

Understanding Data Warehousing

A Foundation for Data-Driven Decision Making

Introduction

What is Data Warehousing?

At its core, a **Data Warehouse (DW)** is a centralized repository of integrated, structured, and historical data. Unlike operational databases, it's designed specifically to support **Business Intelligence (BI)**, analytics, and complex decision-making processes.

Its primary purpose is to enable organizations to **consolidate data from multiple disparate sources** (like ERP, CRM, and other transaction systems) into a single, consistent format. This consolidation allows for comprehensive reporting, deep analysis, and the extraction of strategic insights that would be difficult or impossible with fragmented operational data.



Data Consolidation

Combines data from various sources into one unified view.

2

Strategic Insights

Supports long-term planning and decision-making.

3

Historical Analysis

Retains data over time for trend analysis and forecasting.

Ultimately, data warehousing provides the robust analytical environment necessary for data-driven strategies and operational efficiencies.

Data Lakes vs. Data Warehouses

While both are storage repositories, Data Lakes and Data Warehouses serve distinct purposes and handle data differently. Understanding their nuances is key to modern data architecture.

Data Warehouse

- **Data Type:** Structured, processed data.
- **Schema:** Schema-on-write (predefined structure).
- **Purpose:** Optimized for analytical queries, reporting, OLAP.
- **Users:** Business professionals, analysts.
- **Cost:** Can be more expensive per gigabyte due to processing.



Data Lake

- **Data Type:** Raw, unstructured, semi-structured, or structured data.
- **Schema:** Schema-on-read (flexible, applied during analysis).
- **Purpose:** Suited for exploratory analytics, machine learning, data science.
- **Users:** Data scientists, engineers.
- **Cost:** Low-cost storage, scalable.

They are complementary: Data Lakes are ideal for raw data exploration and machine learning, while Data Warehouses excel at structured reporting and BI.

OLTP vs. OLAP: Transactional vs. Analytical

Understanding the distinction between OLTP and OLAP is fundamental to grasping data warehousing's role. They represent two fundamentally different approaches to data processing.

OLTP (Online Transaction Processing)

- **Purpose:** Manages daily business operations and transactions (e.g., order entry, ATM withdrawals).
- **Data:** Current, highly volatile data in normalized databases.
- **Operations:** High volume of simple inserts, updates, and deletes.
- **Optimization:** Optimized for fast write/read of individual records.
- **Example:** E-commerce checkout system, banking transactions.

OLAP (Online Analytical Processing)

- **Purpose:** Supports complex, ad-hoc queries for strategic decision-making and analysis.
- **Data:** Historical, aggregated data in denormalized data warehouses.
- **Operations:** Lower volume of complex reads involving large data sets.
- **Optimization:** Optimized for fast data retrieval for analytical purposes.
- **Example:** Sales trend analysis, customer behavior insights.

Relationship: OLTP systems are the source of raw data that, after being processed and transformed, feeds into OLAP systems for insightful analysis. They are interdependent parts of a complete data ecosystem.

Normalization vs. Denormalization

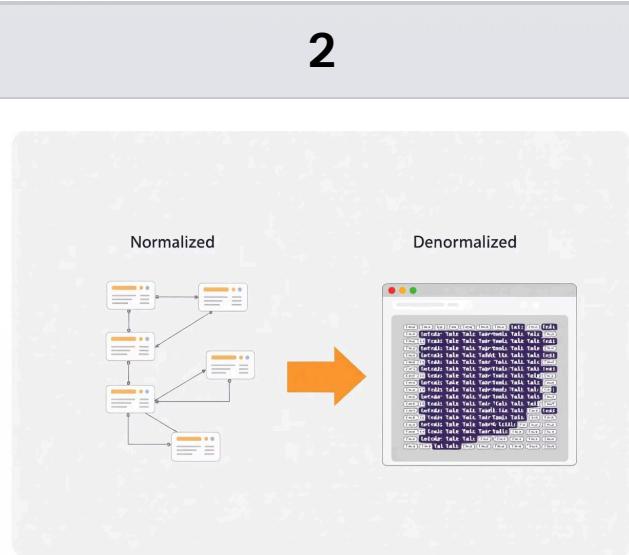
These two data organization strategies are critical for database design, each optimized for different objectives within a data ecosystem.

1

Normalization

- Goal:** Reduce data redundancy and improve data integrity.
- Method:** Breaks down data into multiple, smaller, related tables.
- Use Case:** Primarily used in **OLTP systems** where data updates are frequent.
- Impact:** More joins needed for queries, which can slow down complex analytical reporting.

2



3

Denormalization

- Goal:** Improve query performance by reducing the number of table joins.
- Method:** Combines data into fewer tables, introducing intentional redundancy.
- Use Case:** Crucial for **OLAP systems** and data warehouses where analytical queries are paramount.
- Impact:** Faster analytical query execution, but requires more storage space and careful management of redundancy.

The choice between them is a trade-off between **data consistency** and **query performance**, tailored to the system's primary function.

Data Warehouse Characteristics & Types

Data warehouses possess unique characteristics that differentiate them from operational databases, and they come in various types to serve different organizational needs.

