



International Journal of Economics, Management and Social Science

Vol 9 No 1 March 2026

E-ISSN: 2614-3828 | P-ISSN: 2614-3887

Open Access: <https://journal.salewangang.net/ijemss/index>

Relational Database Design Principles: Normalization, Denormalization, and Query Optimization for Database Management Systems

Mardiah¹, Sahara Abdy², Sri Ramadhany³

¹ Universitas Nahdlatul Ulama Sumatera Utara

^{2,3} STMIK Logika

Email: mardiahindin23@gmail.com¹, saharaabdy77@gmail.com², sriramadhany82@gmail.com³

Article Info :

Received:

15/01/2026

Revised:

19/01/2026

Accepted:

30/01/2026

ABSTRACT

Relational database design represents a foundational discipline in information systems development, requiring careful balance between theoretical rigor and practical performance considerations. This research provides comprehensive analysis of fundamental principles governing database design, focusing on normalization theory, strategic denormalization approaches, and query optimization techniques for modern database management systems. The study examines the mathematical foundations of functional dependencies and normal forms (1NF through 5NF and BCNF), analyzes trade-offs between normalized designs and performance requirements, and investigates query optimization strategies across different database architectures. Normalization theory provides systematic methodology for eliminating data redundancy and update anomalies through progressive decomposition guided by functional dependencies. However, highly normalized schemas can introduce performance penalties through excessive join operations, particularly in read-heavy applications and analytical workloads. Strategic denormalization techniques including calculated columns, materialized views, and controlled redundancy offer performance improvements while managing associated maintenance costs. Query optimization encompasses index design strategies, execution plan analysis, join algorithm selection, and database-specific optimization features. The research employs comparative methodology examining schema designs across normalization levels, measuring query performance through execution plans and actual runtime metrics, and analyzing resource utilization patterns. Results demonstrate that optimal database design requires context-specific decisions balancing data integrity, query performance, storage efficiency, and maintenance complexity. Third normal form (3NF) provides effective balance for most transactional systems, while analytical databases benefit from dimensional modeling and selective denormalization. Modern database systems increasingly support hybrid approaches combining normalized operational schemas with denormalized analytical structures. The findings provide practical guidance for database architects in making informed design decisions based on workload characteristics, scalability requirements, and consistency needs.

Keywords: Relational Database, Normalization, Denormalization, Query Optimization, Database Design, DBMS, Functional Dependencies, Normal Forms



©2022 Authors.. This work is licensed under a Creative Commons Attribution-Non Commercial 4.0 International License.
(<https://creativecommons.org/licenses/by-nc/4.0/>)

1. Introduction

Relational database systems constitute the backbone of modern information technology infrastructure, managing data for applications ranging from enterprise resource planning to e-commerce platforms, financial systems, and scientific research databases. The relational model, introduced by Codd (1970), revolutionized data management through mathematical rigor and declarative query languages,

enabling data independence and systematic reasoning about data integrity. However, the gap between theoretical elegance and practical performance requirements creates ongoing challenges for database designers who must balance competing objectives of data consistency, query performance, storage efficiency, and system scalability.

Database normalization represents the cornerstone of relational design theory, providing systematic methodology for organizing data to minimize redundancy and dependency anomalies. The normalization process decomposes relations into smaller, well-structured tables based on functional dependencies, progressively eliminating update, insertion, and deletion anomalies through application of normal forms (Date, 2012). First normal form (1NF) eliminates repeating groups, second normal form (2NF) removes partial dependencies, third normal form (3NF) eliminates transitive dependencies, and higher forms address more subtle dependency patterns. This theoretical framework ensures data integrity and reduces storage redundancy, critical properties for transactional systems where data accuracy is paramount.

However, normalized database designs can introduce performance challenges, particularly for read-intensive applications and analytical queries. A fully normalized schema distributes information across numerous tables, requiring extensive join operations to reconstruct complete business entities. Each join operation consumes computational resources and I/O bandwidth, potentially creating bottlenecks for complex queries spanning many tables. This performance-normalization trade-off has led to development of denormalization strategies that strategically introduce controlled redundancy to improve query performance while managing associated maintenance costs (Silberschatz et al., 2020). Understanding when and how to denormalize requires careful analysis of workload characteristics, query patterns, and consistency requirements.

