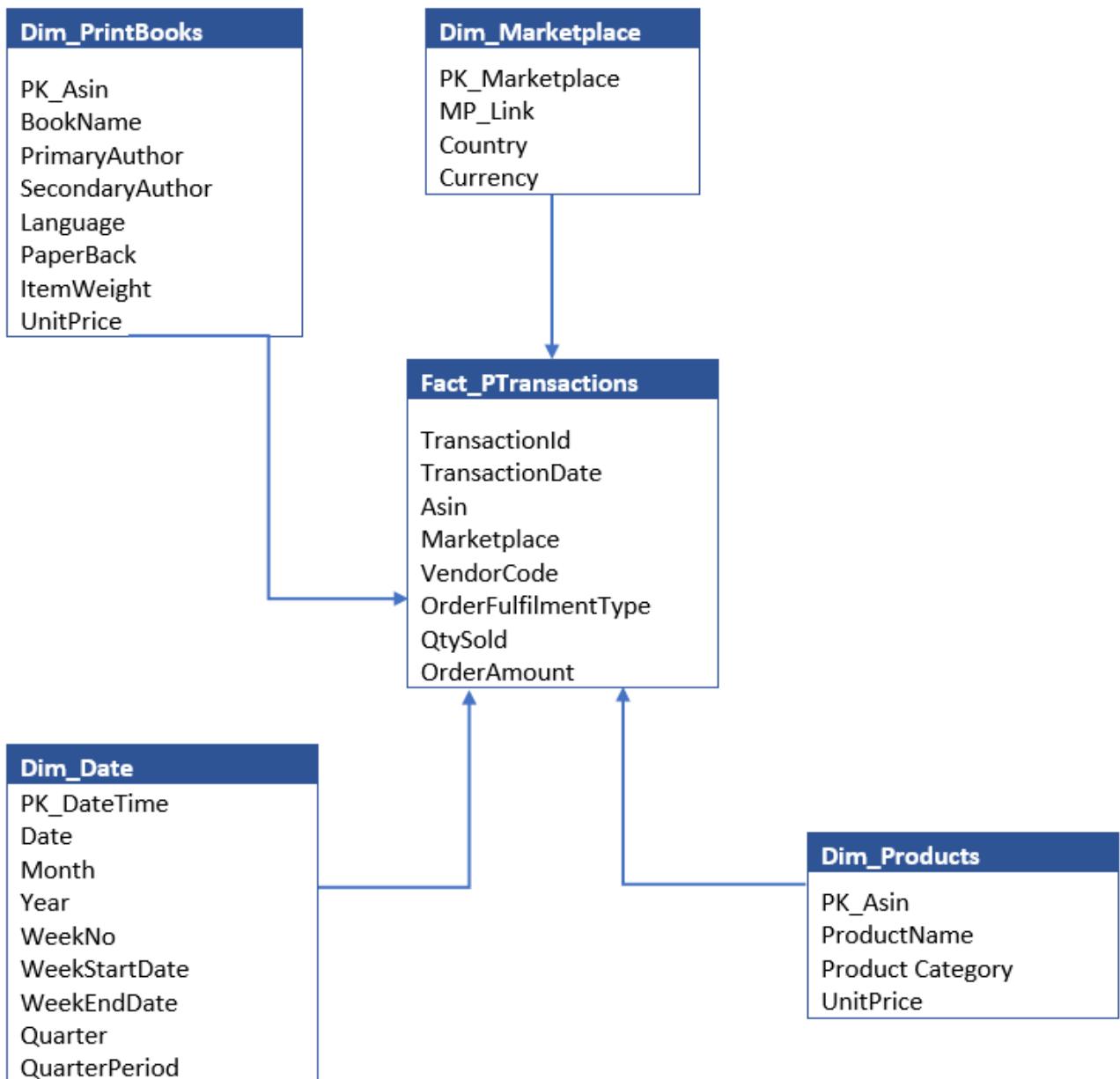


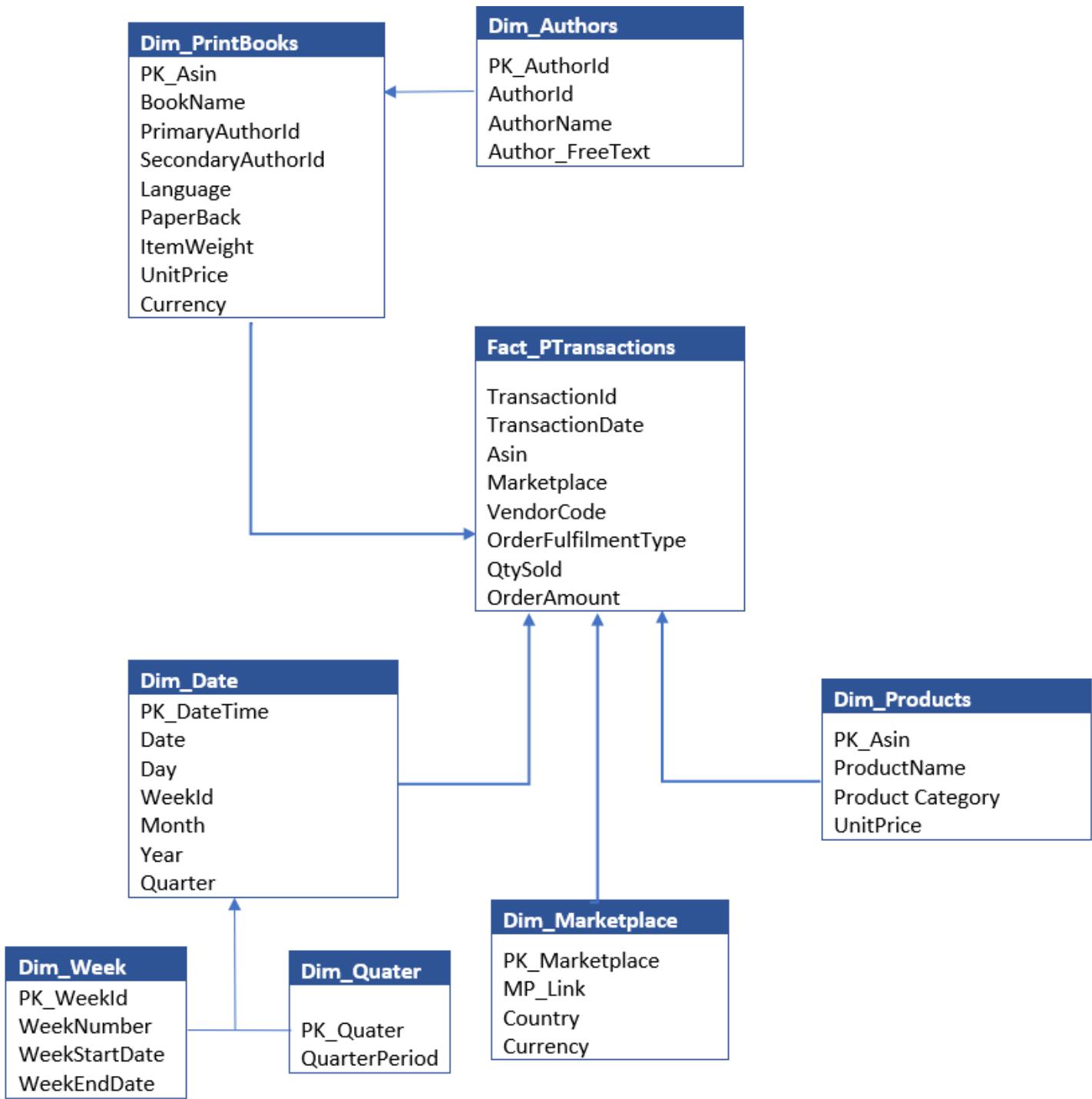
# Agenda

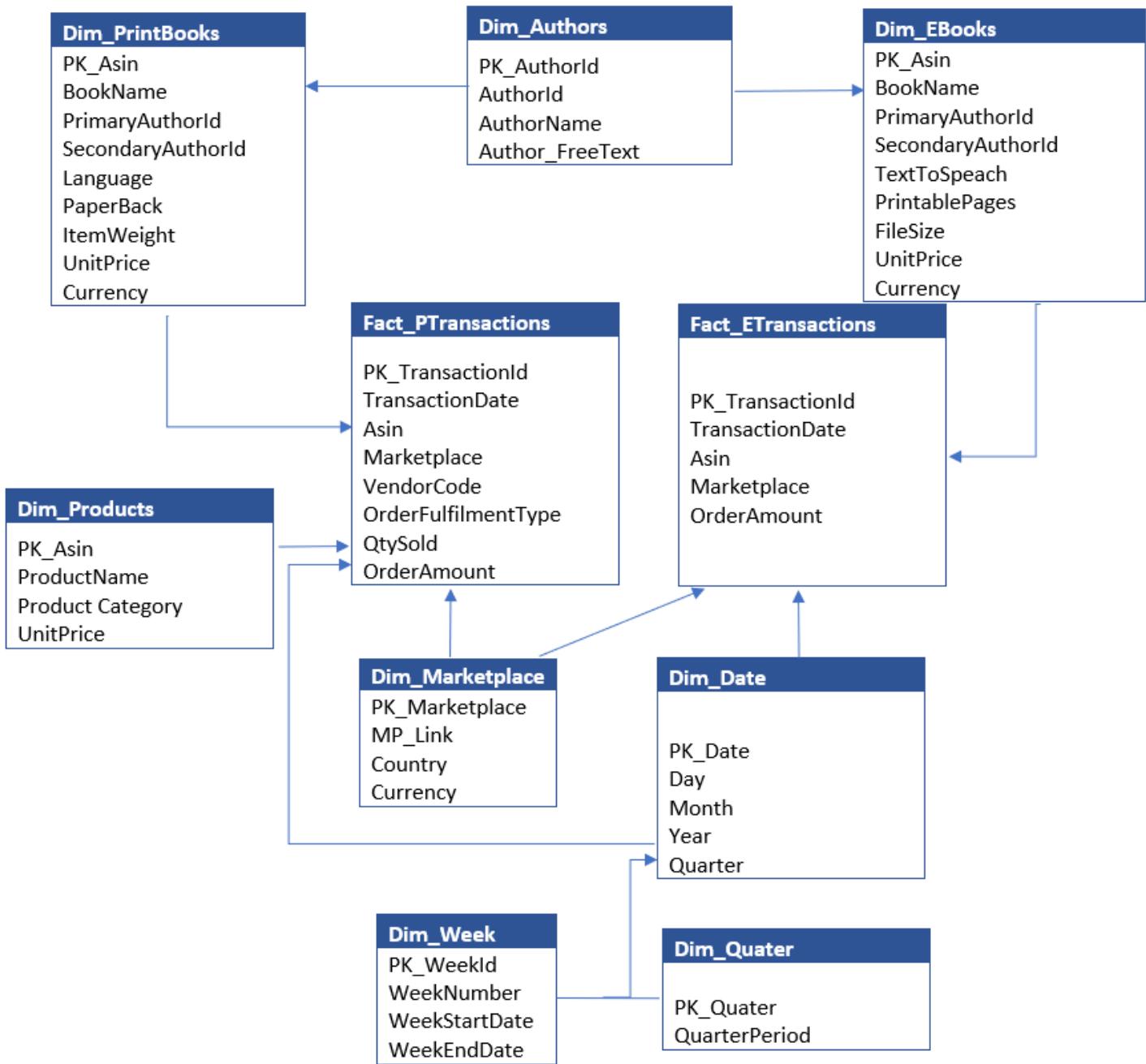
Dimensional Modeling Approach

- Problem statement
- Existing Solution
- What are data warehouses and what is OLAP?
- How does Apache hive act as a data warehouse?
- What makes hive as Hive? (Architecture)
- How to interact with Hive?
- What are different HQL commands that are used commonly?
- What are the optimization techniques in Hive ?
- Use cases of hive/ where to use hive/why and why not

