

Agenda

- a. Star Schema
- b. Snowflake Schema
- c. Galaxy Fact Schema
- d. Slowly Changing Dimensions
- e. Surrogate Keys
- f. Problem statement
- g. Existing Solution
- h. What are data warehouses and what is OLAP?
- i. How does Apache hive act as a data warehouse?
- j. What makes hive as Hive? (Architecture)

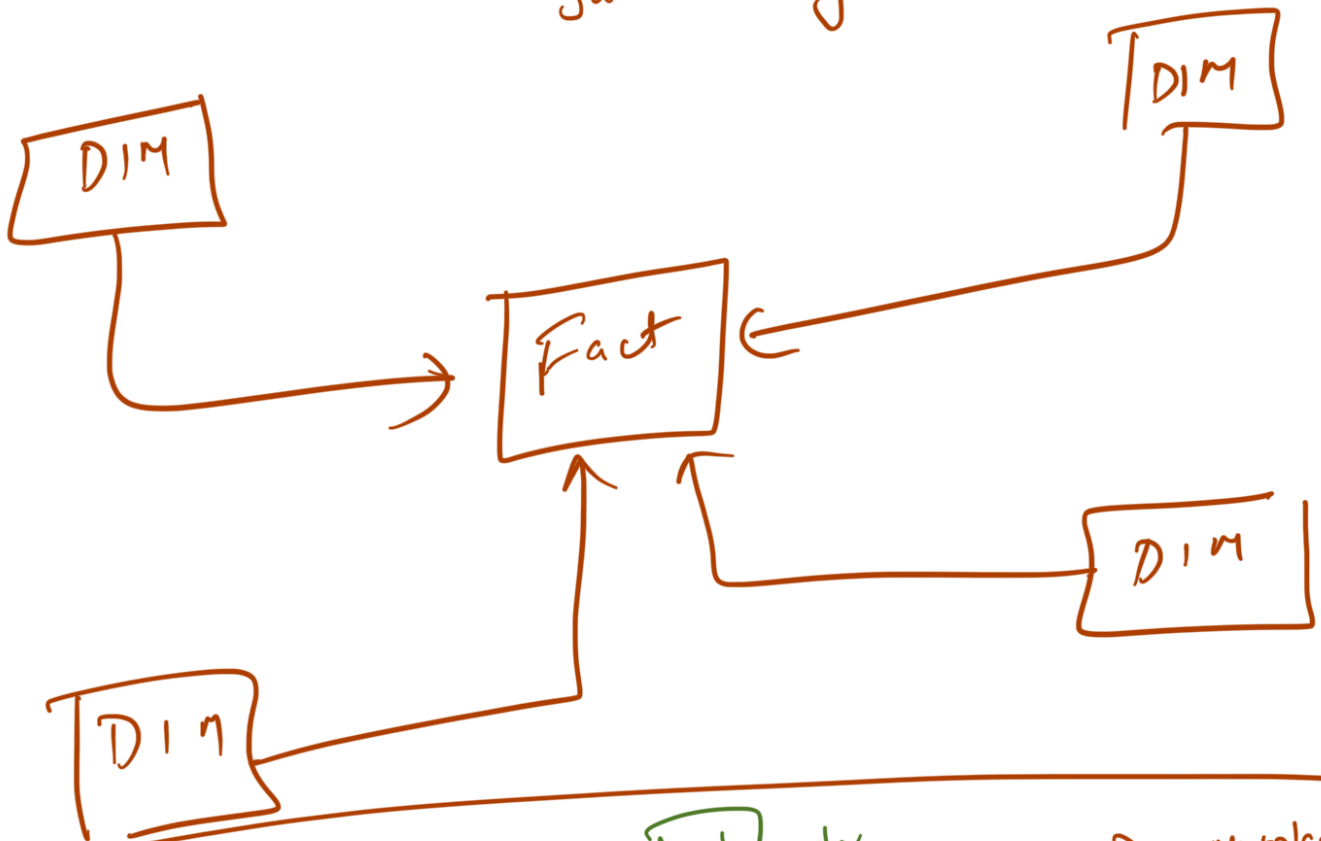


Amazon Books

- Fact table (Numerical + FK)
- Dimension table (PK + Contentual)

① Star Schema

- ↳ simplest form of data relationship structure.
- ↳ Fact table is at center, surrounded by a series of dimensions.



