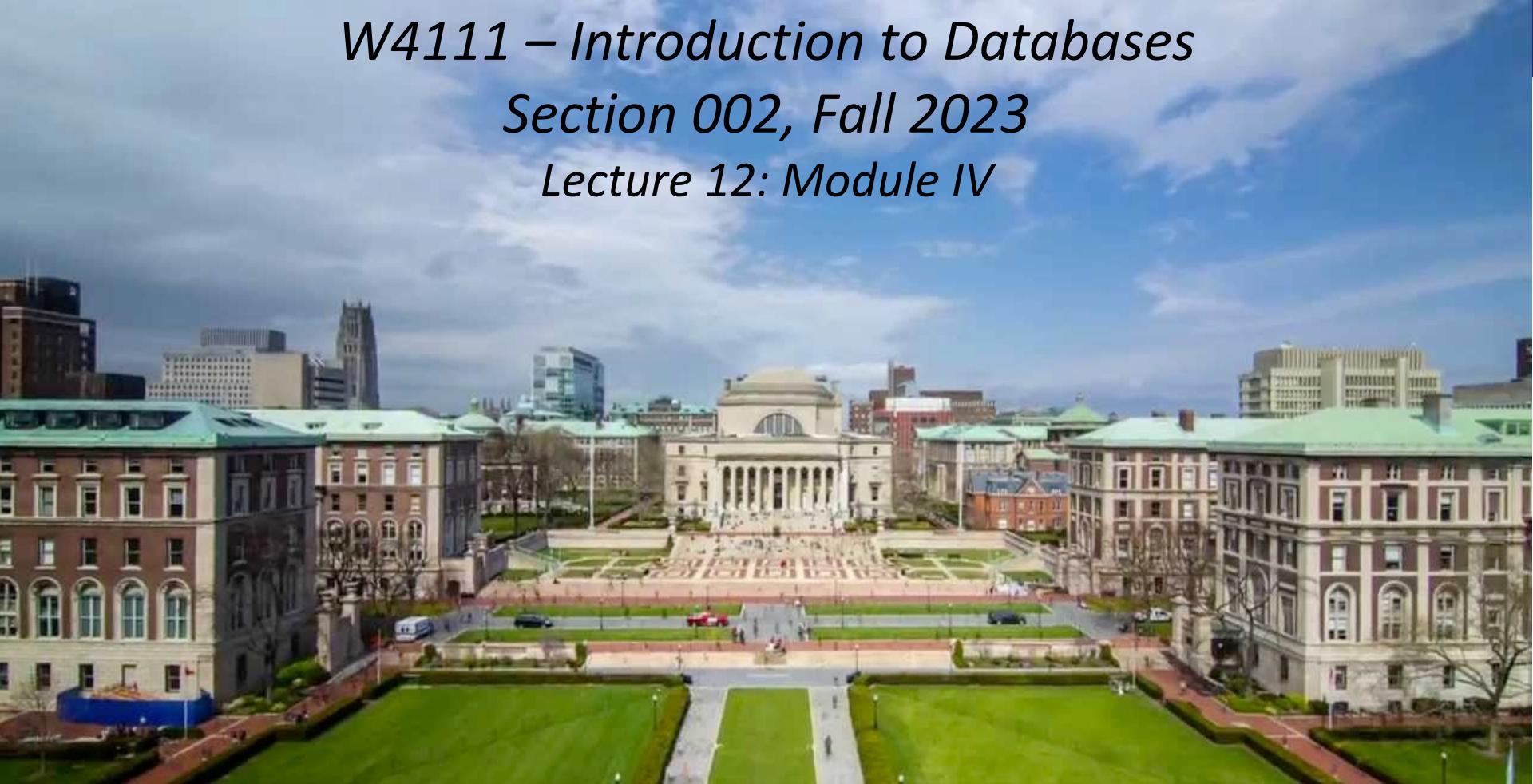


W4111 – Introduction to Databases
Section 002, Fall 2023
Lecture 12: Module IV



Contents

Relational Database Design Normalization



Overview of Normalization



Features of Good Relational Designs

- Hypothetically, assume our schema originally had one relation that provided information about faculty and departments.
- Suppose we combine *instructor* and *department* into *in_dep*, which represents the natural join on the relations *instructor* and *department*

| <i>ID</i> | <i>name</i> | <i>salary</i> | <i>dept_name</i> | <i>building</i> | <i>budget</i> |
|-----------|-------------|---------------|------------------|-----------------|---------------|
| 22222 | Einstein | 95000 | Physics | Watson | 70000 |
| 12121 | Wu | 90000 | Finance | Painter | 120000 |
| 32343 | El Said | 60000 | History | Painter | 50000 |
| 45565 | Katz | 75000 | Comp. Sci. | Taylor | 100000 |
| 98345 | Kim | 80000 | Elec. Eng. | Taylor | 85000 |
| 76766 | Crick | 72000 | Biology | Watson | 90000 |
| 10101 | Srinivasan | 65000 | Comp. Sci. | Taylor | 100000 |
| 58583 | Califieri | 62000 | History | Painter | 50000 |
| 83821 | Brandt | 92000 | Comp. Sci. | Taylor | 100000 |
| 15151 | Mozart | 40000 | Music | Packard | 80000 |
| 33456 | Gold | 87000 | Physics | Watson | 70000 |
| 76543 | Singh | 80000 | Finance | Painter | 120000 |

$X \rightarrow Y$

$X = (\text{dept_name})$

$Y = (\text{building}, \text{budget})$

Pix Physics

Pix Physics

Piy (Watson, 70000)

- There is repetition of information
- Need to use null values (if we add a new department with no instructors)



Decomposition

- The only way to avoid the repetition-of-information problem in the *in_dep* schema is to decompose it into two schemas – instructor and *department* schemas.
- Not all decompositions are good. Suppose we decompose

employee(*ID*, *name*, *street*, *city*, *salary*)

into

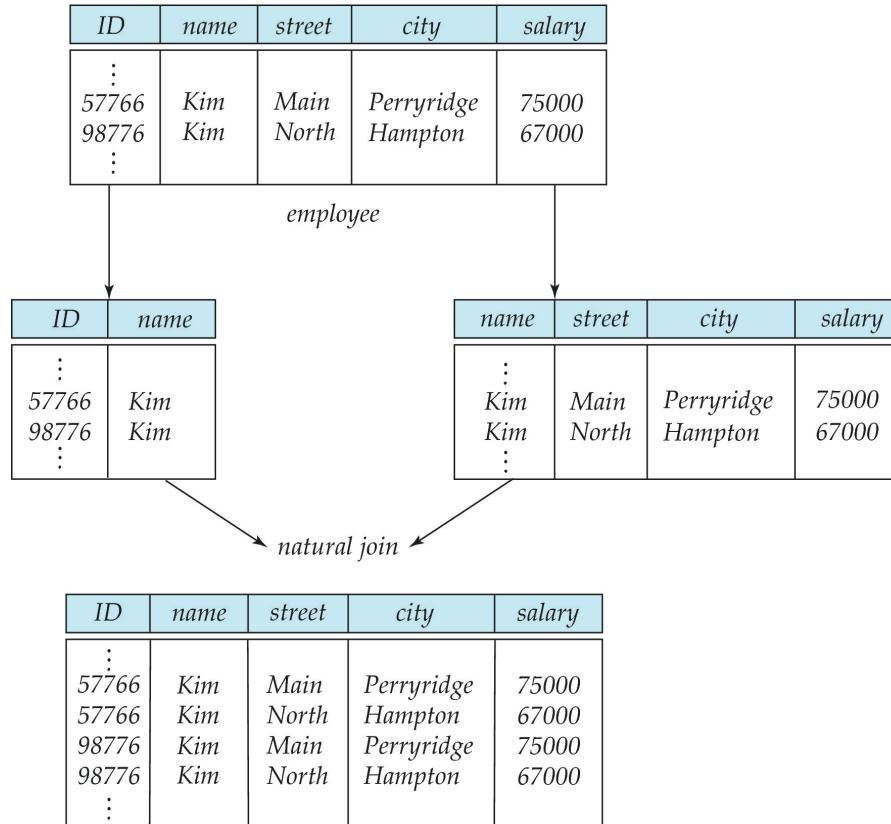
employee1 (*ID*, *name*)

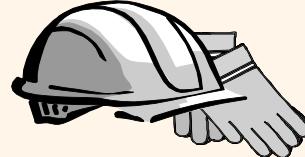
employee2 (*name*, *street*, *city*, *salary*)

- The problem arises when we have two employees with the same name
- The next slide shows how we lose information -- we cannot reconstruct the original *employee* relation -- and so, this is a **lossy decomposition**.



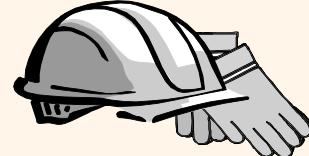
Lossy Decomposition





The Evils of Redundancy

- ❖ *Redundancy* is at the root of several problems associated with relational schemas:
 - redundant storage, insert/delete/update anomalies
- ❖ Integrity constraints, in particular *functional dependencies*, can be used to identify schemas with such problems and to suggest refinements.
- ❖ Main refinement technique: *decomposition* (replacing ABCD with, say, AB and BCD, or ACD and ABD).
- ❖ Decomposition should be used judiciously:
 - Is there reason to decompose a relation?
 - What problems (if any) does the decomposition cause?



Example (Contd.)

- ❖ Problems due to $R \rightarrow W$:
 - Update anomaly: Can we change W in just the 1st tuple of SNLRWH?
 - Insertion anomaly: What if we want to insert an employee and don't know the hourly wage for his rating?
 - Deletion anomaly: If we delete all employees with rating 5, we lose the information about the wage for rating 5!

Will 2 smaller tables be better?

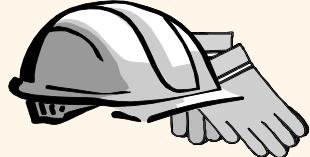
Wages

| R | W |
|---|----|
| 8 | 10 |
| 5 | 7 |

Hourly_Emps2

| S | N | L | R | H |
|-------------|-----------|----|---|----|
| 123-22-3666 | Attishoo | 48 | 8 | 40 |
| 231-31-5368 | Smiley | 22 | 8 | 30 |
| 131-24-3650 | Smethurst | 35 | 5 | 30 |
| 434-26-3751 | Guldu | 35 | 5 | 32 |
| 612-67-4134 | Madayan | 35 | 8 | 40 |

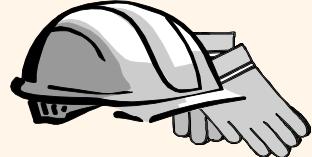
| S | N | L | R | W | H |
|-------------|-----------|----|---|----|----|
| 123-22-3666 | Attishoo | 48 | 8 | 10 | 40 |
| 231-31-5368 | Smiley | 22 | 8 | 10 | 30 |
| 131-24-3650 | Smethurst | 35 | 5 | 7 | 30 |
| 434-26-3751 | Guldu | 35 | 5 | 7 | 32 |
| 612-67-4134 | Madayan | 35 | 8 | 10 | 40 |



Functional Dependencies (FDs)

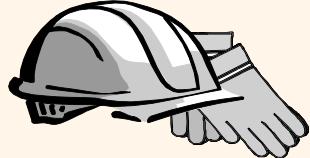
- ❖ A functional dependency $X \rightarrow Y$ holds over relation R if, for every allowable instance r of R:
 - $t1 \in r, t2 \in r, \pi_X(t1) = \pi_X(t2)$ implies $\pi_Y(t1) = \pi_Y(t2)$
 - i.e., given two tuples in r , if the X values agree, then the Y values must also agree. (X and Y are *sets* of attributes.)
- ❖ An FD is a statement about *all* allowable relations.
 - Must be identified based on semantics of application.
 - Given some allowable instance $r1$ of R, we can check if it violates some FD f , but we cannot tell if f holds over R!
- ❖ K is a candidate key for R means that $K \rightarrow R$
 - However, $K \rightarrow R$ does not require K to be *minimal*!

