

1] DATA MODELING

① INTRO

DB: Structure repo that is stored & retrieved electronically for use in apps

DBMS: DB Management System Software used to access the DB

→ Relational DB: A DB based on the Relational Model

↳ Relational model

Types (common):

- Oracle
- PostgreSQL
- Teradata
- SQLite
- MySQL

Tuple {						
Attribute						
Relation						

Basics:

- DB / Schema: collection of tables
- Table / Relation: group of rows w/ the same labeled elements
- Column / Attribute: labeled element
- Rows / Tuple: single item

② When to use a Relational DB?

- ✓ Flexibility for writing SQL queries
- ✓ Modeling data, not queries
- ✓ Smaller data volumes
- ✓ Ability to do joins
- ✓ ACID Transactions
- ✓ Ability to do aggregations & analytic
- ✓ Easier to change to less requirements
- ✓ Secondary indexes available

↑ ACID Transactions

③ When not to use a Relational DB?

- ✗ Large amounts of data
- ✗ Need to be able to store diff. data type
- ✗ Need ↑ throughput - fast reads
- ✗ Flexible schema
- ✗ High availability
- ✗ Need horizontal scalability

◦ Atomicity: a whole tx is processed or nothing is processed

◦ Consistency: only tx that abide by constraints are written in the DB

◦ Isolation: Tx are processed independently and serially. Order does not matter.

◦ Durability: completed tx are saved to DB even in cases of syst failure

→ NoSQL DB: Not only SQL or NonRelational DB

Types (common):

- Apache Cassandra (Partition Row store)
- MongoDB (Document store)
- DynamoDB (Dynamo Key-value store)
- Apache HBase (Wide Column store)
- Neo4J (Graph DB)

Basics [APACHE CASSANDRA]: → Scalability + High availability

- Keyspace: collection of tables
- Table: group of partitions
- Row: single item
- Partition: fundamental unit of access (collection of rows)
- Primary Key: Made up of a partition key (+) clustering column
- Columns: Clustering & data. Labeled element

① When to use a NoSQL DB? → (⊖ when NOT to use Relational DBs)

+ ② Users are distributed → -low latency

② When not to use a NoSQL DB? → Scale up linearly

→ (⊖ when to use a Rel. DB).

② RELATIONAL DATA MODELS

DB: A set of related data and the way it's organised

DBMS: Computer software that allows users to interact w/ DB.

Relational
Importance

- Standardization of data model: Once the data is transformed into rows & columns format \rightarrow queryable
- Flexibility in adding and altering tables:
- Data integrity
- Structured Query Language (SQL): Lang. to access data
- Simplicity: Data is systematic stored and modelled in tables
- Intuitive Organization: as a spreadsheet.

OLAP: Online Analytical Processing

- DBs optimized for these workloads allow for complex analytical and ad hoc queries. OPTIMIZED for READS

OLTP: Online Transactional Processing

DBs optim. for these workloads allow for less complex queries in large volumes. The types of queries of these DB: **READ, INSERT, UPDATE and DELETE** \rightarrow Very fast query processing

STRUCTURING your DB:

- Normalization: Reduce data redundancy and increase data integrity
- Denormalization: Must be done in read-heavy workloads to ↑ performance

Objectives of NORMAL FORM

1. To free the DB from unwanted insertions, updates and deletions dependencies.
2. To reduce the need for refactoring the DB as new types of data are introduced.
3. To make the relational model more informative to users
4. To make the DB neutral to the query statistics

\hookrightarrow opposite of NoSQL

\rightarrow Process of normalization

1NF \rightarrow 2NF \rightarrow 3NF

[NORMALIZATION]

- 1NF:**
- Atomic values: Each cell contains unique and single values
 - Be able to add data w/o altering tables
 - Separate different relations into different tables
 - Keep relationships btw tables together w/ foreign keys

- 2NF:**
- Have reached 1NF
 - All columns in the table must rely on the Primary Key

Ex. [Store Customer Table]

Store ID	Customer ID	Store ID	Customer ID	Name
001	004	001	004	...
001	005	001	005	...
003	006	003	006	...

