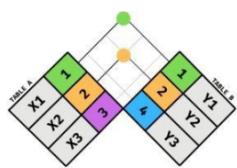


# Data Eng. Roadmap

## Advanced Joins SQL:

pkey: primary key  
fkey: Foreign key

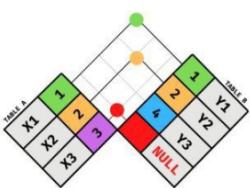
PKeys in tables that are FKeys on another ones.



INNER JOIN

```
SELECT
<SELECT LIST>
FROM TABLE_A A
INNER JOIN TABLE_B B
ON A.KEY = B.KEY
```

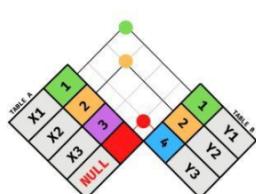
KEY	VAL_X	VAL_Y
1	X1	Y1
2	X2	Y2



LEFT JOIN ( kinda inner join! )

```
SELECT
<SELECT LIST>
FROM TABLE_A A
LEFT JOIN TABLE_B B
ON A.KEY = B.KEY
```

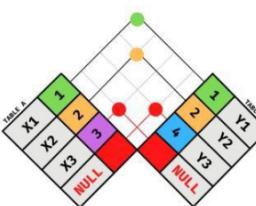
KEY	VAL_X	VAL_Y
1	X1	Y1
2	X2	Y2
3	X3	NULL



RIGHT JOIN ( kinda outer join! )

```
SELECT
<SELECT LIST>
FROM TABLE_A A
RIGHT JOIN TABLE_B B
ON A.KEY = B.KEY
```

KEY	VAL_X	VAL_Y
1	X1	Y1
2	X2	Y2



FULL OUTER JOIN

```
SELECT
<SELECT LIST>
FROM TABLE_A A
FULL OUTER JOIN TABLE_B B
ON A.KEY = B.KEY
```

KEY	VAL_X	VAL_Y
1	X1	Y1
2	X2	Y2
3	X3	NULL
4	NULL	Y3

1.

① Alias: apelido de tablas

E.g.: FROM BOOK B WHERE B.BOOK-ID ...



- Foreign Keys : point to primary keys on another table .
- Primary Keys : uniquely identifies an entity / table

## Data Modelling

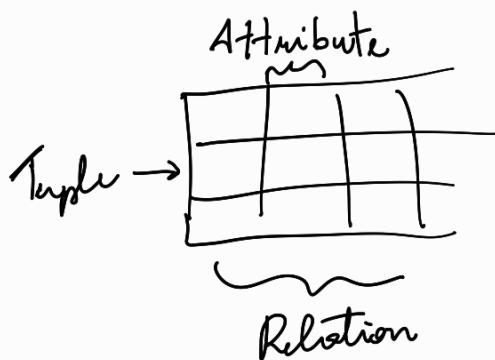
Data Eng with AWS

- Data modelling = Database modelling
- An "abstract process"

Process :

1. Gather requirements
2. Conceptual modeling (entity mapping)
3. Logical " " (tables, schemas, columns)
4. Physical "  
    ↳ DDL (Data Def. Language)

### ② Relational Model



SQL - Structured  
Query  
language

- Scheme : collection of Tables .

- The RDB and relational db's were created in 60' / 70's.

Attribute → Column  
(feature)

Tuple → Row

When to use Relational DB ?

- Joins só podem ser feitos em RDBs
- Smaller data volumes
- Secondary Index for quick searching  
(Primary index is the primary key)

ACID Transactions  
(Transactions são operações ACID feitas  
e percebidas como uma  
única operação)

ACID  
Atomicity: each transaction or a single "unit", which either succeeds or fails completely. In fail, the db is left unchanged; rolls back to prior state.

Isolation: Concurrent transactions have db in the same state that would have been obtained if the transactions were executed sequentially; not impacted by other transactions.

Consistency: Data integrity; recall constraints, rules, etc.

OBS : Join só podem ser feitos se existirem  
valores que são match nos duas  
tabelas!

## • When NOT to use RDBs ?

- Large amounts of data
- Quando precisar armazenar dados de diferentes formatos
  - ↳ RDBs só escalam verticalmente, i.e., adicionando mais memória na máquina; NÃO É DISTRIBUÍDO!
- Qdo necessitar de leitura rápida
  - ↳ ACID temps mais lento ↓
- Qdo precisar escalar horizontalmente
  - ↳ add more nodes/machines (<sup>melhor</sup>  
<sup>mais</sup>  
<sup>NoSQL</sup>)
- " " " de alta disponibilidade

## ④ NoSQL Databases

- Simpler horizontal scaling
  - Different types
- more suitable when:
    - need high availability
    - " to scale out quickly
    - large amounts of data
    - different data types
    - high throughput: faster reading/writing

- Partition Row store
- Document store
- key-value
- wide column ("flexible schema")

