

INTRO TO DATA WAREHOUSES

15

"DW is a copy of transactions data specifically structured for query and analysis"

- Kimball

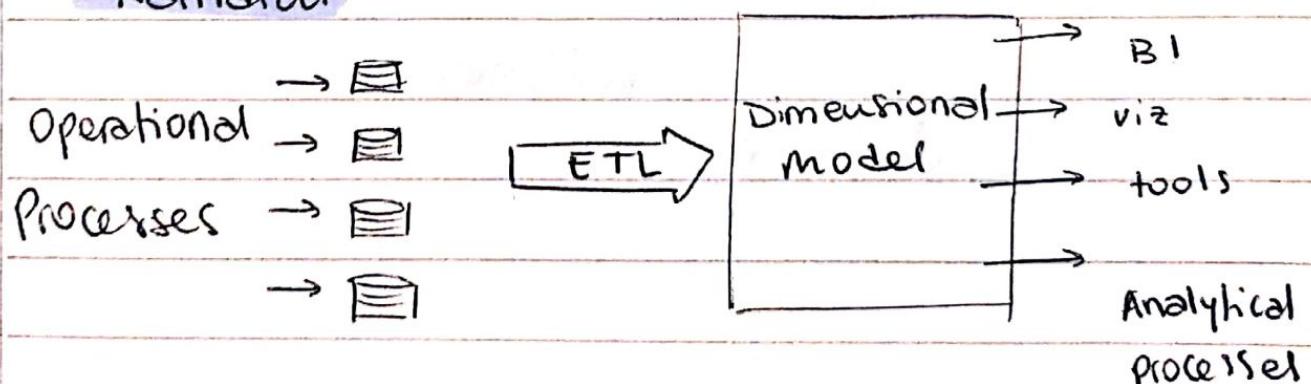
"DW is a subject-oriented, integrated, non-volatile, time-variant collection of data in support of management's decisions"

- Inmon

"DW is a system that retrieves and consolidates data periodically from source systems into dimensional or normalized data store."

It keeps years of history and is queried for business intelligence. Typically updated in batches'

- Rainardi



ETL

pip install ipython-sql

% load -ext sql

load

% .sql < one lines every > \$variables

% % .sql < multiline every >

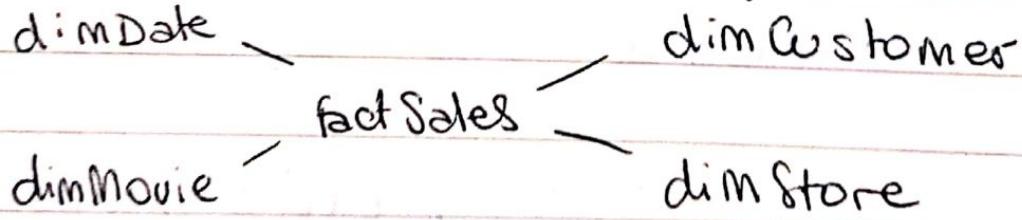
→ = "postgressql://{{...}"

% .sql \$conn-string

1º) Defino la estructura del modelo dimensional

La idea es extraer data de relational DB
en 3NF donde tengo que hacer
muchos Joins para obtener la info que
quiero a

→ Dim - Fact Tables - por ejemplo :



2º) Creo las tablas dimensiones y hechos 16

Para crear las tablas:

%% SQL

```
CREATE TABLE dimDate
{
    date-key SERIAL PRIMARY KEY,
    date date NOT NULL,
    year integer NOT NULL,
}
```

%.% SQL

```
CREATE TABLE FactSales
{
    sale-key SERIAL PRIMARY KEY,
    date-key integer REFERENCES
    dimDate (date-key),
    ...
    sales.amount numeric NOTNULL
}
```

Referencias {
dimensions {
Factos /
metricas {

3º) ETL data de tablas 3NF a Facts y Dims

% % SOL

L \leftarrow INSERT INTO dimDate (date-key, date, year, ..., is-weekend)

