

SUBJECT-Relational database

[bca-IIIrd yr]

A **relational database (RDB^[1])** is a [database](#) based on the [relational model](#) of data, as proposed by [E. F. Codd](#) in 1970.^[2] A [database management system](#) used to maintain relational databases is a **relational database management system (RDBMS)**. Many relational database systems are equipped with the option of using [SQL](#) (Structured Query Language) for querying and updating the database.^[3]

History

The concept of relational database was defined by [E. F. Codd](#) at [IBM](#) in 1970. Codd introduced the term *relational* in his research paper "A Relational Model of Data for Large Shared Data Banks".^[2] In this paper and later papers, he defined what he meant by *relation*. One well-known definition of what constitutes a relational database system is composed of [Codd's 12 rules](#). However, no commercial implementations of the relational model conform to all of Codd's rules,^[4] so the term has gradually come to describe a broader class of database systems, which at a minimum:

1. Present the data to the user as [relations](#) (a presentation in tabular form, i.e. as a *collection* of [tables](#) with each table consisting of a set of [rows](#) and [columns](#));
2. Provide relational operators to manipulate the data in tabular form.

In 1974, IBM began developing [System R](#), a research project to develop a prototype RDBMS.^{[5][6]} The first system sold as an RDBMS was [Multics Relational Data Store](#) (June 1976).^[citation needed] [Oracle](#) was released in 1979 by Relational Software, now [Oracle Corporation](#).^[7] [Ingres](#) and [IBM BS12](#) followed. Other examples of an RDBMS include [IBM Db2](#), [SAP Sybase ASE](#), and [Informix](#). In 1984, the first RDBMS for [Macintosh](#) began being developed, code-named Silver Surfer, and was released in 1987 as [4th Dimension](#) and known today as 4D.^[8]

The first systems that were relatively faithful implementations of the relational model were from:

- University of Michigan – [Micro DBMS](#) (1969)^[9]
- Massachusetts Institute of Technology (1971)^[10]
- IBM UK Scientific Centre at Peterlee – [IS1](#) (1970–72),^[11] and its successor, [PRTV](#) (1973–79).^[12]

The most common definition of an RDBMS is a product that presents a view of data as a collection of rows and columns, even if it is not based strictly upon [relational theory](#). By this definition, RDBMS products typically implement some but not all of Codd's 12 rules.

A second school of thought argues that if a database does not implement all of Codd's rules (or the current understanding on the relational model, as expressed by [Christopher J. Date](#), [Hugh Darwen](#) and others), it is not relational. This view, shared by many theorists and other strict adherents to Codd's principles, would disqualify most DBMSs as not relational. For clarification, they often refer to some RDBMSs as *truly-relational database management systems* (TRDBMS), naming others *pseudo-relational database management systems* (PRDBMS)

As of 2009, most commercial relational DBMSs employ [SQL](#) as their [query language](#).^[13]

Alternative query languages have been proposed and implemented, notably the pre-1996 implementation of [Ingres QUEL](#).

Relational model

A relational model organizes data into one or more [tables](#) (or "relations") of [columns](#) and [rows](#), with a unique key identifying each row. Rows are also called [records](#) or [tuples](#).^[14] Columns are also called attributes. Generally, each table/relation represents one "entity type" (such as customer or product). The rows represent instances of that type of [entity](#) (such as "Lee" or "chair") and the columns represent values attributed to that instance (such as address or price).

For example, each row of a class table corresponds to a class, and a class corresponds to multiple students, so the relationship between the class table and the student table is "one to many"^[15]

Keys

Each row in a table has its own unique key. Rows in a table can be linked to rows in other tables by adding a column for the unique key of the linked row (such columns are known as [foreign keys](#)). Codd showed that data relationships of arbitrary complexity can be represented by a simple set of concepts.^[2]

Part of this processing involves consistently being able to select or modify one and only one row in a table. Therefore, most physical implementations have a unique [primary key](#) (PK) for each row in a table. When a new row is written to the table, a new unique value for the primary key is generated; this is the key that the system uses primarily for accessing the table. System performance is optimized for PKs. Other, more [natural keys](#) may also be identified and defined as [alternate keys](#) (AK). Often several columns are needed to form an AK (this is one reason why a single integer column is usually made the PK). Both PKs and AKs have the ability to uniquely identify a row within a table. Additional technology may be applied to ensure a unique ID across the world, a [globally unique identifier](#), when there are broader system requirements.