## → Apache Cassandra

- Keyspace
  - collection of tables
- tables
  - group of partitions
- rows
- Uses its own query language : CQL

~~~~~

## Relational Data Models

- OLAP (analytical processing)
  - ↳ more aggregations
  - ↳ read-only
- OLTP (transactional processing)

## ⇒ Normalization & Denormalization

### • Normalization

- ↳ reduce data redundancy and increase data integrity
- ↳ reduce copies of data

### • Denorm

- ↳ must be done to increase performance in mad heavy workloads

## ⇒ Objectives of Normal Form

1. Free db's from unwanted insertions, updates and deletion dependencies
2. Reduzir necessidade de misturas e de novos tipos de dados no intuito de
3. Make more information
4. Formar o DB neutro às queries  
(é fazer o DB pl. a query, como é no caso do NoSQL)

1NF

- Unique values on each cell  
(atomicity)

2NF

- All columns in the **normalized** table must depend on a **Primary Key**.

- No "partial dependencies": só com composite keys!

- Partial depend. é qdo uma non-key column depende da chave composta parcialmente; depende só de 1 campo da mesma.

3NF

- Eliminate columns that don't depend on a Primary Key
- No transitive dependencies ( $A \rightarrow C$ , avoiding going through B)

- Transitive Dependency:  
Um campo C depende de A (chave primária), mas esse mesmo campo C também depende de B (pode ser determinado por) outro B (que NÃO é chave primária)



## NoSQL Data Models

NoSQL = Not only SQL

- High availability → copies of data

Distributed Databases

### VOCAB!

- \* stale data = dados obsoletos
- \* retailer = varejista  
retail = varejo

availability:  
make requests and get a response.

## CAP Theorem

Impossible for a distributed data store to simultaneously provide more than 2 of the following guarantees:

- Consistency      - Availability      - Partition Tolerance



Every read from the db gets the latest piece of data or an error.



System continues to work regardless of losing network between nodes.

Apache Cassandra is an AP (Availability, Partition Tolerance) db.

• Eventual Consistency in NoSQL vs Consistency in ACID SQL dbs

→ ACID: only transactions that obey (abide) by constraints and db rules are committed into the db, otherwise db keeps previous state.

→ CAP: cf. definition above!

- Availability and Part. Tolerance on the biggest Mgs.

- Denormalization is a MUST in Cassandra!

- ↳ Think queries first!
- ↳ No JOINs!
- ↳ 1 query per table!

- Primary Keys (PK)

- First element of the PK is the PARTITION KEY
- PK made up of either just PART. key or with the addition of CLUSTERING COLUMNS
- PART. KEY will determine distribution of data across the system.

| year                  | artist-name               | album-name |
|-----------------------|---------------------------|------------|
| ↑<br>Partition<br>Key | ↑<br>Clustering<br>Column |            |

Primary Key

```
CREATE TABLE music-lib
(year int, artist-name text,
album-name text,
PRIMARY KEY ((year),
artist-name, album-name))
```

OBS: If PK are wrapped in "()", then the PK is composite!  
Otherwise, the 1st field is the Partition key.

- Clustering columns will sort data in ASC order
- More than one column can be added.
- Clustering columns will sort in order of how they were added to the primary key.

→ See example of query in table with PART.KEY and CLUSTERING COL. below:

```

Lesson 3 Exercise 2 Primary X Lesson 3 Demo 2 Primary | X Lesson 3 Exercise 3 Cluster X Lesson 3 Exercise 3 Cluster X
+ X C Markdown ▾
1966 The Monkees The Monkees Los Angeles
[18]: import pandas as pd

# CREATE TABLE IF NOT EXISTS music_library
# (year int, city text, artist_name text, album_name text, PRIMARY KEY ((year), artist_name, album_name))
#
# The query above has "year" as PARTITION KEY and artist_name & album_name as CLUSTERING COLUMNS.
# The result set is below, with the data in ascending order.

