

# DATA ENGINEERING

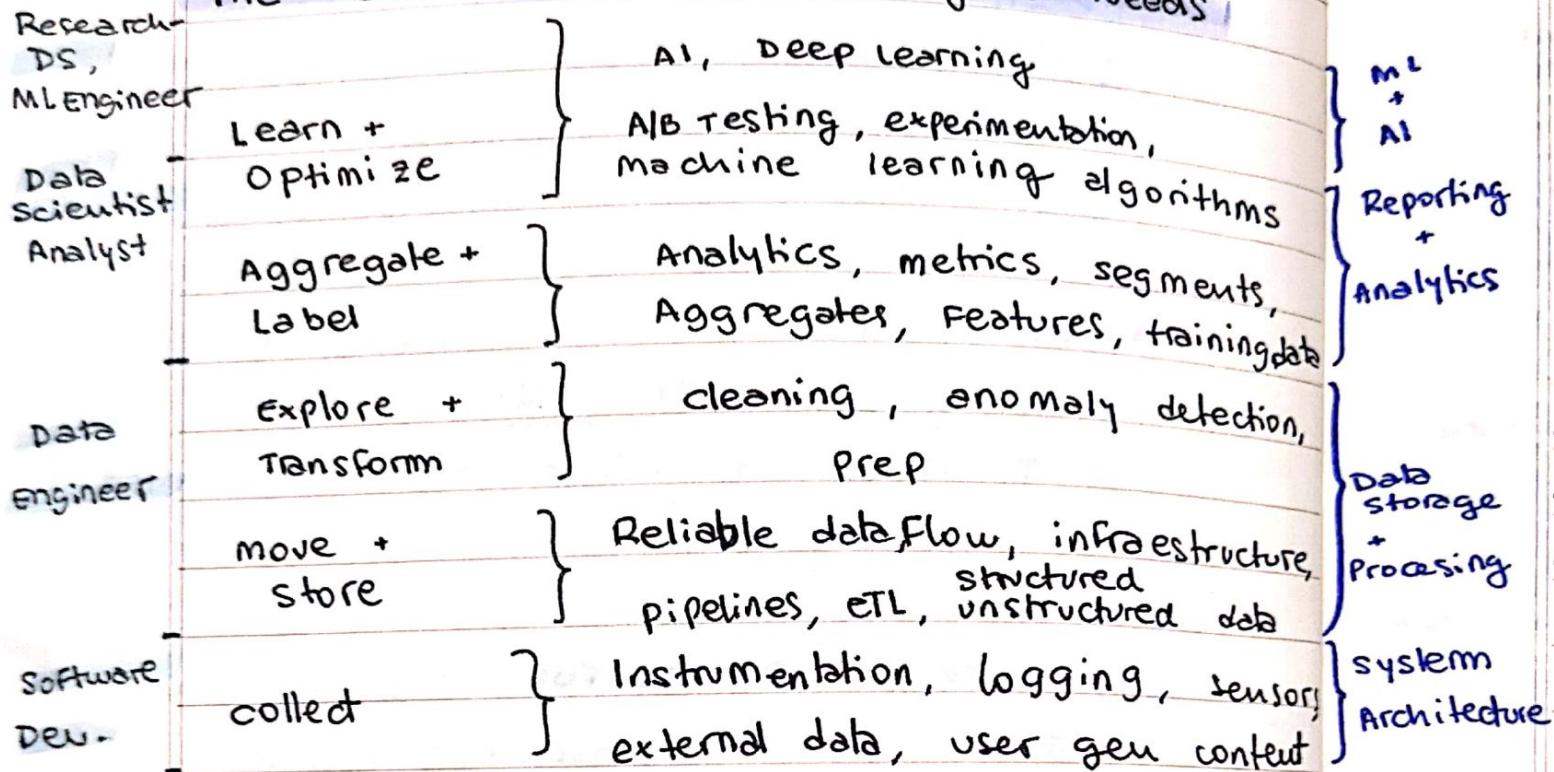
Learn to design DATA MODELS,  
build DATA WAREHOUSES and  
DATA LAKES, automate DATA PIPELINES  
and work with massive datasets.

