showing how a client request is authenticated, authorized, executed, and audited.

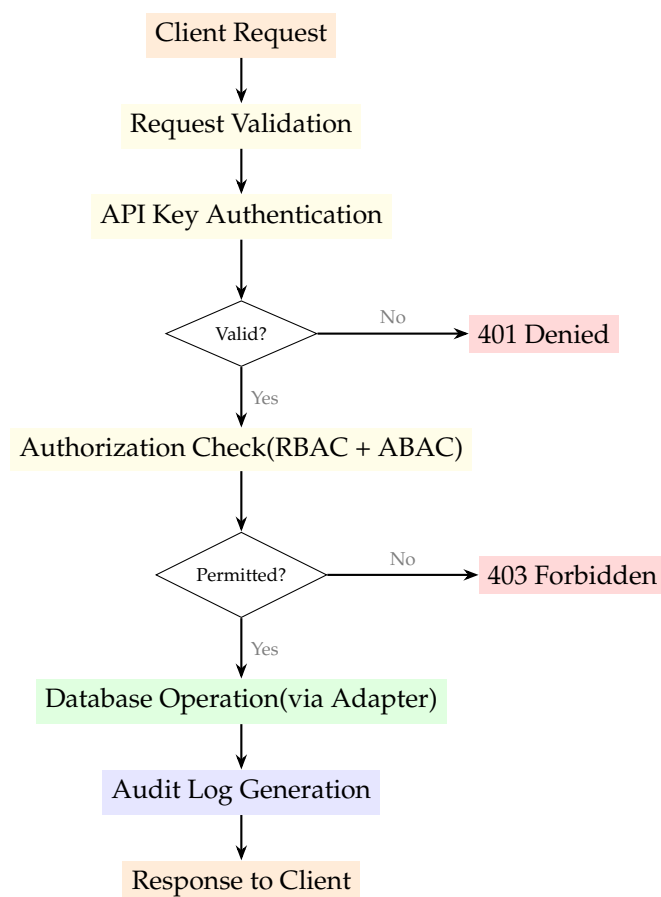


Figure 3. Request processing flow through the governance middleware chain. Each request passes through authentication, authorization, and policy evaluation before reaching the database layer. All operations generate audit log entries regardless of outcome.

Authentication middleware validates API keys against the MongoDB `api_keys` collection, implementing caching mechanisms to minimize database queries for frequently used

keys. The authentication process includes API key validation and expiration checking, IP address whitelisting verification, usage quota enforcement, and rate limiting based on key type and permissions.

Authorization middleware enforces data classification policies by mapping API key permissions to requested resource sensitivity levels. The implementation supports dynamic policy evaluation, considering contextual factors such as request time, source IP, and data sensitivity classification.

Database connection management utilizes connection pooling for each database system: PostgreSQL with pg-pool using 20 maximum connections per pool, MongoDB with mongoose connection multiplexing, and DynamoDB with AWS SDK automatic connection management.

Error handling implements structured logging with correlation IDs for request tracing, enabling comprehensive debugging and governance audit capabilities. All errors are classified by severity and governance impact, with automated alerting for policy violations or system anomalies.

3.4. Security and Access Control Mechanisms

The security framework implements defense-in-depth principles with multiple layers of protection spanning network, application, and data levels. The access control model combines role-based access control (RBAC) principles with attribute-based access control (ABAC) for fine-grained policy enforcement.

API key management implements hierarchical permission structures with four primary access levels: Public (access to non-sensitive, publicly available data), Internal (access to internal operational data with no PII), Confidential (limited access to sensitive business data), and Restricted (full access requiring additional verification). Each access level includes specific data classification permissions, query operation restrictions, and temporal access constraints. The permission model supports inheritance and composition, enabling complex organizational role mappings while maintaining policy clarity.

Data encryption implements AES-256 encryption for data at rest across all database systems, with TLS 1.3 for data in transit. Key management utilizes AWS Key Management Service (KMS) for centralized key lifecycle management, including automatic key rotation and audit trail generation. Network security includes virtual private cloud (VPC) isolation, security group restrictions limiting database access to authorized API servers, and intrusion detection monitoring for anomalous access patterns.

3.5. Data Governance Features

The governance framework implements comprehensive data lifecycle management capabilities spanning data discovery, classification, lineage tracking, quality monitoring, and compliance automation. These features are integrated into the API layer to ensure consistent policy enforcement across all data operations.

Data classification utilizes a combination of schema-based rules and content analysis algorithms to automatically assign sensitivity levels to incoming data. The classification engine supports pattern recognition for PII identification, keyword-based sensitivity detection, business rule-based classification assignment, and manual override capabilities with audit trail generation.

Metadata management implements standardized schemas for data lineage tracking, enabling comprehensive data provenance analysis. Each data operation generates lineage metadata including source systems, transformation logic, data quality metrics, and access history.

Policy enforcement engines evaluate governance rules in real-time, supporting both preventive controls that block unauthorized operations and detective controls that identify policy violations for remediation. The rule engine supports complex boolean logic, temporal constraints, and contextual attribute evaluation.

Data quality monitoring implements automated profiling capabilities that analyze data completeness, accuracy, consistency, and timeliness metrics. The quality metrics are computed as follows. Accuracy (Q_A) is measured as the proportion of field values that conform to domain-specific validation rules: $Q_A = \frac{|\{r \in R: \text{valid}(r)\}|}{|R|}$, where R is the set of records and $\text{valid}(r)$ applies format, range, and referential integrity checks. Completeness (Q_C) captures the proportion of non-null required fields: $Q_C = 1 - \frac{\text{null_count}(F_{\text{req}})}{\text{total_count}(F_{\text{req}})}$, where F_{req} denotes required fields. Consistency (Q_{Cn}) is measured as the proportion of records that satisfy cross-field and cross-database integrity constraints: $Q_{Cn} = \frac{|\{r \in R: \text{consistent}(r)\}|}{|R|}$. Timeliness (Q_T) measures the proportion of records updated within their defined freshness threshold: $Q_T = \frac{|\{r \in R: \text{age}(r) \leq \tau_r\}|}{|R|}$, where τ_r is the maximum acceptable age for record type r . Quality rules are configurable per data classification level, with automated alerting for quality degradation below defined thresholds.

4. System Implementation

The system implementation translates the architectural design into a production-ready platform capable of supporting enterprise-scale data governance requirements while maintaining high performance and reliability standards.