df = pd.DataFrame(list(session.execute(query)))
df

```

|   | artist_name    | album_name    | city        | year |
|---|----------------|---------------|-------------|------|
| 0 | The Carpenters | Close To You  | San Diego   | 1970 |
| 1 | The Beatles    | Let it Be     | Liverpool   | 1970 |
| 2 | The Beatles    | Rubber Soul   | Oxford      | 1965 |
| 3 | The Who        | My Generation | London      | 1965 |
| 4 | The Monkees    | The Monkees   | Los Angeles | 1966 |

- Na query, a PK DEVE ser utilizada!  
Além disso, as colunas de clustering (CLUSTERING COLS.) devem vir depois na ordem da query  
foram colados na PRIMARY KEY.

**Important!**

- CREATE and INSERT statements must abide the order of the PK fields. Example below. FOR CASSANDRA!

Ex.: If in CREATE the PK is "PRIMARY KEY(id, orderNum)"  
the INSERT must be:

> `INSERT INTO <table> (id, orderNum, orderBuid, ...)`  
`VALUES (1, 999, 11, ...)`

ordered according  
to the PK

## Using SELECT \* in an SQL Query Is a Bad Practice

Be specific with your selects

### SELECT \* FROM TABLE

- 1. Unnecessary I/O
- 2. Increased Network traffic
- 3. More Application memory
- 4. Depends on Column Order
- 5. Fragile Views
- 6. Conflict in JOIN Query
- 7. Risky while copying data

) Why it's a bad  
idea to use  
"SELECT \*"

# Data Warehouses in the Cloud

17/04/23

## Data Warehouse Architecture

A data warehouse is a copy of transaction data specifically structured for query and analysis. - Kimball

A data warehouse is a subject-oriented, integrated, nonvolatile, and time-variant collection of data in support of management's decisions. - Inmon

A data warehouse is a system that retrieves and consolidates data periodically from the source systems into a dimensional or normalized data store. It usually keeps years of history and is queried for business intelligence or other analytical activities. It is typically updated in batches, not every time a transaction happens in the source system. - Rainard

## OLTP vs OLAP

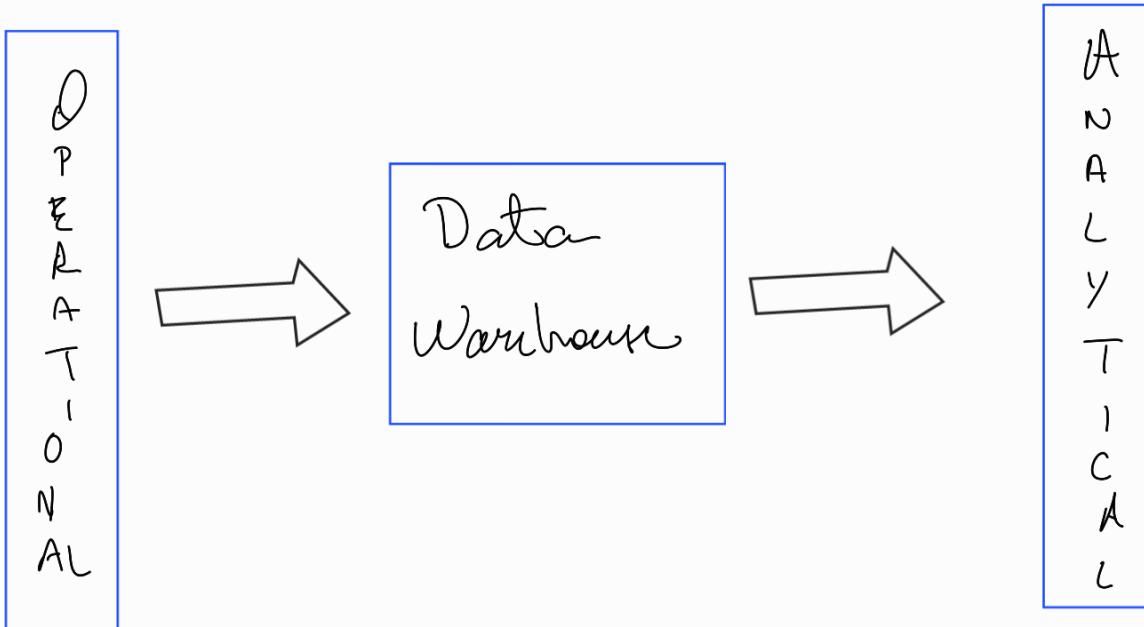
Online Transactional P.  
or Analytical P.

