

Agenda

Map reduce optimisations

What is OLTP?

Why it is not suitable for Analytics

Introduction to Data warehousing

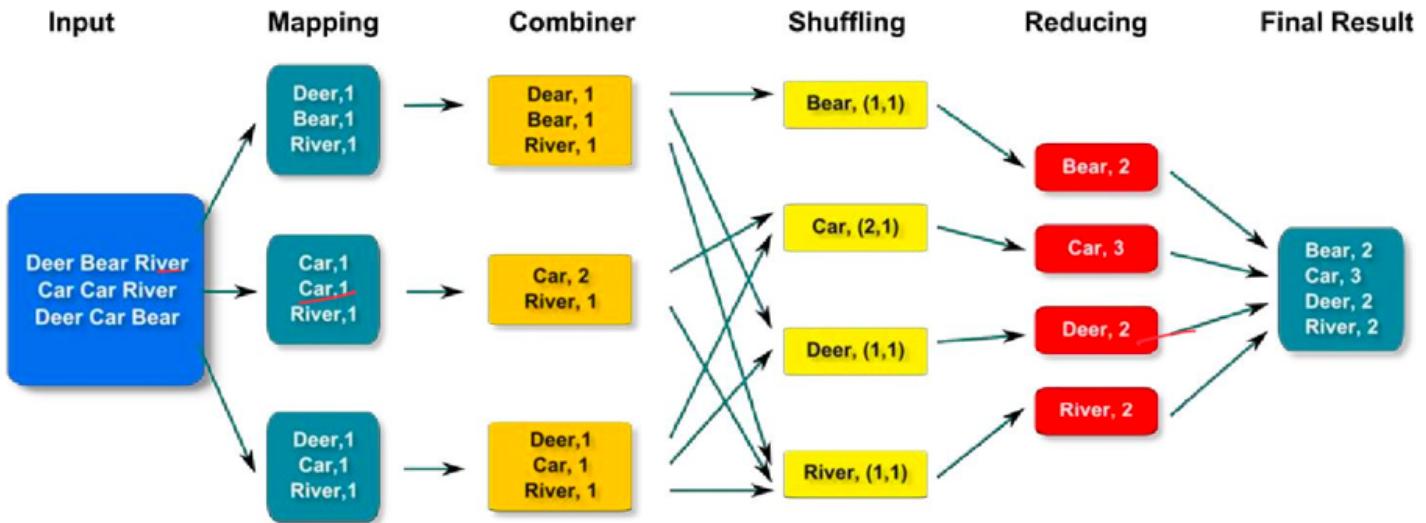
Why is data warehouse a better option?