E \leftarrow SELECT DISTINCT (TO-CHAR(payment-date::DATE
'YYYYMMDD')::integer) AS date-key,
EXTRACT (year FROM payment-date),
AS year,
date(payment-date) AS date,
...
CASE WHEN EXTRACT (ISODOW FROM
payment-date) IN (6,7) THEN TRUE
ELSE FALSE
FROM payment;

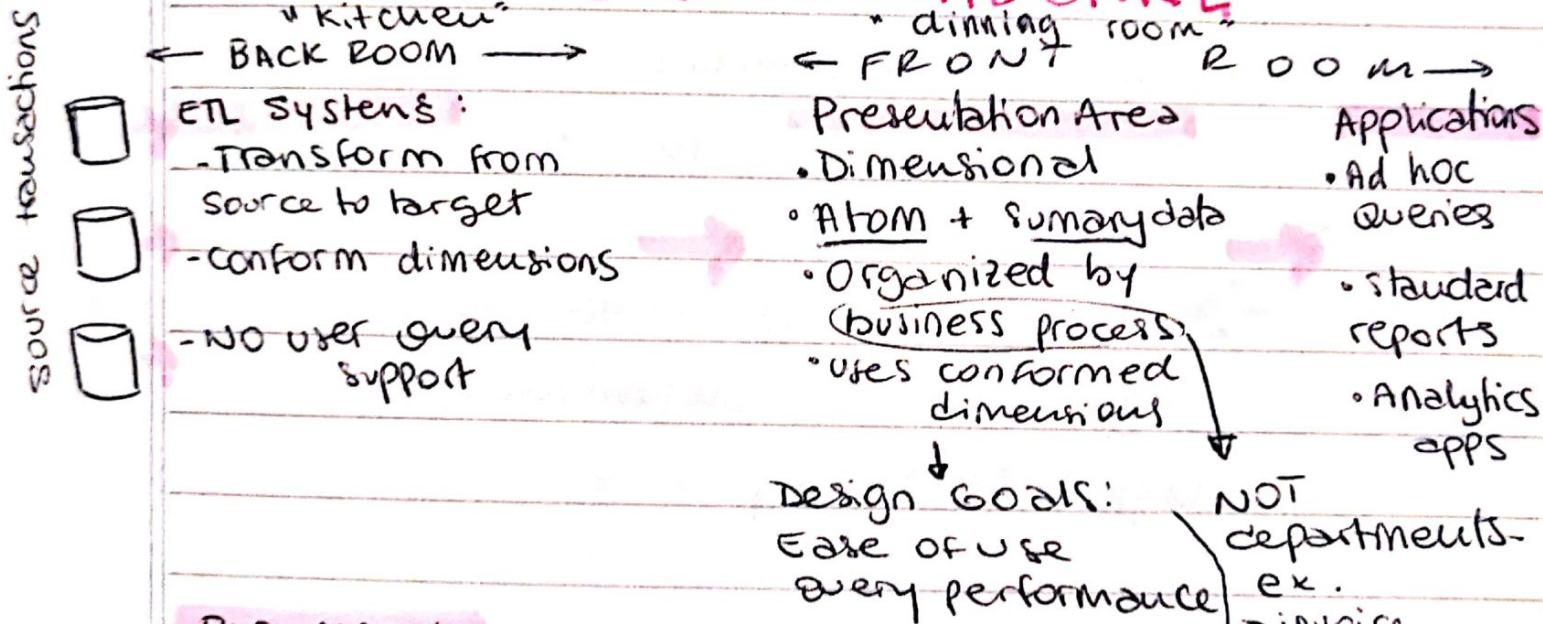
T {

we denormalize tables (put joins on
our dim tables) and transform the
data.