Example: Constraints on Entity Set



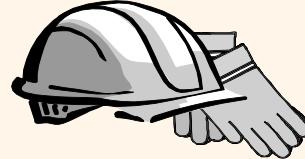
- ❖ Consider relation obtained from Hourly_Emps:
 - Hourly_Emps (ssn, name, lot, rating, hrly_wages, hrs_worked)
- ❖ Notation: We will denote this relation schema by listing the attributes: **SNLRWH**
 - This is really the *set* of attributes {S,N,L,R,W,H}.
 - Sometimes, we will refer to all attributes of a relation by using the relation name. (e.g., Hourly_Emps for SNLRWH)
- ❖ Some FDs on Hourly_Emps:
 - *ssn* is the key: $S \rightarrow \text{SNLRWH}$
 - *rating* determines *hrly_wages*: $R \rightarrow W$

- (uni, email, last_name, first_name)
uni and email are candidate (not NULL, unique)
 - Uni → email
 - Email → (last_name, first_name)
 - ➔ Uni → (last_name, first_name)
- Person:(uni, email, last_name, first_name),
ContactInfo: (email, phone_number, dorm_name)
 - Uni → email
 - Email → phone_number, dorm_name
 - Uni → phone_number, dorm_name



Reasoning About FDs

- ❖ Given some FDs, we can usually infer additional FDs:
 - $ssn \rightarrow did$, $did \rightarrow lot$ implies $ssn \rightarrow lot$
- ❖ An FD f is implied by a set of FDs F if f holds whenever all FDs in F hold.
 - $F^+ = \text{closure of } F$ is the set of all FDs that are implied by F .
- ❖ Armstrong's Axioms (X, Y, Z are sets of attributes):
 - Reflexivity: If $X \subseteq Y$, then $Y \rightarrow X$
 - Augmentation: If $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z
 - Transitivity: If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$
- ❖ These are sound and complete inference rules for FDs!



Lossless Decomposition

- ❖ Let R be a relation schema and let R_1 and R_2 form a decomposition of R .
 - That is $R = R_1 \cup R_2$
 - *Note: This notation is confusing. This is a statement about the schema, not about the data in relations.*
- ❖ We say that the decomposition is a **lossless decomposition** if there is no loss of information by replacing R with the two relation schemas $R_1 \cup R_2$
- ❖ Formally,

$$\Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$$

- ❖ And, conversely a decomposition is lossy if

$$r \subset \Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$$



Normalization Theory

- Decide whether a particular relation R is in “good” form.
- In the case that a relation R is not in “good” form, decompose it into set of relations $\{R_1, R_2, \dots, R_n\}$ such that
 - Each relation is in good form
 - The decomposition is a lossless decomposition
- Our theory is based on:
 - Functional dependencies
 - Multivalued dependencies



Functional Dependencies

- There are usually a variety of constraints (rules) on the data in the real world.
- For example, some of the constraints that are expected to hold in a university database are:
 - Students and instructors are uniquely identified by their ID.
 - Each student and instructor has only one name.
 - Each instructor and student is (primarily) associated with only one department.
 - Each department has only one value for its budget, and only one associated building.



Functional Dependencies (Cont.)

- An instance of a relation that satisfies all such real-world constraints is called a **legal instance** of the relation;
- A legal instance of a database is one where all the relation instances are legal instances
- Constraints on the set of legal relations.
- Require that the value for a certain set of attributes determines uniquely the value for another set of attributes.
- A functional dependency is a generalization of the notion of a *key*.



Functional Dependencies Definition

- Let R be a relation schema

$$\alpha \subseteq R \text{ and } \beta \subseteq R$$

- The **functional dependency**

$$\alpha \rightarrow \beta$$

holds on R if and only if for any legal relations $r(R)$, whenever any two tuples t_1 and t_2 of r agree on the attributes α , they also agree on the attributes β . That is,

$$t_1[\alpha] = t_2[\alpha] \Rightarrow t_1[\beta] = t_2[\beta]$$

- Example: Consider $r(A,B)$ with the following instance of r .

| | |
|---|---|
| 1 | 4 |
| 1 | 5 |
| 3 | 7 |

- On this instance, $B \rightarrow A$ hold; $A \rightarrow B$ does **NOT** hold,



Closure of a Set of Functional Dependencies

- Given a set F set of functional dependencies, there are certain other functional dependencies that are logically implied by F .
 - If $A \rightarrow B$ and $B \rightarrow C$, then we can infer that $A \rightarrow C$
 - etc.
- The set of **all** functional dependencies logically implied by F is the **closure** of F .
- We denote the *closure* of F by F^+ .



Keys and Functional Dependencies

- K is a superkey for relation schema R if and only if $K \rightarrow R$
- K is a candidate key for R if and only if
 - $K \rightarrow R$, and
 - for no $\alpha \subset K$, $\alpha \rightarrow R$
- Functional dependencies allow us to express constraints that cannot be expressed using superkeys. Consider the schema:

in_dep (ID, name, salary, dept_name, building, budget).

We expect these functional dependencies to hold:

$dept_name \rightarrow building$

$ID \rightarrow building$

but would not expect the following to hold:

$dept_name \rightarrow salary$

DFF Note:

- In the current data: $ID \rightarrow dept_name \rightarrow building$
- But
- In the schema, $dept_name$ may be NULL..



Use of Functional Dependencies

- We use functional dependencies to:
 - To test relations to see if they are legal under a given set of functional dependencies.
 - ▶ If a relation r is legal under a set F of functional dependencies, we say that r **satisfies** F .
 - To specify constraints on the set of legal relations
 - ▶ We say that F **holds on** R if all legal relations on R satisfy the set of functional dependencies F .
- Note: A specific instance of a relation schema may satisfy a functional dependency even if the functional dependency does not hold on all legal instances.
 - For example, a specific instance of *instructor* may, by chance, satisfy $name \rightarrow ID$.



Normal Forms



Boyce-Codd Normal Form

- A relation schema R is in BCNF with respect to a set F of functional dependencies if for all functional dependencies in F^+ of the form

$$\alpha \rightarrow \beta$$

where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following holds:

- $\alpha \rightarrow \beta$ is trivial (i.e., $\beta \subseteq \alpha$)
- α is a superkey for R

DFF Note:

- The theoretical treatment is no conveying the practical intent.
- If α is a superkey, I can set a primary key/unique constraint on α .
- Consider tables with address info in the rows.



Decomposing a Schema into BCNF

- Let R be a schema R that is not in BCNF. Let $\alpha \rightarrow \beta$ be the FD that causes a violation of BCNF.
- We decompose R into:
 - $(\alpha \cup \beta)$
 - $(R - (\beta - \alpha))$
- In our example of *in_dep*,
 - $\alpha = \text{dept_name}$
 - $\beta = \text{building}, \text{budget}$and *in_dep* is replaced by
 - $(\alpha \cup \beta) = (\text{dept_name}, \text{building}, \text{budget})$
 - $(R - (\beta - \alpha)) = (\text{ID}, \text{name}, \text{dept_name}, \text{salary})$

DFF Note – again this is baffling

- $R = (\text{id}, \text{name_last}, \text{name_first}, \text{street}, \text{city}, \text{state}, \text{zipcode})$
- $\alpha \rightarrow \beta$ means $\text{zipcode} \rightarrow (\text{city}, \text{state})$
- $(\alpha \cup \beta) = (\text{zipcode}, \text{city}, \text{state})$, which we call Address
- $R - (\beta - \alpha) = R - \beta + \alpha = (\text{id}, \text{name_last}, \text{name_first}, \text{zipcode})$, which we call Person
- Setting primary key ID in Address and zipcode on Address preserves the dependency and is lossless.

Note:

- This is an example only.
- Zip code does not imply city or state in real world.



BCNF and Dependency Preservation

- It is not always possible to achieve both BCNF and dependency preservation
- Consider a schema:
$$\text{dept_advisor}(s_ID, i_ID, \text{department_name})$$
- With function dependencies:
$$i_ID \rightarrow \text{dept_name}$$

$$s_ID, \text{dept_name} \rightarrow i_ID$$
- dept_advisor is not in BCNF
 - i_ID is not a superkey.
- Any decomposition of dept_advisor will not include all the attributes in
$$s_ID, \text{dept_name} \rightarrow i_ID$$
- Thus, the composition is NOT be dependency preserving



Third Normal Form

- A relation schema R is in **third normal form (3NF)** if for all:

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- $\alpha \rightarrow \beta$ is trivial (i.e., $\beta \in \alpha$)
- α is a superkey for R
- Each attribute A in $\beta - \alpha$ is contained in a candidate key for R .

(**NOTE**: each attribute may be in a different candidate key)

- If a relation is in BCNF it is in 3NF (since in BCNF one of the first two conditions above must hold).
- Third condition is a minimal relaxation of BCNF to ensure dependency preservation (will see why later).



3NF Example

- Consider a schema:

$\text{dept_advisor}(s_ID, i_ID, \text{dept_name})$

- With function dependencies:

$i_ID \rightarrow \text{dept_name}$

$s_ID, \text{dept_name} \rightarrow i_ID$

- Two candidate keys = $\{s_ID, \text{dept_name}\}, \{s_ID, i_ID\}$
- We have seen before that dept_advisor is not in BCNF
- R , however, is in 3NF
 - $s_ID, \text{dept_name}$ is a superkey
 - $i_ID \rightarrow \text{dept_name}$ and i_ID is NOT a superkey, but:
 - ▶ $\{ \text{dept_name} \} - \{i_ID\} = \{\text{dept_name}\}$ and
 - ▶ dept_name is contained in a candidate key



Comparison of BCNF and 3NF

- Advantages to 3NF over BCNF. It is always possible to obtain a 3NF design without sacrificing losslessness or dependency preservation.
- Disadvantages to 3NF.
 - We may have to use null values to represent some of the possible meaningful relationships among data items.
 - There is the problem of repetition of information.



Goals of Normalization

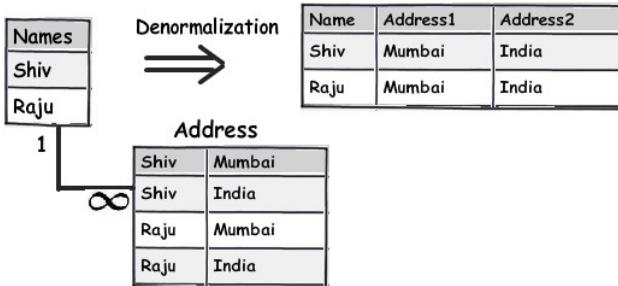
- Let R be a relation scheme with a set F of functional dependencies.
- Decide whether a relation scheme R is in “good” form.
- In the case that a relation scheme R is not in “good” form, need to decompose it into a set of relation scheme $\{R_1, R_2, \dots, R_n\}$ such that:
 - Each relation scheme is in good form
 - The decomposition is a lossless decomposition
 - Preferably, the decomposition should be dependency preserving.