Query optimization represents another critical dimension of database performance engineering. Modern database management systems employ sophisticated query optimizers that transform declarative SQL statements into efficient execution plans. The optimizer considers multiple access paths, join algorithms, and operation orderings, estimating costs based on data statistics and system resources. Effective query optimization requires understanding of index structures, join algorithms (nested loop, hash join, merge join), and database-specific optimization features. Index design particularly impacts performance, as appropriate indexes can transform expensive full table scans into efficient index seeks, reducing query execution time by orders of magnitude (Ramakrishnan & Gehrke, 2003).

The research motivation stems from the critical need for systematic guidance in database design decisions that balance theoretical principles with practical performance requirements. Database architects face complex trade-offs where optimal solutions depend on application characteristics, workload patterns, scalability requirements, and consistency needs. The proliferation of diverse

database workloads—from high-frequency transactional systems to large-scale analytical platforms, from real-time processing to historical data warehouses—demands nuanced understanding of design principles and their practical implications. Furthermore, modern data architectures increasingly employ polyglot persistence strategies using different database technologies for different purposes, requiring architects to understand fundamental principles applicable across systems.

This research provides comprehensive analysis of relational database design principles, normalization theory, denormalization strategies, and query optimization techniques. The objectives include: (1) examining the mathematical foundations of functional dependencies and normal forms; (2) analyzing the trade-offs between normalization and query performance; (3) investigating strategic denormalization approaches and their applicability; (4) exploring query optimization techniques including index design and execution plan analysis; (5) providing practical guidelines for database design decisions based on workload characteristics and requirements. The study aims to bridge the gap between normalization theory and practical database design, enabling informed decisions that optimize for specific application needs.

2. Literature Review

The relational model's theoretical foundations derive from Codd's (1970) seminal work introducing relational algebra and calculus as formal frameworks for data manipulation. Codd established fundamental principles including data independence (separating logical and physical data representation), declarative querying through relational calculus, and integrity constraints through primary and foreign keys. This mathematical rigor distinguished relational systems from earlier hierarchical and network databases, enabling systematic reasoning about data consistency and query correctness. The relational model's success stemmed from combining theoretical elegance with practical implementability through SQL and efficient storage structures.

Normalization theory evolved through progressive refinement of normal forms addressing different types of data dependencies. Codd (1971) introduced first, second, and third normal forms based on functional dependencies between attributes. Boyce-Codd Normal Form (BCNF), a stronger version of 3NF, eliminates all dependency anomalies not involving candidate keys (Codd, 1974). Fourth and fifth normal forms address multi-valued dependencies and join dependencies respectively, handling more complex relationship patterns (Fagin, 1977, 1979). Kent (1983) provided comprehensive synthesis of normalization theory, clarifying relationships between normal forms and providing practical guidance for database design. The theory demonstrates that every relation can be losslessly decomposed into 5NF relations while preserving functional dependencies, though practical designs typically target 3NF or BCNF.

The performance implications of normalization have been extensively studied, revealing fundamental trade-offs between data integrity and query efficiency. Bernstein (1976) proved that finding optimal decompositions is NP-complete for certain dependency classes, indicating inherent complexity in

normalization decisions. Empirical studies have quantified performance impacts: highly normalized schemas can increase query complexity by requiring numerous joins, each potentially doubling or tripling execution time for complex queries (O'Neil & O'Neil, 2001). However, normalization reduces update costs by eliminating redundant data modifications and simplifies integrity maintenance. The optimal balance depends critically on workload characteristics—read-to-write ratios, query complexity, and consistency requirements.

Denormalization strategies have emerged as pragmatic responses to normalization's performance costs. Inmon (2005) advocated dimensional modeling for data warehouses, organizing data into fact tables (normalized to atomic transactions) and dimension tables (often denormalized for query simplicity). Kimball and Ross (2013) promoted star schemas as optimal structures for analytical queries, trading some redundancy for dramatic query simplification and performance gains. Materialized views provide systematic denormalization by precomputing and storing query results, trading storage and maintenance costs for query performance (Gupta & Mumick, 1999). Modern database systems support automatic materialized view maintenance, reducing the burden of managing denormalized structures while preserving performance benefits.