- vital role in storing, organizing +  
making use of data.
  - ↑ speed + volume of info generated  
every day

## HOW TO BE SUCCESSFUL IN THE PROGRAM

- ❖ BE an ACTIVE STUDENT : TAKE notes
  - ❖ write the CODE YOURSELF
  - ❖ time management : go to class as a routine
- ② how to address BLOCKERS
- Research
  - Read
  - Ask 4 help

## The Data Science Hierarchy of Needs



### Roles + Skills

	SOFT DEV	DE	DS
HTML	web services	Hadoop	SQL
CSS	sw dev.	Data wareh.	Big data
XML	Agile Meth.	ETL	Statistics
JQuery		MySQL	Data analysis
PHP		Hive	ML
.NET		Linux	BI
		Java	SAS
		Python	Algorithms
		Oracle	Data Mining
	JS	Databases	
		C	
		SQL	

## common data engineering Activities:

- Ingest data from a data source
- Build + maintain a data warehouse
- Create a data pipeline
- Create an analytics table for a use case
- Migrate data to cloud
- Schedule + automate pipelines
- Backfill data
- Debug data quality issues
- Optimize queries
- Design a Database.

# DATA MODELING

→ Abstraction that organizes elements of data and how they will relate to each other.

→ The process of creating db models for an information system.

Organize data into db systems to ensure that data is persisted and easily usable.

n.

T

4

I

3

- Process to support business and user needs
- Gather requirements from apps team to users

- ① • Conceptual data modeling w/ entity mapping
  - Map concepts of data that you have/will have and relate db concepts

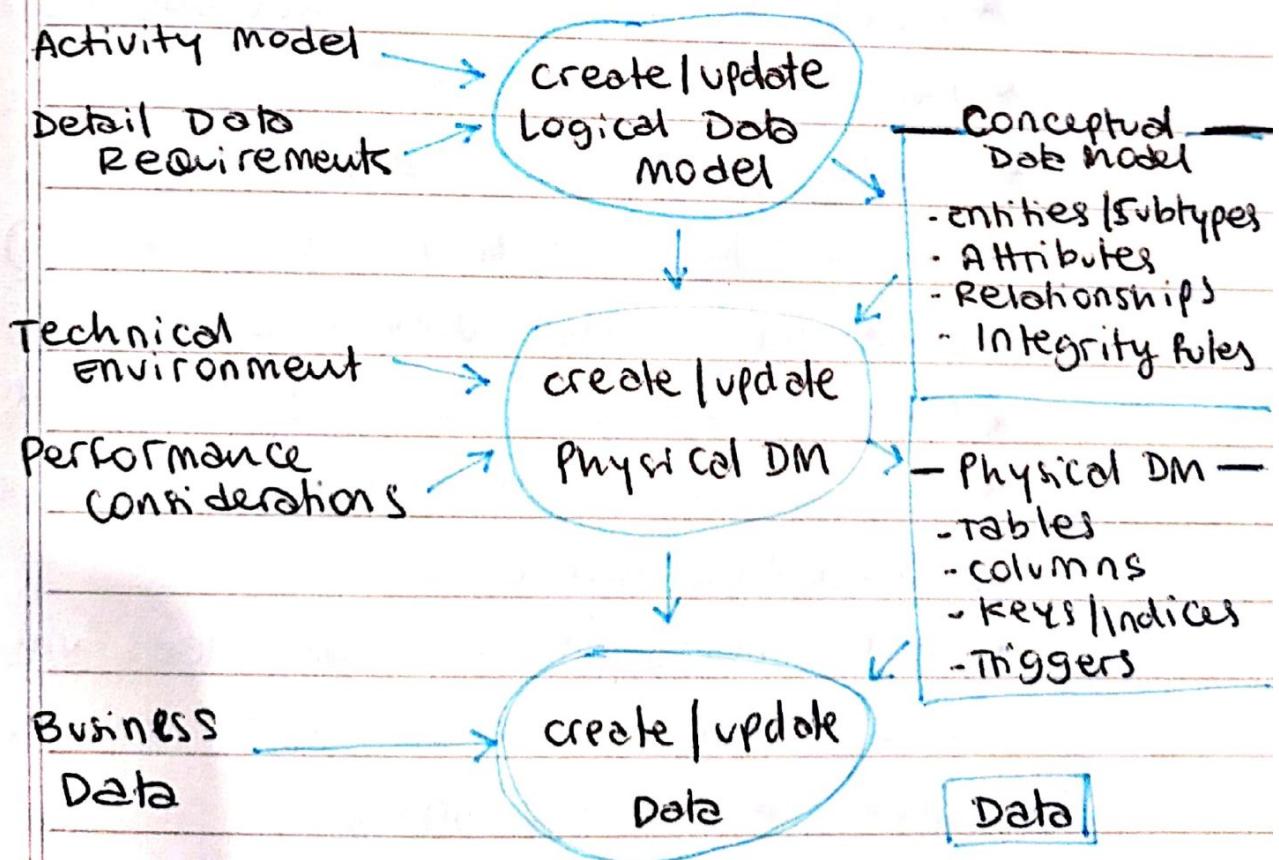
- ② • Logical data modeling
  - Map concepts to logical models using schemas, tables, columns ...