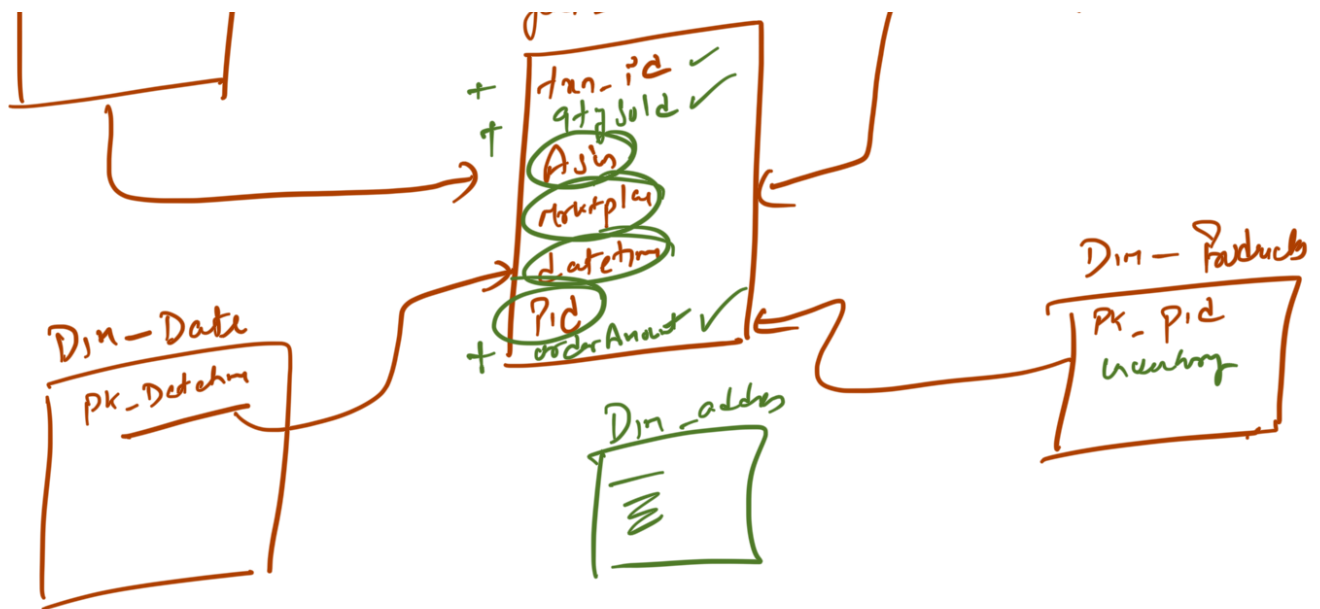
Dim - Product Books

PK_Asis
Product
descend

Fact order
fact - transaction

Dim - Multiplex

PK_Multiplex



② Snowflake Schema !!

↳ extended form of star schema by normalizing dimension table.