D.W. ARCHITECTURE

- Kimball's Bus Architecture
- Independent Data Marts
- Inmon's Corporate Information Factory
- Hybrid Bus + CIF

KIMBALL'S BUS ARCHITECTURE



BUS MATRIX

Business processes x Dimensions table

	Date	Product	Warehouse	...
BP1	x		x	
BP2	x	x		x
:				

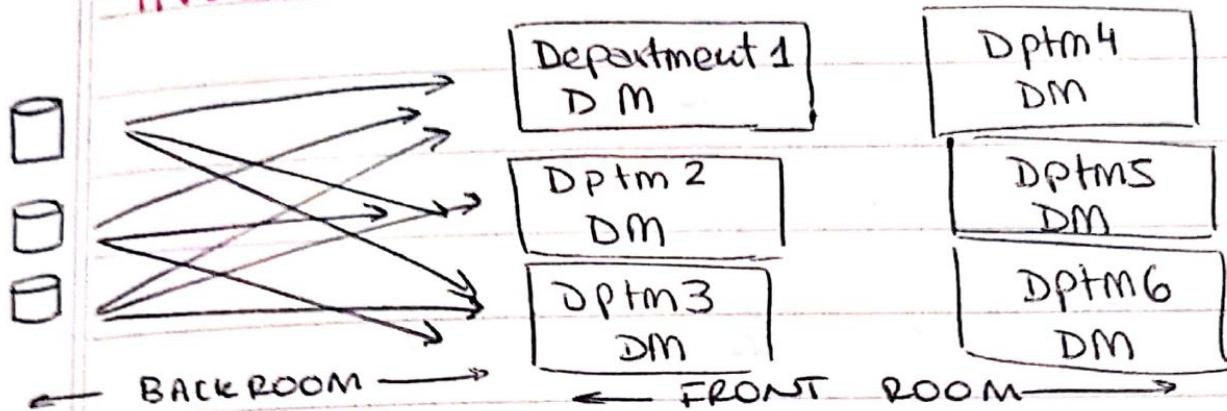
• Dimensions are across business departments

ETL

- EXTRACT
 - Get data from <> sources , OLTP
 - Delete stale/old data from sources
- TRANSFORMATION
 - Integrate <> sources together
 - Data cleansing : inconsistencies, duplications, missing values, etc.
 - Produce diagnostic metadata (this source has % of unreliable data)
- LOADING
 - Structure and load data to dimensional data model.

→ 1 Dimensional data model shared by all the departments, organized by business processes.

INDEPENDENT DATA MARTS



Departments have independent ETL processes and dimensional tables

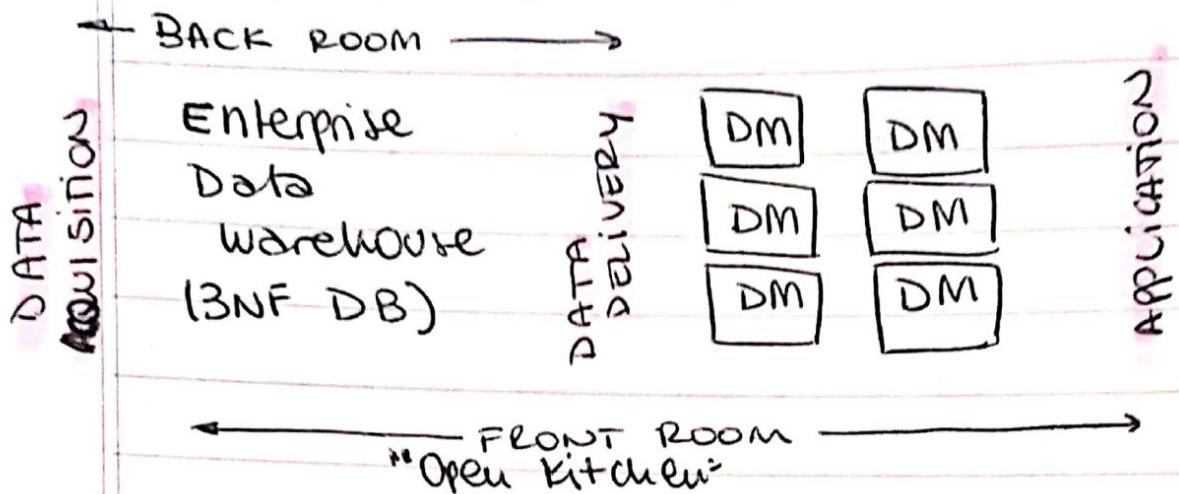
These separated dim tables are called DM.