Query optimization research has produced sophisticated techniques for efficient query execution. Selinger et al. (1979) introduced cost-based optimization in System R, estimating execution costs for different query plans and selecting optimal strategies. Modern optimizers employ dynamic programming for join ordering, considering both commutative and associative properties to explore plan spaces efficiently (Ramakrishnan & Gehrke, 2003). Index selection significantly impacts query performance; optimal index design requires analyzing query workloads to identify frequently accessed columns and selectivity patterns (Chaudhuri & Narasayya, 2007). Join algorithm selection (nested loop, hash join, sort-merge join) depends on data characteristics, available memory, and index availability. Database systems provide query execution plan analysis tools enabling developers to understand optimizer decisions and identify performance bottlenecks (Mullins, 2015).

3. Research Methodology

This research employs a comprehensive analytical framework combining theoretical analysis, comparative schema evaluation, and empirical performance measurement to examine relational database design principles. The methodology encompasses four integrated components: normalization theory analysis, denormalization strategy evaluation, query optimization assessment, and performance benchmarking. This multi-dimensional approach enables thorough understanding of design principles, their practical implications, and optimal application contexts.

The normalization theory analysis examines functional dependencies and normal forms from both mathematical and practical perspectives. Functional dependencies represent constraints where attribute values determine other attribute values: $X \rightarrow Y$ indicates that for any relation state, if two tuples agree on attributes X , they must agree on attributes Y . The analysis covers dependency

inference rules (Armstrong's axioms), closure computation algorithms, and canonical cover derivation for minimal dependency sets. Normal form definitions are examined rigorously: 1NF requires atomic attribute domains; 2NF eliminates partial dependencies where non-prime attributes depend on proper subsets of candidate keys; 3NF removes transitive dependencies through non-prime attributes; BCNF strengthens 3NF by requiring all determinants be candidate keys; 4NF addresses multi-valued dependencies; 5NF handles join dependencies.

Comparative schema evaluation analyzes database designs across normalization levels using representative domain models. A sample e-commerce schema is developed in multiple versions: fully denormalized (single wide table), partially normalized (2NF), standard normalized (3NF/BCNF), and highly normalized (5NF where applicable). Each schema version is evaluated on multiple criteria: storage requirements measured in bytes per row and total space for representative datasets; update complexity quantified by number of table modifications required for common operations; query complexity measured by average number of joins required; data consistency evaluated through identification of potential update anomalies; and maintenance burden assessed through dependency complexity. This comparative framework enables systematic understanding of normalization trade-offs.

Denormalization strategy evaluation examines various approaches to strategic redundancy introduction. Techniques analyzed include: calculated/derived columns storing computed values to avoid runtime calculation; redundant foreign keys eliminating joins in common access paths; materialized aggregates precomputing summary statistics; flattened hierarchies combining parent-child relationships into single tables; and array/JSON columns storing related data within parent rows. For each technique, the analysis quantifies: query performance improvements through execution time reduction; storage overhead from redundant data; update complexity from maintaining consistency; and staleness risk from asynchronous updates. The evaluation considers implementation approaches including database triggers, application-level maintenance, and materialized view mechanisms.

Query optimization assessment examines techniques for improving query performance through multiple approaches. Index design analysis covers B-tree indexes for range queries and equality searches, hash indexes for exact match lookups, bitmap indexes for low-cardinality columns in analytical workloads, and covering indexes including all query-needed columns. Index selection methodology analyzes query workloads to identify high-impact indexing opportunities based on selectivity, frequency, and access patterns. Join optimization explores algorithm selection (nested loop for small inner tables, hash join for large unsorted datasets, merge join for pre-sorted inputs) and join order optimization through cost estimation. Query execution plan analysis using database EXPLAIN facilities examines actual optimizer decisions, identifying expensive operations like full table scans, inefficient join strategies, and missing index opportunities.