How to design a Data Warehouse

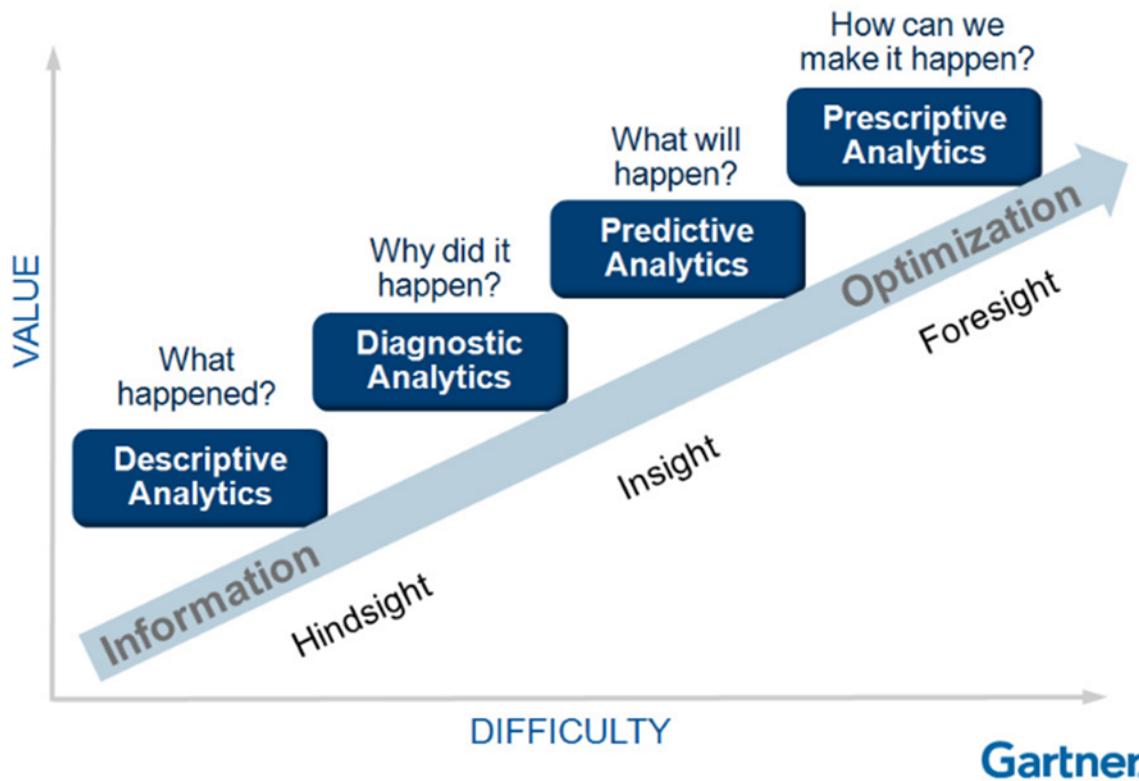
Dimensional Modeling Approach



Combiner - Local Reduce



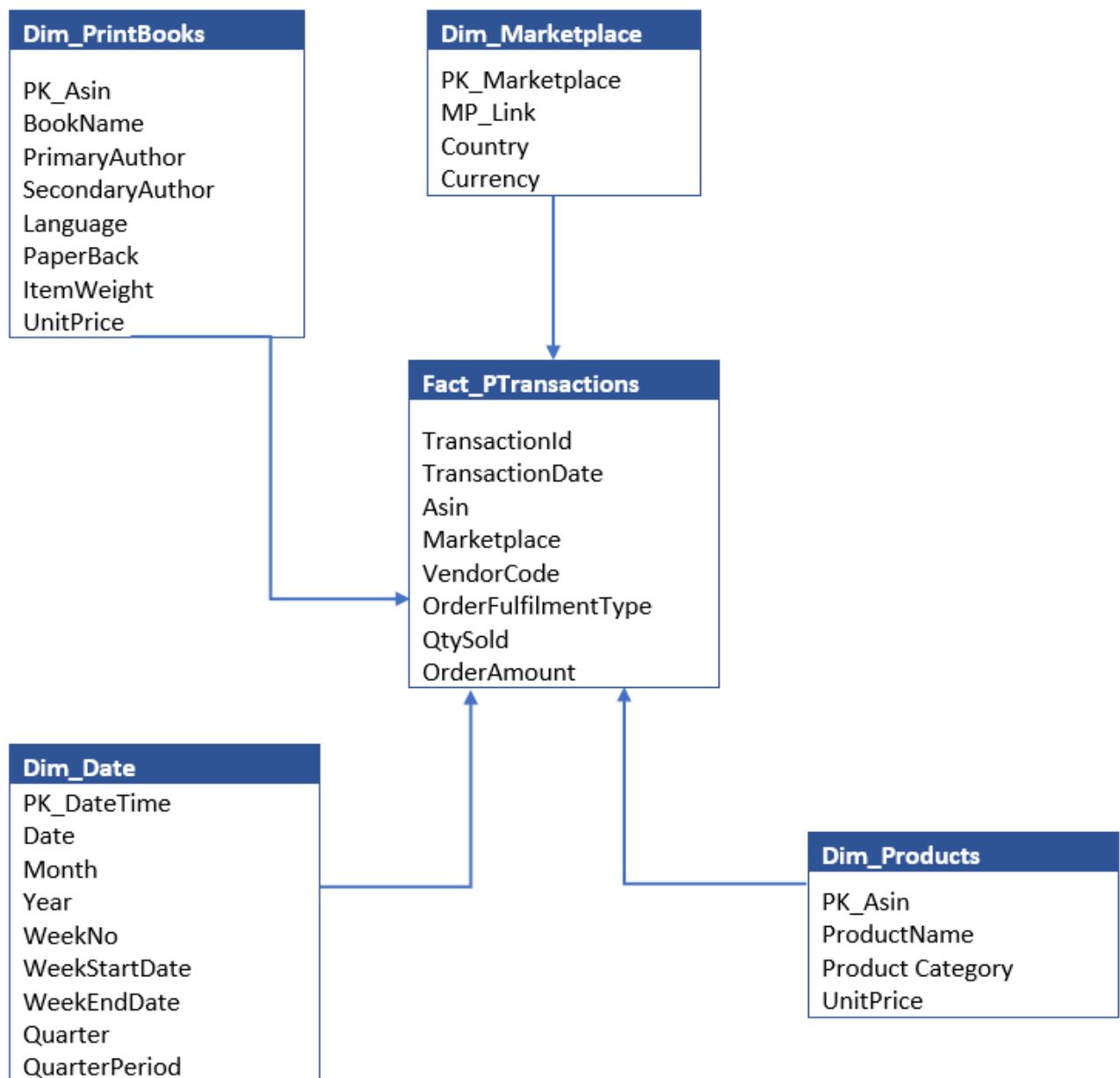
Analytic Value Escalator

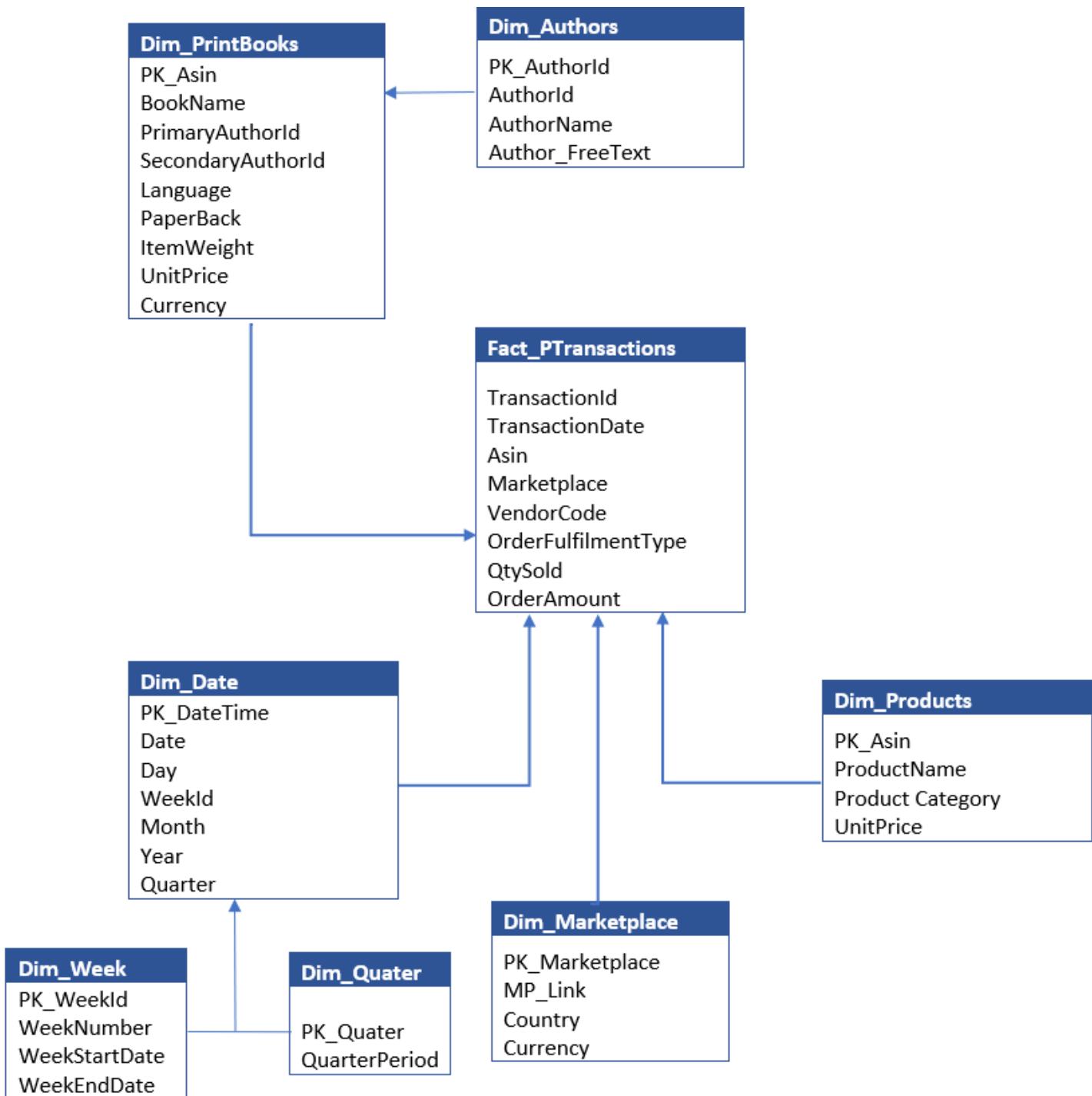


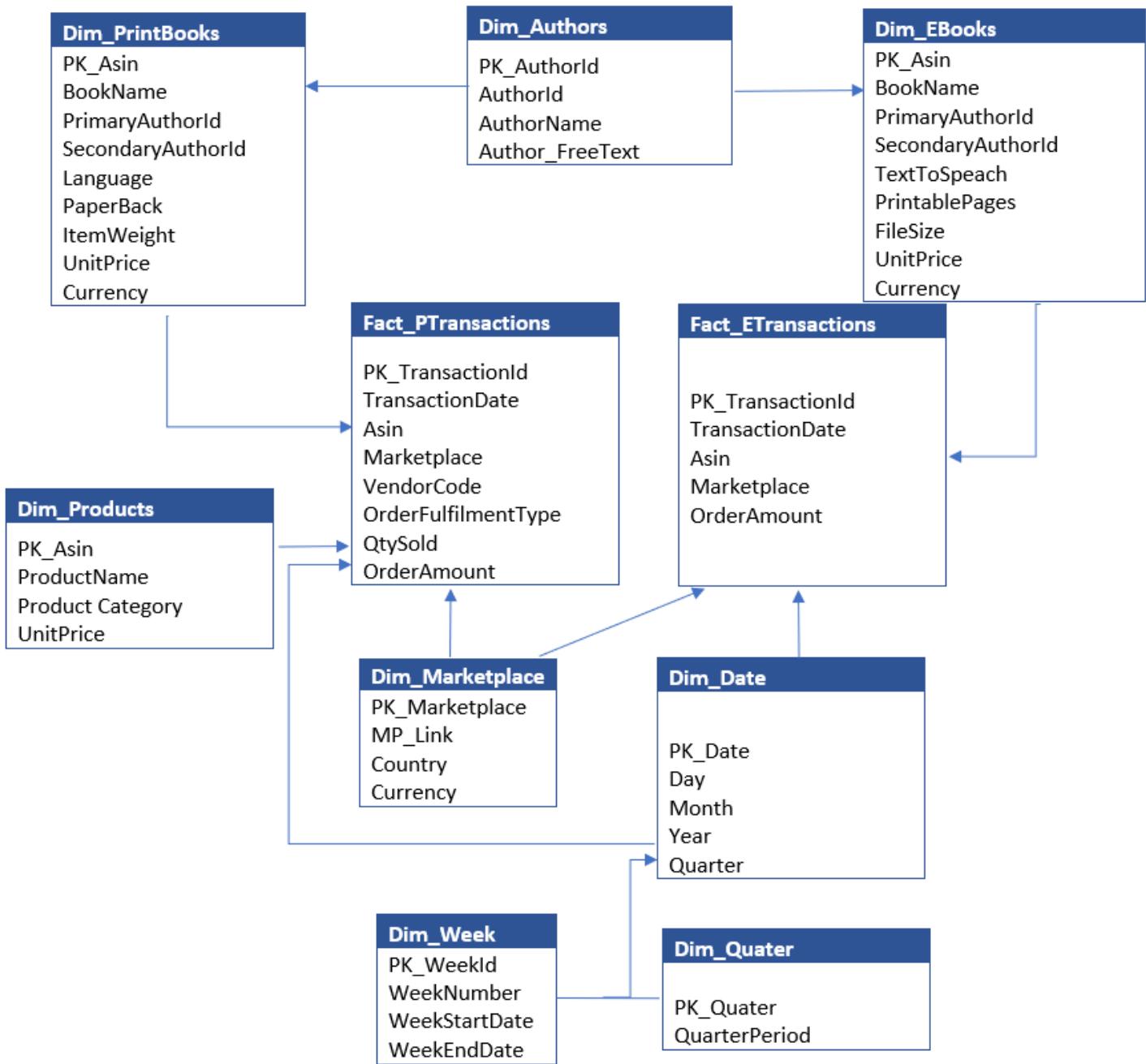


OrderId	CustomerId	InvoiceNo	Invoice Date	ShippingAddress	PaymentMethod	SoldBy	BillingAddress	PlaceOfSupply	PlaceOfDelivery	ASIN	Item.No	Description	UnitPrice	Qty	Net Amount	Tax Amount	Total Amount	Shipping
7146	101	0990	23 August 2022	Bangalore	COD	CT	Bangalore	Chennai	Bangalore	B0963RYQSW	1	Adapter	100	1	100	18	118	0

OrderId	CustomerId	InvoiceNo	Invoice Date	ShippingAddress	PaymentMethod	SoldBy	BillingAddress	PlaceOfSupply	PlaceOfDelivery	ASIN	Item.No	Description	UnitPrice	Qty	Net Amount	Tax Amount	Total Amount	Shipping
7146	101	0990	23 August 2022	sikgfjsvifsvn	COD	CT	sikgfjsvifsvn	Chennai	sikgfjsvifsvn	B0963RYQSW	1	Adapter	100	1	100	18	118	0