Key Characteristics



Subject-Oriented

Organized by business subjects (e.g., sales, customers), not by applications.



Integrated

Data from disparate sources is made consistent and uniform.



Time-Variant

Stores historical data, allowing for trend analysis over time.



Nonvolatile

Once data is in, it generally remains static, preserving history.

Types of Data Warehouses

- **Enterprise Data Warehouse (EDW):** A central, organization-wide repository integrating all data sources for a holistic view.
- **Data Marts:** Smaller, subject-oriented data warehouses for specific departments (e.g., marketing, finance). Can be dependent (from EDW) or independent.
- **Operational Data Store (ODS):** A database that integrates data from various sources for operational reporting and analysis, often with near real-time data.



These characteristics and types ensure data warehouses provide the right data, at the right time, for the right level of analysis across the organization.

Dimensional Modeling for BI

Dimensional modeling is the cornerstone of data warehouse design, optimizing databases specifically for analytical querying and reporting, not transactional processing.

Key Components

- **Facts:** Numeric measures of a business process (e.g., sales amount, units sold, profit). They are stored in **fact tables**, which contain foreign keys to dimension tables and the measures themselves. The "grain" defines the level of detail (e.g., per transaction, per day).
- **Dimensions:** Descriptive attributes related to the facts, providing context for analysis (e.g., Time, Product, Customer, Store). Dimensions often include hierarchies (e.g., Year → Quarter → Month).

Common Schemas

- **Star Schema:** The most common and recommended. Features a central fact table directly linked to denormalized dimension tables. It's fast, simple to understand, and widely supported by BI tools.
- **Snowflake Schema:** An extension of the star schema where dimensions are normalized into multiple related tables. Reduces redundancy but increases complexity and query join overhead.



Dimensional modeling ensures business users can easily understand and query data for their analytical needs, providing a flexible framework for BI.

ETL: Extract, Transform, Load

ETL is the foundational process for populating a data warehouse, turning raw, disparate data into clean, consistent, and usable information.

01

Extract

Pulling data from various source systems. This can include transactional databases (ERP, CRM), flat files, APIs, web logs, or streaming data. Challenges include identifying relevant data, handling different data formats, and ensuring efficient extraction frequency (daily, hourly, real-time).

02

Transform

Cleaning, standardizing, and enhancing the extracted data. This is the most critical step for data quality. Tasks include:

- Data cleansing (handling missing values, correcting errors)
- Format conversion (e.g., dates, currencies)
- Deduplication and validation
- Aggregation and summarization
- Data enrichment (adding external data for context)

03

Load

Writing the transformed data into the data warehouse.

- **Initial Load:** Full historical data load when the warehouse is first set up.
- **Incremental Load:** Only new or changed data is added, typically on a recurring schedule. This is crucial for efficiency and maintaining up-to-date information.

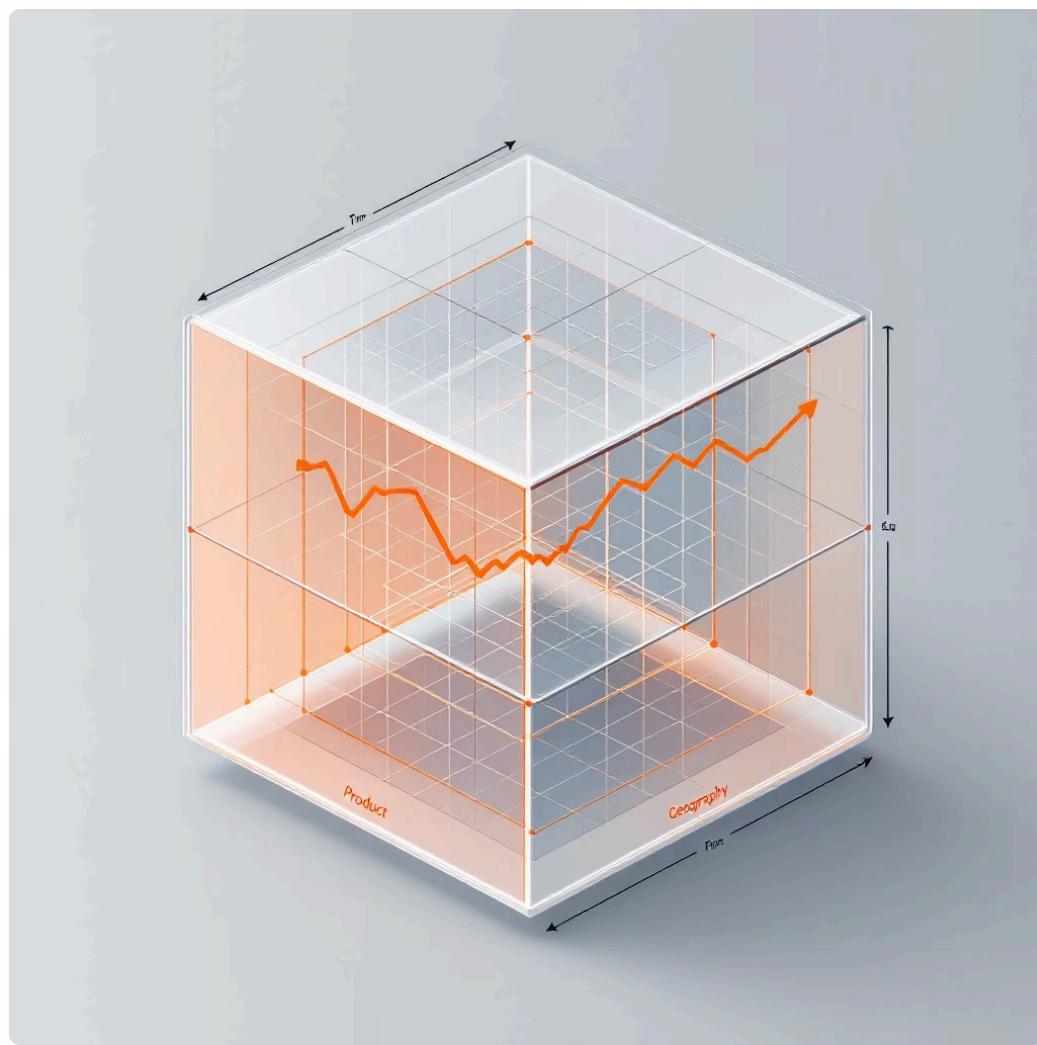
ETL ensures data quality and consistency, making the data reliable for accurate business intelligence and analytical insights.