Each have its own fact tables for the same event → no conformed dimensions.

Uncoordinated efforts might lead to inconsistent views

It happens ad-hoc because it's simpler and autonomous, but it's highly discouraged.

INMON'S CORP. INFORMATION FACTORY (CIF)

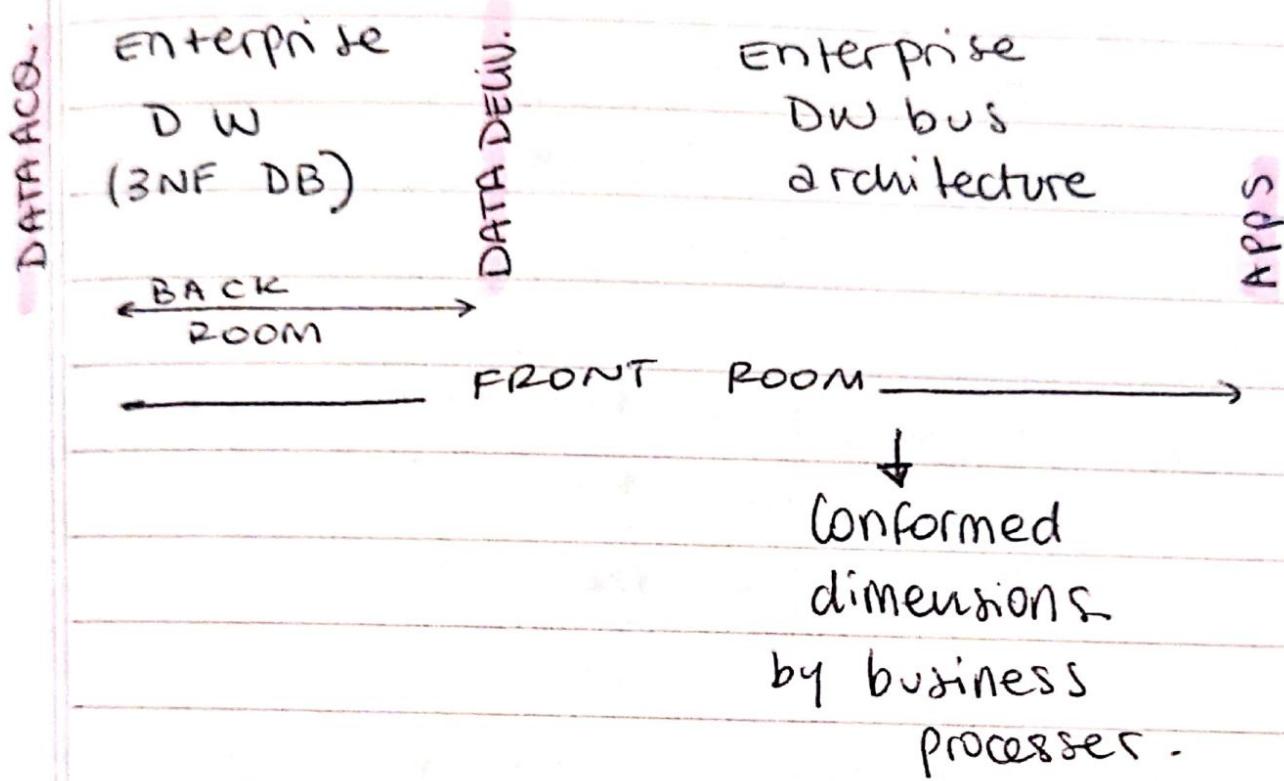


2 ETL processes → 1^o: Source → DW (3NF DB)
2^o: DW → DM's

3NF DB acts as an enterprise wide data store
→ single integrated source of truth for users
→ Accessible by end users

DM's are dimensionally modelled, and contain mostly aggregated data.