→ [Customer Table]

- 3NF:**
- Have reached 2NF

- No transitive dependencies

When we want to update data, we want to be able to do it in 1 place.

Customer ID	Name
004	...
005	...
006	...

[Awards Table]

Ex:

Music Award	Year	Winner Rec.	Lead Singer	Music Award	Year	Winner Rec.
Grammy	1965	The Beatles	John Lennon	Grammy	1965	The B.
CMA	2000	Faith Hill	Faith Hill	CMA	2000	Faith Hill
Grammy	1970	The Beatles	John Lennon	Grammy	1970	The B.

[Lead Singer]

Band Name	Lead Singer
The Beatles	John Lennon
Faith Hill	Faith Hill

DENORMALIZATION

- Process to improve the read performance of a DB at the expense of losing some write performance by adding redundant copies of data.

JOINS allow for outstanding flexibility BUT extremely slow.

Logical design change:

- The designer is in charge of keeping data consistency
- Reads will be faster (select)
- Writes will be slower (insert, update, delete)

FACT TABLES: Measurements, metrics or facts of a business process. (Transactions)

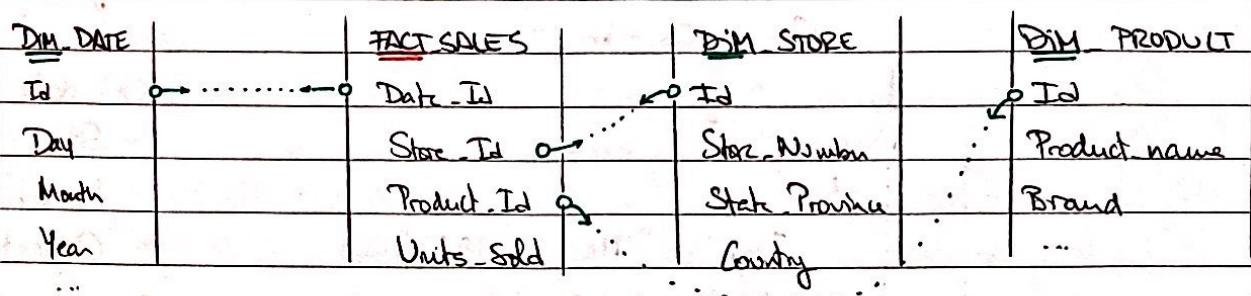
DIMENSION TABLES: Categorizes facts & measures. Enable users to answer less specific questions, people, products, place and time

[DATA MART] SCHEMAS

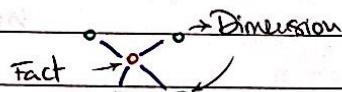
STAR: Simplest style of data mart schema. 1 fact table referring any # dimtbls

SNOWFLAKE: Arrangement of tables in a multidim. dB. Centralized facts conn. to multiple dimensions.

Ex: STAR *



Benefits of * Schema



• Denormalized

• Simplified queries

• Fast aggregations

Drawbacks: • Issues that come w/ denormalization

• query perf. • Data integrity

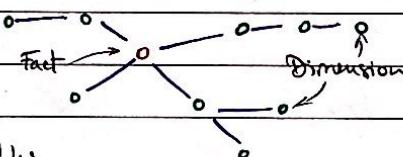
• Many to many relationships

Ex: SNOWFLAKE *

• 1 fact table with

multiple dimension tables

with more child dimension tables



• Star Schema is a simplified case of a snowflake sch.

Note: While using SQL to create STAR/SNOWFLAKE schemas, there are some data definitions and constraints to take into account. → CREATE statement

- NOT NULL: The column cannot contain a Null value. This occurs w/ a composite key (numerically).
- UNIQUE: The data across all rows in 1 column is unique.
- PRIMARY KEY: Defined on a single column. The values in this column identify the table. If a group of columns are defined as a PK, they are called COMPOSITE KEY.