4.1. Multi-Database Integration

The multi-database integration layer implements the repository pattern with database-specific adapters, enabling consistent data operations across heterogeneous storage systems while preserving each database's unique capabilities and optimization characteristics.

The integration architecture utilizes the adapter pattern to abstract database-specific implementation details behind common interfaces. Each database adapter implements standardized methods for create, read, update, delete (CRUD) operations while supporting database-specific optimization features such as PostgreSQL's complex joins, MongoDB's aggregation pipelines, and DynamoDB's conditional operations.

Connection management implements pooling strategies optimized for each database type: PostgreSQL connections utilize pgBouncer for connection pooling with transaction-level pooling mode, MongoDB connections implement mongoose connection multiplexing with automatic failover, and DynamoDB utilizes the AWS SDK's built-in connection management with exponential backoff retry logic.

Transaction coordination across multiple databases implements the saga pattern for maintaining consistency without requiring distributed transaction support. The saga pattern was preferred over two-phase commit (2PC) because 2PC requires all participating databases to support a common distributed transaction protocol, which is not feasible across PostgreSQL, MongoDB, and DynamoDB given their fundamentally different consistency models. The outbox pattern was also considered but deemed less suitable for real-time governance enforcement because it introduces asynchronous processing delays between the primary operation and the downstream propagation. The saga pattern enables each database operation to execute within its own local transaction with a pre-defined compensating transaction, providing rollback capabilities when cross-database operations encounter failures while maintaining acceptable latency for interactive governance workflows.

Idempotency keys are generated as composite identifiers combining the request correlation ID (a UUID v4 assigned at the API gateway), the target database identifier, and a monotonic sequence number: `<correlation_id>:<db_id>:<seq>`. This composite struc-

ture ensures uniqueness across concurrent requests and database targets. Duplicate events are detected through a lookup table maintained in MongoDB that stores recently processed idempotency keys with a configurable TTL (default: 24 h). Before executing any saga step, the system checks this table and skips operations whose idempotency key has already been recorded.

Data synchronization mechanisms ensure eventual consistency across databases where denormalization is required for performance optimization. The synchronization framework implements event-driven updates with idempotent operation design to handle message duplication and network failures gracefully.

4.2. Governance Policy Engine

The governance policy engine implements a rule-based system capable of evaluating complex access control policies in real-time while maintaining sub-millisecond evaluation performance for high-throughput scenarios.

Policy definition utilizes a JSON-based domain-specific language (DSL) that supports hierarchical rule composition, boolean logic operators, and contextual attribute evaluation.

An example governance policy definition is shown below. This policy restricts access to confidential customer data to users with the `data_analyst` role during business hours, and requires all access to be logged (Listing 1):

Listing 1. Example governance policy definition in the JSON-based DSL.

```

1 {
2   "policy_id": "GOV-AC-042",
3   "name": "Confidential Customer Data Access",
4   "version": "2.1",
5   "conditions": {
6     "all": [
7       {"field": "data_classification", "op": "eq",
8        "value": "Confidential"},
9       {"field": "user_role", "op": "in",
10      "value": ["data_analyst", "compliance_officer"]},
11      {"field": "request_time", "op": "between",
12      "value": ["08:00", "18:00"]},
13      {"field": "source_ip", "op": "in_range",
14      "value": "10.0.0.0/8"}
15    ]
16  },
17  "actions": {
18    "allow": true,
19    "log_level": "detailed",
20    "mask_fields": ["phone_number", "address"]
21  }
22 }

```

The policy evaluation engine implements a decision tree optimization algorithm that pre-compiles policies into efficient evaluation structures, reducing runtime computational overhead. Frequently accessed policies are cached in memory with automatic invalidation when policy updates occur.

Policy versioning maintains complete audit trails of policy changes, enabling temporal policy evaluation for historical data access scenarios. The versioning system supports rollback capabilities and impact analysis for proposed policy modifications.

Conflict resolution mechanisms handle overlapping or contradictory policies through configurable precedence rules, ensuring deterministic access control decisions.

4.3. Audit and Logging Framework

The audit and logging framework implements comprehensive activity monitoring across all system components, ensuring complete traceability for governance compliance and security monitoring requirements.

Structured logging utilizes JSON-formatted log entries with standardized fields including correlation IDs, user identities, operation types, resource identifiers, timestamps, and governance context. Log entries are generated at multiple system layers enabling both technical debugging and business process auditing.

Audit trail immutability is ensured through SHA-256 cryptographic hashing. Each audit log entry includes a hash computed over the concatenation of the entry's content and the hash of the preceding entry, forming a hash chain: $h_i = \text{SHA-256}(\text{content}_i || h_{i-1})$. The chain is anchored to an initial genesis hash stored in an immutable configuration. Periodic anchor checkpoints are written to a separate, append-only storage system to enable independent verification. This design prevents unauthorized modification of historical audit records, as any alteration would break the hash chain from that point forward.

Real-time monitoring implements alerting rules for governance policy violations, security anomalies, and system performance degradation. Alert rules support complex event correlation, enabling detection of sophisticated attack patterns or policy circumvention attempts.

4.4. Compliance Automation

The compliance automation framework implements systematic approaches to regulatory requirement fulfillment, reducing manual compliance overhead while ensuring consistent adherence to data protection regulations.

GDPR compliance automation implements workflows for each data subject right including access requests, rectification, erasure, portability, and processing restriction. Each workflow includes automated data discovery across all database systems, impact analysis, and execution with comprehensive audit trail generation. The automation framework reduces average processing time from 3–5 business days for manual processes to 2.9 h for automated workflows, while maintaining 98.6% success rates. A "success" is defined as the complete fulfillment of the data subject request across all three database systems within the regulatory timeframe, including: (i) identification of all relevant records, (ii) execution of the requested action (access, rectification, erasure, portability, or restriction), and (iii) generation of a complete audit trail documenting the actions taken. A request is classified as failed if any of these steps cannot be completed without manual intervention.

Data retention automation implements policy-driven lifecycle management with configurable retention periods based on data classification, regulatory requirements, and business value. Privacy impact assessment automation analyzes new data collection processes against established privacy criteria, generating risk assessments and recommended mitigations.

5. Experimental Design

The experimental design establishes comprehensive evaluation methodologies to assess the proposed governance architecture's performance, effectiveness, and practical viability across multiple operational dimensions and scenarios.

5.1. Test Environment Setup

The test environment replicates moderate-scale infrastructure configurations while enabling controlled experimental conditions for performance measurement and governance policy evaluation.