HYBRID KIMBALL BUS & INMON CIF



A bit of both words

When you feel like quitting ...
Think about why you started.

OLAP CUBES

Aggregation of fact metric on a number of dimensions.

Eg. Movie, Branch, Month

	MAR	NY	Paris	SF
FEB	NY	Paris	SF	
Avatar	\$25k	\$5k	\$15k	
Starwars	\$15k	\$7k	\$10k	
Batman	\$3k	\$2k	\$3k	

Easy to communicate to end user

Roll-up: Sums the sales of each city by country (less columns in branch dim)

eg.

	FEB	US	FR	ARG	...
Avatar					
Starwars					
Batman					

Groups up a dimension, and sums the fact

Drill-down: Decompose sales of each city into smaller districts (more ~~dimensions~~ columns in branch dimension).

⇒ OLAP cubes should store the finest grain of data (atomic data) in case we need to drill-down to the lowest level.

eg. Country → City → District → Street.

Slice: Reduce N dimensions to N-1 dimensions by restricting one dimension to a single value.

Eg. month = "MAR"

→ I will be seeing only a 2dim table.

Dice: Same dimensions but computing a sub-cube by restricting some values of the dimensions

Eg. month in ['FEB', 'MAR'] and movie in ['Batman', 'Avatar'] and branch ['NY']

→ 'smaller sub cubes'

OLAP CUBES Query Optimization

Business users will typically want to slice, dice, rollup, drill down all the time
→ Ad hoc queries

Each combination / operation will potentially go through all the fact table
→ sub optimal

* GROUP BY CUBE (movie, branch, month)
will make ① pass through the facts table and aggregate all possible combs.

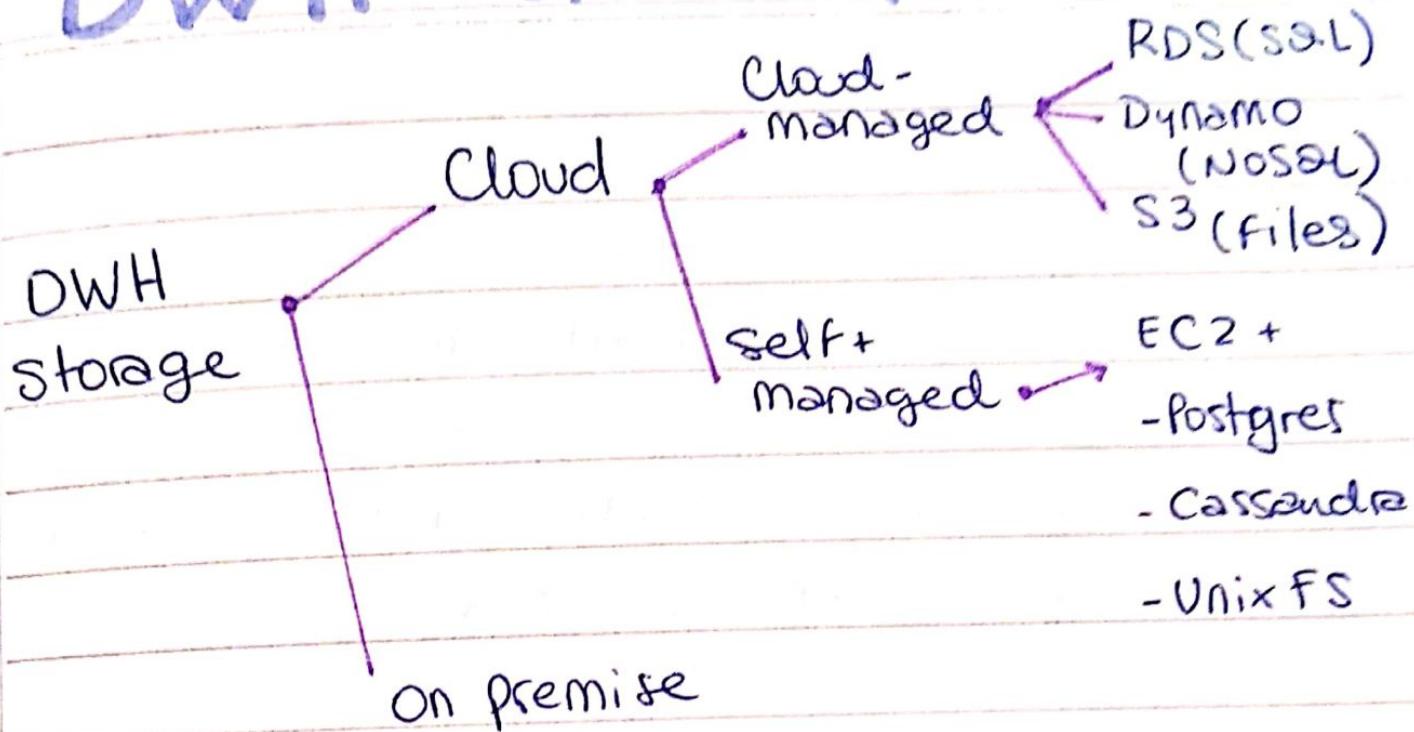
of groupings, of length 0, 1, 2, 3, etc.