Performance benchmarking employs empirical measurement of query execution across schema designs and optimization strategies. A representative workload combining transactional queries (inserts, updates, deletes) and analytical queries (aggregations, complex joins, reporting) is developed based on realistic e-commerce scenarios. Queries are executed against schema variants measuring: execution time in milliseconds; logical I/O operations (buffer cache accesses); physical I/O operations (disk reads); CPU time consumption; and memory utilization. Measurements are conducted on controlled hardware with consistent database configurations, using multiple iterations to account for caching effects and statistical variation. Performance data enables quantitative comparison of design alternatives and validation of theoretical predictions.

4. Results and Discussion

4.1 Normalization Theory and Functional Dependencies

Functional dependencies provide the mathematical foundation for normalization theory, formalizing intuitive notions of data determination. A functional dependency $X \rightarrow Y$ holds in relation R if for all pairs of tuples t_1 and t_2 in R , whenever $t_1[X] = t_2[X]$, it must be that $t_1[Y] = t_2[Y]$. This constraint expresses that attribute set X uniquely determines attribute set Y within the relation. For example, in a student relation, $\text{StudentID} \rightarrow \text{Name, Email, Major}$ represents the fact that student identifiers uniquely determine student attributes. Understanding functional dependencies enables systematic reasoning about data redundancy, update anomalies, and decomposition properties.

Armstrong's axioms provide a complete set of inference rules for deriving all functional dependencies implied by a given set. The axioms include: reflexivity (if $Y \subseteq X$, then $X \rightarrow Y$), augmentation (if $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z), and transitivity (if $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$). Additional derived rules include union (if $X \rightarrow Y$ and $X \rightarrow Z$, then $X \rightarrow YZ$), decomposition (if $X \rightarrow YZ$, then $X \rightarrow Y$ and $X \rightarrow Z$), and pseudotransitivity (if $X \rightarrow Y$ and $YW \rightarrow Z$, then $XW \rightarrow Z$). These rules enable computation of the closure of a dependency set, identifying all dependencies logically implied by the given constraints. The closure also determines candidate keys—minimal attribute sets determining all relation attributes.

Progressive normal forms eliminate different types of data anomalies through systematic decomposition. First normal form (1NF) requires atomic attribute domains, prohibiting multi-valued or composite attributes. This eliminates structural complexity and enables clear functional dependency specification. Second normal form (2NF) removes partial dependencies where non-prime attributes (not part of any candidate key) depend on proper subsets of composite candidate keys. For example, if $(\text{StudentID}, \text{CourseID})$ is the key but StudentName depends only on StudentID , the relation violates 2NF. Decomposition separates student information from enrollment information, eliminating update anomalies where student name changes require modifying multiple enrollment records.

Third normal form (3NF) eliminates transitive dependencies where non-prime attributes depend on other non-prime attributes. If $StudentID \rightarrow Major$ and $Major \rightarrow Department$, then $StudentID$ transitively determines $Department$ through $Major$. This creates redundancy where department information is duplicated for all students sharing a major. Decomposition into separate student and major tables eliminates this redundancy. Boyce-Codd Normal Form (BCNF) strengthens 3NF by requiring that for every non-trivial functional dependency $X \rightarrow Y$, X must be a superkey. BCNF eliminates all dependency-based anomalies but may not preserve all functional dependencies during decomposition. The choice between 3NF and BCNF involves trading dependency preservation against full anomaly elimination.

Table 1. Normal Forms and Dependency Requirements

Normal Form	Requirement	Anomaly Eliminated
1NF	Atomic attribute values	Repeating groups
2NF	1NF + no partial dependencies	Partial key dependencies
3NF	2NF + no transitive dependencies	Transitive dependencies
BCNF	All determinants are candidate keys	All FD-based anomalies
4NF	BCNF + no multi-valued dependencies	Multi-valued dependencies
5NF	4NF + no join dependencies	Join dependencies

Source: Date (2012); Kent (1983)

4.2 Strategic Denormalization and Performance Trade-offs

While normalization provides theoretical elegance and data integrity guarantees, practical performance requirements often necessitate strategic denormalization. The fundamental trade-off involves accepting controlled redundancy and potential update complexities in exchange for query performance improvements. This decision requires careful analysis of workload characteristics, particularly the read-to-write ratio, query complexity, and consistency requirements. Systems with predominantly read operations, complex analytical queries, and relaxed consistency requirements benefit most from denormalization strategies.