• Operational processes: make it work

- I. Find goods and make orders
- II. Stock and find goods
- III. Pick up and deliver goods

• Analytical n.: what is going on?

- I. Assess performance of sales staff
- II. See the effect of different sales channels
- III. Monitor sales growth.



Ex.: track something in a retail is OLTP, but examine a market segmentation over a period of time is OLAP.

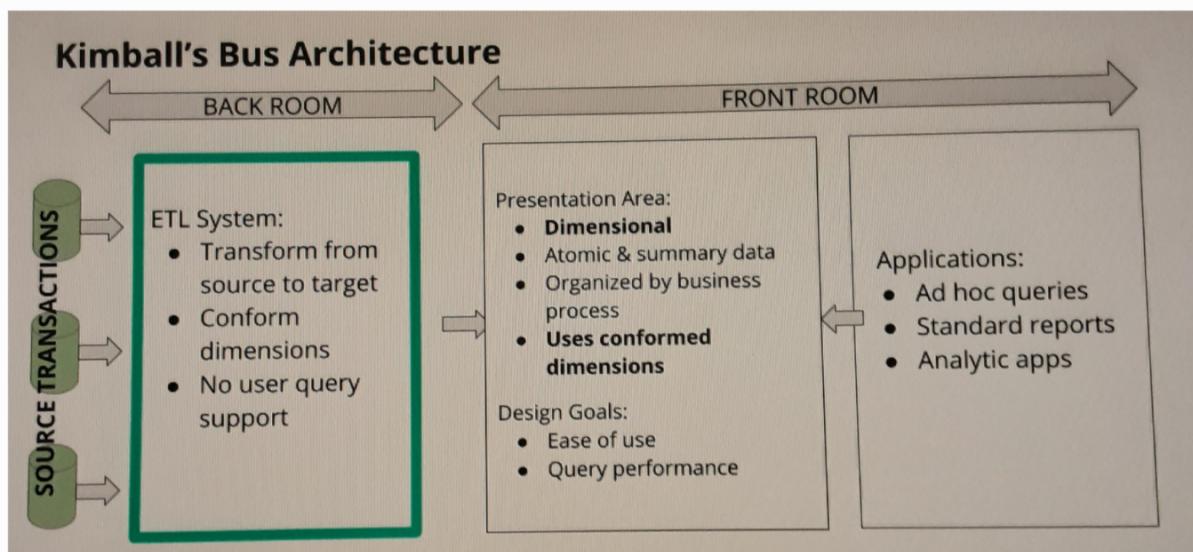
(OBS): If a database is small, it may be fine for both OLTP and OLAP.

If big: queries slow; schema hard to understand!

→ Data is loaded in a DW in batches, not after every operation of a transactional system.

② No streaming?

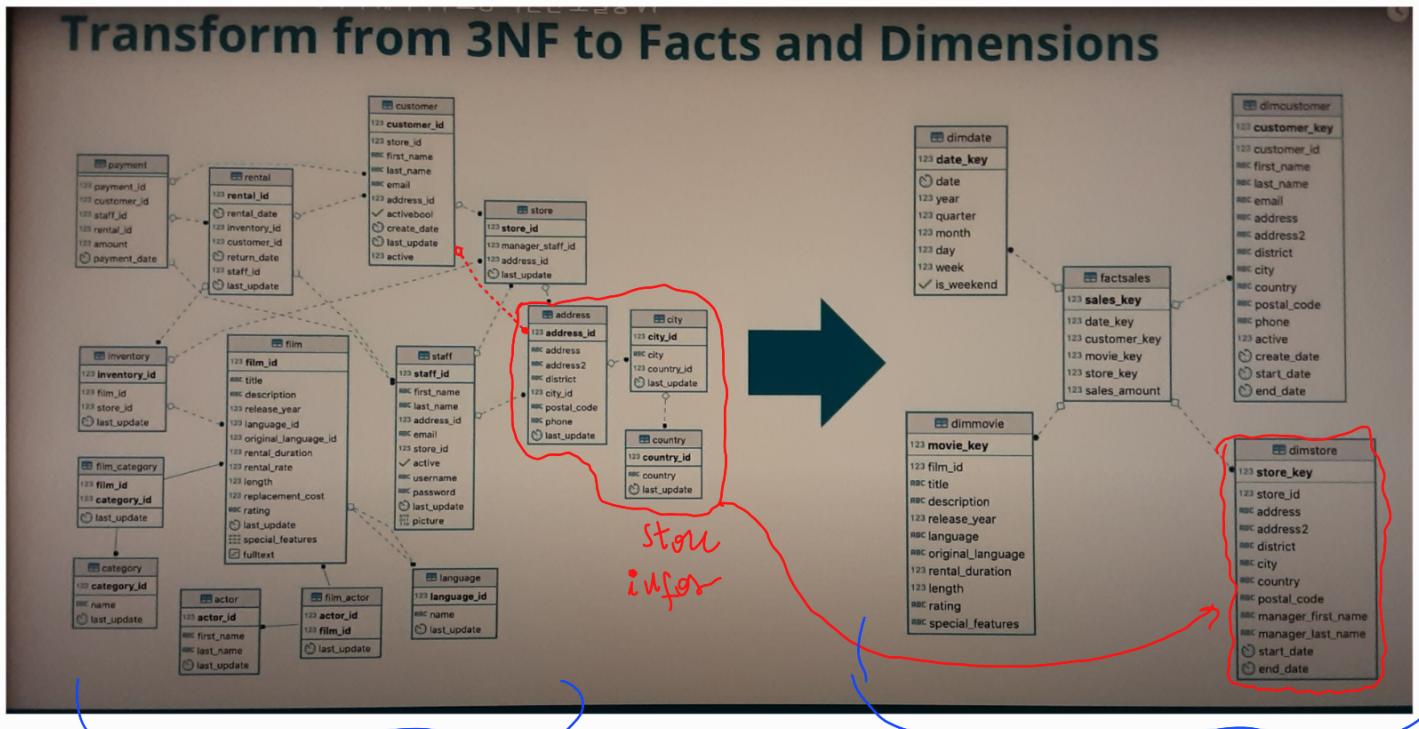
→ Data structured in a Dimensional Modeling, easy to query.



Kimball's Bus Arch:

1. Results in common data model shared by different departments
2. Data not kept at aggregated level, rather at atomic level
3. Organized by business processes, and used by different departments.

# Identifying Facts and Dimensions



difficult for a business user to understand

easier for a business user to understand

~ From 3NF to ETL :

Extract

- query 3NF db

Transform

- join tables
- change types
- add new columns

Load

- insert into facts | dimensions

## Facts

- Business units
- Quantifiable data
- Numeric data
- "Fato"

Examples: nº items bought by a customer;  
instances of positive feedback;  
sales amount;  
length of phone call.

## Dimensional

- Context of business units (dates, time, etc)
- Ex.: store where item was purchased, the customer, etc.

Examples: tipo de resposta;  
tipos de móveis buscados online;  
local loja;  
employee ID, etc.

## OBS

"Dimensional Modeling" exercise and "ETL"

exercise → using "SERIAL" while creating  
fact/dimension tables in PostgreSQL will generate

an auto-incrementing number once records are inserted!

## Grouping Sets [SQL]

SELECT ... FROM ... JOIN ...

GROUP BY grouping sets ( f1,f2,(f1,f2) );

See the example below using the F/DIM tables already shown before (cf. "Identifying Fact & Dimension tables"):