• Total revenue	• Revenue by movie	• Revenue by branch, month	• Revenue by movie, branch, month
	- movie	- branch	
	- branch	- month	
	- month		

Saving / materializing the output of the CUBE op and using it → enough to answer all forthcoming agg. from business users w/o processing the whole Facts table

DWH ON AWS

21



AMAZON REDSHIFT

MPP: massively parallel processing DB.

→ parallelizes query execution on multiple CPU machines

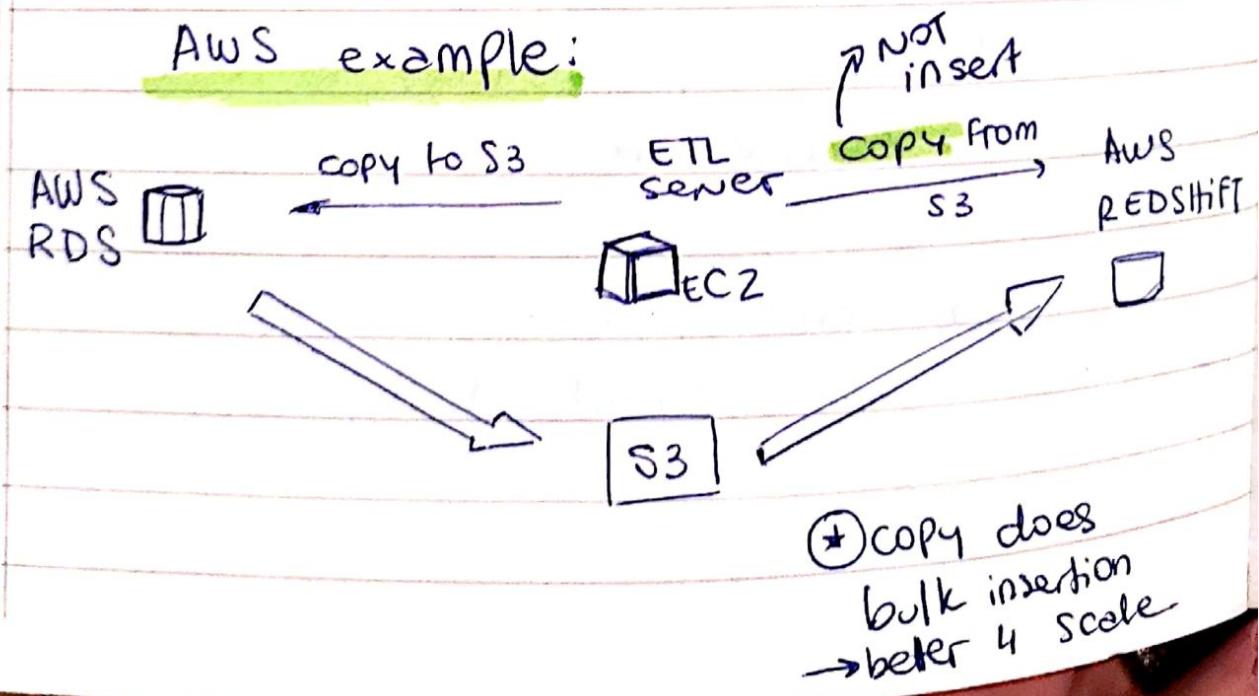
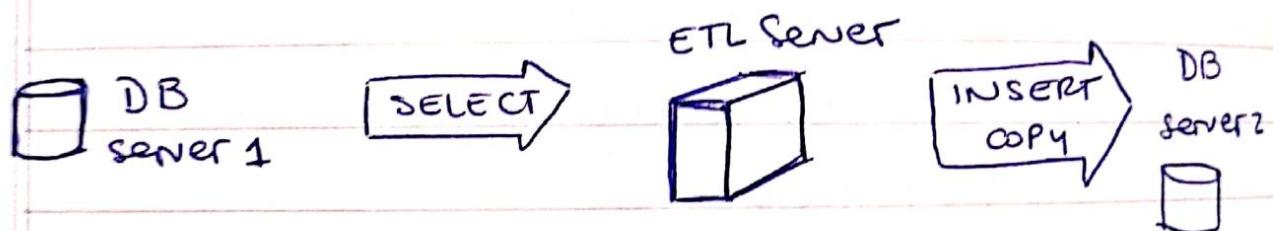
Tables are partitioned, and each partition is processed in parallel