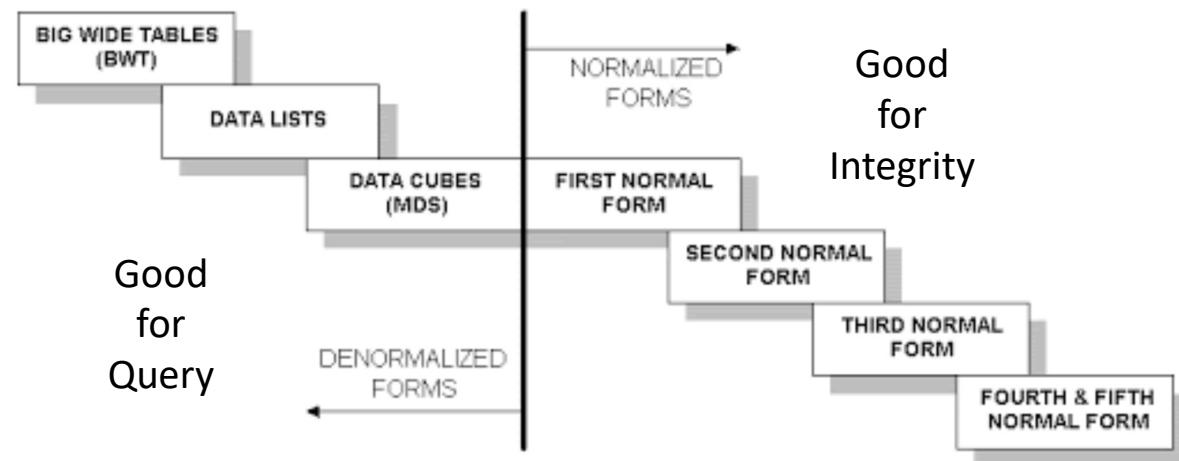
Denormalization for Performance

- May want to use non-normalized schema for performance
- For example, displaying *prereqs* along with *course_id*, and *title* requires join of *course* with *prereq*
- Alternative 1: Use denormalized relation containing attributes of *course* as well as *prereq* with all above attributes
 - faster lookup
 - extra space and extra execution time for updates
 - extra coding work for programmer and possibility of error in extra code
- Alternative 2: use a materialized view defined a $course \bowtie prereq$
 - Benefits and drawbacks same as above, except no extra coding work for programmer and avoids possible errors

Wide Flat Tables



- Improve query performance by precomputing and saving:
 - JOINs
 - Aggregation
 - Derived/computed columns
- One of the primary strength of the relational model is maintaining “integrity” when applications create, update and delete data. This relies on:
 - The core capabilities of the relational model, e.g. constraints.
 - A well-designed database (We will cover a formal definition – “normalization” in more detail later.)
- Data models that are well designed for integrity are very bad for read only analysis queries.
We will build and analyze wide flat tables as part of the analysis tasks in HW3, HW4 as projects.





First Normal Form

- Domain is **atomic** if its elements are considered to be indivisible units
 - Examples of non-atomic domains:
 - ▶ Set of names, composite attributes
 - ▶ Identification numbers like CS 101 that can be broken up into parts
- A relational schema R is in **first normal form** if the domains of all attributes of R are atomic
- Non-atomic values complicate storage and encourage redundant (repeated) storage of data
 - Example: Set of accounts stored with each customer, and set of owners stored with each account
 - We assume all relations are in first normal form (and revisit this in Chapter 22: Object Based Databases)



First Normal Form (Cont.)

- Atomicity is actually a property of how the elements of the domain are used.
 - Example: Strings would normally be considered indivisible
 - Suppose that students are given roll numbers which are strings of the form *CS0012* or *EE1127*
 - If the first two characters are extracted to find the department, the domain of roll numbers is not atomic.
 - Doing so is a bad idea: leads to encoding of information in application program rather than in the database.

Big Data
Business Intelligence
Decision Support

Core Concepts



Chapter 20: Data Analysis

- Decision Support Systems
- Data Warehousing
- ~~Data Mining~~
- ~~Classification~~
- ~~Association Rules~~
- ~~Clustering~~



Decision Support Systems

- **Decision-support systems** are used to make business decisions, often based on data collected by on-line transaction-processing systems.
- Examples of business decisions:
 - What items to stock?
 - What insurance premium to change?
 - To whom to send advertisements?
- Examples of data used for making decisions
 - Retail sales transaction details
 - Customer profiles (income, age, gender, etc.)



Decision-Support Systems: Overview

- **Data analysis** tasks are simplified by specialized tools and SQL extensions
 - Example tasks
 - ▶ For each product category and each region, what were the total sales in the last quarter and how do they compare with the same quarter last year
 - ▶ As above, for each product category and each customer category
- **Statistical analysis** packages (e.g., : S++) can be interfaced with databases
 - Statistical analysis is a large field, but not covered here
- **Data mining** seeks to discover knowledge automatically in the form of statistical rules and patterns from large databases.
- A **data warehouse** archives information gathered from multiple sources, and stores it under a unified schema, at a single site.
 - Important for large businesses that generate data from multiple divisions, possibly at multiple sites
 - Data may also be purchased externally

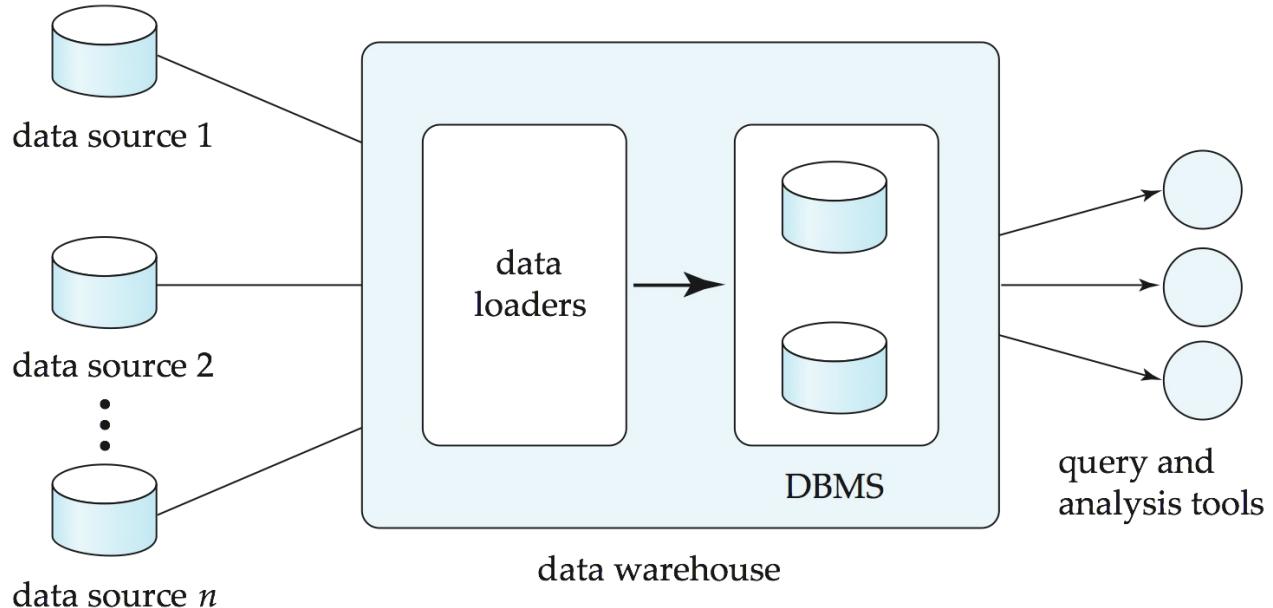


Data Warehousing

- Data sources often store only current data, not historical data
- Corporate decision making requires a unified view of all organizational data, including historical data
- A **data warehouse** is a repository (archive) of information gathered from multiple sources, stored under a unified schema, at a single site
 - Greatly simplifies querying, permits study of historical trends
 - Shifts decision support query load away from transaction processing systems



Data Warehousing



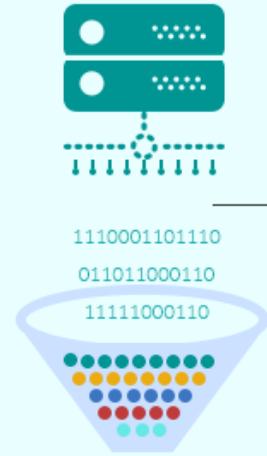
Data Warehouse vs Data Lake

<https://www.grazitti.com/blog/data-lake-vs-data-warehouse-which-one-should-you-go-for/>

DATA WAREHOUSE

VS

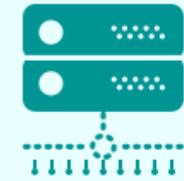
DATA LAKE



- Data is processed and organized into a single schema before being put into the warehouse



- The analysis is done on the cleansed data in the warehouse



- Raw and unstructured data goes into a data lake

- Data is selected and organized as and when needed



Simplistic: Data Warehouse vs Data Lake

<https://panoply.io/data-warehouse-guide/data-warehouse-vs-data-lake/>

| DATA WAREHOUSE | vs. | DATA LAKE |
|----------------------------------|------------|---|
| structured, processed | DATA | structured / semi-structured / unstructured, raw |
| schema-on-write | PROCESSING | schema-on-read |
| expensive for large data volumes | STORAGE | designed for low-cost storage |
| less agile, fixed configuration | AGILITY | highly agile, configure and reconfigure as needed |
| mature | SECURITY | maturing |
| business professionals | USERS | data scientists et. al. |



Design Issues

■ When and how to gather data

- **Source driven architecture**: data sources transmit new information to warehouse, either continuously or periodically (e.g., at night)
- **Destination driven architecture**: warehouse periodically requests new information from data sources
- Keeping warehouse exactly synchronized with data sources (e.g., using two-phase commit) is too expensive
 - ▶ Usually OK to have slightly out-of-date data at warehouse
 - ▶ Data/updates are periodically downloaded from online transaction processing (OLTP) systems.

■ What schema to use

- Schema integration



More Warehouse Design Issues

■ *Data cleansing*

- E.g., correct mistakes in addresses (misspellings, zip code errors)
- **Merge** address lists from different sources and **purge** duplicates

■ *How to propagate updates*

- Warehouse schema may be a (materialized) view of schema from data sources

■ *What data to summarize*

- Raw data may be too large to store on-line
- Aggregate values (totals/subtotals) often suffice
- Queries on raw data can often be transformed by query optimizer to use aggregate values

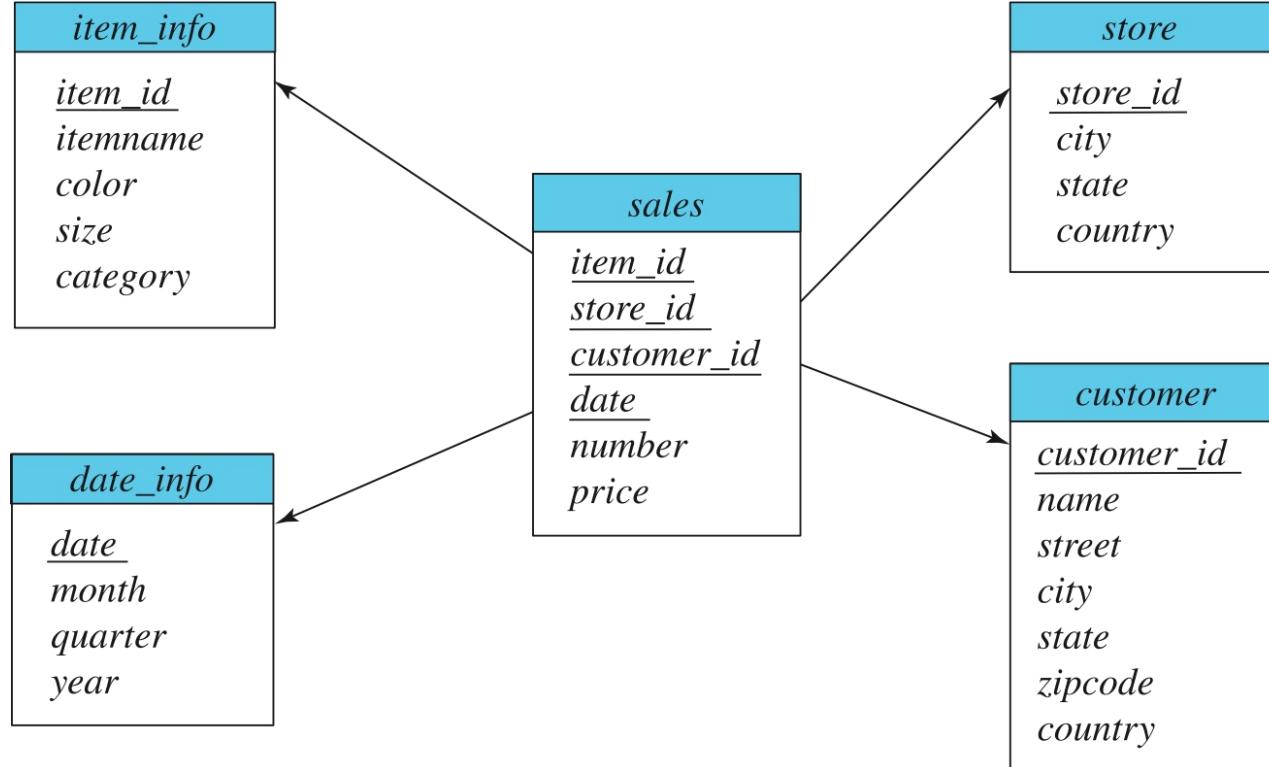


Warehouse Schemas

- Dimension values are usually encoded using small integers and mapped to full values via dimension tables
- Resultant schema is called a **star schema**
 - More complicated schema structures
 - ▶ **Snowflake schema**: multiple levels of dimension tables
 - ▶ **Constellation**: multiple fact tables



Data Warehouse Schema



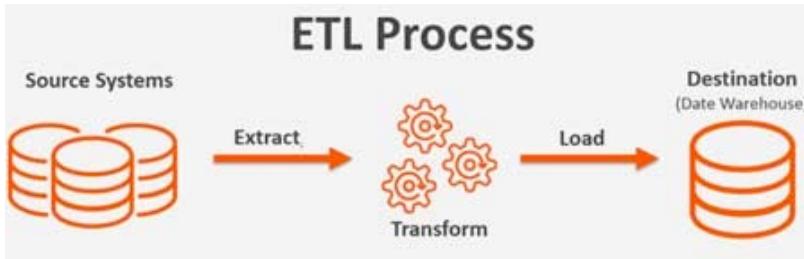


Overview (Cont.)