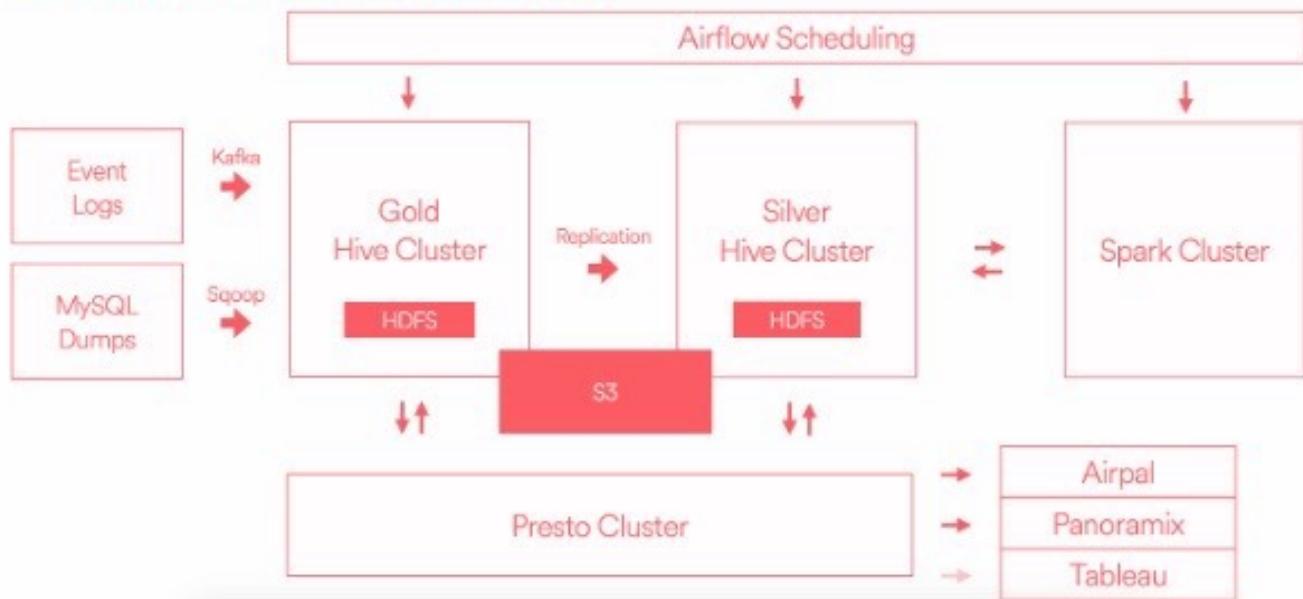




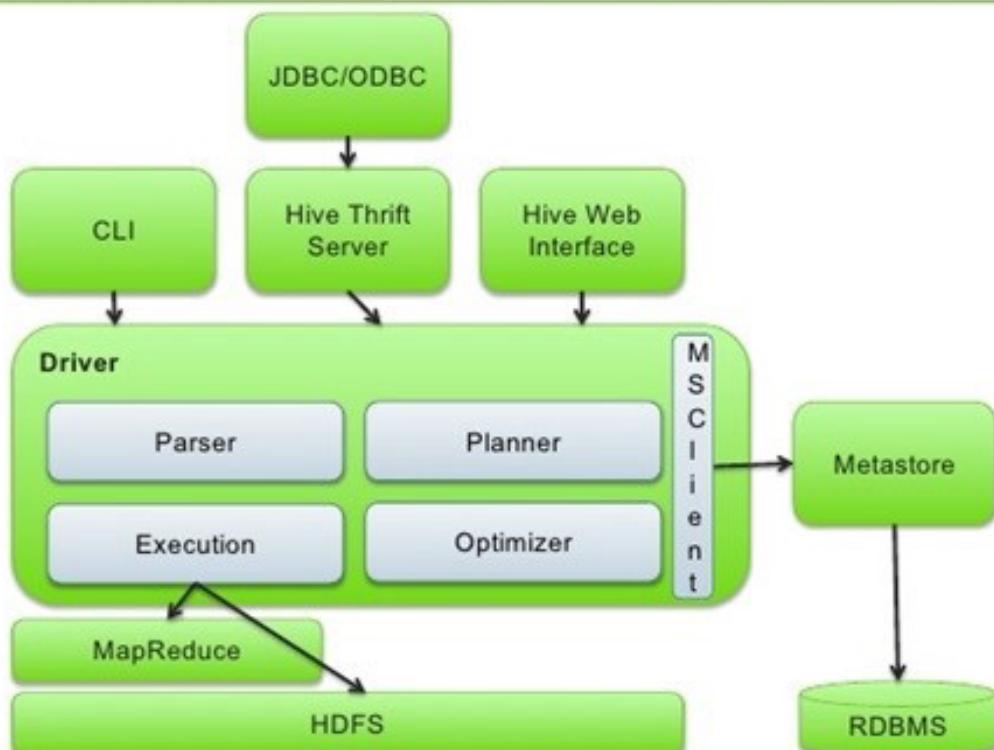




# AIRBNB DATA INFRA



## Apache Hive Architecture



## **Receive SQL query**

- 1** Parse HiveQL
- 2** Make optimizations
- 3** Plan execution
- 4** Submit job(s) to cluster
- 5** Monitor progress
- 6** Process data in MapReduce or Apache Spark
- 7** Store the data in HDFS

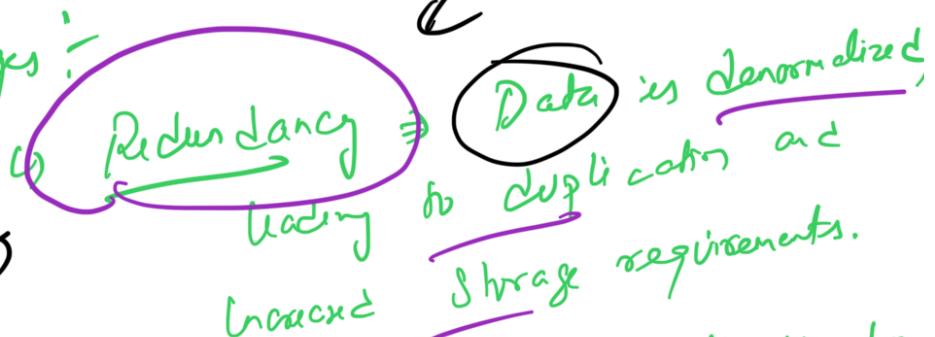
# Star Schema

## Advantages

- ↳ Star schema is easy to understand and provides optimal disk usage.
- ↳ Schema is widely supported by BI tools.

## Disadvantages:

1	HP	Vol	120
1	HP	Vol	120



↳ limited flexibility  $\Rightarrow$  less flexible when dealing with complex queries.

↳ Performance issues  $\Rightarrow$  large dimension tables can slow down performance.

↳ Maintenance challenges: updates/changes in this schema can be more complex & error-prone.

↳ Migration challenges:

② Snowflake Schema:- extended from of star Schema, by normalizing dimension table.

- ↳ advantages:-
- (a) Low Data redundancy
  - (b) Consumes less storage

Disadvantages:-

- (a) Complexity
- (b) Query Performance is lower
- (c) Designing is complex than star schema.