- Column oriented storage.
- It's a modified Postgres.

SQL TO SQL

Select fact1, fact2
INTO newFactTable
FROM table X, Y
WHERE X.id = Y.id, x,y > NULL
GROUP BY Y.d

IF (because vendors) we cannot insert directly into new table (if we are using a different server);



```

COPY sporting-event-ticket FROM
's3://udacity-labs/tickets|split/part'
CREDENTIALS '...';
gzip DELIMITER ';' REGION 'us-west-2';

```

Va a copiar todas las particiones que
empiezan con ese key.
Para ser más concreta:

cust.
manifest

```

{ "entries": [
    { "url": "...", "mandatory": true },
    ...
  ]
}

```

Copy customer

FROM 'S3://my bucket/cust.manifest'

IAM-ROLE '...'

manifest ;

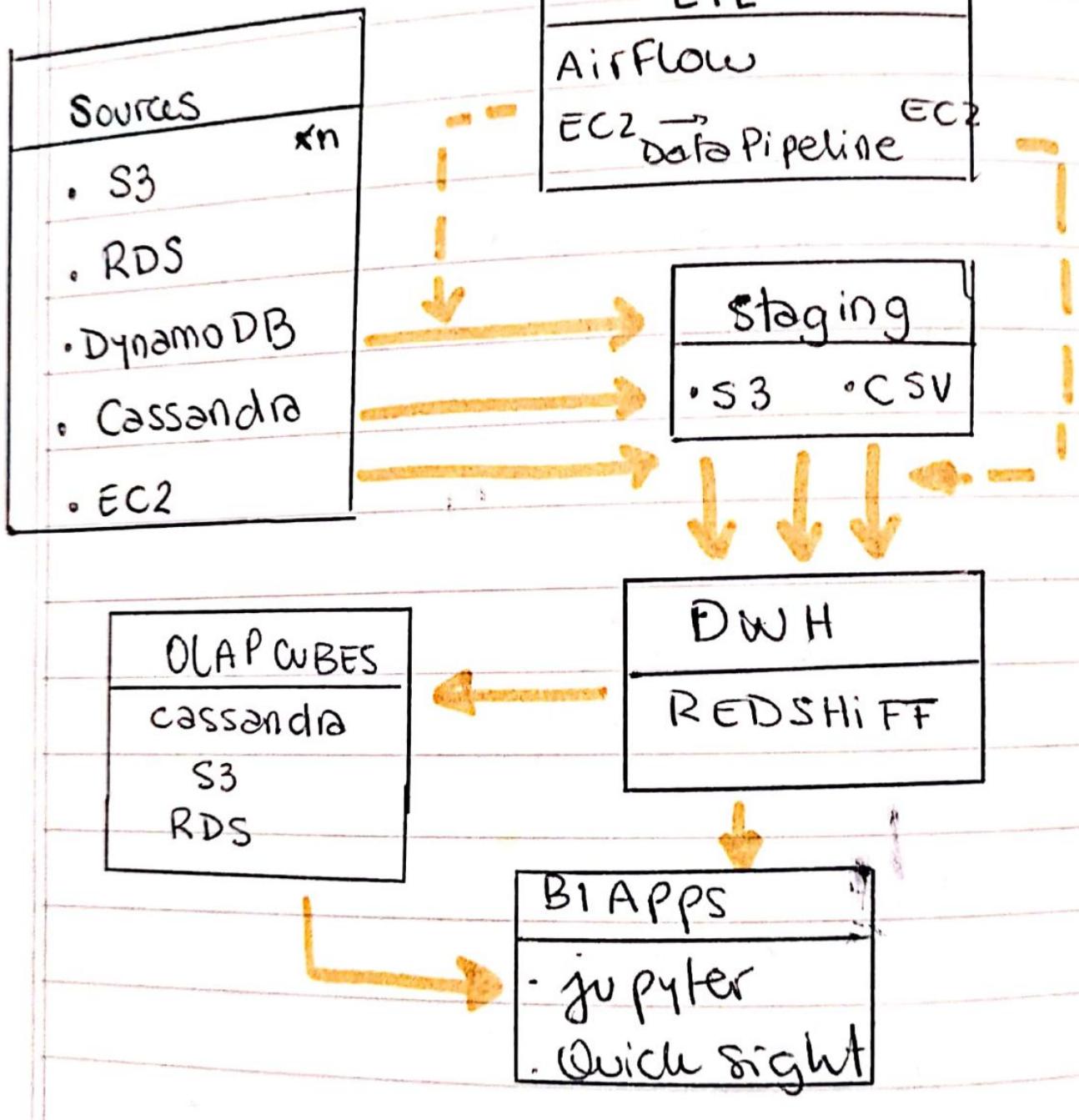
ETL
out

```

} UNLOAD ('select * from venue limit 20')
to 's3://my bucket/venue-pipe'
iam-role '...';

```

Redshift + ETL Context



Infrastructure as a Code (IaC)

- aws-cli scripts
- AWS SDK (ex. Python → boto3)
- Amazon Cloud formation

Benefits.

- sharing
- reproducibility
- multiple deploy meets
- maintainability

- 1° Create user w/ admin privilege
- 2° with boto3 + credentials create clients for IAM, EC2, S3, Redshift.