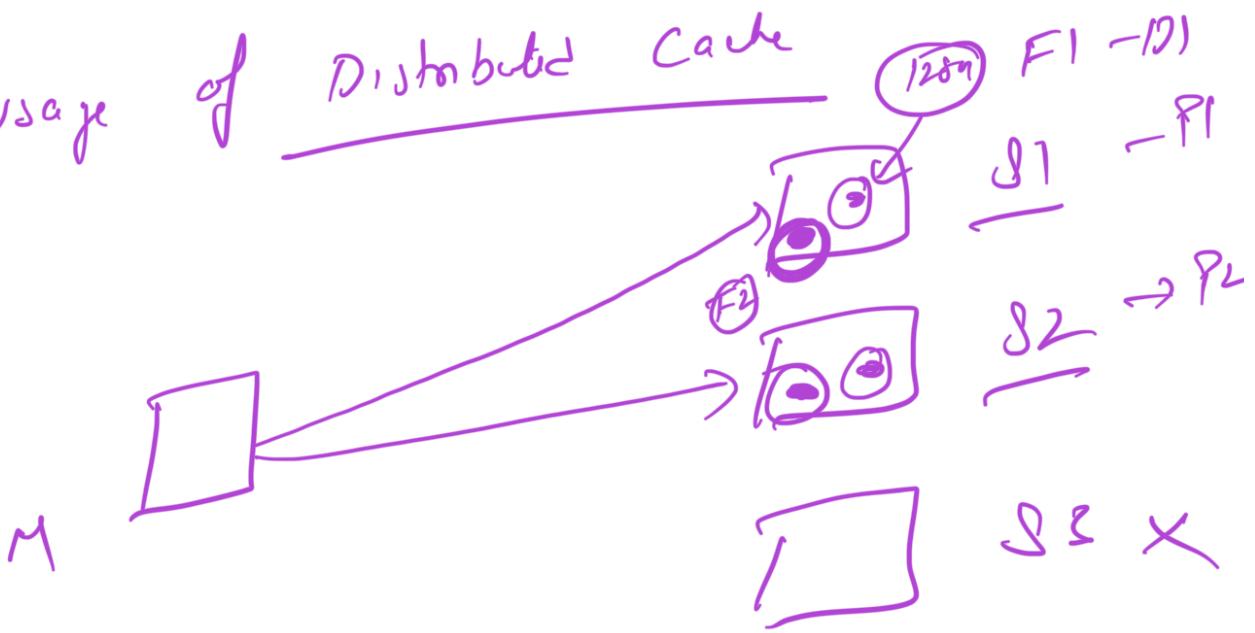




Map Reduce optimizations

① Usage of

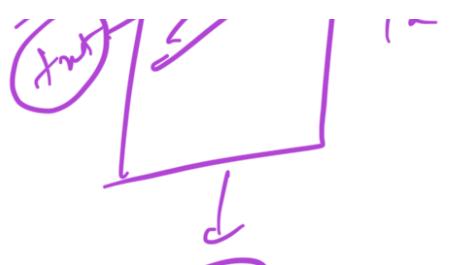
Distributed Cache



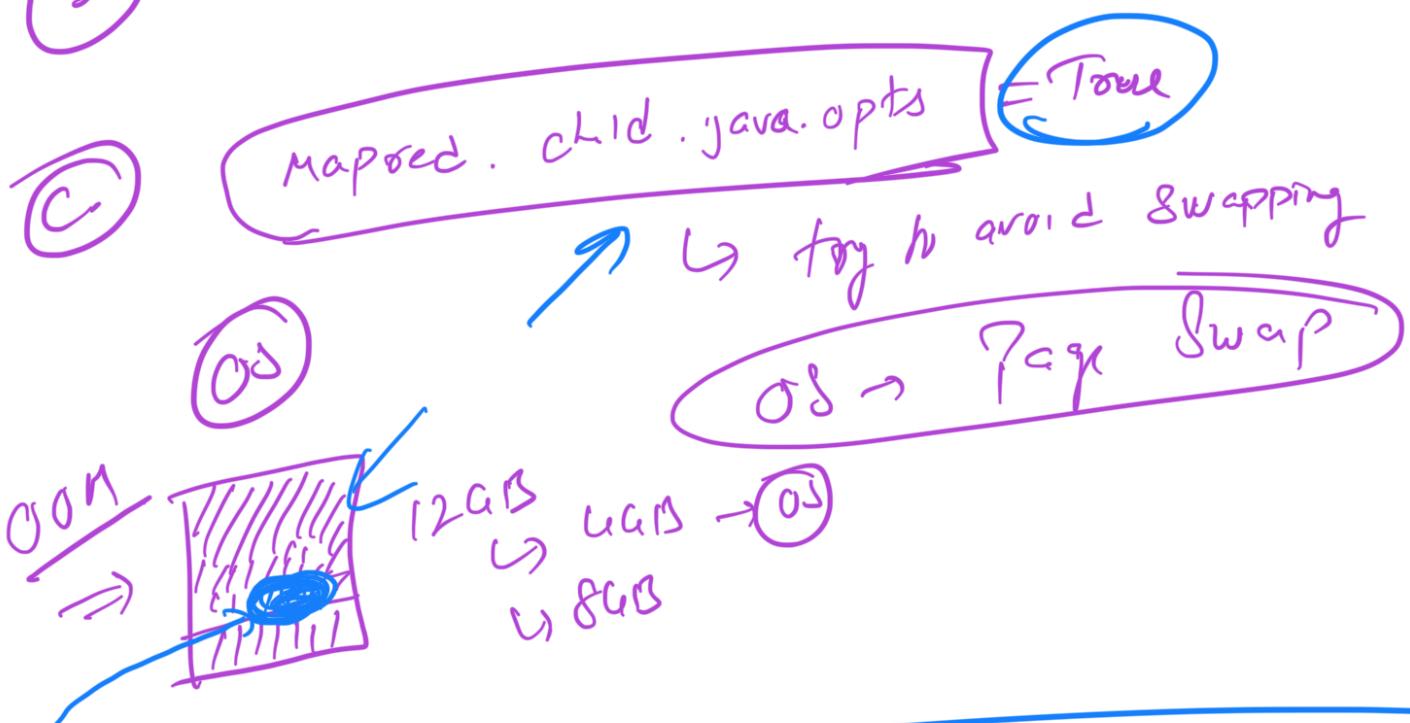
→ MR provides facility to Cache Small files which are used only by files like `fout`, `zip`, `jar`.

→ Each DataNode gets a copy of the file which is sent through DC when job gets finished, this can be deleted from DN.