③ Galaxy | Fact Constellation Schema

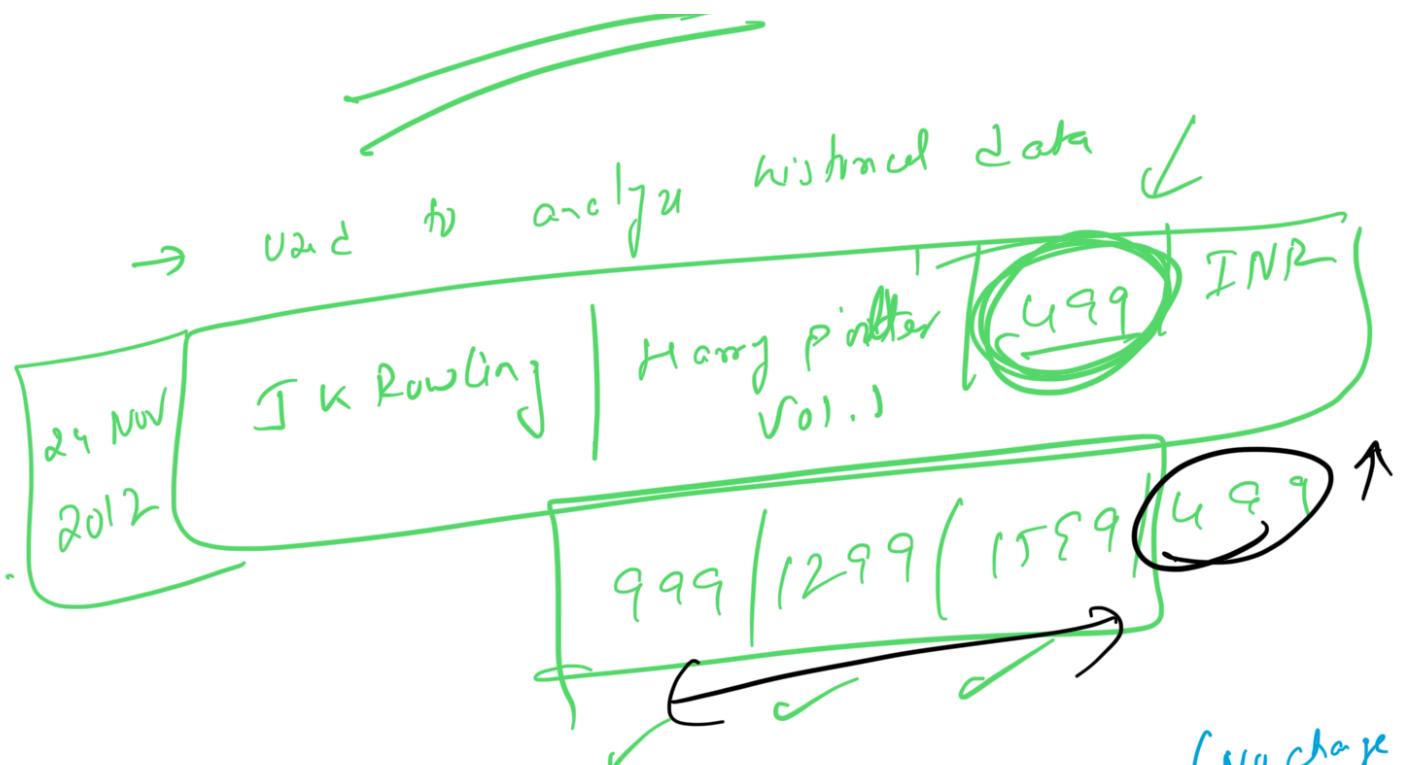
↳ Compound Schema which contains more than 1 fact and normalized Dimension.