The primary keys within a database are used to define the relationships among the tables. When a PK migrates to another table, it becomes a foreign key in the other table. When each cell can contain only one value and the PK migrates into a regular entity table, this design pattern can represent either a [one-to-one](#) or [one-to-many](#) relationship. Most relational database designs resolve [many-to-many](#) relationships by creating an additional table that contains the PKs from both of the other entity tables – the relationship becomes an entity; the resolution table is then named appropriately and the two FKs are combined to form a PK. The migration of PKs to other tables is the second major reason why system-assigned integers are used normally as PKs; there is usually neither efficiency nor clarity in migrating a bunch of other types of columns.

Relationships

Relationships are a logical connection between different tables (entities), established on the basis of interaction among these tables. These relationships can be modelled as an [entity-relationship model](#).

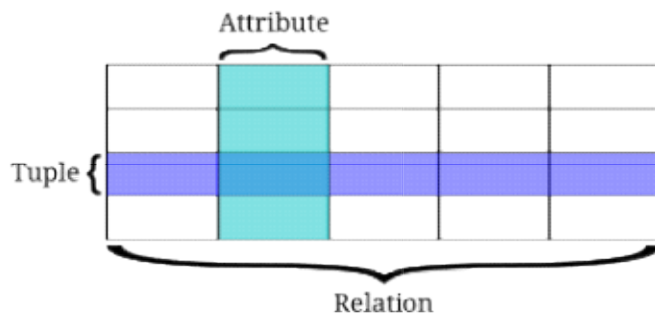
Transactions

In order for a database management system (DBMS) to operate efficiently and accurately, it must use [ACID transactions](#).^{[16][17][18]}

Stored procedures

Part of the programming within a RDBMS is accomplished using [stored procedures](#) (SPs). Often procedures can be used to greatly reduce the amount of information transferred within and outside of a system. For increased security, the system design may grant access to only the stored procedures and not directly to the tables. Fundamental stored procedures contain the logic needed to insert new and update existing data. More complex procedures may be written to implement additional rules and logic related to processing or selecting the data.

Terminology



Relational database terminology

The relational database was first defined in June 1970 by [Edgar Codd](#), of IBM's [San Jose Research Laboratory](#).^[2] Codd's view of what qualifies as an RDBMS is summarized in [Codd's 12 rules](#). A relational database has become the predominant type of database. Other models besides the *relational model* include the [hierarchical database model](#) and the [network model](#).

The table below summarizes some of the most important relational database terms and the corresponding [SQL](#) term:

SQL term	Relational database term	Description
Row	Tuple or record	A data set representing a single item
Column	Attribute or field	A labeled element of a tuple, e.g. "Address" or "Date of birth"
Table	Relation or Base relvar	A set of tuples sharing the same attributes; a set of columns and rows
View or result set	Derived relvar	Any set of tuples; a data report from the RDBMS in response to a query

Relations or tables

In a relational database, a [relation](#) is a set of [tuples](#) that have the same [attributes](#). A tuple usually represents an object and information about that object. Objects are typically physical objects or concepts. A relation is usually described as a [table](#), which is organized into [rows](#) and [columns](#). All the data referenced by an attribute are in the same [domain](#) and conform to the same constraints.

The relational model specifies that the tuples of a relation have no specific order and that the tuples, in turn, impose no order on the attributes. Applications access data by specifying queries, which use operations such as *select* to identify tuples, *project* to identify attributes, and *join* to combine relations. Relations can be modified using the *insert*, *delete*, and *update* operators. New tuples can supply explicit values or be derived from a query. Similarly, queries identify tuples for updating or deleting.

Tuples by definition are unique. If the tuple contains a [candidate](#) or primary key then obviously it is unique; however, a primary key need not be defined for a row or record to be a tuple. The definition of a tuple requires that it be unique, but does not require a primary key to be defined. Because a tuple is unique, its attributes by definition constitute a [superkey](#).

Base and derived relations

All data are stored and accessed via [relations](#). Relations that store data are called "base relations", and in implementations are called "tables". Other relations do not store data, but are computed by applying relational operations to other relations. These relations are sometimes called "derived relations". In implementations these are called "[views](#)" or "queries". Derived relations are convenient in that they act as a single relation, even though they may grab information from several relations. Also, derived relations can be used as an [abstraction layer](#).

Domain

A domain describes the set of possible values for a given attribute, and can be considered a constraint on the value of the attribute. Mathematically, attaching a domain to an attribute means that any value for the attribute must be an element of the specified set. The character string "ABC", for instance, is not in the integer domain, but the integer value 123 is. Another example of domain describes the possible values for the field "CoinFace" as ("Heads", "Tails"). So, the field "CoinFace" will not accept input values like (0,1) or (H,T).