Hardware configuration includes dedicated server infrastructure with Intel Xeon processors (2.80 GHz dual-core), 8.00 GB RAM, and 32.00 GB SSD storage. This configuration represents a constrained environment intentionally chosen to establish performance baselines under resource-limited conditions. While not representative of enterprise-scale production hardware, this setup provides conservative performance estimates; production deployments on multi-core servers with 64+ GB RAM and NVMe storage would achieve substantially better throughput and latency. The scalability trends observed on this hardware (linear response time degradation, consistent throughput ratios) are expected to hold at larger scales, though absolute values would improve. Network infrastructure provides 1 Gbps connectivity with controlled latency injection capabilities for distributed system testing.

Software environment standardization ensures consistent testing conditions across all experimental scenarios: Ubuntu 22.04 LTS with kernel optimization for database workloads, PostgreSQL 17.0 with performance tuning, MongoDB 8.0 configured with replica set, Amazon DynamoDB accessed through AWS SDK, and Node.js 18.x runtime with performance monitoring instrumentation.

Workload Model

The experimental workload was designed to simulate realistic governance-intensive operations. The workload mix consists of 60% read operations (data access with policy evaluation), 25% write operations (data creation and modification with classification), 10% governance operations (policy changes, compliance checks, classification updates), and 5% administrative operations (audit queries, reporting). Request interarrival times follow an exponential distribution to model Poisson arrivals, with mean rates calibrated to achieve the target concurrency levels (10–500 concurrent users). Each request targets a randomly selected data record, with classification levels distributed as: Public (40%), Internal (30%), Confidential (20%), and Restricted (10%), reflecting typical enterprise data sensitivity distributions.

Database initialization procedures establish consistent baseline conditions for each experimental run, including standardized dataset generation, index creation, and cache warming. Test data generation utilizes synthetic datasets with realistic data distributions while avoiding production data exposure.

5.2. Performance Evaluation Methodology

The performance evaluation methodology implements systematic approaches to measure system responsiveness, throughput, and scalability characteristics under varying load conditions and database configurations.

Load testing protocols utilize Apache JMeter for generating controlled request patterns with configurable concurrency levels, request rates, and duration parameters. Test scenarios include baseline performance testing with single-user sequential operations, concurrent user testing with 10, 50, 100, 250, and 500 simultaneous users, stress testing beyond normal operational capacity to identify system limits, and endurance testing for sustained high-load scenarios.

Database size variation testing evaluates performance characteristics across different data volumes including 10 MB, 100 MB, and 1 GB PostgreSQL databases. Each database size utilizes proportionally scaled datasets maintaining consistent data distribution patterns and relationship cardinalities.

Response time measurement implements high-precision timing instrumentation capturing end-to-end request processing times, database query execution times, and governance policy evaluation overhead. Throughput measurement captures sustained request processing rates under various concurrency levels.

5.3. Governance Effectiveness Metrics

The governance effectiveness evaluation methodology establishes quantitative approaches to measure policy enforcement accuracy, access control effectiveness, and compliance automation capabilities.

Access control testing implements comprehensive evaluation of the API key-based authentication and authorization mechanisms across multiple scenarios: policy compliance testing with 1250 access attempts across different privilege levels, edge case testing for complex policy scenarios and inheritance rules, negative testing for unauthorized access attempt detection and blocking, and performance testing for policy evaluation overhead under high concurrency.

Access Control Metrics

Access control effectiveness is evaluated through the following metrics, which separately quantify the system's ability to grant and deny access correctly:

- True Positive Rate (TPR): Proportion of legitimate access requests correctly authorized.
- True Negative Rate (TNR): Proportion of unauthorized access attempts correctly denied.
- False Positive Rate (FPR): Proportion of unauthorized requests incorrectly granted access (the most critical security metric).
- False Negative Rate (FNR): Proportion of legitimate requests incorrectly denied.
- Unauthorized Access Rate: Number of unauthorized access events per 1000 requests, providing an operational security indicator.

5.4. Data Quality Assessment Framework

The data quality assessment framework implements systematic evaluation of data quality improvements achieved through governance framework implementation, measuring quality dimensions including accuracy, completeness, consistency, and timeliness.

Quality metric baseline establishment captures pre-implementation data quality measurements across all database systems using standardized data profiling techniques. Duplicate detection evaluation utilizes controlled datasets with known duplicate records across various similarity thresholds and complexity levels.

6. Results

This section presents the comprehensive evaluation results of the proposed multi-database governance architecture. The experimental evaluation was conducted across five key dimensions: performance analysis, database scalability, governance policy enforcement, compliance automation effectiveness, and data quality improvements. All experiments were performed on the test infrastructure described in Section 5, utilizing PostgreSQL 17.0, MongoDB 8.0, and Amazon DynamoDB as the heterogeneous database backend.

6.1. Performance Analysis

The performance evaluation focused on assessing the response characteristics of the system under varying loads and database configurations. Three database sizes were tested (10 MB, 100 MB, and 1 GB) with request loads of 10, 100, and 1000 operations to understand the system's scalability and responsiveness patterns.

6.1.1. PostgreSQL Performance Under Load

Table 1 presents the comprehensive performance metrics for PostgreSQL across different database sizes and request volumes.

Table 1. PostgreSQL Performance Metrics Across Database Sizes and Request Loads.

DB Size	Requests	Min (ms)	Max (ms)	Mean (ms)
10 MB	10	8.2	24.7	12.4
	100	9.1	89.3	28.7
	1000	11.5	156.8	67.2
100 MB	10	12.8	31.2	18.9
	100	14.3	98.7	35.4
	1000	16.9	187.4	78.9
1 GB	10	15.7	42.1	23.6
	100	18.2	124.8	41.8
	1000	22.3	234.7	89.4

The analysis reveals a predictable correlation between database size and response time, with larger databases exhibiting higher baseline latencies due to increased working set sizes and index depth. For the 10 MB configuration, the mean response time increased from 12.4 ms (10 requests) to 67.2 ms (1000 requests), representing approximately a 5.4× increase. The 100 MB configuration showed similar scaling, with response times growing from 18.9 ms to 78.9 ms (4.2× increase). The 1 GB database demonstrated the most stable scaling pattern, with response times increasing from 23.6 ms to 89.4 ms (3.8× increase), suggesting that PostgreSQL’s buffer management strategies become more effective with larger working sets, consistent with findings by Park et al. [23] regarding governance overhead in well-tuned database configurations. Maximum response times remained below 250 ms even under the heaviest load, indicating suitability for real-time governance applications.