```
total revenue
1: %%sql
SELECT sum(sales_amount) as revenue
FROM factSales

revenue by country
1: %%sql
SELECT dimStore.country,sum(sales_amount) as revenue
FROM factSales
JOIN dimStore on (dimStore.store_key = factSales.store_key)
GROUP by dimStore.country
order by dimStore.country, revenue desc;

* postgresql://student:***@127.0.0.1:5432/pagila_star
2 rows affected.

1: country revenue
Australia 33726.77
Canada 33689.74

revenue by month
2: %%sql
SELECT dimDate.month,sum(sales_amount) as revenue
FROM factSales
JOIN dimDate on (dimDate.date_key = factSales.date_key)
GROUP by dimDate.month
order by dimDate.month, revenue desc;

* postgresql://student:***@127.0.0.1:5432/pagila_star
5 rows affected.

2: month revenue
1 4824.43
2 9631.88
3 23886.56
4 28559.46
5 514.18
```

```
revenue by month & country
1: %%sql
SELECT dimDate.month,dimStore.country,sum(sales_amount) as revenue
FROM factSales
JOIN dimDate on (dimDate.date_key = factSales.date_key)
JOIN dimStore on (dimStore.store_key = factSales.store_key)
GROUP by (dimDate.month, dimStore.country)
order by dimDate.month, dimStore.country, revenue desc;

* postgresql://student:***@127.0.0.1:5432/pagila_star
10 rows affected.

1: month country revenue
1 Australia 2364.19
1 Canada 2460.24
2 Australia 4895.10
2 Canada 4736.78
3 Australia 12060.33
3 Canada 11826.23
4 Australia 14136.07
4 Canada 14423.39
5 Australia 271.08
5 Canada 243.10
```

## revenue total, by month, by country, by month & country All in one shot

- watch the nones