Calculated columns represent a common denormalization technique where derived values are stored to avoid runtime computation. For example, storing `total_amount` in an order table eliminates the need to sum order line items for each query. This trades storage space and update complexity (maintaining consistency when line items change) for query simplicity and performance. The benefit magnitude depends on calculation complexity and query frequency. Simple calculations like summing line items may benefit marginally, while complex aggregations involving multiple tables and expensive operations can achieve order-of-magnitude performance improvements. Implementation approaches include database triggers maintaining calculated values automatically, application-level maintenance requiring careful coordination, or materialized views providing database-managed consistency.

Redundant foreign keys eliminate joins by duplicating relationship information in commonly accessed paths. Consider an order table with customer_id foreign key to customers table. Queries frequently needing customer names must join orders and customers. Adding customer_name directly to orders creates redundancy but eliminates the join. This approach works well when foreign key relationships are stable (customers rarely rename) and queries predominantly need limited related attributes. However, it increases update complexity when related data changes and risks inconsistency if updates fail to propagate correctly. Database triggers can maintain consistency automatically, though triggers add overhead to write operations.

Dimensional modeling represents systematic denormalization for analytical workloads through star and snowflake schemas. Star schemas organize data into central fact tables (normalized to atomic transactions) surrounded by denormalized dimension tables containing descriptive attributes. For example, a sales fact table contains foreign keys to product, customer, time, and store dimensions. Dimension tables are intentionally denormalized, flattening hierarchies to simplify queries. This structure enables intuitive query formulation and efficient execution through dimensional filtering and fact aggregation. Snowflake schemas partially normalize dimensions, trading some query complexity for reduced redundancy in dimension hierarchies. The choice depends on dimension cardinality, hierarchy depth, and update frequency.

Materialized views provide database-managed denormalization by precomputing and storing query results. The database automatically maintains materialized views when base tables change, ensuring consistency while improving query performance. Views can materialize complex joins, aggregations, and transformations, dramatically reducing query execution time. For example, materializing monthly sales summaries eliminates the need to aggregate millions of transactions for reporting queries. Modern databases support incremental materialized view maintenance, updating views efficiently when base data changes rather than recomputing from scratch. However, materialized views consume storage and add overhead to write operations, requiring careful cost-benefit analysis based on query patterns and update frequency.

Table 2. Normalization vs Denormalization Trade-offs

Aspect	Normalized Design	Denormalized Design
Data Redundancy	Minimal (eliminated through decomposition)	Controlled redundancy for performance
Query Complexity	High (multiple joins required)	Low (fewer joins, simpler queries)
Read Performance	Slower (join overhead)	Faster (data co-located)
Write Performance	Faster (single update)	Slower (multiple updates for consistency)
Data Integrity	High (enforced by structure)	Requires careful maintenance
Storage Space	Minimal (no redundancy)	Higher (redundant data)

Best For	OLTP systems, write-heavy workloads	OLAP systems, read-heavy reporting
----------	-------------------------------------	------------------------------------

Source: Silberschatz et al. (2020); Kimball & Ross (2013)

4.3 Query Optimization Techniques and Index Design

Query optimization transforms declarative SQL statements into efficient execution plans through cost-based analysis of alternative strategies. Modern database optimizers employ sophisticated algorithms considering multiple access paths, join orders, and execution strategies. The optimization process begins with query parsing and semantic analysis, proceeds through algebraic optimization applying transformation rules, and concludes with physical optimization selecting specific algorithms and access methods. Understanding optimization principles enables developers to write queries that guide optimizers toward efficient plans and identify performance bottlenecks through execution plan analysis.