- Advantages
- (1) Reusability
  - (2) Supports Complex Business Processing
  - (3) Scalability
- ..... Complexity

- Disadvantages :-
- ① Increase - Challenges
  - ② Maintenance Challenges
  - ③ Steeper Learning Curve

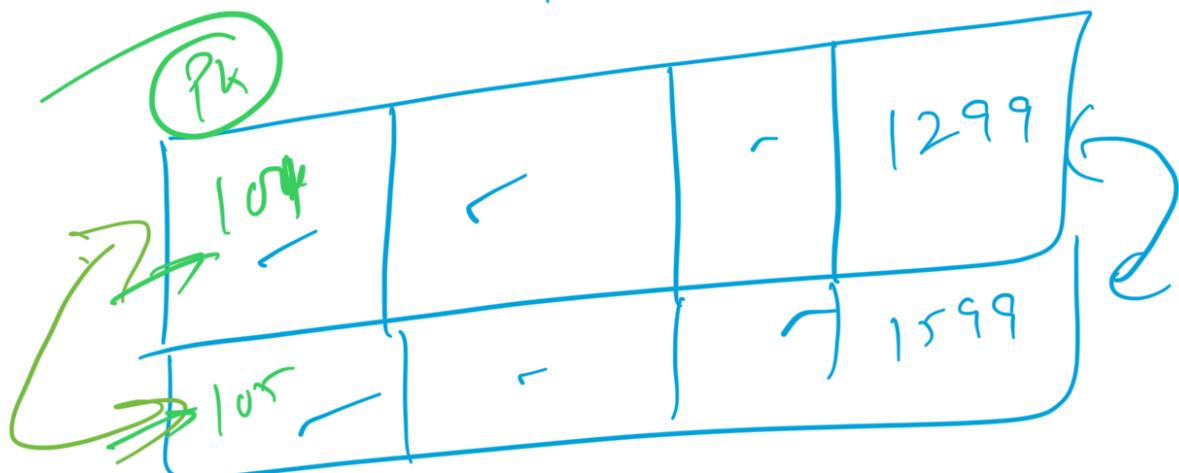
	Snowflake	Star
① Ease of Maintenance	<ul style="list-style-type: none"> <li>→ No redundancy and hence more easy to maintain &amp; storage.</li> </ul>	<ul style="list-style-type: none"> <li>→ has redundant data and hence less easy to maintain / change.</li> </ul>
② Ease of J2E	<ul style="list-style-type: none"> <li>→ More complex queries and hence less easy to understand</li> </ul>	<ul style="list-style-type: none"> <li>→ less complex queries and easy to understand</li> </ul>
③ Query Performance	<ul style="list-style-type: none"> <li>→ More no of table, More no of FKs, More execution time.</li> </ul>	<ul style="list-style-type: none"> <li>→ less no of table</li> <li>→ less no of FKs</li> <li>→ less execution time</li> </ul>
④ Joins	Higher no of joins	<ul style="list-style-type: none"> <li>→ fewer joins</li> </ul>
⑤ When to use J2E	<ul style="list-style-type: none"> <li>when the dimension table is relatively big in size, it reduces space</li> </ul>	<ul style="list-style-type: none"> <li>→ when the dimension table contains less no. of rows, use star schema</li> </ul>

SCD → slowly changing Dimension



① SCD type 0 : Passive Method (No change is logical)

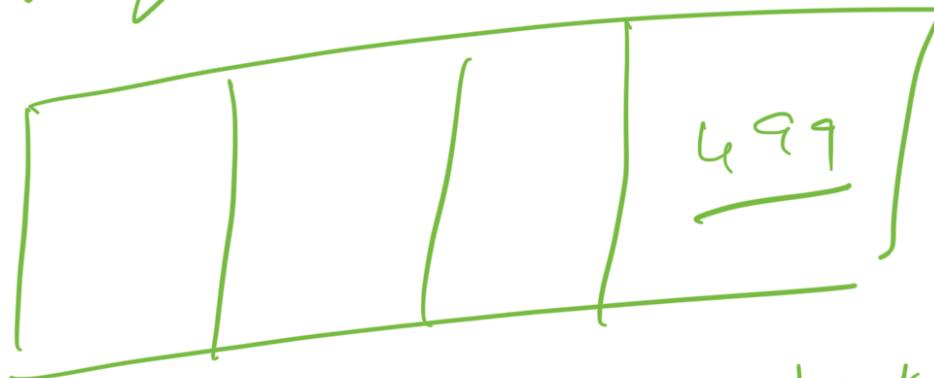
↳ In this method, no special action is performed upon dimensional changes



Pros : no extra consideration for designing DDL.