- ③ • Physical data modeling
  - DDL: Data definition language code to persist our data

- \* Data organization is critical - makes everyone's life easier.
- \* Organized data determines later data use
- \* Prior step to building out app, business logic, and analytical models. - start early.
- \* Iterative process - as new reqs + data come in.

W H O ?

→ Everyone who works with data.



# Relational Model

This model organizes data into one or more tables of columns + rows, with a unique key identifying each row. Each table represents an entity type. (General  
entity type)

Columns = Attributes

Rows = single items

## WHY USE RDBMS?

Ease of use:

- \* Flexibility for writing SQL queries
- \* Ability to use JOINS
- \* Ability to do agg and analytics (Group by, sum, min, MAX, AVG...)
- \* Great for smaller data volumes
- \* Easy to change business requirements
- \* You model the data, independent of queries
- \* Secondary Indexes
- \* ACID Transactions = guarantees validity even in the event of errors, power failures, → ⇒ Data integrity

4

Atomicity: The whole transaction is processed or nothing is processed  
→ All or nothing.

Consistency: Only transactions that abide by constraints and rules is written into the db, otherwise DB keeps previous state. → Data should be correct across all rows and tables

Isolation: Transactions are processed independently and securely, order does not matter. → Reduces changes of concurrency effects.

Durability: Completed transactions are saved to DB even of cases of system failure → Transactions are recorded on non-volatile memory

## WHY NOT USE RDBMS ?

- \* Large amount of data: can only scale vertically  $\Rightarrow$  They are not distributed. Limited by the host server.
- \* Need to store diff data types formats.
- \* Need high throughput = fast reads. ACID transactions slow down. Flw process.
- \* Need flexible schema := columns added to tables, but not used on every row.
- \* Need high availability = AS they are not distributed  $\rightarrow$  SPOF.  
IF DB goes down, failover backup sys. must happen  $\Rightarrow$  takes time.
- \* Need horizontal scalability