Constraints

Constraints are often used to make it possible to further restrict the domain of an attribute. For instance, a constraint can restrict a given integer attribute to values between 1 and 10. Constraints provide one method of implementing [business rules](#) in the database and support subsequent data use within the application layer. SQL implements constraint functionality in the form of [check constraints](#). Constraints restrict the data that can be stored in [relations](#). These are usually defined using expressions that result in a [Boolean](#) value, indicating whether or not the data satisfies the constraint. Constraints can apply to single attributes, to a tuple (restricting combinations of attributes) or to an entire relation. Since every attribute has an associated domain, there are constraints (**domain constraints**). The two principal rules for the relational model are known as [entity integrity](#) and [referential integrity](#).

Primary key

Every [relation](#)/table has a primary key, this being a consequence of a relation being a [set](#).^[19] A primary key uniquely specifies a tuple within a table. While natural attributes (attributes used to describe the data being entered) are sometimes good primary keys, [surrogate keys](#) are often used instead. A surrogate key is an artificial attribute assigned to an object which uniquely identifies it (for instance, in a table of information about students at a school they might all be assigned a student ID in order to differentiate them). The surrogate key has no intrinsic (inherent) meaning, but rather is useful through its ability to uniquely identify a tuple. Another common occurrence, especially in regard to N:M cardinality is the [composite key](#). A composite key is a key made up of two or more attributes within a table that (together) uniquely identify a record.^[20]

Foreign key

Foreign key refers to a field in a relational table that matches the primary key column of another table. It relates the two keys. Foreign keys need not have unique values in the referencing relation. A foreign key can be used to [cross-reference](#) tables, and it effectively uses the values of attributes in the referenced relation to restrict the domain of one or more attributes in the referencing relation. The concept is described formally as: "For all tuples in the referencing relation projected over the referencing attributes, there must exist a tuple in the referenced relation projected over those same attributes such that the values in each of the referencing attributes match the corresponding values in the referenced attributes."

Stored procedures

A stored procedure is executable code that is associated with, and generally stored in, the database. Stored procedures usually collect and customize common operations, like inserting a [tuple](#) into a [relation](#), gathering statistical information about usage patterns, or encapsulating complex [business logic](#) and calculations. Frequently they are used as an [application programming interface](#) (API) for security or simplicity. Implementations of stored procedures on SQL RDBMS's often allow developers to take advantage of [procedural](#) extensions (often vendor-specific) to the standard [declarative](#) SQL syntax. Stored procedures are not part of the relational database model, but all commercial implementations include them.

Index

An index is one way of providing quicker access to data. Indices can be created on any combination of attributes on a [relation](#). Queries that filter using those attributes can find matching tuples directly using the index (similar to [Hash table](#) lookup), without having to check each tuple in turn. This is analogous to using the [index of a book](#) to go directly to the page on which the information you are looking for is found, so that you do not have to read the entire book to find what you are looking for. Relational databases typically supply multiple indexing techniques, each of which is optimal for some combination of data distribution, relation size, and typical access pattern. Indices are usually implemented via [B+ trees](#), [R-trees](#), and [bitmaps](#). Indices are usually not considered part of the database, as they are considered an implementation detail, though indices are usually maintained by the same group that maintains the other parts of the database. The use of efficient indexes on both primary and foreign keys can dramatically improve query performance. This is because B-tree indexes result in query times proportional to $\log(n)$ where n is the number of rows in a table and hash indexes result in

constant time queries (no size dependency as long as the relevant part of the index fits into memory).

Relational operations

Queries made against the relational database, and the derived [relvars](#) in the database are expressed in a [relational calculus](#) or a [relational algebra](#). In his original relational algebra, Codd introduced eight relational operators in two groups of four operators each. The first four operators were based on the traditional mathematical [set operations](#):

- The [union](#) operator (\cup) combines the tuples of two [relations](#) and removes all duplicate tuples from the result. The relational union operator is equivalent to the [SQL UNION](#) operator.
- The [intersection](#) operator (\cap) produces the set of tuples that two relations share in common. Intersection is implemented in SQL in the form of the [INTERSECT](#) operator.
- The [set difference](#) operator ($-$) acts on two relations and produces the set of tuples from the first relation that do not exist in the second relation. Difference is implemented in SQL in the form of the [EXCEPT](#) or MINUS operator.
- The [cartesian product](#) (\times) of two relations is a join that is not restricted by any criteria, resulting in every tuple of the first relation being matched with every tuple of the second relation. The cartesian product is implemented in SQL as the [Cross join](#) operator.