- Common steps in data analytics
 - Gather data from multiple sources into one location
 - Data warehouses also integrated data into common schema
 - Data often needs to be **extracted** from source formats, **transformed** to common schema, and **loaded** into the data warehouse
 - Can be done as **ETL (extract-transform-load)**, or **ELT (extract-load-transform)**
 - Generate aggregates and reports summarizing data
 - Dashboards showing graphical charts/reports
 - **Online analytical processing (OLAP) systems** allow interactive querying
 - Statistical analysis using tools such as R/SAS/SPSS
 - Including extensions for parallel processing of big data
 - Build **predictive models** and use the models for decision making

ETL Concepts

<https://databricks.com/glossary/extract-transform-load>



Extract

The first step of this process is extracting data from the target sources that could include an ERP, CRM, Streaming sources, and other enterprise systems as well as data from third-party sources. There are different ways to perform the extraction: **Three Data Extraction methods:**

1. Partial Extraction – The easiest way to obtain the data is if the source system notifies you when a record has been changed
2. Partial Extraction- with update notification – Not all systems can provide a notification in case an update has taken place; however, they can point those records that have been changed and provide an extract of such records.
3. Full extract – There are certain systems that cannot identify which data has been changed at all. In this case, a full extract is the only possibility to extract the data out of the system. This method requires having a copy of the last extract in the same format so you can identify the changes that have been made.

Transform

Next, the transform function converts the raw data that has been extracted from the source server. As it cannot be used in its original form in this stage it gets cleansed, mapped and transformed, often to a specific data schema, so it will meet operational needs. This process entails several transformation types that ensure the quality and integrity of data; below are the most common as well as advanced transformation types that prepare data for analysis:

- Basic transformations:
- Cleaning
- Format revision
- Data threshold validation checks
- Restructuring
- Deduplication
- Advanced transformations:
- Filtering
- Merging
- Splitting
- Derivation
- Summarization
- Integration
- Aggregation
- Complex data validation

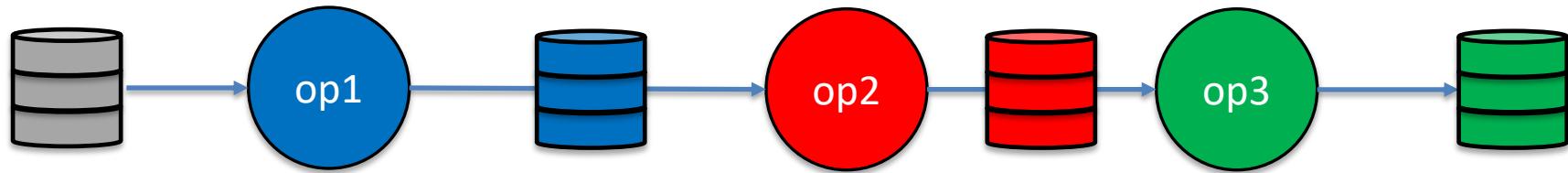
Load

Finally, the load function is the process of writing converted data from a staging area to a target database, which may or may not have previously existed. Depending on the requirements of the application, this process may be either quite simple or intricate.

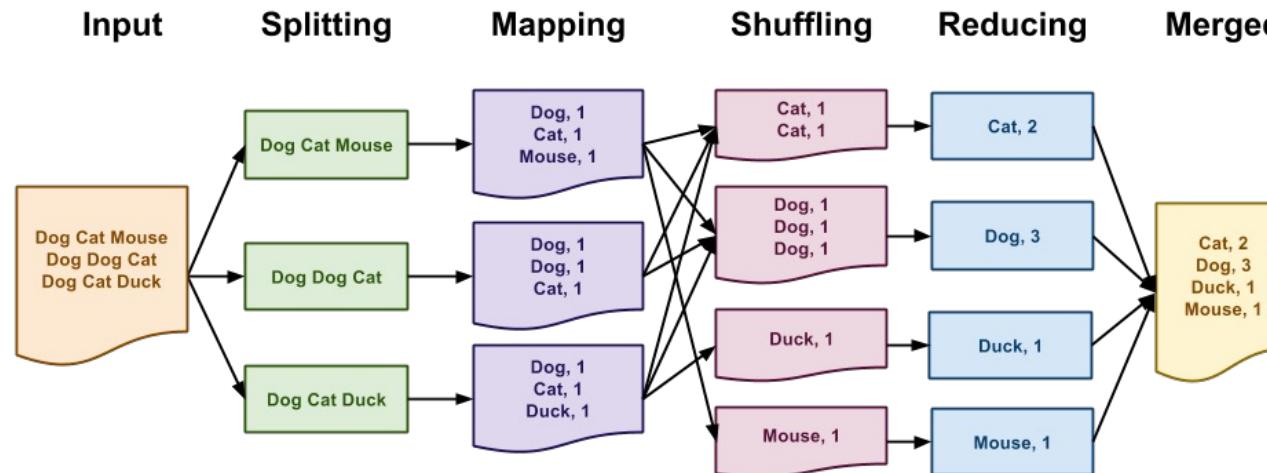
MapReduce, Spark (and ETL)

MapReduce

MapReduce is a data flow program with relatively simple operators on the data set.



With each operator implemented in parallel on multiple nodes for performance.



What if we want more complex “operators?”

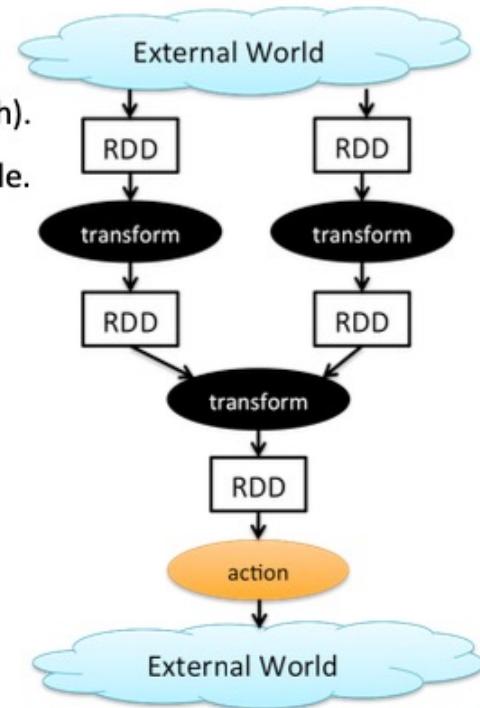
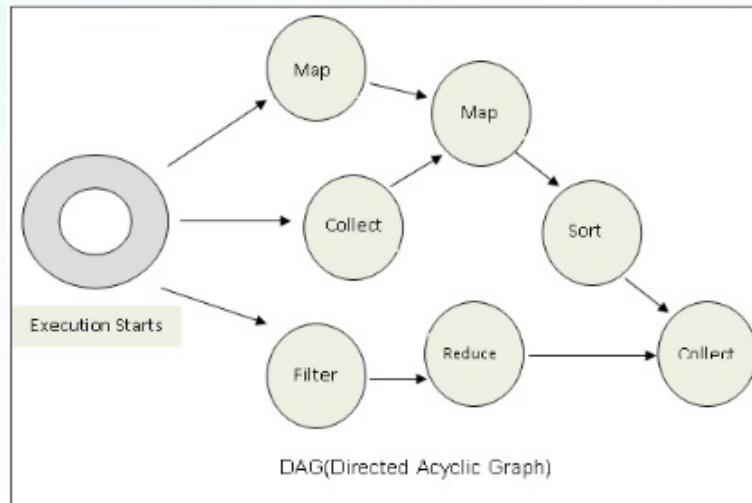


Algebraic Operations

- Current generation execution engines
 - natively support algebraic operations such as joins, aggregation, etc. natively.
 - Allow users to create their **own algebraic operators**
 - Support trees of algebraic operators that can be executed on **multiple nodes in parallel**
- E.g. Apache Tez, Spark
 - Tez provides low level API; Hive on Tez compiles SQL to Tez
 - Spark provides more user-friendly API

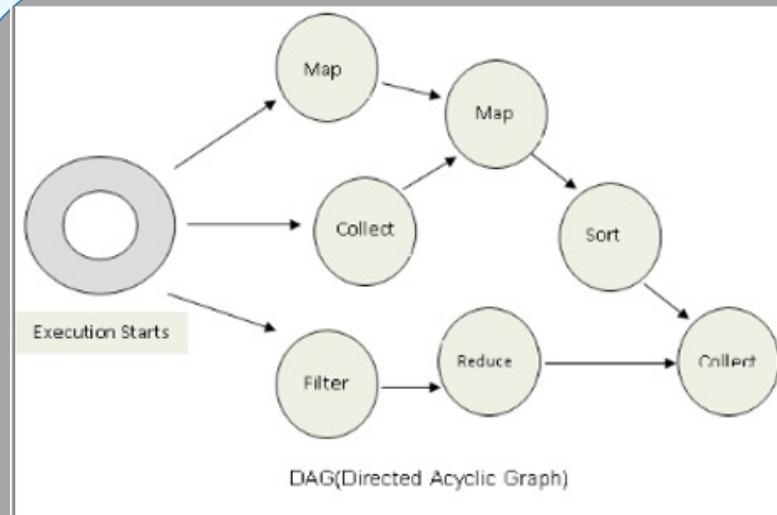
- In the relational model,
 - The are relations.
 - A fixed set of operators that produce relations from relations (closed).
 - Are declarative languages and compute the execution graph/plan.
- The Spark/... ... engines:
 - Enable programmers to develop and add new operators.
 - Operators convert reliable distributed datasets/data frames to new RDDs/data frames.
 - Data engineer explicitly defines the execution graph (data flow).

- All jobs in spark comprise a series of operators and run on a set of data.
- All the operators in a job are used to construct a DAG (Directed Acyclic Graph).
- The DAG is optimized by rearranging and combining operators where possible.