OLAP: Powering Multidimensional Analysis

OLAP technology goes beyond traditional SQL queries, enabling fast, flexible, and interactive multidimensional analysis of business data.

How OLAP Works

OLAP uses a conceptual data structure known as an **OLAP cube**. This cube organizes data by multiple dimensions (e.g., Time, Product, Region) and contains numerical facts (measures) at their intersections (e.g., Sales Quantity, Revenue).



This structure allows users to explore data from different angles, providing deep insights into business performance.

Key Analytical Operations

- **Roll-Up:** Summarizing data by aggregating to a higher level of hierarchy (e.g., daily sales to monthly sales).
- **Drill-Down:** Navigating from summarized data to more detailed data (e.g., monthly sales to daily sales, then to individual transactions).
- **Slice:** Selecting a specific dimension value to view data for that particular slice (e.g., sales for a specific product category).
- **Dice:** Creating a sub-cube by selecting specific values across multiple dimensions (e.g., sales of specific products in a particular region during a given quarter).
- **Pivot (Rotate):** Changing the dimensional orientation of the cube view (e.g., changing the rows and columns in a cross-tabulated report to view sales by product instead of by region).

OLAP tools empower business users with intuitive, interactive analysis capabilities, crucial for quick insights and flexible reporting.

Analysis & Dimension Types

Effective data warehousing supports diverse analytical needs and leverages specialized dimension types to provide comprehensive data context.

Analysis Types



Ad Hoc Analysis

On-demand, custom queries for specific, often unplanned, business questions (e.g., "Why did sales drop in Region X last week?").



Interactive Analysis

Structured exploration of data, typically through dashboards, reports, and visualizations, allowing users to drill down or filter within predefined contexts.

Key Dimension Types

- **Conformed Dimension:** A dimension (e.g., Date, Customer) that is shared across multiple fact tables, ensuring consistent analysis across different business processes. This is key for enterprise-wide BI.
- **Degenerate Dimension:** An operational transaction identifier (e.g., Order ID, Ticket Number) that has no corresponding dimension table and is stored directly in the fact table.
- **Slowly Changing Dimension (SCD):** Handles changes to dimension attributes over time (e.g., a customer's address changing).
 - **Type 1:** Overwrite the old value (no history).
 - **Type 2:** Create a new row for each change, preserving full history (most common for analytical purposes).

The right combination of analysis techniques and carefully modeled dimensions ensures that data warehouses deliver comprehensive, accurate, and adaptable insights for all business needs.

Key Takeaways: Empowering Data-Driven Decisions

Data warehousing is a critical discipline for organizations aiming to leverage their data assets for strategic advantage. It transforms disparate, raw data into a cohesive, analytical resource.

Centralized Data Hub

A data warehouse serves as a unified, structured repository, consolidating data from various operational systems for analytical purposes.

Dimensional Clarity

Dimensional modeling (Star/Snowflake schemas) provides an intuitive and efficient framework for querying and understanding business data, crucial for BI.

ETL & OLAP Power

ETL processes ensure data quality and readiness, while OLAP capabilities enable fast, multidimensional analysis and interactive data exploration.

Strategic Insights

Ultimately, data warehousing empowers organizations with reliable, accessible data, leading to deeper insights, better decision-making, and improved business performance.

By systematically collecting, transforming, and organizing data, data warehouses provide the foundation for powerful analytics and a competitive edge in today's data-rich environment.

Thank You!

We hope this presentation has shed light on the immense value of data warehousing in driving informed, strategic decisions.

Embracing a robust data architecture empowers your organization to unlock deeper insights and maintain a competitive edge.

I'm happy to answer any questions you may have.