b. Take the help from Combiner



Data Modelling

\hookrightarrow SQL

\hookrightarrow Datawarehouses

Amazon Books

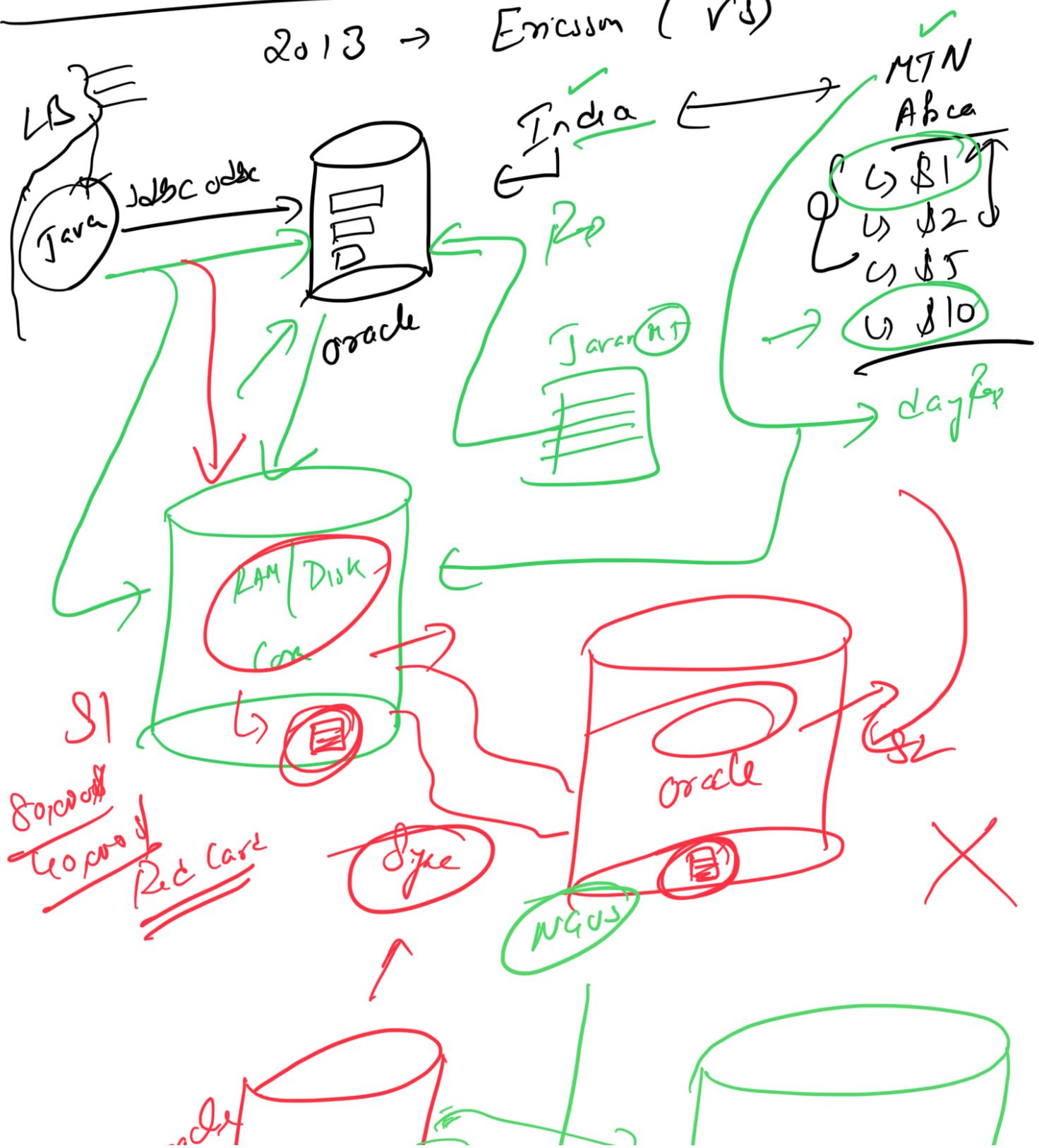
\hookrightarrow physical

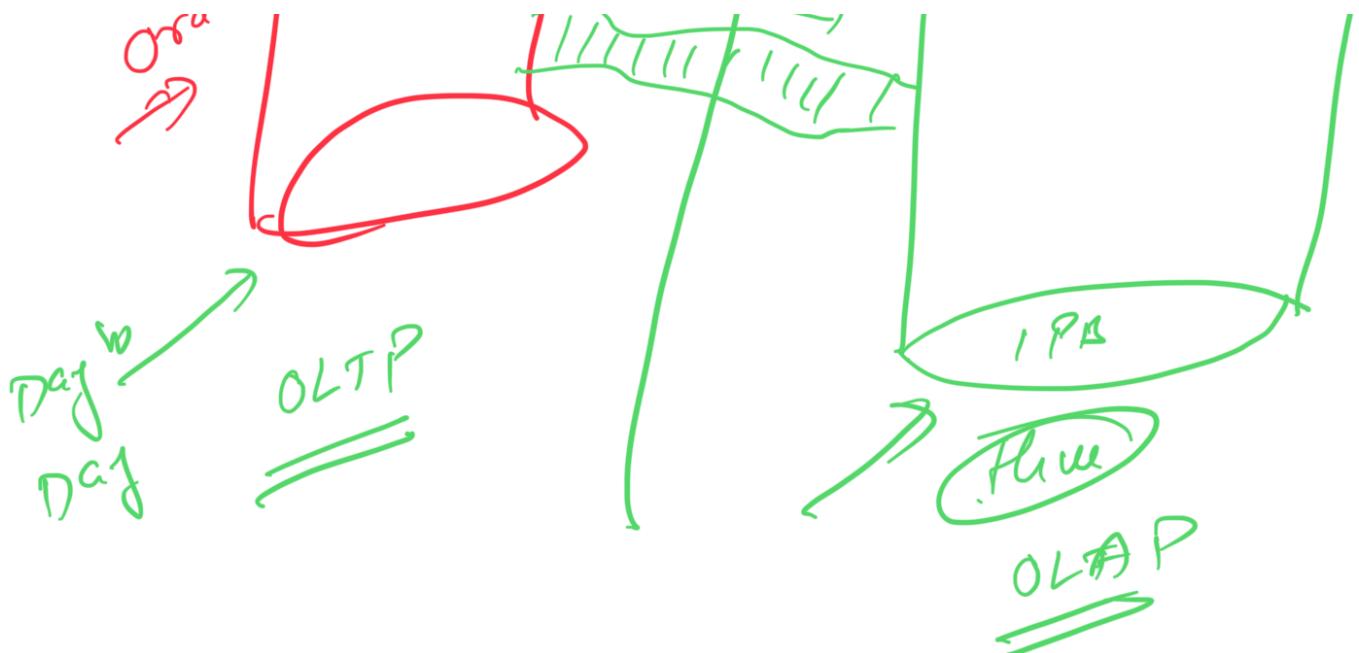
\hookrightarrow e-books

... weekly 9

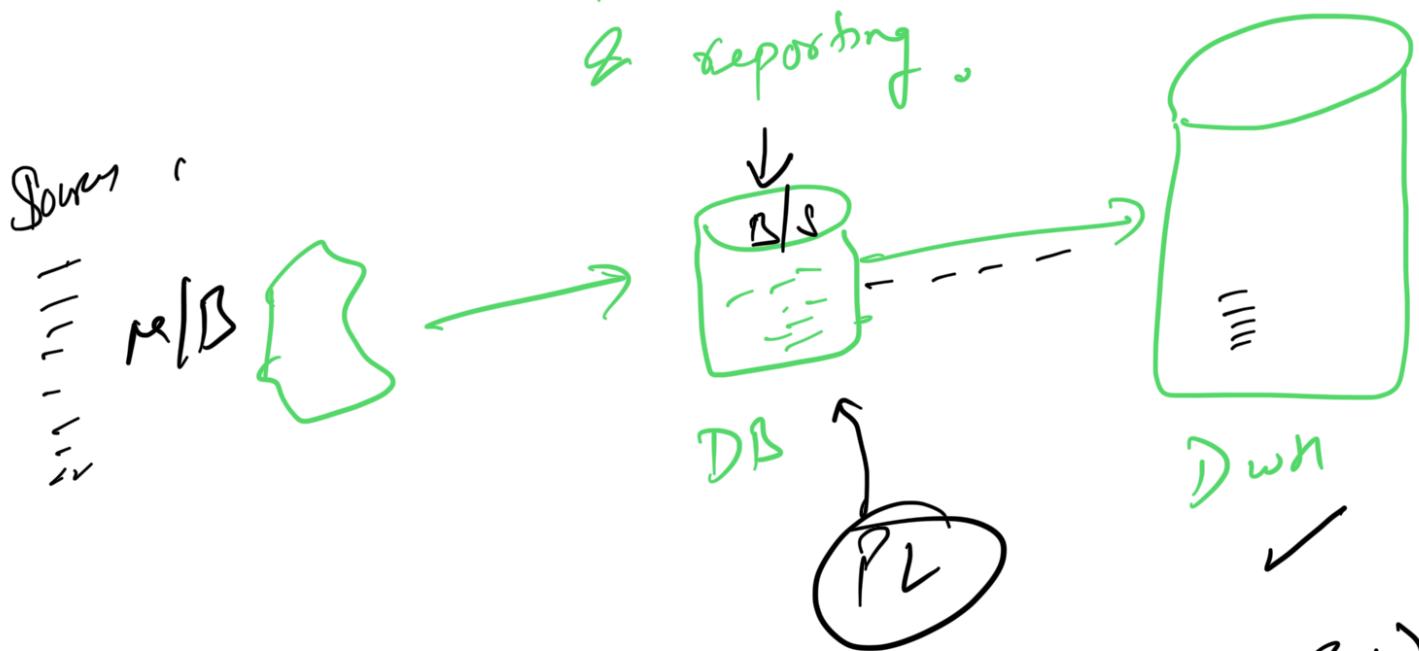
