

What is Data Warehousing?

The background is a solid orange color. It features several decorative elements: a large teal circle with a white pie chart inside it, positioned in the upper right; several smaller orange circles of varying sizes scattered around; and a bar chart in the bottom right corner with four bars of increasing height.

Data warehousing is a process of collecting, storing, and managing large volumes of data from different sources within an organization. The purpose of data warehousing is to provide a centralized and unified view of the organization's data, making it easier for decision-makers to analyze and derive insights. Here are key components and characteristics of data warehousing

Characteristics of Data Warehousing

- Centralized Repository
- Integration of Data
- Historical Data Storage
- Optimized for Query and Analysis
- Data Cleansing and Transformation
- Support for Decision-Making
- Separation from Operational Systems
- Scalability
- Data Marts
- Business Intelligence Tools

Implementing a data warehouse involves careful planning, data modeling, and the use of technologies and methodologies to ensure the effectiveness of the stored data for analytical purposes. The insights derived from a data warehouse can contribute significantly to strategic planning, forecasting, and overall business performance.

Keys in Data Warehousing

In data warehousing, keys play a crucial role in organizing and structuring data. Keys are used to establish relationships between tables and ensure data integrity. Here are some key types commonly used in data warehousing:

Key	Description
● Primary Key:	It ensures that each row in the table can be uniquely identified, and no two rows have the same primary key value.
● Foreign Key:	A foreign key is a field in a table that refers to the primary key in another table. It establishes a link between two tables, enforcing referential integrity.
● Candidate Key:	A candidate key is a set of one or more columns that can uniquely identify a record. It is a potential candidate for becoming the primary key.
● Composite Key:	A composite key is a key that consists of two or more columns in a table. Together, these columns uniquely identify a record in the table.

Keys in Data Warehousing Cont'd

In data warehousing, keys play a crucial role in organizing and structuring data. Keys are used to establish relationships between tables and ensure data integrity. Here are some key types commonly used in data warehousing:

Key	Description
• Surrogate Key:	A surrogate key is an artificially generated key used as the primary key. It does not have any business meaning and is often an auto-incremented number.
• Natural Key:	A natural key is a key that has a meaningful business significance. It is derived from the natural characteristics of the data, such as a product code or employee ID.
• Super Key:	A super key is a set of one or more keys that, taken together, can uniquely identify a record. It is a broader concept that includes candidate keys and other combinations.

Primary Key

A Primary Key is a fundamental concept in database management and data warehousing. It is a unique identifier for a record in a database table, ensuring that each row can be distinctly identified. The primary key plays a crucial role in maintaining data integrity, as it ensures that there are no duplicate records in the table and provides a means to establish relationships between tables.

Here's an example of a simple table with a primary key: In this example, the "EmployeeID" column serves as the primary key for uniquely identifying each employee record in the table.

EmployeeID	FirstName	LastName	Department
1	John	Smith	HR
2	Jane	Doe	IT
3	Bob	Johnson	Finance

Key Characteristics of Primary Key

● Uniqueness:	Each value in the primary key column must be unique. This uniqueness ensures that no two rows in the table have the same identifier.
● Non-null Values:	The primary key column cannot contain null (empty) values. Every record must have a valid and non-null identifier.
● Indexed:	The primary key column is often indexed for faster data retrieval. This indexing improves the performance of queries that involve searching for specific records based on their primary key.
● Single Column or Composite Key:	A primary key can be a single column or a combination of multiple columns, forming a composite key. In the case of a composite key, the combination of values must be unique.

Key Characteristics of Primary Key Cont'd

<ul style="list-style-type: none">• Used for Referential Integrity:	<p>In relational databases, the primary key of one table can be referenced as a foreign key in another table. This establishes relationships between tables, ensuring referential integrity.</p>
<ul style="list-style-type: none">• Auto-incrementing:	<p>In many cases, primary keys are auto-incremented or generated automatically by the database system. This eliminates the need for users or applications to manually assign unique identifiers.</p>
<ul style="list-style-type: none">• Immutable:	<p>The value of a primary key, once assigned, should ideally be immutable. Changing the primary key value for a record can have implications on referential integrity and data consistency.</p>

Foreign Key

A Foreign Key is a relational database concept that establishes a link between two tables by connecting a field in one table to the primary key in another table. The purpose of a foreign key is to enforce referential integrity, ensuring that relationships between tables are maintained consistently. Foreign keys are a fundamental component of relational database systems, including those used in data warehousing.