↳ NOT NULL + UNIQUE constraints

→ UPSERT: inserting a new row or updating an existing one

INSERT INTO customer_address (cust_id, cust_street, cust_city, cust_state)
VALUES (432, '123 Ab Street', 'Albany', 'NY')

ON CONFLICT (cust_id)

DO UPDATE

SET cust_street = EXCLUDED.customer_street

table: customer_address

- customer_id
- customer_street
- customer_city
- customer_state

Useful Jupyter Notebooks / Python scripts:

JSON : Great to represent data w/ complex data hierarchy

object
string
Values for both JSON objects and arrays can be any valid JSON data types

- String
- Number
- Object
- Array
- Boolean
- null

```
{  
    "Directed by": "Steven Spielberg",  
    "Produced by": [  
        "Kathleen Kennedy",  
        "Steven Spielberg"  
    ],  
    "Release": [  
        {"Date": "May 26, 1982",  
         "Location": "Cannes"},  
        {"Date": "June 11, 1982",  
         "Location": "USA"}  
    ]  
}
```

③ DATA MODELING WITH POSTGRES

Intro: Sparkify: music streaming app. Analyze data: SONGS & USER act.

What songs users are listening | JSON logs

Create Postgres DB to optimize queries | JSON metadata

① DB Schema | ③ Test DB & ETL

② ETL pipeline

Description: • Define fact and dimension tables for a star schema using Python

• Write an ETL pipeline that transfers data: files ↓
DB

Song Dataset: Real data from Million Song Dataset

Each JSON file: metadata about a song and an artist

Files partitioned by the 3 first letters of each song's track ID

e.g. song-data/A/B/C/TRABCE128F474C983.json

→ Ex. { "num_songs": 1
 "artist_id": "ARJIE...B7"
 "artist_latitude": null
 "artist_longitude": null
 "artist_location": ""
 "artist_name": "Line Renaud"
 "song_id": "SOUPI...E1"
 "title": "Der Kleine Dampfaff"
 "duration": 152.92036
 "year": 0 }

Log Dataset: Log files (JSON) by an event simulation based on the songs (1).

Log files partitioned by year & month

Ex: log-data/2018/11/2018-11-12-events.json

→ To look at the JSON data in log-data files → pandas dataframe

`df = pd.read_json(filepath, lines=True)`

↑
 data /log-data /2018/11/2018-11-01-events.json

FACT TABLE : ① songplays : records in log data. Records page NextSong

- col: songplay_id, start_time, user_id, level
song_id, artist_id, session_id, location,
user agent

DIMENSION : ② users :

TABLES

- col: user_id, first_name, last_name, gender, level

③ songs :

- col: song_id, title, artist_id, year, duration

④ artists :

- col: artist_id, name, location, latitude, longitude

⑤ time : timestamps of records in songplays broken
down into specific units

- col: start_time, hour, day, week, month,
year, weekday.

STEPS

- ① Create tables:
- CREATE in 'sql_queries.py'
 - DROP in 'sql_queries.py'
 - 'run' 'create_tables.py'
 - 'run' 'test.ipynb' confirm creation
 - Restart kernel to close connection to DB
after running 'test.ipynb'

- TIPS:
- CREATE statement: CREATE TABLE IF NOT EXISTS
table_name (column data_type, ...);
 - DROP statement: DROP TABLE IF EXISTS table_name;
 - File > New > Terminal # python create_tables
 - Restart kernel w/ [C] button

- ② Build ETL processes - Instructions in 'etl.ipynb'

↳ At the end of each table section: 'run' 'test.ipynb'
↳ Tip! Rerun 'create_tables.py' to reset tables

TIPS:

- 'sql_queries.py':