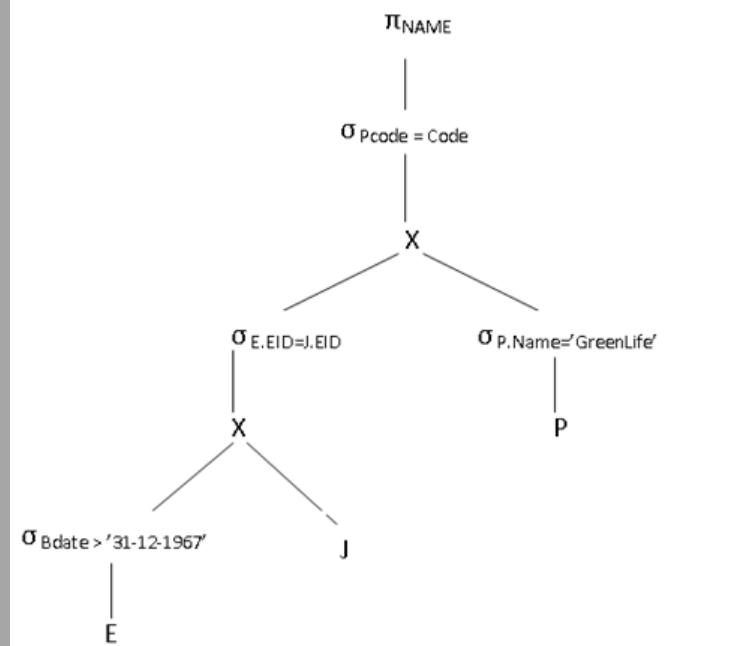
www.edureka.co/apache-spark-scala-training

Spark vs SQL



Conceptually similar to query evaluation graph, but ...

- You explicitly define the graph.
- You can develop your own operators.
- You can define levels of parallelism that the infrastructure manages.



SQL is declarative, and the query engine transparently produces the execution plan.

www.edureka.co/apache-spark-scala-training



Algebraic Operations in Spark

- **Resilient Distributed Dataset (RDD) abstraction**
 - Collection of records that can be stored across multiple machines
- RDDs can be created by applying algebraic operations on other RDDs
- RDDs can be lazily computed when needed
- ~~Spark programs can be written in Java/Scala/R~~
 - ▲ Our examples are in Java
- ~~Spark makes use of Java 8 Lambda expressions; the code~~
 - `s -> Arrays.asList(s.split(" ")).iterator()`
~~defines unnamed function that takes argument s and executes the expression `Arrays.asList(s.split(" ")).iterator()` on the argument~~
 - Lambda functions are particularly convenient as arguments to map, reduce and other functions

Spark/PySpark

DataFrame

edureka!

Inspired by DataFrames in R and Python (Pandas).

DataFrames API is designed to make big data processing on tabular data easier.

DataFrame is a distributed collection of data organized into named columns.

Provides operations to filter, group, or compute aggregates, and can be used with Spark SQL.

Can be constructed from structured data files, existing RDDs, tables in Hive, or external databases.

1. DataFrame in PySpark: Overview

In Apache Spark, a DataFrame is a distributed collection of rows under named columns. In simple terms, it is same as a table in relational database or an Excel sheet with Column headers. It also shares some common characteristics with RDD:

- Immutable in nature** : We can create DataFrame / RDD once but can't change it. And we can transform a DataFrame / RDD after applying transformations.
- Lazy Evaluations**: Which means that a task is not executed until an action is performed.
- Distributed**: RDD and DataFrame both are distributed in nature.

My first exposure to DataFrames was when I learnt about Pandas. Today, it is difficult for me to run my data science workflow with out Pandas DataFrames. So, when I saw similar functionality in Apache Spark, I was excited about the possibilities it opens up!

<https://www.analyticsvidhya.com/blog/2016/10/spark-dataframe-and-operations/>

DataFrame features

edureka!

Ability to scale from KBs to PBs

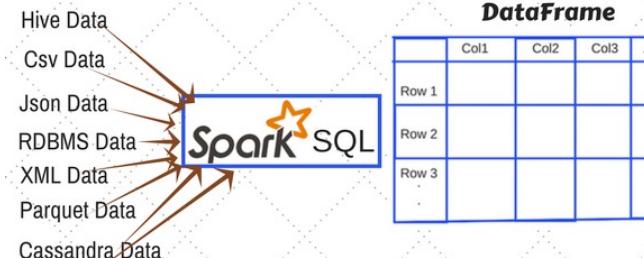
Support for a wide array of data formats and storage systems

State-of-the-art optimization and code generation through the spark SQL catalyst optimizer

Seamless integration with all big data tooling and infrastructure via spark

APIs for Python, Java, Scala, and R

Ways to Create DataFrame in Spark



Conceptual Example

- Remember that really unpleasant set of queries and stored procedure we had to write for the midterm? You had to convert from

| nconst | name | dateOfBirth | dateOfDeath | primaryProfession | knownFor |
|-----------|-----------------|-------------|-------------|-------------------------------------|---|
| nm0000003 | Brigitte Bardot | 1934 | | actress,soundtrack,music_department | tt0057345,tt0059956,tt0054452,tt0049189 |
| nm0000004 | John Belushi | 1949 | 1982 | actor,soundtrack,writer | tt0077975,tt0072562,tt0080455,tt0078723 |
| nm0000005 | Ingmar Bergman | 1918 | 2007 | writer,director,actor | tt0050986,tt0083922,tt0060827,tt0050976 |
| nm0000006 | Ingrid Bergman | 1915 | 1982 | actress,soundtrack,producer | tt0038787,tt0036855,tt0034583,tt0038109 |
| nm0000007 | Humphrey Bogart | 1899 | 1957 | actor,soundtrack,producer | tt0043265,tt0040897,tt0034583,tt0037382 |
| nm0000008 | Marlon Brando | 1924 | 2004 | actor,soundtrack,director | tt0078788,tt0070849,tt0047296,tt0068646 |
| nm0000009 | Richard Burton | 1925 | 1984 | actor,soundtrack,producer | tt0087803,tt0061184,tt0059749,tt0057877 |
| nm0000010 | James Cagney | 1899 | 1986 | actor,soundtrack,director | tt0042041,tt0035575,tt0029870,tt0031867 |

- To

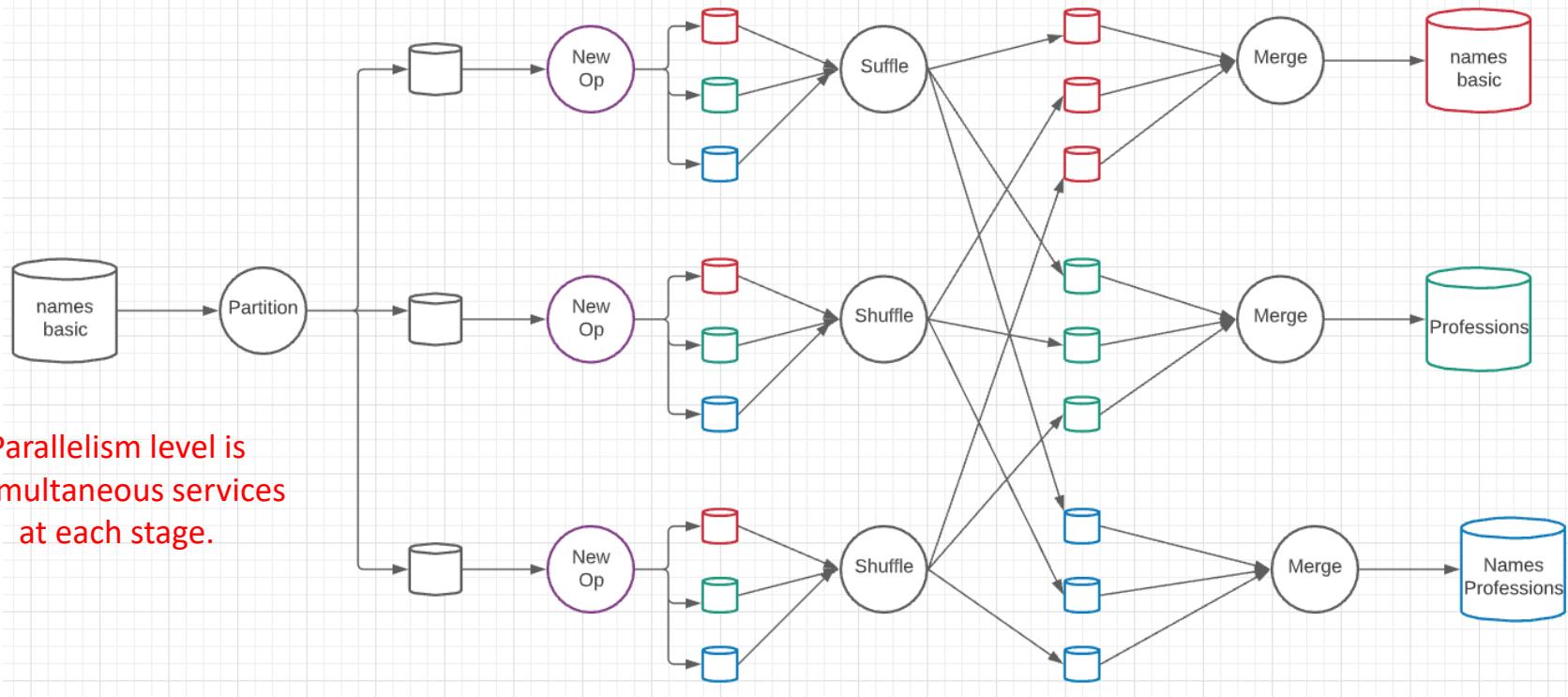
| | | | |
|-----------|-----------------|------|------|
| nm0000001 | Fred Astaire | 1899 | 1987 |
| nm0000002 | Lauren Bacall | 1924 | 2014 |
| nm0000003 | Brigitte Bardot | 1934 | |
| nm0000004 | John Belushi | 1949 | 1982 |
| nm0000005 | Ingmar Bergman | 1918 | 2007 |
| nm0000006 | Ingrid Bergman | 1915 | 1982 |
| nm0000007 | Humphrey Bogart | 1899 | 1957 |
| nm0000008 | Marlon Brando | 1924 | 2004 |
| nm0000009 | Richard Burton | 1925 | 1984 |
| nm0000010 | James Cagney | 1899 | 1986 |

| nconst | profession_id |
|-----------|---------------|
| nm0000001 | 1 |
| nm0000001 | 2 |
| nm0000001 | 3 |
| nm0000002 | 3 |
| nm0000002 | 4 |
| nm0000003 | 3 |
| nm0000003 | 4 |
| nm0000003 | 5 |
| nm0000004 | 1 |
| nm0000004 | 3 |

| Profession | Profession_id |
|--------------------|---------------|
| actor | 1 |
| miscellaneous | 2 |
| soundtrack | 3 |
| actress | 4 |
| music_department | 5 |
| writer | 6 |
| director | 7 |
| producer | 8 |
| make_up_department | 9 |
| composer | 10 |
| assistant_director | 11 |

A New Algebraic Operator

•
•



OLAP



Data Analysis and OLAP

- **Online Analytical Processing (OLAP)**
 - Interactive analysis of data, allowing data to be summarized and viewed in different ways in an online fashion (with negligible delay)
- We use the following relation to illustrate OLAP concepts
 - *sales (item_name, color, clothes_size, quantity)*

This is a simplified version of the *sales* fact table joined with the dimension tables, and many attributes removed (and some renamed)



Example sales relation

| item_name | color | clothes_size | quantity |
|-----------|--------|--------------|----------|
| dress | dark | small | 2 |
| dress | dark | medium | 6 |
| dress | dark | large | 12 |
| dress | pastel | small | 4 |
| dress | pastel | medium | 3 |
| dress | pastel | large | 3 |
| dress | white | small | 2 |
| dress | white | medium | 3 |
| dress | white | large | 0 |
| pants | dark | small | 14 |
| pants | dark | medium | 6 |
| pants | dark | large | 0 |
| pants | pastel | small | 1 |
| pants | pastel | medium | 0 |
| pants | pastel | large | 1 |
| pants | white | small | 3 |
| pants | white | medium | 0 |
| pants | white | large | 2 |
| shirt | dark | small | 2 |
| shirt | dark | medium | 6 |
| shirt | dark | large | 6 |
| shirt | pastel | small | 4 |
| shirt | pastel | medium | 1 |
| shirt | pastel | large | 2 |
| shirt | white | small | 17 |
| shirt | white | medium | 1 |
| shirt | white | large | 10 |
| skirt | dark | small | 2 |
| skirt | dark | medium | 5 |

... | ... | ... | ...