- What top 10 books sold were frequently bought with books?
- What other products were frequently bought with books?
- Who are the top authors based on book sales?

2013 → Encision (v3)





DwH : is a storage unit environment to support BA activities like analysis & reporting.



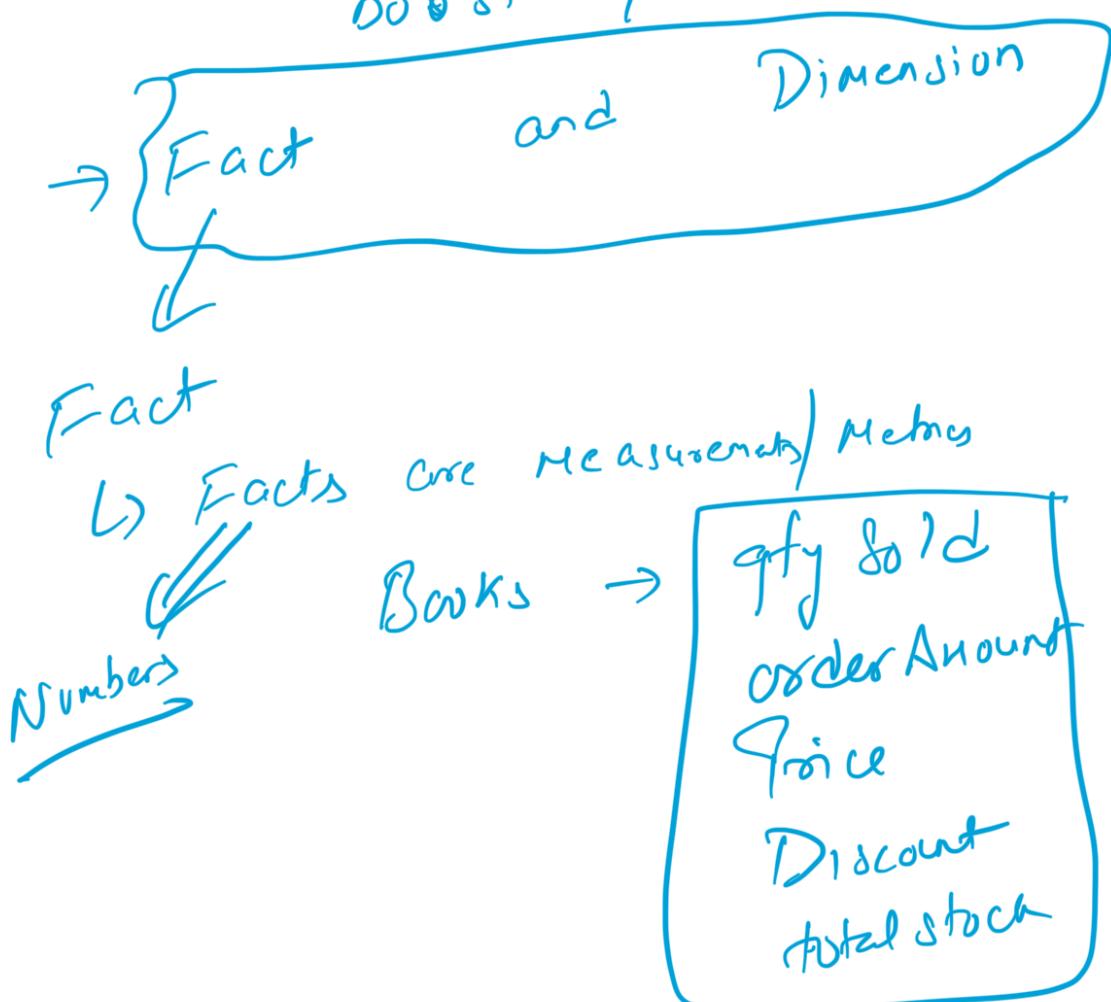
Design a DwH (Amazon Books)