`ec2 = boto3.resource('ec2', region, creds)`
- 3° Get bucket : `s3.Bucket("name")`
- 4° Create IAM Role to allow Redshift read access to S3
- 5° Create Redshift Cluster w/ created role
 Params : cluster Type, Node Type, N° Nodes, DBName, cluster ID, Master Username + Pass.

6° Open incoming TCP port to cluster

Get Redshift cluster's VPC id, endpoint, etc.

① Get VPC of cluster's vpc id

→ Get security groups of VPC

→ Auth. ingress to SG

To and From DWH - Port

7° Connect to DB

!- SQL 'postgres': / DBuser : pass @Endpoint : port)

8° CREATE Table -

DB

9° LOAD data partitioned to cluster

COPY DB-NAME from 'S3://bucket/parts'

Credentials aws-iam-role =

gzip delimiter ' ;'

region ;

Table Design Optimization

Based on Partition strategies based on frequent access pattern.

→ DISTRIBUTION STYLE

- EVEN: Table is partitioned on slices such that each slice would have almost equal N^2 of records from partitioned table
- ALL or BROADCASTING: Replicates a table on all slices
- AUTO: Allow Redshift to determine distrib. style.
- KEY: rows having similar values are placed in the same slice.

→ SORTING KEY

Define a column as sort key.