**HA** : very little downtime of the system, it is always on and functioning .

**H.S** : ability to add servers and nodes to improve performance

# POSTGRE SQL

5

```
import psycopg2
```

```
try:
```

```
    conn = psycopg2.connect("host=127.0.0.1  
                            dbname=studentdb user=student password=1234")  
except psycopg2.Error as e:  
    print("Error caught")  
    print(e)
```

```
cur = conn.cursor()
```

```
conn.set_session(autocommit=True)
```

```
cur.execute("create database <name>")  
conn.close()
```

```
cur.execute('CREATE TABLE IF NOT EXISTS  
<table_name> (<col-name> <col-type>,  
...);')
```

```
('INSERT INTO <table_name> (<col-name>, ..)  
VALUES (%s, ...),  
(<value-1>, ..));
```

# NO SQL

not only

= Non Relational DB

simple design, horizontal scaling,  
finer control of availability.

Data structures  $\leftrightarrow$  RDB

Faster operations.

- Different types of NoSQL  
implementations

Common types =

DB	Type
Apache Cassandra	Partition row store
Mongo DB	Document store
Dynamo DB	Key-Value store
Apache HBase	Wide Column store
Neo 4j	Graph DB

# CASSANDRA

6

Data distributed by partitions across nodes or servers and the data is organized in the columns and rows format. → Own language: CQL

**Keyspace**: Collection of tables [ $\approx$  DB]

**Table**: A group of partitions

**Rows**: A single item

**Partition**: Fundamental unit of access.

Collection of row(s). - How data is distrib.

**Primary Key**: PK is made of a partition key + clustering columns-

**Columns**: Clustering and Data  
Labeled elements-

→ Provides scalability and high availability without compromising performance.

Linear scalability + Fault-tolerance on commodity hw / cloud infra

⇒ Perfect platform for mission-critical data.

Good for

- Transactional logging (retail, healthcare)
- IoT
- Time series data
- Heavy on writes workloads.

Bad for

- Analytics
- Group by's
- Ad-hoc queries
- Joins

## WHEN USE NOSQL

- \* Large amount of data
- \* Need horizontal scale → + nodes  
+ performance
- \* High throughput - reads
- \* Need flexible schema
- \* Need high availability
- \* Different data type formats
- \* Distributed users - low latency

## WHEN NOT USE NOSQL

- \* ACID Transactions
- \* Need to use JOINS
- \* Aggregations + Analytics
- \* Changing business requirements
- \* Small dataset
- \* Queries not available, flexible.

→ These properties are by design.

Some products can be customized  
to add this support.

```
import cassandra  
  
from cassandra.cluster import Cluster  
cluster = cluster(['127.0.0.1'])  
session = cluster.connect()  
session.execute("CREATE KEYSPACE IF NOT EXISTS name  
WITH REPLICATION =  
{'class': 'SimpleStrategy',  
'replication_factor': 1}  
...")  
  
session.set_keyspace('name')
```

# RELATIONAL DATA MODEL

## IMPORTANCE

- ✓ Standardization of data model
- ✓ Flexibility in adding / altering tables
- ✓ Data Integrity
- ✓ SQL
- ✓ Simplicity
- ✓ Intuitive Organization

OLAP - Online Analytical Processing

→ DB optimized for analytical + ad hoc queries.

OLTP - Online Transactional Processing

→ DB optimized for less complex queries  
in large volumes

# Normal Forms

Normalization → to reduce data redundancy  
and increase data integrity.  
correctness of data

copies  
of  
data

Denormalization → must be done in read heavy workloads to increase perform.

1. Free DB of unwanted insertions, updates, deletion dependencies
2. Reduce need of refactoring the DB as new types of data are introduced.
3. Make data model informative to user
4. Make DB neutral to query statistics.

OBJECTIVES

## normal forms

### 1NF

- Atomic values: each cell must contain unique single values (no sets, arrays, etc)
- Be able to add data w/o altering tables.
- Separate different relations into diff tables.
- Keep relation between tables together with FK

## 2NF

### • $\oplus$ 1NF

- All columns in the table must rely on the PK

## 3NF

### • $\oplus$ 2NF

- No transitive dependencies

⇒ to get from  $A \rightarrow C$  you shouldn't need to go through B

: when you want to update data, we want to be able to do in just 1 place.