- (1) Books / Books related to Author
- (2) Sacks
- (3) Marketplace
- (4) Corporate

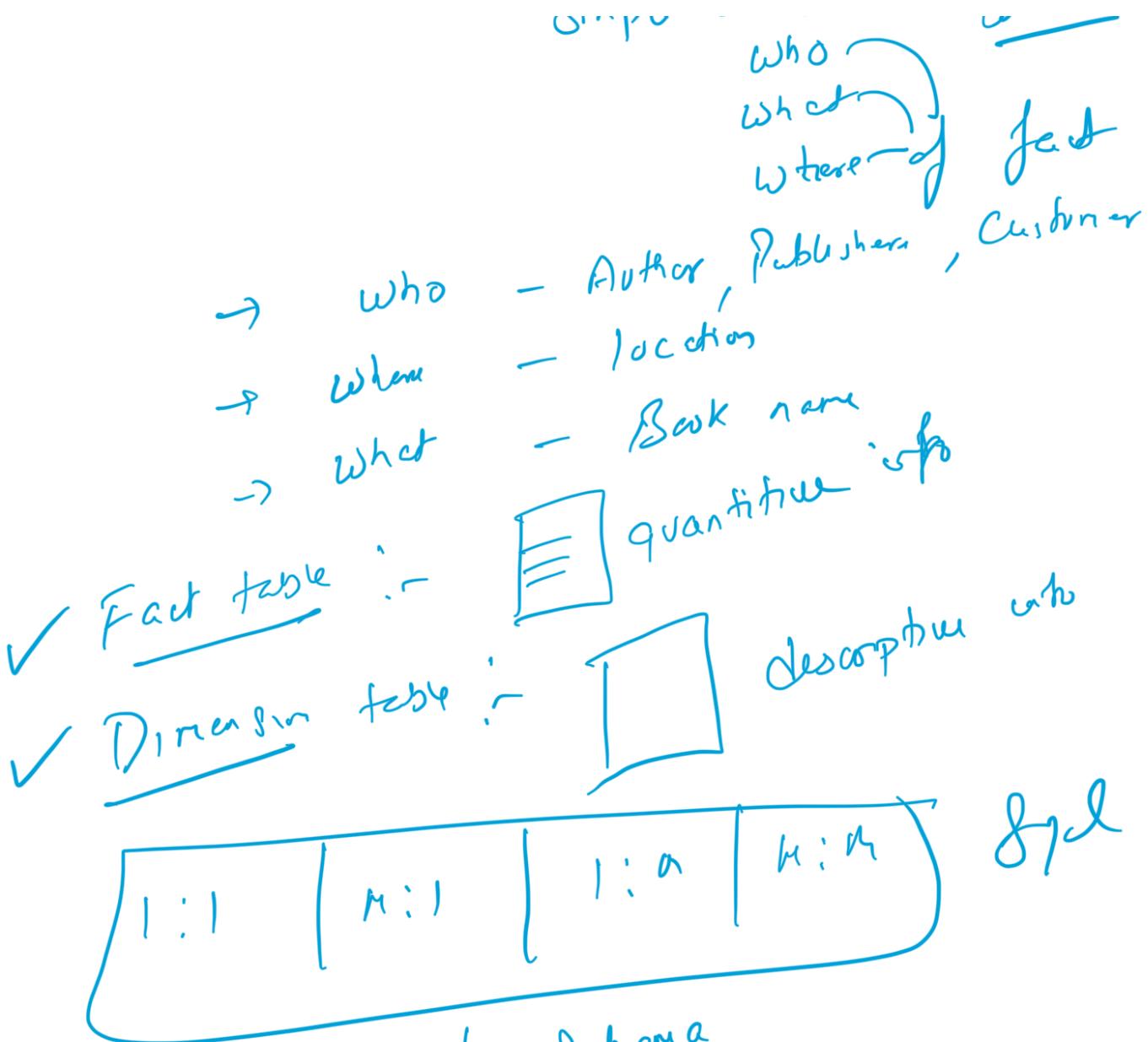
↳ Vendor

↳ Dimensional Modeling

↳ as a **technique** which include a set of measures, processes and concepts to design and store data in a DW for storage optimization & boost performance.



Dimension → provides the context surrounding a business process event, in very simple terms:



↳ Schema

- (1) Star Schema
- (2) Snowflake Schema
- (3) Galaxy Schema

→ Star Schema :- type of schema that utilizes facts & dimensions to create a simpler yet efficient dimension

Model.

- Every dimension in a star Schema is represented with only 1 dimension table.
- Dim. table should contain the set of attributes.
- Dim. table is joined to fact table using a Fk.
- Dim. table are not joined to each other.
- Fact table contains key & measure.