Cons : provide no control over recording 1. Analyzing historical data &

→ total changes.  
 SCD Type 1: overwriting few old values



Pros:

- straightforward to implement
- we ignore old values
- uses less disk space since only 1 record exists.

Cons:

- stores no historical data for analysis
- Hard to find previous state

SCD type 2:

Creating a new addition record.

PK-SK	Pk-Asian	Bookname	Author	UnitPrice	Currency
18n1	101	H Prolly	1001	12.00	INR
18n2	101	N Prolly	1011	99.99	INR
18n3	101			2.99	INR
18n4	101			4.99	INR

Intermediate Keys → It is an artificial key.

Overwriting which is a unique identifier  
a row used by different  
addresses of the same entity.

- it is a system generated
- it is not manipulable by a user.
- it has no business context and is independent of domain.

Pros :- unlimited history  
→ original table schema does not change.

Cons :-  
↳ takes up more disk space.  
↳ introduce change is very expensive.

SCD Type 3 :- Adding a new column

↳ An additional column available to capture the previous value of the attributes and the existing value of the present value of the record by overwriting the old value.

Final answer - Presenting

Log

Age	BN	Author	Year	Word	
101	H P Vro	J K Rau	499	INR	999

Pros :- Suitable for tracking only the recent changes.

↳ Disk Space.

Cons :- It is by design limited both no of columns designated for strong historical values.

3rd track :- Using historical table

Pros :- Provides clear distinct b/w historical & current data.  
↳ It resembles how an audit tool stores the change data.  
↳ Query on non-historical is fast.

∴ A Disk Space (Light)  
.. h. normal

Cons :   
 ② Additional effort w/ ...  
 both tables.

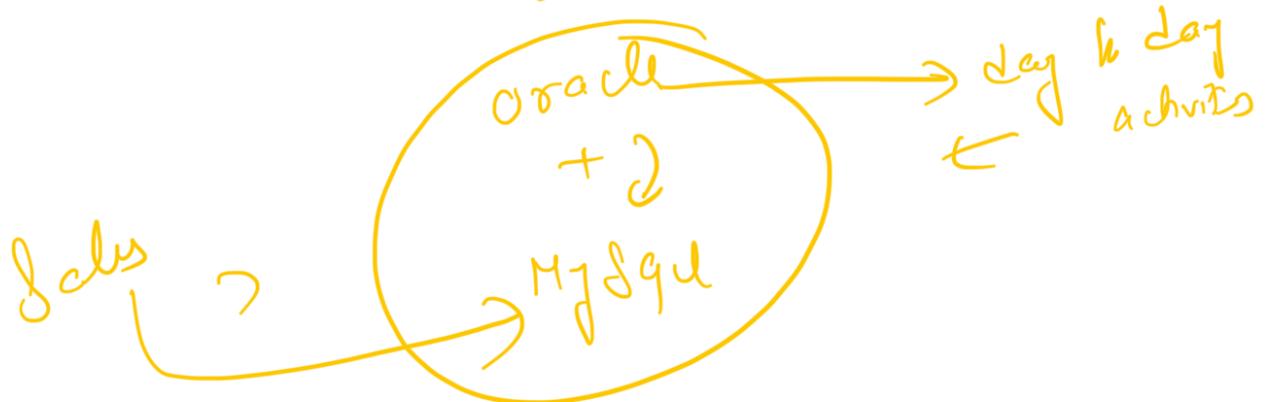
Scd Type 6 : (1, 2, 3)

Pros : Shows complete picture of history in a 1 table.  
 → Allows us to store multiple state changes while using firs

Cons : More complex to implement.

Airbnb

↳ 2014



→ Dashboard/Reporting  
...  
m...

## ↳ Adhoc Queries

- Reports jobs becomes ~~factory~~
- Business Critical jobs
- Adhoc Queries → RTO
- ML models ↗