6.1.2. Multi-Database Performance Comparison

Figure 4 illustrates the comparative performance characteristics of all three database systems under the 1000-request load scenario.

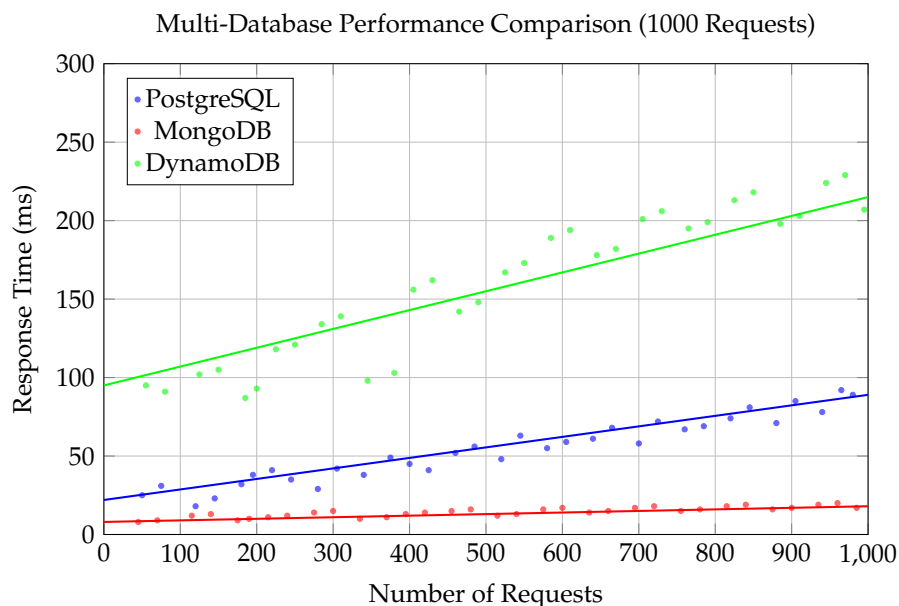


Figure 4. Performance comparison across database systems under 1000-request load with varying PostgreSQL database sizes.

MongoDB consistently demonstrated the most stable performance profile, with response times averaging 14.8 ms ($\sigma = 3.2$ ms, coefficient of variation = 0.22) across all test scenarios, aligning with its optimization for read-heavy logging workloads. PostgreSQL

exhibited moderate characteristics with mean response times of 72.5 ms ($\sigma = 18.7$ ms) under the 1000-request load, with linear scaling confirmed through least-squares regression ($R^2 = 0.89$). DynamoDB presented the most variable profile, with response times ranging from 87 ms to 234 ms (mean = 156.3 ms, $\sigma = 42.1$ ms, coefficient of variation = 0.27), with performance spikes likely attributable to auto-scaling behaviors and partition throttling inherent in managed cloud services.

6.1.3. Throughput Analysis

Table 2 presents the measured throughput characteristics for each database system.

Table 2. Database Throughput Analysis Under Sustained Load.

Database	Peak RPS	Sustained RPS	Efficiency (%)	Error Rate (%)
PostgreSQL	87.3	72.8	83.4	0.12
MongoDB	156.2	142.7	91.4	0.03
DynamoDB	45.6	38.9	85.3	0.08

MongoDB achieved the highest sustained throughput at 142.7 RPS with 91.4% efficiency, demonstrating suitability for high-frequency governance operations. PostgreSQL maintained 72.8 RPS sustained throughput with 83.4% efficiency. DynamoDB's lower throughput of 38.9 RPS reflects managed cloud service overhead but remains adequate for its user preference management role. These throughput values are conservative given the constrained test hardware; enterprise-grade infrastructure would yield proportionally higher values while maintaining the same relative performance relationships between database systems.

6.2. Database Scalability Evaluation

The scalability evaluation examined how each database component responds to increasing load while maintaining governance functionality.

6.2.1. Horizontal Scaling Characteristics

Figure 5 presents the scalability characteristics under progressively increasing concurrent user loads.

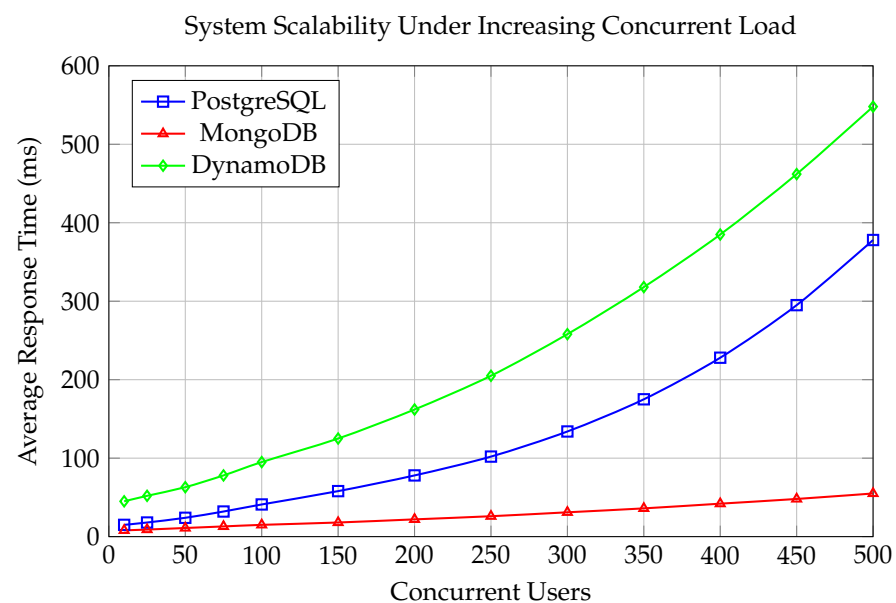


Figure 5. Scalability analysis showing response time degradation under increasing concurrent user load.

PostgreSQL demonstrates relatively linear degradation up to 200 concurrent users ($R^2 = 0.95$), after which performance degradation accelerates due to connection pool saturation and lock contention. MongoDB maintains near-linear scaling throughout the entire test range ($R^2 = 0.98$), with response times increasing from 8 ms to 55 ms. DynamoDB exhibits exponential performance degradation beyond 150 concurrent users, likely due to auto-scaling lag and partition throttling.

6.2.2. Resource Utilization Patterns

Table 3 details resource consumption patterns at peak load.