Here's a simple example to illustrate a foreign key relationship:

Table: Employees

EmployeeID	FirstName	LastName	DepartmentID
1	John	Smith	101
2	Jane	Doe	102
3	Bob	Johnson	101

Table: Departments

DepartmentID	DepartmentName
101	HR
102	IT

In this example, the "DepartmentID" column in the "Employees" table is a foreign key that references the "DepartmentID" column in the "Departments" table. This establishes a relationship between the two tables based on the department ID.

Key Characteristics of Foreign Key

● Reference to Primary Key:	A foreign key in one table refers to the primary key in another table. This establishes a relationship between the two tables.
● Ensures Referential Integrity:	Referential integrity ensures that relationships between tables are consistent. The values in the foreign key column must correspond to existing values in the primary key column of the referenced table.
● Supports Joins:	Foreign keys play a crucial role in enabling the merging of data from different tables through joins. This is essential for querying and analyzing data across related tables.
● Optional or Mandatory:	A foreign key relationship can be optional or mandatory. In an optional relationship, the foreign key column in a table may contain null values. In a mandatory relationship, every record in the referencing table must have a corresponding value in the referenced table.

Key Characteristics of Foreign Key Cont'd

<ul style="list-style-type: none">● Cascading Actions:	<p>Cascading actions define how changes to the referenced primary key are reflected in the referencing foreign key. Common options include:</p> <ul style="list-style-type: none">● CASCADE: Changes in the primary key are automatically reflected in the foreign key.● SET NULL: Foreign key values are set to NULL if the corresponding primary key is deleted or updated.● SET DEFAULT: Foreign key values are set to their default values if the corresponding primary key is deleted or updated.● NO ACTION: The database prevents the deletion or update of the primary key if it is referenced by a foreign key.
<ul style="list-style-type: none">● Indexed for Performance:	<p>To enhance performance, foreign key columns are often indexed, allowing for faster retrieval of data when querying based on the foreign key values.</p>

Candidate Key

A Candidate Key is a set of one or more columns in a database table that can uniquely identify a record within that table. These are potential keys that could be chosen as the primary key of the table. Candidate keys are a concept in database design and normalization, and they serve as a pool of options from which the primary key can be selected.

Here's a simple example to illustrate Candidate Keys:

EmployeeID	SocialSecurityNumber	Email	DepartmentID
1	123-45-6789	john.smith@email.com	101
2	987-65-4321	jane.doe@email.com	102
3	555-11-2222	bob.johnson@email.com	101

In this example, the "EmployeeID," "SocialSecurityNumber," and "Email" columns could all be potential candidate keys because each of them is unique for each employee. The "DepartmentID" column is not a candidate key on its own because it does not uniquely identify individual employees.

Key Characteristics of Candidate Keys

• Uniqueness:	Like a primary key, a candidate key must ensure that the values within the chosen set of columns are unique across all records in the table.
• Irreducibility:	A candidate key should not have any unnecessary columns. It is considered irreducible, meaning that removing any column from the set would result in a loss of uniqueness.
• Minimality:	The set of columns chosen as a candidate key should be minimal, meaning that no proper subset of the candidate key should be a unique identifier for the records.
• Potential Primary Key:	A candidate key is a set of columns that has the potential to become the primary key of the table. The primary key is usually chosen from the set of candidate keys.

Key Characteristics of Candidate Keys Cont'd

<ul style="list-style-type: none">Multiple Candidate Keys:	A table may have more than one candidate key. Each candidate key provides a unique way of identifying records within the table.
<ul style="list-style-type: none">Selection of Primary Key:	The primary key is ultimately selected from the pool of candidate keys based on various factors, including simplicity, stability, and potential for efficient indexing.

The selection of the primary key from these candidate keys would depend on factors such as the stability of the data, ease of use, and the efficiency of indexing for querying purposes.

Composite Key

A Composite Key is a combination of two or more columns in a database table that, when taken together, uniquely identifies each record in the table. Unlike a single-column primary key, a composite key is formed by the combination of multiple columns. Each column within the composite key might not be unique by itself, but the combination of all the columns in the key must be unique across all records in the table.

Here's a simple example to illustrate Composite Keys:

Table: OrderDetails

OrderID	ProductID	LineNumber	Quantity
1	101	1	10
1	102	2	5
2	101	1	8

In this example, the combination of "OrderID" and "LineNumber" forms a composite key, ensuring the uniqueness of each order line.

Key Characteristics of Composite Keys

• Uniqueness:	The combination of columns in a composite key must be unique across all records in the table.
• Complexity:	Composite keys are more complex than single-column primary keys, as they involve multiple columns.
• No Single Column Identity:	No single column within the composite key can be used on its own to uniquely identify records.