INSERT INTO	table_name	(column, ...)
VALUES	(%s, %s, ...)	: (+) ON CONFLICT ...
- song_data = df[['song_id', 'title', 'artist_id', ...]]
song_data_list = song_data.values[0].tolist()
cur.execute(song_table_insert, song_data_list)

- FILTERING w/ Pandas: $df[df.\text{column} == \text{'Value'}]$
- COPY df: $t = df.\text{copy}()$
- DATEREIME: $t['ts'] = pd.\text{to_datetime}(df['ts'], \text{unit}=\text{'ms'})$
- Extract time: $t['ts'].dt.\text{hour}$ or $t['ts'].dt.\text{week}$
- frame = { 'column_1': t['ts'].(...), 'column_2': t['ts'].(...), ... }
 ↳ time_df = pd.DataFrame(frame)

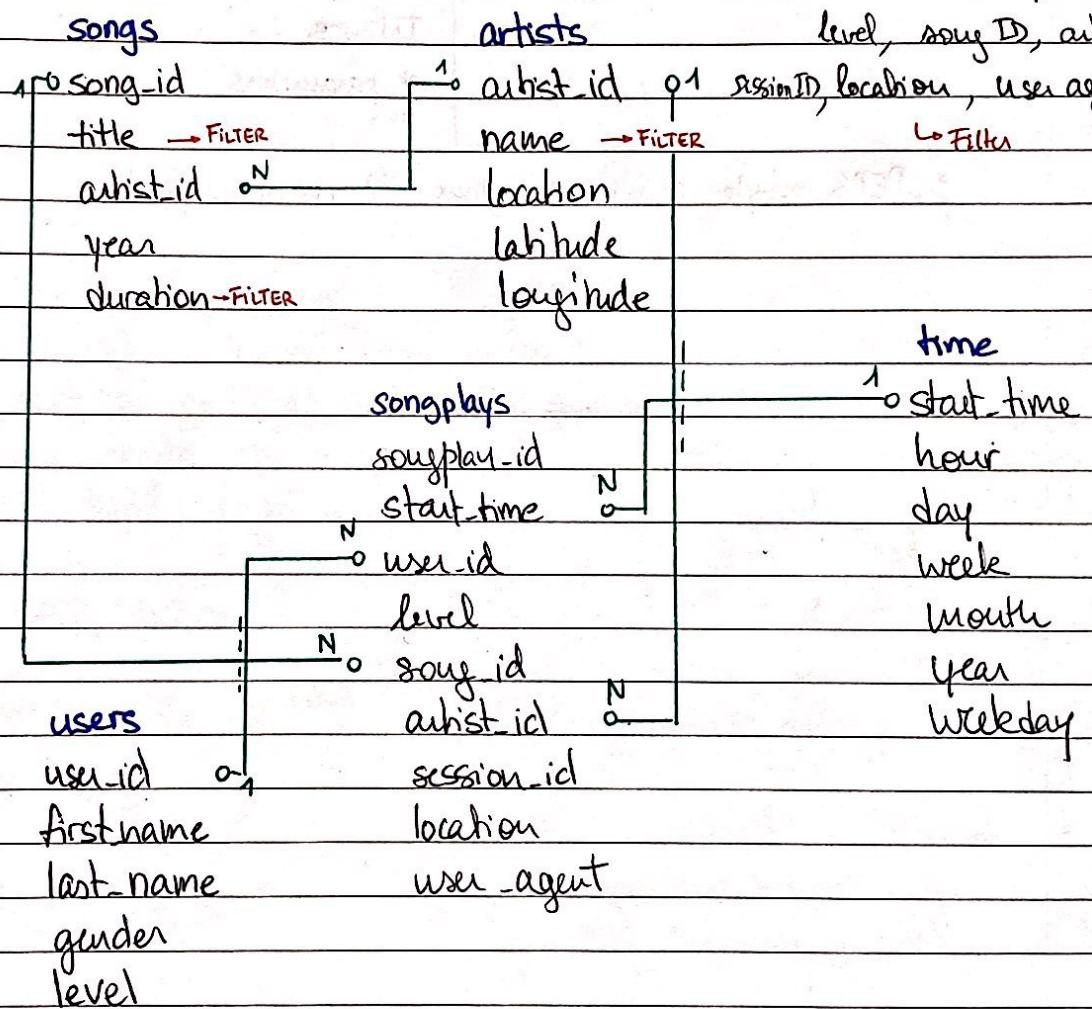
songplays TABLE: song-select query : ① find song_id, artist_id based on the title, artist name & duration.