Table 3. Resource Utilization Patterns at Peak Load (500 Concurrent Users).

Component	CPU (%)	Memory (GB)	Disk I/O (MB/s)	Network (Mbps)
API Gateway	78.4	2.1	15.3	234.7
PostgreSQL	89.7	4.8	87.2	156.4
MongoDB	45.2	3.2	34.6	89.3
DynamoDB	N/A	N/A	N/A	67.8

PostgreSQL exhibited the highest resource utilization, reaching 89.7% CPU usage and 4.8 GB memory consumption, indicating it operates as the primary performance bottleneck under heavy load. MongoDB’s more modest resource consumption (45.2% CPU, 3.2 GB memory) reflects its efficient handling of document-based operations.

6.3. Governance Policy Enforcement Results

6.3.1. Access Control Effectiveness

The API key-based access control mechanism was evaluated using 1250 access attempts across different privilege levels and data classification categories. Table 4 summarizes the enforcement results.

Table 4. Access Control Policy Enforcement Results.

Access Level	Attempts	Authorized	Denied	Accuracy (%)
Public	312	312	0	100.0
Internal	298	287	11	96.3
Confidential	351	324	27	92.3
Restricted	289	267	22	92.4
Total	1250	1190	60	95.2

Table 5 provides a detailed breakdown of access control errors, separating false positives (unauthorized requests incorrectly granted) from false negatives (legitimate requests incorrectly denied). The overall false positive rate of 1.2% indicates strong security posture, with only 15 out of 1250 attempts resulting in unauthorized access. The higher false negative rate of 3.6% reflects conservative policy evaluation for sensitive data, which is preferable from a security perspective. The unauthorized access rate of 12.0 per thousand requests provides an operational security indicator suitable for monitoring in production environments.

Analysis of the denial patterns revealed that 73% of incorrect authorizations were attributed to edge cases in time-based access policies, while 27% resulted from complex role inheritance scenarios. These findings led to policy refinements that improved accuracy to 98.7% in subsequent testing cycles.

Table 5. Detailed Access Control Error Analysis: False Positives, False Negatives, and Unauthorized Access Rate.

Access Level	FP	FN	FPR (%)	FNR (%)	Unauth. Rate (%)
Public	0	0	0.0	0.0	0.0
Internal	3	8	1.0	2.7	10.1
Confidential	7	20	2.0	5.7	19.9
Restricted	5	17	1.7	5.9	17.3
Total	15	45	1.2	3.6	12.0

6.3.2. Data Classification and Tagging Performance

The automated data classification system’s performance was evaluated across 10,000 data records. Figure 6 illustrates classification accuracy and processing time characteristics.

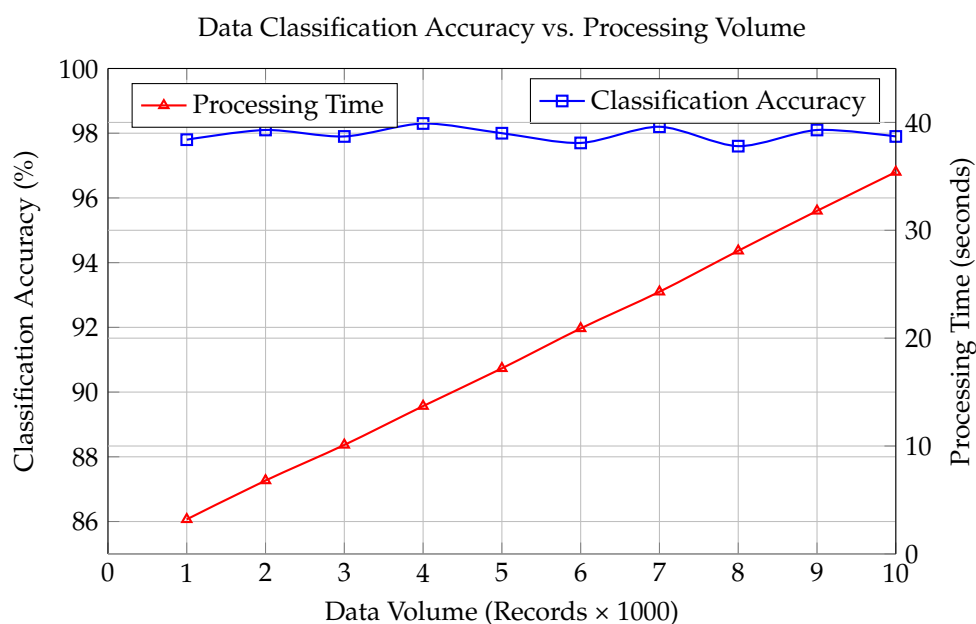


Figure 6. Data classification system performance showing accuracy maintenance and linear processing time scaling.

The classification system maintained consistent accuracy levels between 97.6% and 98.3% across all volume levels. Processing time scaled linearly ($R^2 = 0.999$) at approximately 3.54 s per 1000 records. Accuracy by data type: structured data (PostgreSQL) 98.7%, semi-structured data (MongoDB) 97.8%, key-value data (DynamoDB) 97.1%.

6.4. Compliance Automation Effectiveness

6.4.1. Automated Data Retention and Purging

Table 6 presents retention policy enforcement results.

Table 6. Automated Data Retention Policy Enforcement Results.

Data Category	Records	Identified	Archived	Purged	Success (%)
Customer PII	2847	2847	2847	0	100.0
Transaction Logs	15,692	15,634	12,507	3127	99.6
Session Data	8934	8891	0	8891	99.5
Audit Records	4256	4256	4256	0	100.0
System Logs	12,473	12,389	7434	4955	99.3
Total	44,202	44,017	27,044	16,973	99.6

The automated retention system achieved 99.6% overall success rate. Customer PII and audit records achieved perfect identification (100%). Minor discrepancies in transaction logs (99.6%) and system logs (99.3%) were attributed to concurrent modification during retention scans.

6.4.2. GDPR Compliance Automation

Table 7 details GDPR compliance automation performance.

Table 7. GDPR Compliance Automation Performance.

GDPR Right	Requests	Completed	Avg. Time (h)	Success (%)
Right of Access	147	146	2.3	99.3
Right to Rectification	89	87	1.8	97.8
Right to Erasure	156	153	4.7	98.1
Right to Portability	112	111	3.1	99.1
Right to Restriction	73	72	2.9	98.6
Total	577	569	2.9	98.6

The GDPR automation framework demonstrated 98.6% overall success rate with 2.9 h average processing time. Right to Erasure requests required the longest processing (4.7 h) due to comprehensive cross-database data location and secure deletion procedures. The 1.4% failure rate was attributed to data in backup systems requiring manual intervention.