## Fact & Dimensional Tables

Fact tables → consists of measurement metrics, facts of a business process.

Could contain aggregations of data, but are not meant to be modified in place.

Dimension tables → structures that categorizes facts in order to enable users to answer business questions.

Dimensions are people, products, places, time...

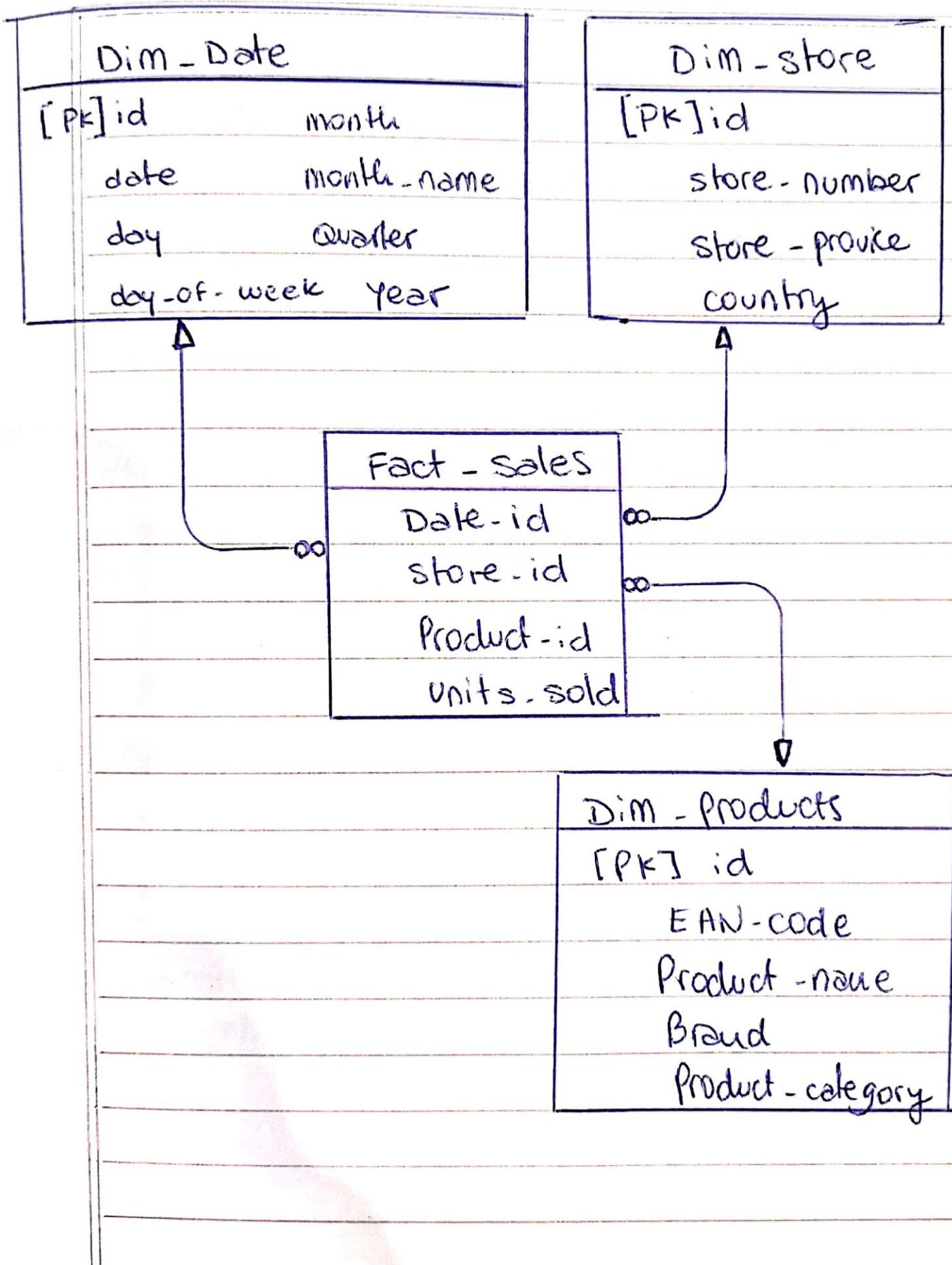
They give context to facts.