Cross Tabulation of sales by *item_name* and *color*

clothes_size all

| | | color | | |
|------------------|-------|-------|--------|-------|
| | | dark | pastel | white |
| <i>item_name</i> | skirt | 8 | 35 | 10 |
| | dress | 20 | 10 | 5 |
| | shirt | 14 | 7 | 28 |
| | pants | 20 | 2 | 5 |
| | total | 62 | 54 | 48 |
| | | total | | |
| | | 53 | 35 | 49 |
| | | 27 | 164 | |

- The table above is an example of a **cross-tabulation (cross-tab)**, also referred to as a **pivot-table**.
 - Values for one of the dimension attributes form the row headers
 - Values for another dimension attribute form the column headers
 - Other dimension attributes are listed on top
 - Values in individual cells are (aggregates of) the values of the dimension attributes that specify the cell.



Data Cube

- A **data cube** is a multidimensional generalization of a cross-tab
- Can have n dimensions; we show 3 below
- Cross-tabs can be used as views on a data cube

| | | item_name | | | | | clothes_size | | | |
|-------|--|-----------|-------|-------|-------|-----|--------------|-------|--------|--------|
| | | skirt | dress | shirt | pants | all | all | large | medium | small |
| color | | dark | 8 | 20 | 14 | 20 | 62 | 34 | 4 | 16 |
| | | pastel | 35 | 10 | 7 | 2 | 54 | 21 | 9 | 18 |
| white | | white | 10 | 5 | 28 | 5 | 48 | 77 | 42 | 45 |
| | | all | 53 | 35 | 49 | 27 | 164 | all | large | medium |



Online Analytical Processing Operations

- **Pivoting:** changing the dimensions used in a cross-tab
 - E.g., moving colors to column names
- **Slicing:** creating a cross-tab for fixed values only
 - E.g., fixing color to white and size to small
 - Sometimes called **dicing**, particularly when values for multiple dimensions are fixed.
- **Rollup:** moving from finer-granularity data to a coarser granularity
 - E.g., aggregating away an attribute
 - E.g., moving from aggregates by day to aggregates by month or year
- **Drill down:** The opposite operation - that of moving from coarser-granularity data to finer-granularity data

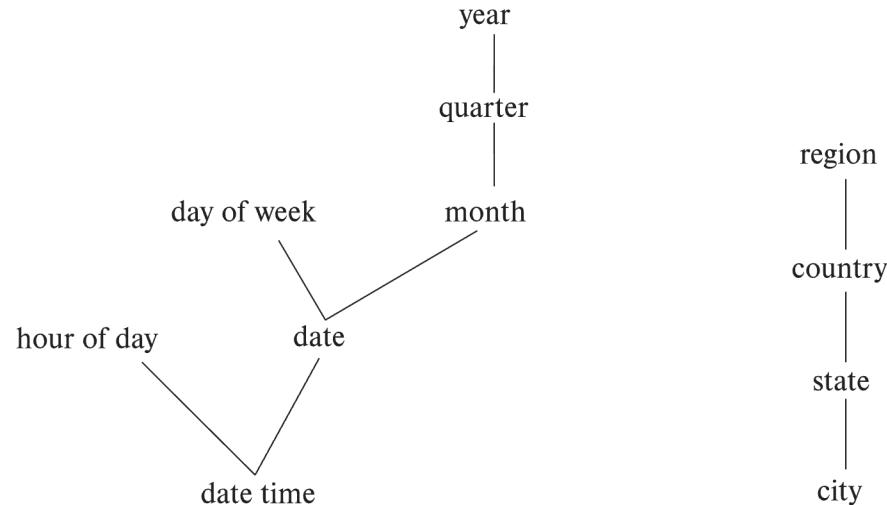
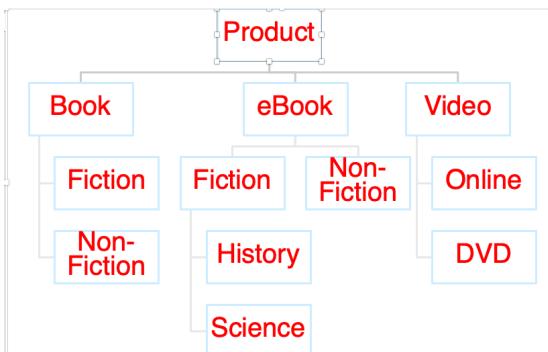


Hierarchies on Dimensions

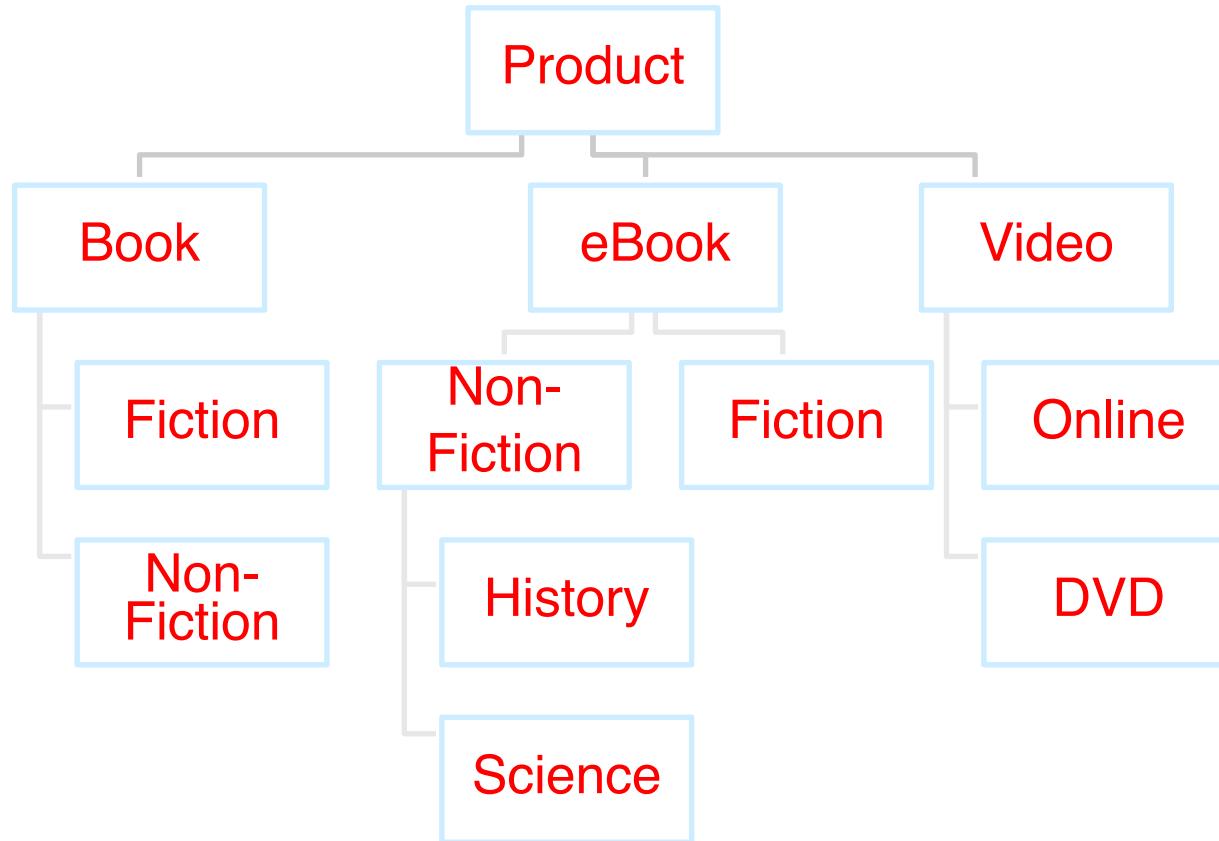
- **Hierarchy** on dimension attributes: lets dimensions be viewed at different levels of detail
- E.g., the dimension *datetime* can be used to aggregate by hour of day, date, day of week, month, quarter or year

Another dimension could be ...

Product category.

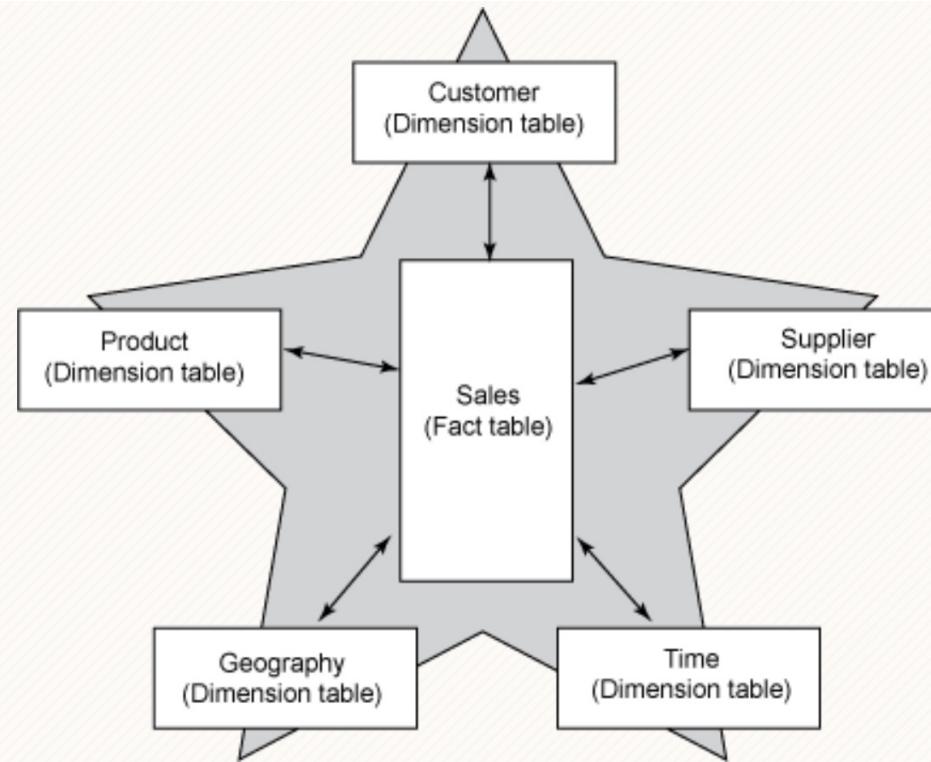


Another Dimension Example – Product Categories





Facts and Dimensions

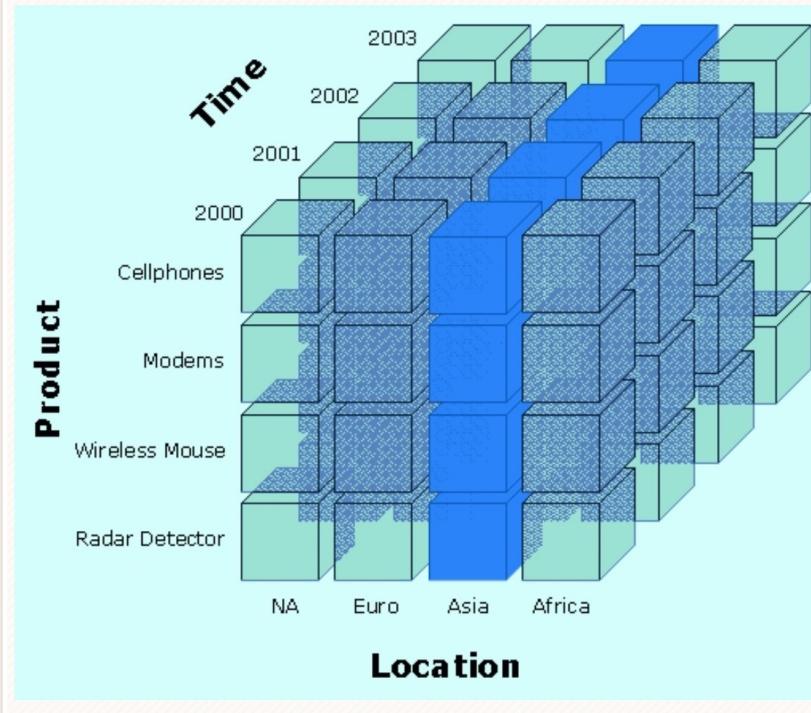




Slice

Slice:

A slice is a subset of a multi-dimensional array corresponding to a single value for one or more members of the dimensions not in the subset.