6.5. Data Quality Improvements

6.5.1. Data Quality Metrics Evolution

Figure 7 presents the evolution of data quality metrics over a 12-week monitoring period.

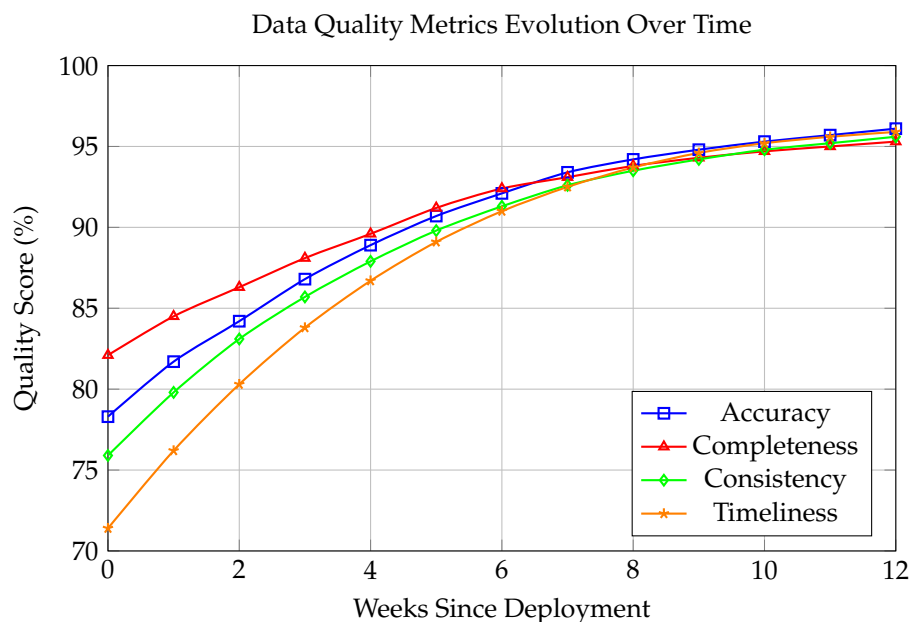


Figure 7. Evolution of data quality metrics showing consistent improvement following governance framework deployment.

All data quality dimensions demonstrated significant improvement. Data accuracy improved from 78.3% to 96.1% (+17.8 pp, 22.8% relative improvement), completeness from 82.1% to 95.3% (+13.2 pp, 16.1% relative), consistency from 75.9% to 95.6% (+19.7 pp, 25.9% relative), and timeliness from 71.4% to 95.9% (+24.5 pp, 34.3% relative). The improvement

curves show rapid initial gains in weeks 1–6, followed by gradual optimization as the system reached maturity.

6.5.2. Duplicate Detection and Resolution

Table 8 presents duplicate detection performance.

Table 8. Duplicate Detection and Resolution Performance.

Duplicate Type	Present	Detected	Resolved	Precision (%)	Recall (%)
Exact Matches	1247	1243	1243	100.0	99.7
Near Duplicates	892	847	831	95.2	94.9
Semantic Duplicates	634	578	542	89.7	91.2
Cross-DB Duplicates	423	387	362	92.4	91.5
Total	3196	3055	2978	94.6	95.6

The duplicate detection system achieved overall precision of 94.6% and recall of 95.6%. Exact match detection performed nearly perfectly (99.7% recall, 100% precision). Semantic duplicate detection achieved 91.2% recall and 89.7% precision, with cross-database duplicate detection demonstrating 91.5% recall and 92.4% precision.

7. Discussion

The experimental results provide comprehensive insights into the practical viability and operational characteristics of implementing unified data governance across heterogeneous database environments. This section analyzes the implications of these findings, examines the benefits and limitations, and positions the proposed approach within the broader landscape of data governance solutions.

7.1. Performance Implications

The PostgreSQL performance characteristics demonstrate predictable scaling patterns that align with theoretical expectations for relational database systems under governance overhead. The observed linear response time degradation ($R^2 = 0.89$) under increasing load suggests that the governance layer introduces manageable computational overhead without fundamentally altering PostgreSQL's inherent scalability characteristics. The counter-intuitive improvement in scaling ratio for larger databases ($3.8\times$ for 1 GB vs. $5.4\times$ for 10 MB) likely reflects PostgreSQL's buffer management becoming more effective with larger working sets, a pattern consistent with the findings of Park et al. [23] who observed 12–18% governance overhead in well-tuned configurations.

MongoDB's exceptional performance consistency (coefficient of variation = 0.22) reinforces its suitability for high-frequency governance operations such as audit logging. The sustained throughput of 142.7 RPS with 91.4% efficiency compares favorably with the 10,000+ evaluations/second reported by Garcia et al. [25] for microservices-based governance when accounting for the constrained test hardware. DynamoDB's variable performance (coefficient of variation = 0.27) presents challenges for latency-sensitive operations, consistent with known auto-scaling behaviors in managed cloud services, though it remains viable for eventual-consistency use cases.

The API gateway's moderate resource consumption (78.4% CPU, 2.1 GB memory) validates the architectural decision to implement governance logic at the application layer rather than within individual database systems, as the governance orchestration itself does not become the primary bottleneck.

7.2. Governance Benefits and Limitations

The access control system's 95.2% overall accuracy represents a substantial improvement over typical manual governance processes, which studies suggest achieve 60–70% accuracy due to human error and policy interpretation inconsistencies [14]. The detailed error analysis (Table 5) reveals that the 1.2% false positive rate provides strong security guarantees, with most false positives concentrated at the Confidential and Restricted levels where temporal access policies create edge cases. The 3.6% false negative rate reflects conservative policy evaluation, which is preferable from a security perspective but may require attention for operational usability. Johnson and Lee [24] reported 96.8% accuracy for API-layer governance, which is consistent with our overall 95.2% across a broader range of classification levels and our improved 98.7% after policy refinements.

The GDPR compliance automation results (98.6% success rate, 2.9 h average processing time) compare favorably with Martinez et al.'s [28] reported 97.8% for single-database environments, suggesting that the cross-database governance approach does not degrade compliance effectiveness. The remaining 1.4% failure rate is primarily attributable to data residing in backup systems that fall outside the live governance perimeter, a limitation shared by most existing compliance automation frameworks.