```
]: %%time
%%sql
SELECT dimDate.month, dimStore.country, sum(sales_amount) as revenue
FROM factSales
JOIN dimDate on (dimDate.date_key = factSales.date_key)
JOIN dimStore on (dimStore.store_key = factSales.store_key)
GROUP BY grouping sets ((), dimDate.month, dimStore.country, (dimDate.month, dimStore.country));
```

\* postgresql://student:\*\*\*@127.0.0.1:5432/pagila\_star  
18 rows affected.  
CPU times: user 0 ns, sys: 0 ns, total: 0 ns  
Wall time: 27.1 ms

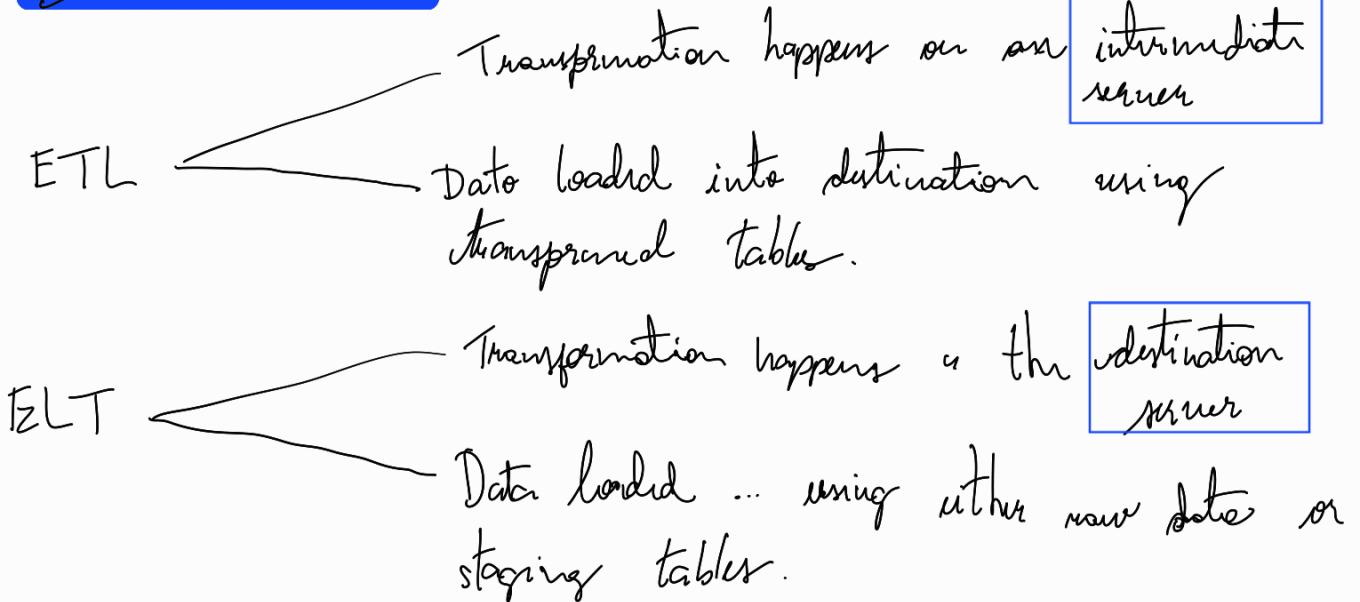
| month | country   | revenue  |
|-------|-----------|----------|
| 1     | Australia | 2364.19  |
| 1     | Canada    | 2460.24  |
| 1     | None      | 4824.43  |
| 2     | Australia | 4895.10  |
| 2     | Canada    | 4736.78  |
| 2     | None      | 9631.88  |
| 3     | Australia | 12060.33 |
| 3     | Canada    | 11826.23 |
| 3     | None      | 23886.56 |
| 4     | Australia | 14136.07 |
| 4     | Canada    | 14423.39 |
| 4     | None      | 28559.46 |
| 5     | Australia | 271.08   |
| 5     | Canada    | 243.10   |
| 5     | None      | 514.18   |
| None  | None      | 67416.51 |
| None  | Australia | 33726.77 |
| None  | Canada    | 33689.74 |

→ "month & country" grouping  
→ "month" grouping  
→ likewise  
→ likewise  
→ likewise  
→ likewise  
→ "( ), i.e., the total revenue regardless month & country  
→ "country" grouping

## ☒ Data Warehouses in the CLOUD

- Why moving to cloud DW's ? = why ETL to ELT?
  - Scalability: large amounts of data.
  - Flexibility: adding/removing data is easier.
  - Cost shifting: "transform" step (more expensive/costly) is done last, no DEngs seen performance. Just-in-time transformations to meet the highest priority business needs first.

## ETL vs ELT



## → No SQL Databases in AWS

- DynamoDB - key value
- Document DB - document
- Keyspaces - column oriented
- Neptune - graph
- Time stream - time series

# AWS Data Warehouse Technologies

- All AWS services require Identity and Access Management (IAM) policies, roles, and their associated permissions
- **IAM User**: is an identity with long-term credentials that is used to interact with AWS in an account.
- **Security Groups**: act as a firewall rules for e.g. a Redshift cluster to control inbound and outbound traffic.
  - Ps.: Users are different from Roles.
    - Person / Application
    - Uniquely associated
    - Long-term credentials
    - Assumable by anyone who needs
    - Temporary credentials
- ~→ **Federated Access**: external access to AWS resources

**ARN**  
Amazon Resource Name

## 💡 Benefits of Infrastructure as Code

- IaC é o provisionamento de infraestrutura da TI via linguagem de codificação descritiva de alto nível (e.g. json, yaml, etc).

### Benefícios da IaC

**Agilidade:** A capacidade de automação da IaC acelera o processo de provisionamento de uma infraestrutura para desenvolver, testar e produzir aplicações. Permitindo a configuração de uma infraestrutura completa executando apenas um script. Assim possuindo os mesmos princípios seguidos na cultura DevOps no quesito da velocidade e consistência do ciclo de vida de entrega do projeto.

**Consistência:** Com a automatização de processos oriundos da IaC, falhas e discrepâncias que seriam criadas por um processo manual são praticamente eliminadas. A infraestrutura como código tem a capacidade de evitar esses possíveis problemas, pois seus arquivos de configuração possuem apenas uma única fonte de informação, assim garantindo a possibilidade de realizar repetidas implantações de forma consistente e sem disparidade de informações.