The remaining operators proposed by Codd involve special operations specific to relational databases:

- The selection, or restriction, operation (σ) retrieves tuples from a relation, limiting the results to only those that meet a specific criterion, i.e. a [subset](#) in terms of set theory. The SQL equivalent of selection is the [SELECT](#) query statement with a [WHERE](#) clause.
- The [projection operation](#) (π) extracts only the specified attributes from a tuple or set of tuples.
- The join operation defined for relational databases is often referred to as a natural join (\bowtie). In this type of join, two relations are connected by their common attributes. MySQL's approximation of a natural join is the [Inner join](#) operator. In SQL, an INNER JOIN prevents a cartesian product from occurring when there are two tables in a query. For each table added to a SQL Query, one additional INNER JOIN is added to prevent a cartesian product. Thus, for N tables in an SQL query, there must be N-1 INNER JOINS to prevent a cartesian product.
- The [relational division](#) (\div) operation is a slightly more complex operation and essentially involves using the tuples of one relation (the dividend) to partition a second relation (the divisor). The relational division operator is effectively the opposite of the cartesian product operator (hence the name).

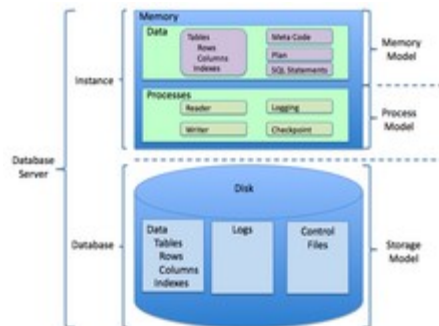
Other operators have been introduced or proposed since Codd's introduction of the original eight including relational comparison operators and extensions that offer support for nesting and hierarchical data, among others.

Normalization

Normalization was first proposed by Codd as an integral part of the relational model. It encompasses a set of procedures designed to eliminate non-simple domains (non-atomic values) and the redundancy (duplication) of data, which in turn prevents data manipulation

anomalies and loss of data integrity. The most common forms of normalization applied to databases are called the [normal forms](#).

RDBMS



The general structure of a relational database

Connolly and Begg define database management system (DBMS) as a "software system that enables users to define, create, maintain and control access to the database".^[21] RDBMS is an extension of that initialism that is sometimes used when the underlying database is relational.

An alternative definition for a *relational database management system* is a database management system (DBMS) based on the [relational model](#). Most databases in widespread use today are based on this model.^[22]

RDBMSs have been a common option for the storage of information in databases used for financial records, manufacturing and logistical information, personnel data, and other applications since the 1980s. Relational databases have often replaced legacy [hierarchical databases](#) and [network databases](#), because RDBMS were easier to implement and administer. Nonetheless, relational stored data received continued, unsuccessful challenges by [object database](#) management systems in the 1980s and 1990s, (which were introduced in an attempt to address the so-called [object-relational impedance mismatch](#) between relational databases and object-oriented application programs), as well as by [XML database](#) management systems in the 1990s.^[23] However, due to the expanse of technologies, such as [horizontal scaling](#) of [computer clusters](#), [NoSQL](#) databases have recently become popular as an alternative to RDBMS databases.^[24]

Distributed relational databases

[Distributed Relational Database Architecture](#) (DRDA) was designed by a workgroup within IBM in the period 1988 to 1994. DRDA enables network connected relational databases to cooperate to fulfill SQL requests.^{[25][26]} The messages, protocols, and structural components of DRDA are defined by the [Distributed Data Management Architecture](#).

List of database engines

^[27]

1. [Oracle Database](#)
2. [MySQL](#)
3. [Microsoft SQL Server](#)
4. [PostgreSQL](#) (free software)

5. [IBM Db2](#)
6. [Microsoft Access](#)
7. [SQLite](#) (free software)
8. [MariaDB](#) (free software)
9. [Snowflake](#)
10. [Microsoft Azure SQL Database](#)
11. [Apache Hive](#) (free software)
12. [Teradata Vantage](#)

According to research company [Gartner](#), in 2011, the five leading [proprietary software](#) relational database vendors by revenue were [Oracle](#) (48.8%), [IBM](#) (20.2%), [Microsoft](#) (17.0%), [SAP](#) including [Sybase](#) (4.6%), and [Teradata](#) (3.7%).^[28]