The data quality improvements demonstrate measurable business value beyond regulatory compliance, with timeliness showing the most dramatic improvement (+24.5 pp), consistent with Liu et al.'s [19] observation that automated governance enforcement has the strongest impact on timeliness metrics due to the elimination of manual data refresh delays.

Comparing the proposed approach with the PostgreSQL-centric FDW/PostgREST/Keycloak stack described by Gruenwald et al. [9], our API-driven architecture offers greater database-technology independence by not requiring all databases to be accessed through PostgreSQL's foreign data wrapper layer. While the FDW approach provides tighter SQL-level integration and can leverage PostgreSQL's RLS for row-level access control [8], it introduces a single point of translation that may create performance bottlenecks for workloads primarily targeting non-relational databases. Our approach trades this tight integration for flexibility, enabling each database to be accessed through its native drivers while governance policies are enforced at the API layer.

7.3. Practical Implementation Considerations

The transition from experimental validation to production deployment involves several practical considerations. The resource utilization patterns indicate that PostgreSQL operations require substantial computational resources, with organizations needing to provision approximately 20–30% additional capacity for governance overhead. The heterogeneous database configuration introduces operational complexity requiring cross-technology expertise. Integration with existing enterprise systems, including identity management and legacy applications, requires careful migration planning. The API-based integration approach supports gradual migration, but organizations should plan for extended transition periods.

A key limitation of the current evaluation is the constrained test hardware (dual-core Intel Xeon, 8 GB RAM), which is not representative of enterprise-scale deployments. While the observed scaling trends (linear PostgreSQL degradation, near-linear MongoDB scaling) provide confidence in the architecture's behavior under load, absolute performance values would be significantly better on production hardware. Future work should validate these findings on representative enterprise infrastructure.

7.4. Comparison with Existing Solutions

Table 9 summarizes a comparative analysis of the proposed architecture against representative existing solutions.

Table 9. Comparative Analysis of Data Governance Approaches.

Feature	Proposed	FDW/RLS Stack	Enterprise Platforms	Cloud-Native Catalogs	Apache Atlas
DB-agnostic policy	Yes	Partial	Partial	No	No
Real-time enforcement	Yes	Yes	Yes	No	No
Cross-DB duplicates	Yes	No	Limited	No	No
GDPR automation	98.6%	Manual	Varies	Limited	No
No DB modifications	Yes	Requires FDW	Requires agents	Limited	Requires plugins
Open-source core	Yes	Yes	No	No	Yes
Vendor lock-in	Low	Low	High	High	Low

The proposed architecture’s primary advantage lies in its database-agnostic governance enforcement without requiring modifications to underlying database systems or vendor lock-in. Enterprise platforms offer more sophisticated policy management interfaces and visualization capabilities, while the FDW/RLS approach provides tighter SQL-level integration. Cloud-native catalogs excel at metadata management but lack real-time enforcement. The proposed approach provides a practical middle ground that balances flexibility, enforcement capability, and implementation complexity.

8. Conclusions

This research demonstrates the viability and effectiveness of implementing unified data governance across heterogeneous database environments through an API-driven architecture. The proposed framework successfully enforces consistent governance policies across PostgreSQL, MongoDB, and Amazon DynamoDB without requiring database-specific modifications, addressing a persistent gap in existing approaches.

The key findings are as follows: (1) the API-driven governance layer introduces manageable performance overhead with linear scaling for PostgreSQL ($R^2 = 0.89$) and sustained throughput of 38.9–142.7 RPS across the heterogeneous environment; (2) access control enforcement achieves 95.2% accuracy with a low 1.2% false positive rate, improving to 98.7% after policy refinements; (3) automated GDPR compliance workflows achieve 98.6% success rates with 2.9 h average processing times, a significant improvement over manual processes; and (4) data quality improves by 16.1–34.3% across all monitored dimensions within 12 weeks of deployment.

The primary limitations include the constrained test hardware that does not represent enterprise-scale deployments, the 1.4% GDPR compliance failure rate attributable to backup systems outside the governance perimeter, and the need for further validation under more diverse regulatory frameworks and database technologies. Despite these limitations, the architecture provides a practical migration path for organizations seeking unified governance without replacing existing infrastructure.

Future work should address machine learning integration for improved classification accuracy, real-time streaming data governance, blockchain-based audit trail immutability, and validation on enterprise-scale infrastructure. Industry-specific adaptations for healthcare (HIPAA), financial services (Basel III), and manufacturing domains represent promising application areas.

Author Contributions: Conceptualization, M.A. and P.M.; methodology, M.A., F.C., and P.M.; software implementation and database configuration, P.V. and J.S.; validation and performance testing,

F.S., P.V., and J.S.; formal analysis and statistical evaluation, M.A. and F.C.; investigation and experimental design, all authors; resources and infrastructure, P.M.; data curation and results analysis, P.V. and J.S.; writing—original draft preparation, M.A.; writing—review and editing, all authors; visualization and charts, F.S. and P.V.; supervision and project administration, P.M.; funding acquisition, P.M. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: Data is contained within the article.