**Segurança:** Assim como todos os serviços de cloud computing, a IaC traz mais segurança para o ambiente de TI da sua empresa. As ferramentas de infraestrutura como código permitem a correção rápida de erros e a solução de problemas de forma automatizada, deixando assim uma infraestrutura mais segura e gerenciada pela empresa.

### Quiz Question

Which of the following are advantages of Infrastructure-as-Code over creating infrastructure by clicking-around?

- Sharing: One can share all the steps with others easily
- Reproducibility: One can be sure that no steps are forgotten
- Multiple deployments: One can create a test environment identical to the production environment
- Maintainability: If a change is needed, one can keep track of the changes by comparing the code

↙ all are correct!

## Common Terms in AWS

### Terms

In the previous illustration we used specific terminology to describe how to obtain access to resources. These IAM terms are commonly used when working with AWS:

#### IAM Resources

The user, group, role, policy, and identity provider objects that are stored in IAM. As with other AWS services, you can add, edit, and remove resources from IAM.

#### IAM Identities

The IAM resource objects that are used to identify and group. You can attach a policy to an IAM identity. These include users, groups, and roles.

#### Principals

A person or application that uses the AWS account root user, an IAM user, or an IAM role to sign in and make requests to AWS. Principals include federated users and assumed roles.

#### Human users

Also known as *human identities*; the people, administrators, developers, operators, and consumers of your applications.

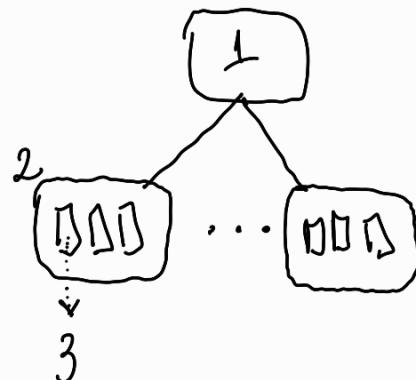
#### Workload

A collection of resources and code that delivers business value, such as an application or backend process. Can include applications, operational tools, and components.

# Implementing a DW on AWS

## ② Redshift Arch.

1. Leader Node
2. Compute nodes
3. Node slices



Important!

- Total # **nodes** in Redshift = # AWS EC2 instances used in the cluster
- Slices have dedicated storage and memory
- Total # **slices** in a **cluster** is the unit of parallelism (max # of partitions/table)

### Redshift Node Types and Slices

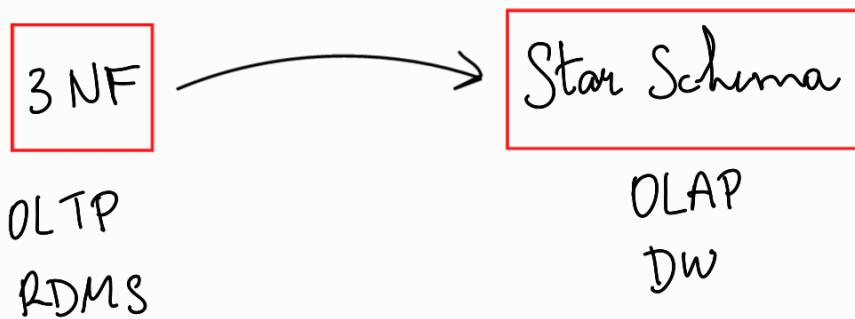
#### Compute Optimized Nodes:

- Start with these for lower costs and a capacity of 5 terabytes

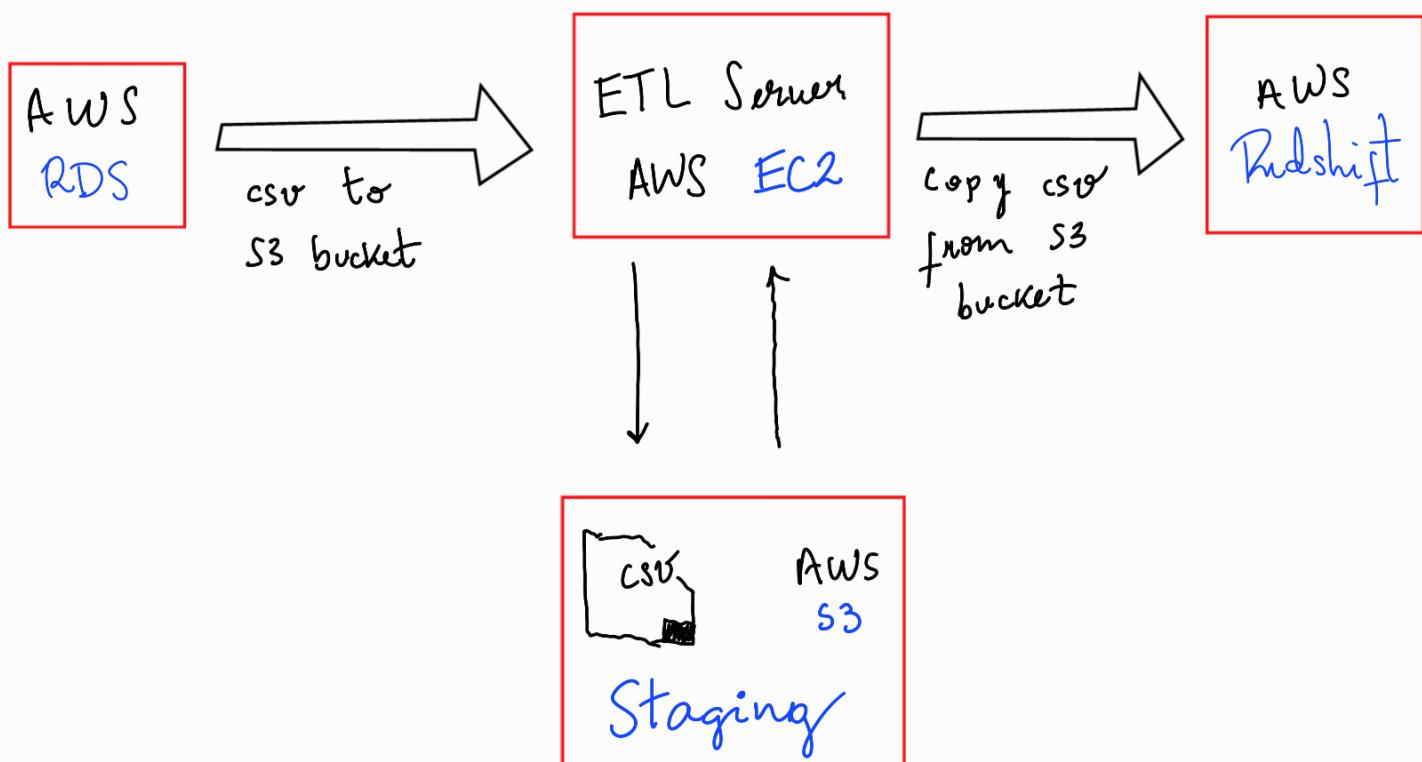
#### Storage Optimized Nodes:

- Higher costs, not as fast, but higher capacity

## B) Ingesting Data at Scale

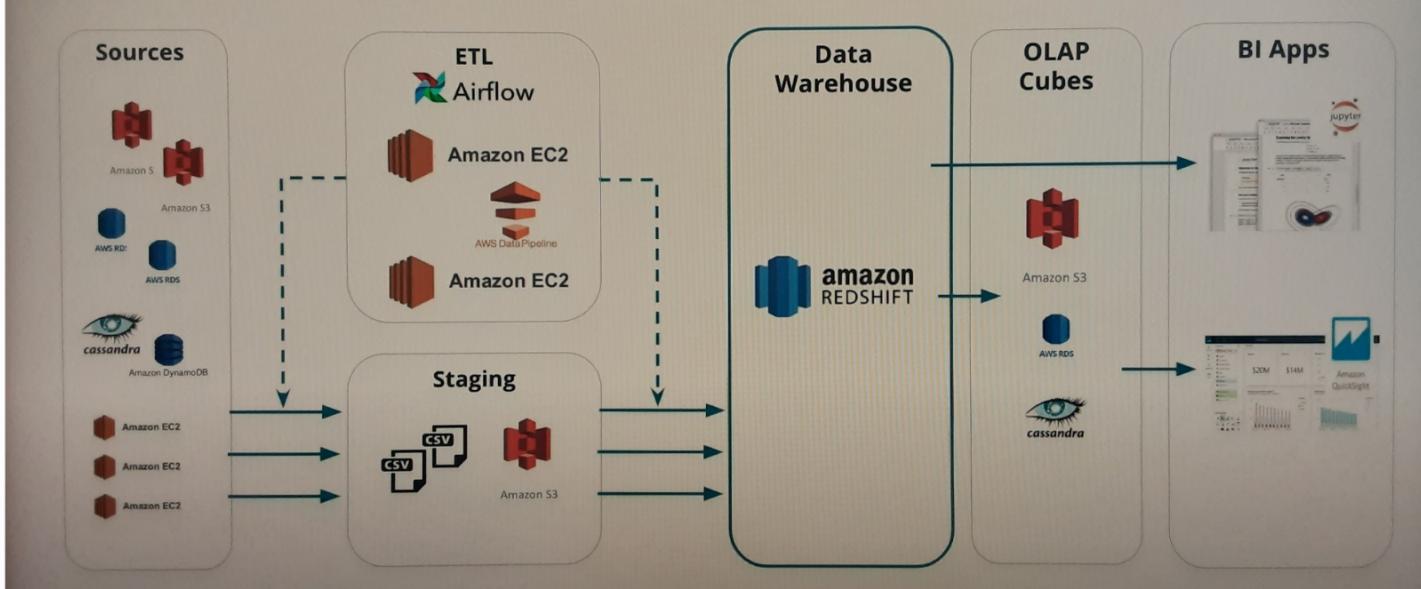


### 1. Redshift SQL to SQL ETL



- Most of AWS products are able to communicate with S3 storage.
- Staging Area: tabela temporária; materialized view no banco ("área temporária") mapeia os dados brutos p/ serem transformados!

# Redshift Architecture & Dataflow



## Ingesting at Scale

- S3 staging → Redshift with COPY command
- Ingest files in parallel ↳ if file large
  - common prefix
  - manifest file (?)
- Big files? Use compression
- Use same AWS region to avoid network traffic

→ INSERT will be very slow