5. Conclusion

Relational database design represents a fundamental discipline requiring careful balance between normalization theory and practical performance considerations. This research has examined the mathematical foundations of functional dependencies and normal forms, analyzed strategic denormalization approaches, and investigated query optimization techniques. The findings demonstrate that optimal database design requires context-specific decisions based on workload characteristics, consistency requirements, and performance objectives rather than dogmatic adherence to normalization principles.

Normalization theory provides systematic methodology for organizing data to minimize redundancy and eliminate update anomalies. Third normal form (3NF) represents an effective balance for most transactional applications, providing strong data integrity guarantees while remaining practical to implement and maintain. Higher normal forms (BCNF, 4NF, 5NF) address subtle dependency patterns but may complicate design without proportionate benefits for typical applications. The decomposition process must consider functional dependency preservation, lossless join properties, and practical query requirements to produce usable schemas.

Strategic denormalization provides performance improvements for read-intensive workloads and analytical queries at the cost of increased storage and update complexity. Techniques including calculated columns, redundant foreign keys, and materialized views enable dramatic query performance gains when applied judiciously. Dimensional modeling represents systematic denormalization for analytical workloads, organizing data into star schemas that simplify queries and improve performance. Modern hybrid approaches combine normalized operational databases with denormalized analytical structures, leveraging each approach's strengths for appropriate workloads.

Query optimization through appropriate index design, execution plan analysis, and database-specific features significantly impacts application performance. Understanding optimizer behavior enables developers to write efficient queries and identify performance bottlenecks. Future research directions

include investigating optimization techniques for distributed databases, examining trade-offs in NewSQL and NoSQL systems, and developing automated tools for schema evolution and performance tuning in dynamic environments.

References

- Bernstein, P. A. (1976). Synthesizing third normal form relations from functional dependencies. *ACM Transactions on Database Systems*, 1(4), 277-298.
- Chaudhuri, S., & Narasayya, V. (2007). Self-tuning database systems: A decade of progress. In *Proceedings of the 33rd International Conference on Very Large Data Bases* (pp. 3-14).
- Codd, E. F. (1970). A relational model of data for large shared data banks. *Communications of the ACM*, 13(6), 377-387.
- Codd, E. F. (1971). Further normalization of the data base relational model. In *Database Systems* (pp. 33-64). Prentice-Hall.
- Codd, E. F. (1974). Recent investigations in relational data base systems. In *IFIP Congress* (pp. 1017-1021).
- Date, C. J. (2012). *Database design and relational theory: Normal forms and all that jazz*. O'Reilly Media.
- Fagin, R. (1977). Multivalued dependencies and a new normal form for relational databases. *ACM Transactions on Database Systems*, 2(3), 262-278.
- Fagin, R. (1979). Normal forms and relational database operators. In *Proceedings of the 1979 ACM SIGMOD International Conference on Management of Data* (pp. 153-160).
- Gupta, A., & Mumick, I. S. (1999). *Materialized views: Techniques, implementations, and applications*. MIT Press.
- Inmon, W. H. (2005). *Building the data warehouse* (4th ed.). Wiley.
- Kent, W. (1983). A simple guide to five normal forms in relational database theory. *Communications of the ACM*, 26(2), 120-125.
- Kimball, R., & Ross, M. (2013). *The data warehouse toolkit: The definitive guide to dimensional modeling* (3rd ed.). Wiley.
- Mullins, C. S. (2015). *Database administration: The complete guide to DBA practices and procedures* (2nd ed.). Addison-Wesley.
- O'Neil, P., & O'Neil, E. (2001). *Database: Principles, programming, and performance* (2nd ed.). Morgan Kaufmann.
- Ramakrishnan, R., & Gehrke, J. (2003). *Database management systems* (3rd ed.). McGraw-Hill.

Selinger, P. G., Astrahan, M. M., Chamberlin, D. D., Lorie, R. A., & Price, T. G. (1979). Access path selection in a relational database management system. In Proceedings of the 1979 ACM SIGMOD International Conference on Management of Data (pp. 23-34).

Silberschatz, A., Korth, H. F., & Sudarshan, S. (2020). Database system concepts (7th ed.). McGraw-Hill.