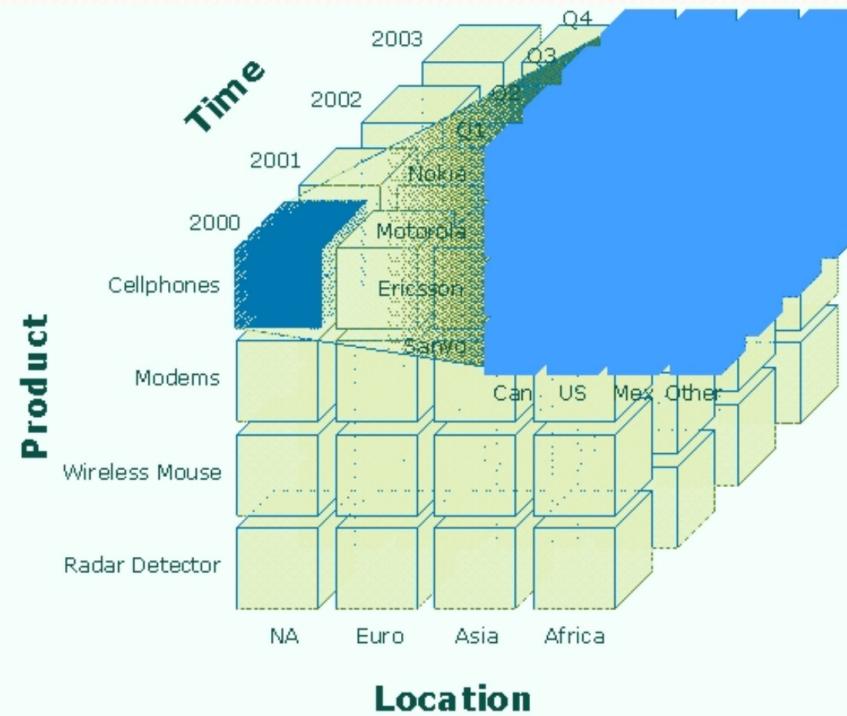




Dice

Dice:

The dice operation is a slice on more than two dimensions of a data cube (or more than two consecutive slices).

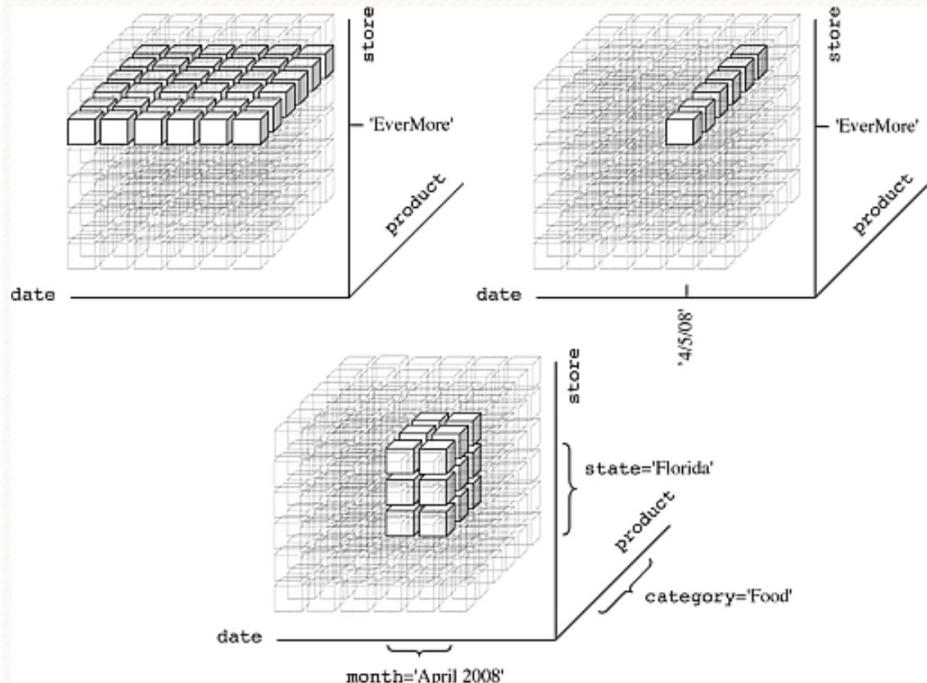




Drilling

Drill Down/Up:

Drilling down or up is a specific analytical technique whereby the user navigates among levels of data ranging from the most summarized (up) to the most detailed (down).

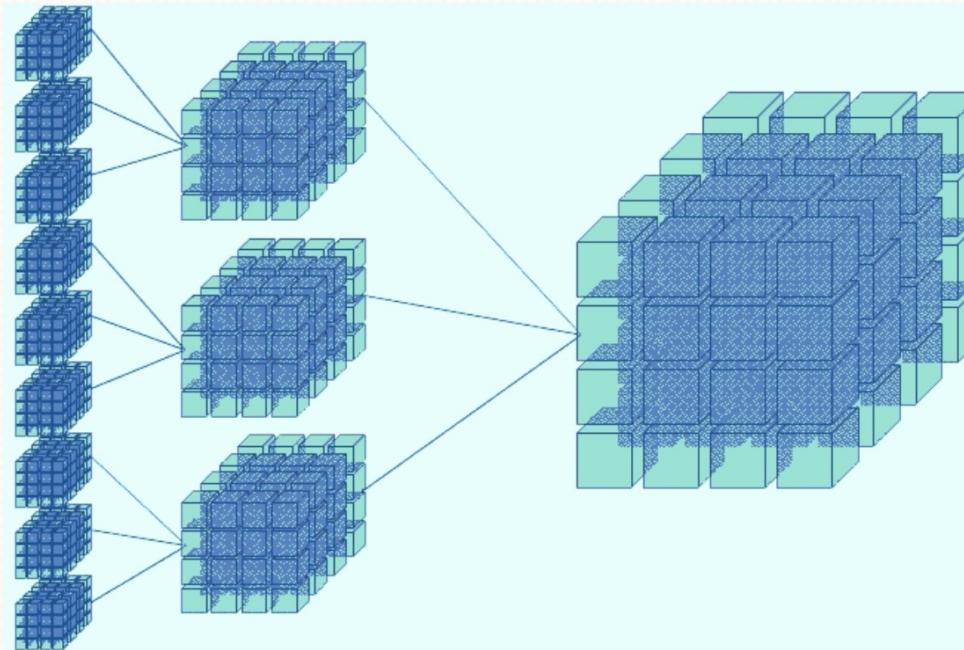




Roll-up

Roll-up:

(Aggregate, Consolidate) A roll-up involves computing all of the data relationships for one or more dimensions. To do this, a computational relationship or formula might be defined.

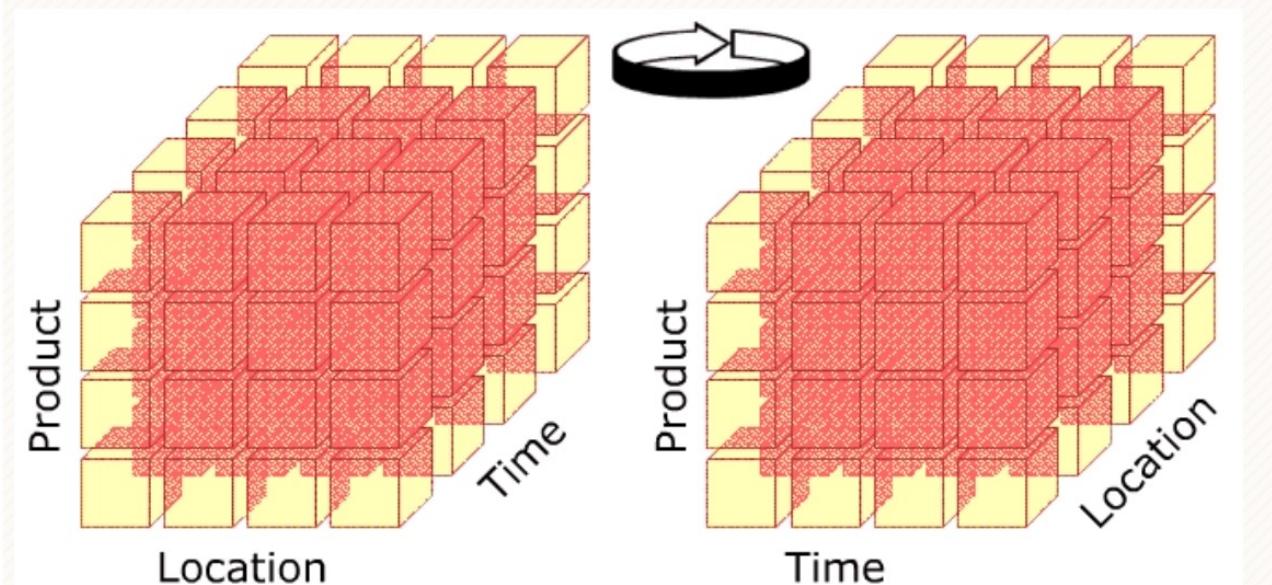




Pivot

Pivot:

This operation is also called rotate operation. It rotates the data in order to provide an alternative presentation of data – the report or page display takes a different dimensional orientation.





Cross Tabulation With Hierarchy

- Cross-tabs can be easily extended to deal with hierarchies
- Can drill down or roll up on a hierarchy
- E.g. hierarchy: *item_name* → *category*

clothes_size: all

| | <i>category</i> | <i>item_name</i> | <i>color</i> | | | |
|------------|-----------------|------------------|--------------|--------|-------|-------|
| | | | dark | pastel | white | total |
| womenswear | skirt | 8 | 8 | 10 | 53 | 88 |
| | dress | 20 | 20 | 5 | 35 | |
| | subtotal | 28 | 28 | 15 | | |
| menswear | pants | 14 | 14 | 28 | 49 | 76 |
| | shirt | 20 | 20 | 5 | 27 | |
| | subtotal | 34 | 34 | 33 | | |
| total | | 62 | 62 | 48 | | 164 |