## QUIZ

- Purpose of EC2 instance in ETL  
Acts as a client to RDS & Redshift to issue COPY commands

- Benefits of using S3 for ETL storage
  - Reliable, scalable and storage free of worry
  - Don't provide processing
- Small data
  - Possible to use EC2 for small data
- 

## ☞ Distribution Styles

- Optimizing the way a table is partitioned up into many pieces and distributed across slices on different machines
  - Column or **DISTKEY** sets distribution style to "KEY".
  - One can use a table attribute to set the dist. style.
- ```
CREATE TABLE <table> (...)
  DISTSTYLE ALL;
```

## EVEN Distribution

- Round-robin method, regardless of values.
- When a table is not used in queries with joins
- Each slice with almost same # records

## KEY "

- Rows with similar values are placed in the same slice
- Used for tables frequently joined together.

## ALL "

- for small tables (e.g. lookups tables)
- Record a copy of the table/column in all slices.

## AUTO "

- Redshift is responsible for distributing

# Redshift Cluster

- [REDACTED]
- [REDACTED]
- [REDACTED]



## Final Project Reviews

### About primary keys in staging tables

Staging tables are properly defined.

#### Reviewer Note

In data warehousing, staging tables serve as temporary storage for raw data being imported from source systems, before it's transformed and loaded into the final warehouse tables. They are a crucial part of the Extract, Transform, Load (ETL) process, as they allow for efficient data extraction and loading, while minimizing the impact on operational systems.

When defining staging tables, we typically avoid creating IDENTITY() columns or any kind of synthetic primary key. This is for several reasons:

1. **Preservation of Raw Data:** The main purpose of a staging table is to be a faithful reflection of the source data. Adding an IDENTITY() column or any other computed columns could distort this raw view. We want to ensure that the data in the staging tables is an exact copy of the source data, allowing us to perform transformations, cleaning, and conforming operations in subsequent stages of the ETL process.
2. **Avoidance of Unnecessary Computation:** Assigning IDENTITY() to a column means SQL Server (or any other relational database management system) would automatically generate a unique number for each row. This computation is often unnecessary in a staging table, as we are more focused on moving the data quickly and efficiently from source to staging, rather than creating new data.
3. **Preparation for Transformation:** In many cases, primary keys or unique identifiers will be generated or transformed in the process of loading data from the staging table into the final tables of the data warehouse. These identifiers are often based on the business logic and data requirements of the final tables, rather than the raw data.

Therefore, keeping staging tables as simple and as close to the source data as possible allows for more flexibility and control in the ETL process. IDENTITY() columns and similar constructs are generally more relevant in the design of the final data warehouse tables, where they can help to enforce data integrity and facilitate querying.

CREATE statements in sql\_queries.py specify all columns for both the songs and logs staging tables with the right data types and conditions.