Surrogate Key

A Surrogate Key is a unique identifier for a record in a database table that is generated independently of the actual data. Unlike a natural key, which is derived from the attributes of the data (like an employee ID or a product code), a surrogate key has no inherent business meaning. Instead, it is often an auto-incremented number or some other system-generated value.

Here's a simple example to illustrate Surrogate Key:

Table: Customers

CustomerID (Surrogate Key)	FirstName	LastName	Email
1	John	Smith	john.smith@email.com
2	Jane	Doe	jane.doe@email.com
3	Bob	Johnson	bob.johnson@email.com

In this example, the "CustomerID" is a surrogate key, providing a unique identifier for each customer record.

Key Characteristics of Surrogate Key

● Artificial Generation:	Surrogate keys are typically generated automatically by the database system and have no meaning outside the database.
● Uniqueness:	Surrogate keys are designed to be unique for each record, ensuring a simple and effective means of identification.
● Stability:	Unlike natural keys, surrogate keys remain stable even if the underlying data changes. This stability is advantageous in systems where the primary key should not change.

Natural Key

A Natural Key is a type of key in a database that is derived from the natural characteristics of the data it represents. Unlike a surrogate key, which is artificially generated and lacks inherent business meaning, a natural key is based on attributes that have a direct and meaningful connection to the real-world entities being represented in the database.

Here's a simple example to illustrate Natural Key:

Table: Employees

EmployeeID (Natural Key)	FirstName	LastName	SocialSecurityNumber
1	John	Smith	123-45-6789
2	Jane	Doe	987-65-4321
3	Bob	Johnson	555-11-2222

In this example, the "SocialSecurityNumber" is a natural key for identifying employees.

Key Characteristics of Natural Key

● Inherent Business Meaning:	A natural key has a business-related meaning or significance. It is often based on attributes like names, codes, or identifiers that are meaningful in the context of the data.
● Derived from Data Attributes:	The natural key is derived directly from the attributes of the data being stored, and it reflects the way information is naturally identified in the real world.
● Stability:	The values of a natural key are subject to the stability of the underlying data. If the business meaning of an attribute changes, the natural key may also change.

Super Key

A Super Key is a set of one or more columns in a database table whose combined values can uniquely identify a record. A super key may include more columns than necessary to uniquely identify a record, and it can contain candidate keys as subsets.

Here's a simple example to illustrate Natural Key:

Table: Students

StudentID	FirstName	LastName	DateOfBirth	PhoneNumber
1	Alice	Johnson	1990-05-15	555-123-4567
2	Bob	Smith	1991-08-22	555-987-6543
3	Carol	Davis	1990-03-10	555-789-0123

In this example, the combination of "StudentID" and "PhoneNumber" forms a super key. Both columns together uniquely identify each student.

Key Characteristics of Super Key

• Uniqueness	The combination of columns in a super key must be unique across all records in the table.
• Redundancy:	A super key can include more columns than necessary to uniquely identify records. It may contain additional attributes that are not required for uniqueness.
• Candidate Keys Subset:	A super key can contain candidate keys as subsets. Candidate keys are minimal super keys.

Keys Comparison

Key type	Uniqueness	Artificial Generation
Primary Key	Unique across the table	No
Foreign Key	Unique in referenced table	No
Candidate Key	Unique across the table	No
Composite Key	Unique across the table	No
Surrogate Key	Unique across the table	Yes
Natural Key	Unique across the table	No
Super Key	Unique across the table	Depends

Keys Comparison Cont'd

- **Uniqueness:** Indicates whether the key values must be unique across the entire table.
- **Artificial Generation:** Whether the key is artificially generated or has inherent business meaning.
- **Meaningful Business Context:** Whether the key has meaningful business significance.
- **Complexity:** Indicates the complexity of the key, considering the number of columns and their relationships.
- **Indexing:** Whether the key is typically indexed for performance.
- **Usage:** Brief description of the primary purpose or usage of each key type.

This table provides a comparative overview of the key characteristics of different types of keys commonly used in database design. Keep in mind that the choice of a specific key type depends on the requirements of the database schema and the nature of the data being modeled.

Conclusion

Each type of key serves a specific purpose in data warehousing, and their proper use is essential for maintaining data integrity and supporting efficient data retrieval. The choice of keys depends on the nature of the data and the relationships between tables in the data warehouse schema. The use of keys helps in organizing data in a way that facilitates effective querying and analysis.