② Select timestamp, user ID,

level, song ID, artist ID,

name → FILTER

→ Filter



Project Review

3 specifications:

- o Add NOT NULL constraints:

Ex: Songplay table : start_time NOT NULL

- o Add FOREIGN keys.

→ Added songplay.start_time NOT NULL REFERENCES time (start_time)

songs.artist_id NOT NULL REFERENCES artists(artist_id)

- o Add docstrings → def function (parameters):

'''

Description

INPUTS

* parameters

'''

- o PEPS style guidelines: max 72 characters per line

④ NoSQL DATA MODELS

NoSQL and Non-Relational are interchangeable terms

↳ Not Only SQL

APACHE CASSANDRA

↓
Distributed DBs

- In order to have high availability,

You will need copies of your data

① A DB made up of multiple machines

② High avail.: (almost) no downtime → with 2 nodes

③ Because 1 node can go down ⇒ need copy

- Open Source NoSQL DB

- Masterless Arch

- High Availability

- Linearly Scalable (adding nodes)

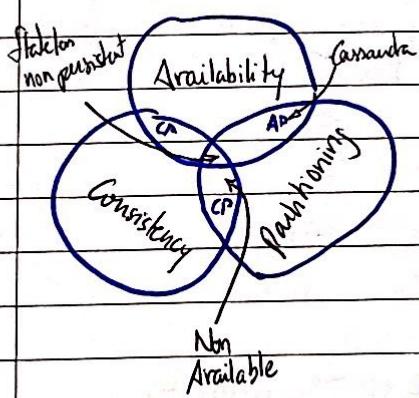
UBER NETFLIX

Hulu TWITTER

▷ Eventual Consistency

Consistency model used in distrib. computing to achieve HIGH AVAILABILITY that informally guarantees that, if no updates are made to a given data item, eventually all accesses to that item will return the last updated value.

CAP Theorem



• Consistency: Every read from the DB gets the latest (and correct) piece of data or an error.

• Availability: Every request is received and a response is given - w/o the guarantee that the data is latest update

• Partition Tolerance: The system continues to work regardless of losing network connectivity between nodes

Apache Cassandra is...

- Highly available DB
- Linearly scalable
- An open source project (Apache Foundation)

DENORMALIZATION

→ It's a must w/ Apache Cassandra

- Done for fast reads

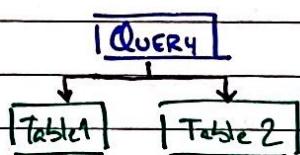
- A. Cassandra has been optimized for fast writes

ALWAYS think queries first

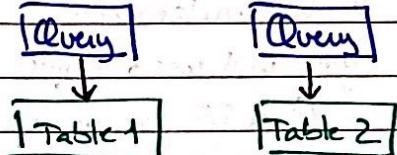
- Good strategy: 1 table per query

→ A. Cassandra does NOT allow for joins b/w tables.