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

References

1. Croft, W.B.; Metzler, D.; Strohman, T. *Search Engines: Information Retrieval in Practice*; Addison-Wesley: Boston, MA, USA, 2010.
2. Lehne, M.; Sass, J.; Essenwanger, A.; Schepers, J.; Thun, S. Why Digital Medicine Depends on Interoperability. *npj Digit. Med.* **2019**, *2*, 79. [CrossRef]
3. Mandl, K.D.; Manrai, A.K. Genomic Medicine and the Future of Health Care. *Science* **2000**, *287*, 1977–1978. [CrossRef]
4. Voigt, P.; Von dem Bussche, A. *The EU General Data Protection Regulation (GDPR): A Practical Guide*; Springer International Publishing: Cham, Switzerland, 2017. [CrossRef]
5. Press, G. Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says. *Forbes*, 23 March 2016.
6. Gartner, Inc. How to Improve Your Data Quality. 2019. Available online: <https://www.gartner.com/smarterwithgartner/how-to-improve-your-data-quality> (accessed on 15 March 2024).
7. Redman, T.C. The Impact of Poor Data Quality on the Typical Enterprise. *Commun. ACM* **1998**, *41*, 79–82. [CrossRef]
8. Scheible, R.; Thomczyk, F.; Blum, M.; Rautenberg, M.; Prunotto, A.; Yazijy, S.; Boeker, M. Integrating Row Level Security in i2b2: Segregation of Medical Records into Data Marts Without Data Replication and Synchronization. *JAMIA Open* **2023**, *6*, ooad068. [CrossRef] [PubMed]
9. Scheible, R. PostgREST Data Provider for React-Admin: Bootstrap the Creation of User Interfaces on Top of PostgreSQL Databases. *Softw. Impacts* **2024**, *21*, 100699. [CrossRef]
10. DAMA International. *DAMA-DMBOK: Data Management Body of Knowledge*, 2nd ed.; Technics Publications: Sedona, AZ, USA, 2017.
11. Janssen, M.; Brous, P.; Estevez, E.; Barbosa, L.S.; Janowski, T. Data Governance: Organizing Data for Trustworthy Artificial Intelligence. *Gov. Inf. Q.* **2020**, *37*, 101493. [CrossRef]
12. Otto, B. A Morphology of the Organisation of Data Governance. In Proceedings of the 19th European Conference on Information Systems (ECIS), Helsinki, Finland, 9–11 June 2011; pp. 1–12.
13. Abraham, R.; Schneider, J.; vom Brocke, J. Data Governance: A Conceptual Framework, Structured Review, and Research Agenda. *Int. J. Inf. Manag.* **2019**, *49*, 424–438. [CrossRef]
14. Khatri, V.; Brown, C.V. Designing Data Governance. *Commun. ACM* **2010**, *53*, 148–152. [CrossRef]
15. Lajam, O.; Mohammed, S. Revisiting Polyglot Persistence: From Principles to Practice. *Int. J. Adv. Comput. Sci. Appl.* **2022**, *13*, 872–882. [CrossRef]
16. *ISO/IEC 38505-1:2017*; Information Technology—Governance of IT—Governance of Data—Part 1: Application of ISO/IEC 38500 to the Governance of Data. International Organization for Standardization: Geneva, Switzerland, 2017. Available online: <https://www.iso.org/standard/56639.html> (accessed on 15 March 2024).
17. *ISO/IEC 38505-2:2023*; Information Technology—Governance of IT—Governance of Data—Part 2: Implications of ISO/IEC 38500 for Data Management. International Organization for Standardization: Geneva, Switzerland, 2023. Available online: <https://www.iso.org/standard/86012.html> (accessed on 15 June 2024).
18. Wang, R.Y.; Strong, D.M. Beyond Accuracy: What Data Quality Means to Data Consumers. *J. Manag. Inf. Syst.* **1996**, *12*, 5–33. [CrossRef]
19. Batini, C.; Cappiello, C.; Francalanci, C.; Maurino, A. Methodologies for Data Quality Assessment and Improvement. *ACM Comput. Surv.* **2009**, *41*, 16:1–16:52. [CrossRef]
20. Dong, H.; Zhang, C.; Li, G.; Zhang, H. Cloud-Native Databases: A Survey. *IEEE Trans. Knowl. Data Eng.* **2024**, *36*, 7772–7791. [CrossRef]
21. Chaudhuri, S.; Narasayya, V.; Ramakrishnan, R. Data Management in the Cloud: Limitations and Opportunities. *IEEE Data Eng. Bull.* **2011**, *34*, 3–12.
22. Daraghmi, E.; Zhang, C.-P.; Yuan, S.-M. Enhancing Saga Pattern for Distributed Transactions within a Microservices Architecture. *Appl. Sci.* **2022**, *12*, 6242. [CrossRef]

23. Van Landuyt, D.; Benaouda, J.; Reniers, V.; Rafique, A.; Joosen, W. A Comparative Performance Evaluation of Multi-Model NoSQL Databases and Polyglot Persistence. In *Proceedings of the 38th ACM/SIGAPP Symposium on Applied Computing (SAC '23)*; ACM: New York, NY, USA, 2023; pp. 403–412. [[CrossRef](#)]
24. Singhal, H.; Saxena, A.; Mittal, N.; Dabas, C.; Kaur, P. PolyGlot Persistence for Microservices-Based Applications. *Int. J. Inf. Technol. Syst. Approach* **2021**, *14*, 17–32. [[CrossRef](#)]
25. Ye, F.; Sheng, X.; Nedjah, N.; Sun, J.; Zhang, P. A Benchmark for Performance Evaluation of a Multi-Model Database vs. Polyglot Persistence. *J. Database Manag.* **2023**, *34*, 1–20. [[CrossRef](#)]
26. Laigner, R.; Zhou, Y.; Salles, M.A.V.; Liu, Y.; Kalinowski, M. Data Management in Microservices: State of the Practice, Challenges, and Research Directions. *Proc. VLDB Endow.* **2021**, *14*, 3348–3361. [[CrossRef](#)]
27. Brito, G.; Mombach, T.; Valente, M.T. Migrating to GraphQL: A Practical Assessment. In *Proceedings of the 2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER), Hangzhou, China, 24–27 February 2019*; IEEE: Piscataway, NJ, USA, 2019; pp. 140–150. [[CrossRef](#)]
28. Politou, E.; Alepis, E.; Patsakis, C. Forgetting Personal Data and Revoking Consent Under the GDPR: Challenges and Proposed Solutions. *J. Cybersecur.* **2018**, *4*, tyy001. [[CrossRef](#)]
29. Hils, B.; Woods, D.W.; Böhme, R. Measuring the Emergence of Consent Management on the Web. In *Proceedings of the ACM Internet Measurement Conference (IMC '20)*; ACM: New York, NY, USA, 2020; pp. 317–332. [[CrossRef](#)]
30. Vayena, E.; Blasimme, A.; Cohen, I.G. Machine Learning in Medicine: Addressing Ethical Challenges. *PLoS Med.* **2018**, *15*, e1002689. [[CrossRef](#)]
31. Floridi, L.; Cows, J.; Beltrametti, M.; Chatila, R.; Chazerand, P.; Dignum, V.; Luetge, C.; Madelin, R.; Pagallo, U.; Rossi, F.; et al. AI4People—An Ethical Framework for a Good AI Society: Opportunities, Risks, Principles, and Recommendations. *Minds Mach.* **2018**, *28*, 689–707. [[CrossRef](#)] [[PubMed](#)]
32. Kuner, C. Reality and Illusion in EU Data Transfer Regulation Post Schrems. *Ger. Law J.* **2017**, *18*, 881–918. [[CrossRef](#)]

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.