PKs of dim table are included in fact tables.

2 main schemas.

→ star schema

→ snowflake schema



## STAR SCHEMA

simplest style of data most schema  
→ consists of 1 or more fact tables  
referencing any number of dimension tables.

### BENEFITS

- Denormalized
- Simplified queries
- Fast Aggregation

### DRAWBACK

- Denormalization issues
- Data integrity
- Decrease query flexibility
- Many-to-many relationship

## SNOW SCHEMA

Tables in a multidim DB represented by centralized fact tables which are connected to multiple dimensions

- Allows for 1→many relationships
- more normalized than star

→ Not 3NF though

# NO SQL DATA MODEL

not-only sql

- Need H.A. in the db

- Have large amounts of data

- Need linear scalability

- Low latency

- Faster reads & writes

## CASSANDRA

- Open Source NoSQL DB
- Masterless Architecture - NO SPOF
- High Availability
- Linearly Scalable → + Nodes ⇒ + perf.
- Used by Hulu, Facebook, Twitter, Netflix

# DISTRIBUTED DB

→ To have H.A. we need more copies of your data.

Database runs over several nodes.

HA: no downtime

## Eventual Consistency

Consistency model used in distrib. DB computing to achieve HA.

→ if no new updates are made to a given data item, eventually all access to the item will return the last updated value - Data may be inconsistent for only milliseconds

## CAP Theorem

It is impossible for a distrib. db to store simultaneously provide more than two out of the following three guarantees of consistency, availability and partition tolerance.

**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 guarantee that the data is the latest update.

**PARTITION TOLERANCE:** System continues to work regardless of losing network connectivity between nodes.

## DENORMALIZATION in CASSANDRA

→ Absolutely critical - A MUST!

The biggest take away when doing modeling in Cassandra is to think your queries first.

- There is no joins here.

- must be done for fast reads

- Cassandra has been optimized for fast writes.

- One query per table → Great strategy.

(partition by search item)

```
import cassandra
cluster = Cluster({'':})
session = cluster.connect()
```

session.execute ("""

CREATE KEYSPACE IF NOT EXISTS

udacity WITH REPLICATION =

'class': 'SimpleStrategy',

'replication\_factor': 1 (""")

session.set\_keyspace ('udacity')

I want to run queries:

① select \* from music-library where year=1970

② select \* from artist-library where artist = "..."  
joint key

for ① Table: music-library

columns: year, artist name, album

PK : (year, artist name) | artist name  
is clustering  
key

for ②: Table artist-library

PK: (artist-name, year)

## PRIMARY KEY

→ how each row can be uniquely id.  
and how data is distributed across  
the nodes in our system.

- Made up of only the PARTITION KEY

→ simple PK

or w/ the addition of CLUSTERING KEY

→ composite PK

Determined the distribution of the data.

⇒ The PK will be hashed and stored on the node that holds that range of values

- To choose a PK:

→ We want to spread our data evenly.  
→ Ensure that is unique for each row.

## CLUSTERING COLUMNS

→ determines the sort order within a partition. (DESC order)

More than 1 can be added, and they will sort in the order they were declared. → nested sorted order

... PRIMARY KEY ((year),  
artist-name, album-name)

→ If you want to use cc in select you must use them in order.

## WHERE CLAUSE

DM in Cassandra is query focused

→ Partition key must be in your where

→ Clustering cols. can appear in order.

Select \* from table → can be done,  
but is highly discouraged.