Relational DB



NoSQL DB



Not possible to JOIN

CQL Cassandra Query language

Language to interact w/ DB. Very similar to SQL

NO: JOINS, GROUP BY or subqueries not in CQL, not supported by CQL

Python wrapper / driver : pip install cassandra-driver

import | import cassandra

connection to DB | from cassandra.cluster import Cluster

try:

cluster = Cluster(['127.0.0.1'])

locally installed Ap. (ass. instance)

session = cluster.connect()

except Exception as e:

print(e)

try:

session.execute ("""

CREATE KEYSPACE IF NOT EXISTS udacity

WITH REPLICATION =

{'class': 'SimpleStrategy', 'replication_factor': 1} """)

except Exception as e:

print(e)

connect | (try/except) session.set_keyspace('udacity')

query = "CREATE TABLE IF NOT EXISTS music_library"

query += "(year int, artist_name text, album_name text,
PRIMARY KEY (year, artist_name))"

PRIMARY KEY : ① PARTITION KEY

② PARTITION KEY + CLUSTERING COLUMNS

PARTITION KEY : Determine the distribution of data across the system

↳ The partition key's row value will be hashed ($\rightarrow \#$) and stored on the node in the system that holds that range of values.

PK

◦ SIMPLE: ①

`CREATE TABLE music_library
(year int, artist_name text, album_name text,
PRIMARY KEY (year))`

↑ partition key

◦ Must be **UNIQUE**

◦ Hashing on this value results in placement on a particular node in the system

◦ COMPOSITE: ②

`CREATE TABLE music_library
(year int, artist_name text, album_name text,
PRIMARY KEY (year, artist_name))`

partition key ↑ clustering column

◦ Data distrib. by this partition key

◦ simple or composite

◦ At clustering columns (optionals)

CLUSTERING COLUMNS:

- Will SORT the data in ASC.
- 1+ (or none) can be added
- The columns will be sorted in order of how they were added to the PK

`CREATE TABLE music_library`

`(year int, artist_name text, album_name text,`

`PRIMARY KEY [(year), artist_name, album_name)] KEY`

year	artist_name	album_name	partition key	1st clust. column (ASC)	2nd clust. column
1965	Elvis	Blue Hawaii			
1965	The Beatles	Rubber Soul			
1965	The Beatles	Revolution			
1965	The Monkees	Meet the Monkees			

↳ for each artist_name it orders the album names ASC.

'WHERE' CLAUSE:

- Must be included to execute queries.
- Recommended that 1 partition be query at a time (performance)
- Possible: 'select * from table' if you add a config: ALLOW FILTERING
 - PARTITION KEY must be included ↑
 - CLUSTERING COLUMNS can be used in order they appear RISKY
 - ↳ We cannot use the 3rd clust. col. if we do not add the 1st & 2nd.

(5) DATA MODELING WITH APACHE CASSANDRA

INTRO: (⇒ POSTGRES)

Create A Cassandra DB → song play data

EVENT DATA

Directory of CSV files partitioned by date

Ex: 'event_data/2018-11-08-events.csv'

CSV: Comma Separated Value

- Jupyter notebook:
- ① Process: 'event_datafile-new.csv'
 - ② Model data tables (w/ queries in mind)
 - ③ Queries needed to keep in mind
 - ④ Load data into tables + Run queries

STEPS

① Modeling NoSQL / Apache Cassandra DB

event_datafile-new.csv

1	artist	BC	5	lastName	BC	9	sessionId	Z
2	firstName	AB	6	length	AZ	10	Song	AC
3	gender	AC	7	level	BC	11	userId	Z
4	itemInSession	Z	8	location	AC			
↳ in array [-1]								

QUERY 1: SELECT 1, 10, 6 FROM music_library WHERE 9 AND 4

QUERY 2: SELECT 1, 10, 2, 5 FROM music_library WHERE 11 AND 9
↳ SORTED BY itemInSession [4]

QUERY 3: SELECT 2, 5 FROM music_library WHERE 10

Columns :	1	normal	Q1)	1	Q2)	1
	2	normal		4		12
	4	PRIM. KEY		6 CC		4
	5	normal		9 PK		5
	6	normal		10 EE		9
	9	PRIM. KEY				10
	10					11
	11					