Relational Representation of Cross-tabs

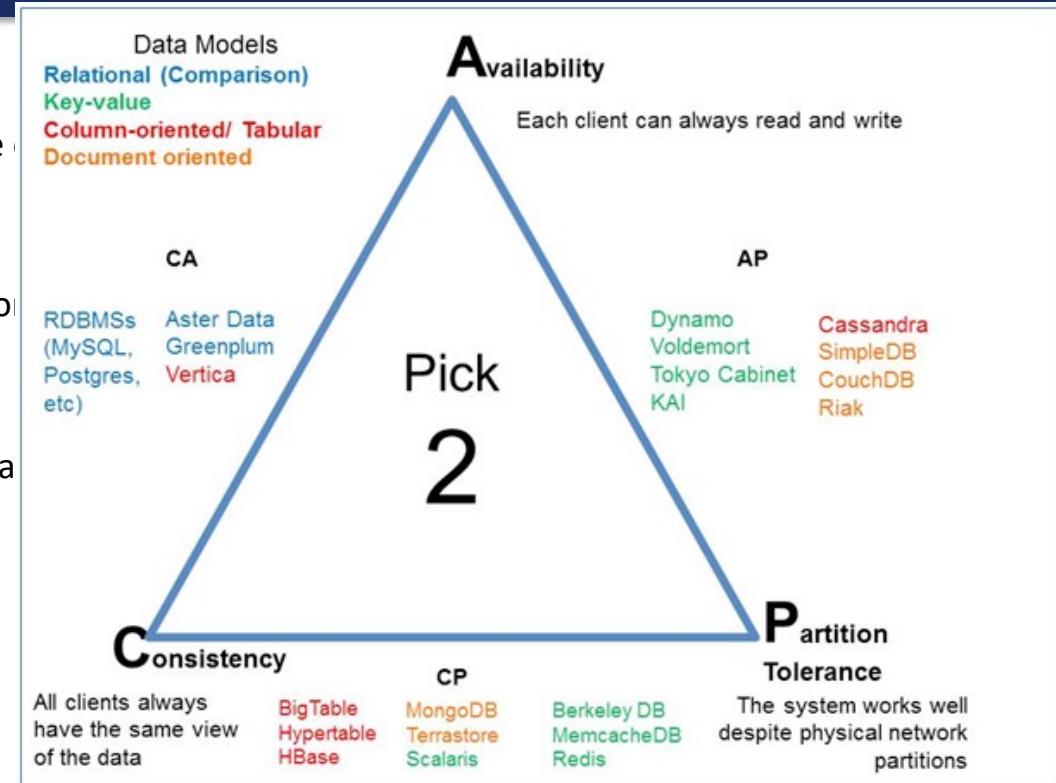
- Cross-tabs can be represented as relations
- We use the value **all** to represent aggregates.
- The SQL standard actually uses *null* values in place of **all**
 - Works with any data type
 - But can cause confusion with regular null values.

| item_name | color | clothes_size | quantity |
|-----------|--------|--------------|----------|
| skirt | dark | all | 8 |
| skirt | pastel | all | 35 |
| skirt | white | all | 10 |
| skirt | all | all | 53 |
| dress | dark | all | 20 |
| dress | pastel | all | 10 |
| dress | white | all | 5 |
| dress | all | all | 35 |
| shirt | dark | all | 14 |
| shirt | pastel | all | 7 |
| shirt | white | all | 28 |
| shirt | all | all | 49 |
| pants | dark | all | 20 |
| pants | pastel | all | 2 |
| pants | white | all | 5 |
| pants | all | all | 27 |
| all | dark | all | 62 |
| all | pastel | all | 54 |
| all | white | all | 48 |
| all | all | all | 164 |

CAP, BASE

CAP Theorem

- **Consistency**
Every read receives the most recent write
- **Availability**
Every request receives a (non-error) response
- **Partition Tolerance**
The system continues to operate despite a network partition between nodes



Consistency Models

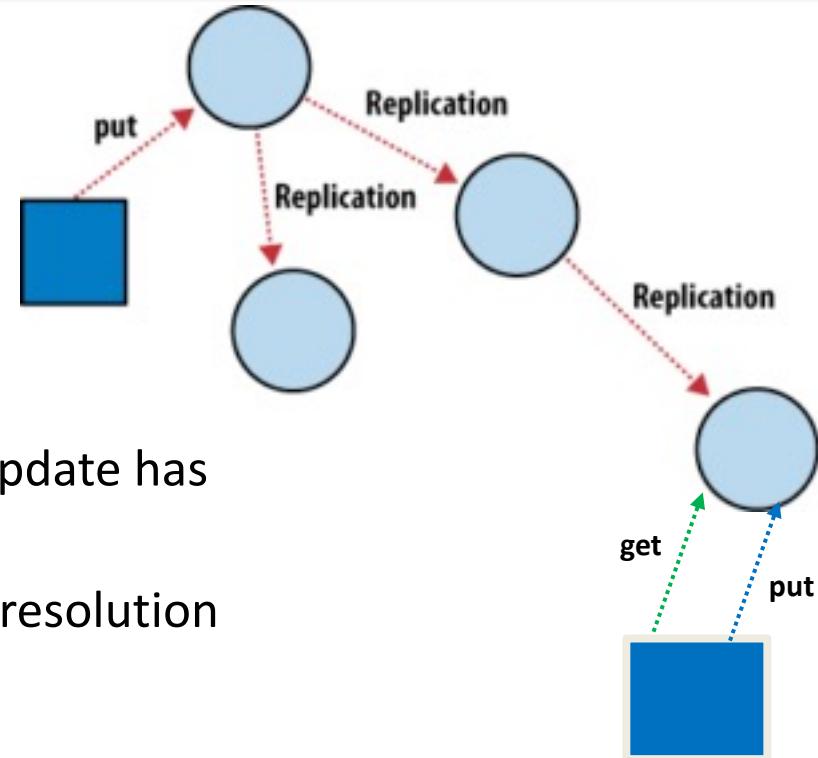
- **STRONG CONSISTENCY:** Strong consistency is a consistency model where all subsequent accesses to a distributed system will always return the updated value after the update.
- **WEAK CONSISTENCY:** It is a consistency model used in distributed computing where subsequent accesses might not always be returning the updated value. There might be inconsistent responses.
- **EVENTUAL CONSISTENCY:** Eventual consistency is a special type of weak consistency method which informally guarantees that, if no new updates are made to a given data item, eventually all accesses to that item will return the last updated value.

Eventually-consistent services are often classified as providing BASE (Basically Available, Soft state, Eventual consistency) semantics, in contrast to traditional ACID (Atomicity, Consistency, Isolation, Durability) guarantees. Rough definitions of each term in BASE:

- **Basically Available:** basic reading and writing operations are available as much as possible (using all nodes of a database cluster), but without any kind of consistency guarantees (the write may not persist after conflicts are reconciled, the read may not get the latest write)
- **Soft state:** without consistency guarantees, after some amount of time, we only have some probability of knowing the state, since it may not yet have converged
- **Eventually consistent:** If the system is functioning and we wait long enough after any given set of inputs, we will eventually be able to know what the state of the database is, and so any further reads will be consistent with our expectations

Eventual Consistency

- Availability and scalability via
 - Multiple, replicated data stores.
 - Read goes to “any” replica.
 - PUT/POST/DELETE
 - Goes to any replica
 - Change propagate asynchronously
- GET may not see the latest value if the update has not propagated to the replica.
- There are several algorithms for conflict resolution
 - Detect and handle in application.
 - Clock/change vectors/version numbers
 -



ACID – BASE (Simplistic Comparison)

| ACID (relational) | BASE (NoSQL) |
|-------------------------------|-------------------------------------|
| Strong consistency | Weak consistency |
| Isolation | Last write wins (Or other strategy) |
| Transaction | Program managed |
| Robust database | Simple database |
| Simple code (SQL) | Complex code |
| Available and consistent | Available and partition-tolerant |
| Scale-up (limited) | Scale-out (unlimited) |
| Shared (disk, mem, proc etc.) | Nothing shared (parallelizable) |

Scalability

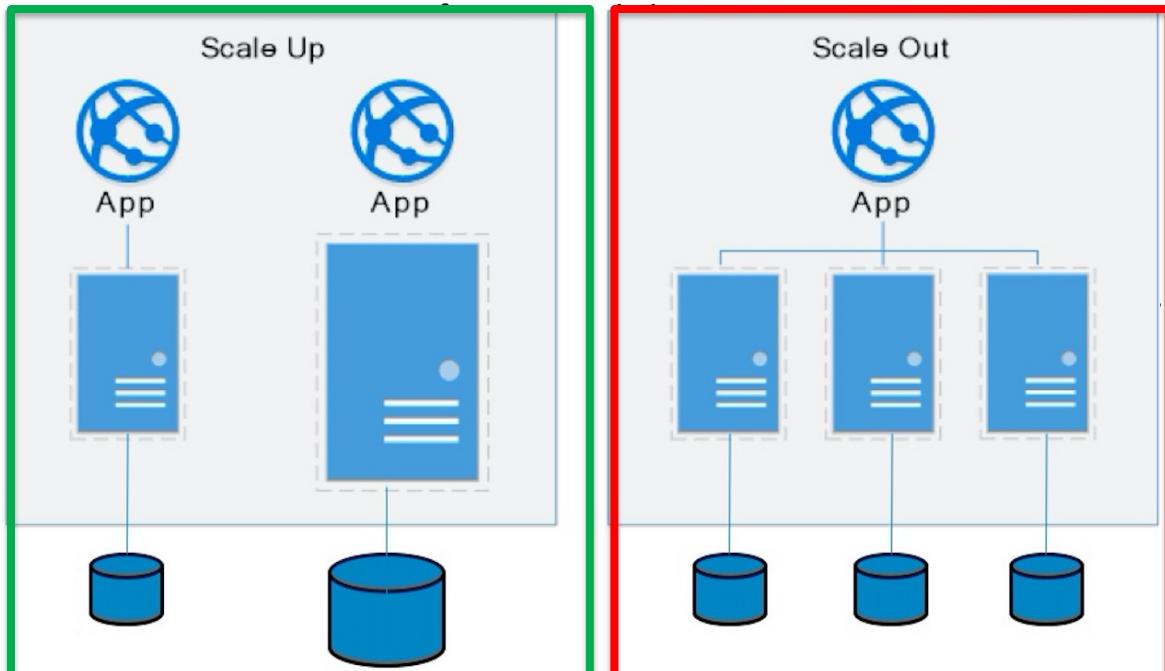
Approaches to Scalability

Scalability is the property of a system to handle a growing amount of work by adding resources to the system.

Scale-up:

Replace system with a bigger machine,
e.g. more memory, CPU, ...
— More disruptive:

Add another system.

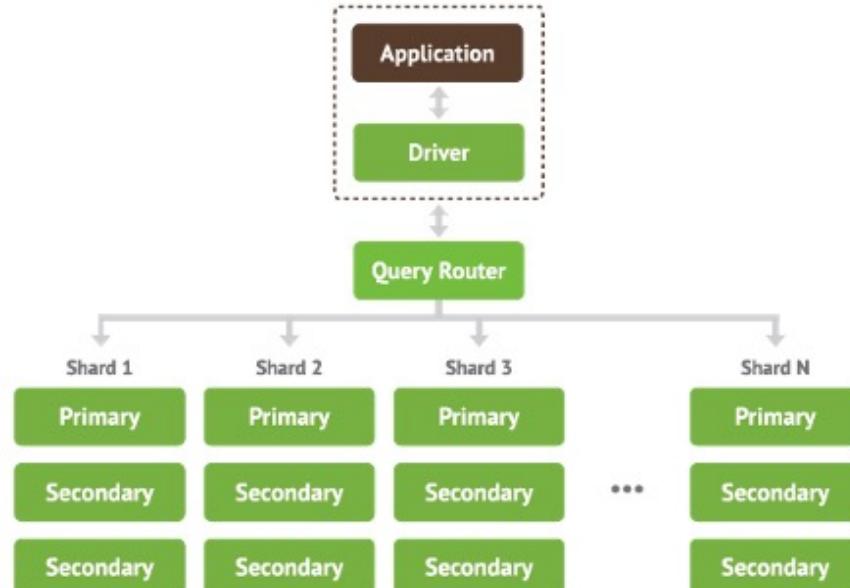


Disk Architecture for Scale-Out

- Share disks:
 - Is basically scale-up for data/disks.
You can use NAS, SAN and RAID.
 - Isolation/Integrity requires distributed locking to control access from multiple database servers.
 - Share nothing:
 - Is basically scale-out for disks.
 - Data is partitioned into *shards* based on a function $f()$ applied to a key.
 - Can improve availability, at the code consistency, with data replication.
 - There is a router that sends requests to the proper shard based on the function.
-
- The diagram illustrates three disk architectures for scale-out:
- Share Everything:** A single database server (DB) is connected to a single disk via an IP network. This is labeled "eg. Unix FS".
 - Share Disks:** Multiple database servers (DB) are connected to a central SAN Disk via an IP network and Fibre Channel (FC). This is labeled "eg. Oracle RAC".
 - Share Nothing:** Multiple database servers (DB) are connected to their own local storage disks via an IP network. This is labeled "eg. HDFS".

Shared Nothing, Scale-Out

MongoDB Sharding



DynamoDB Partitioning