Snowflake (N)	star (D)
① No redundancy, hence more easy to maintain & change.	① has redundant Data and hence less easy to maintain/change.
② More complex queries and less easy to understand	② less complex queries and easy to understand.
③ More foreign keys and hence more execution time	③ less no. of keys and hence lesser execution time
④ Good to use for DWH Core to simplify relationships (M:M)	④ Good for data marts with simple relations (1:1, 1:many)
	⇒ fewer joins

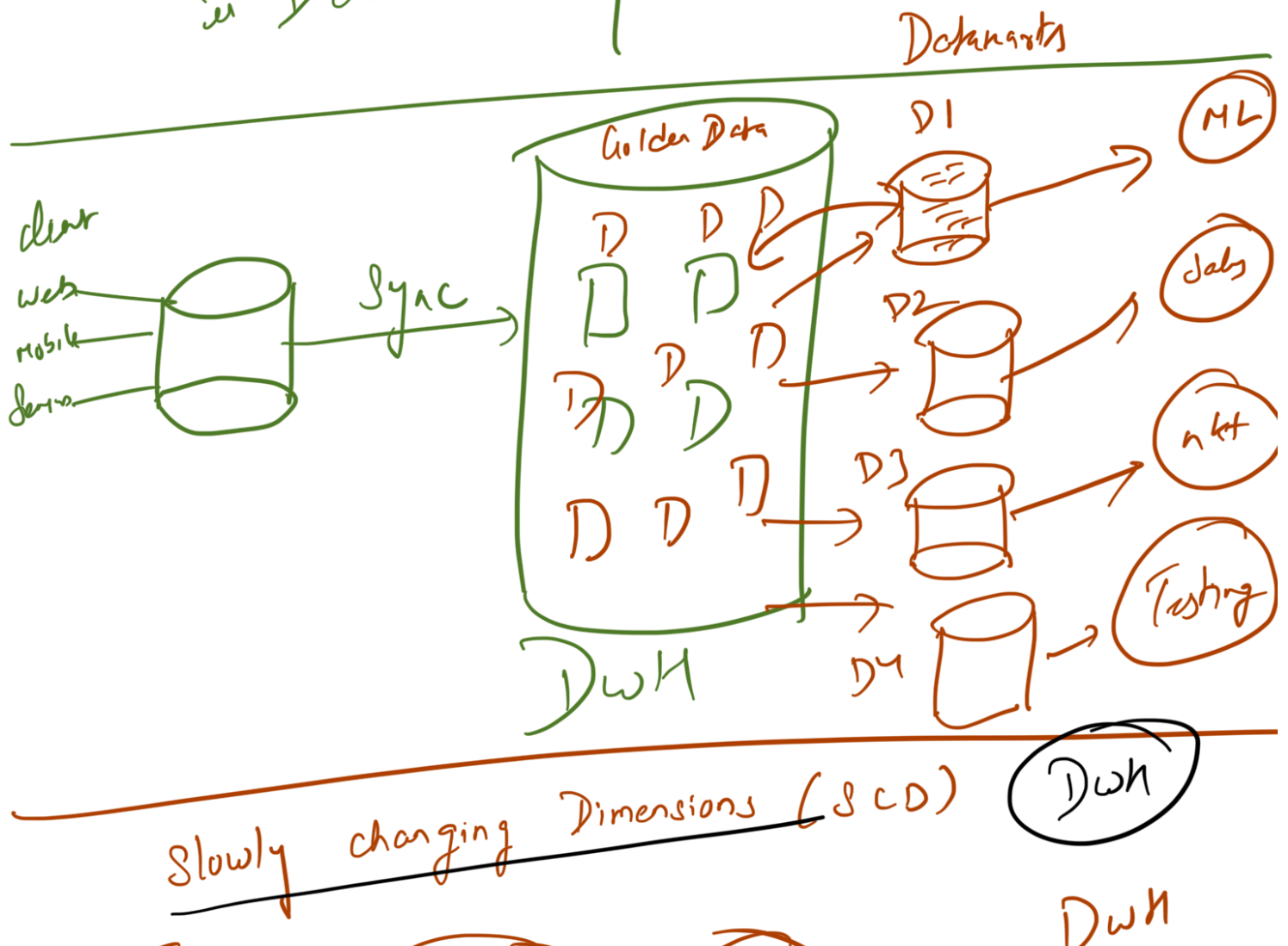
Compare

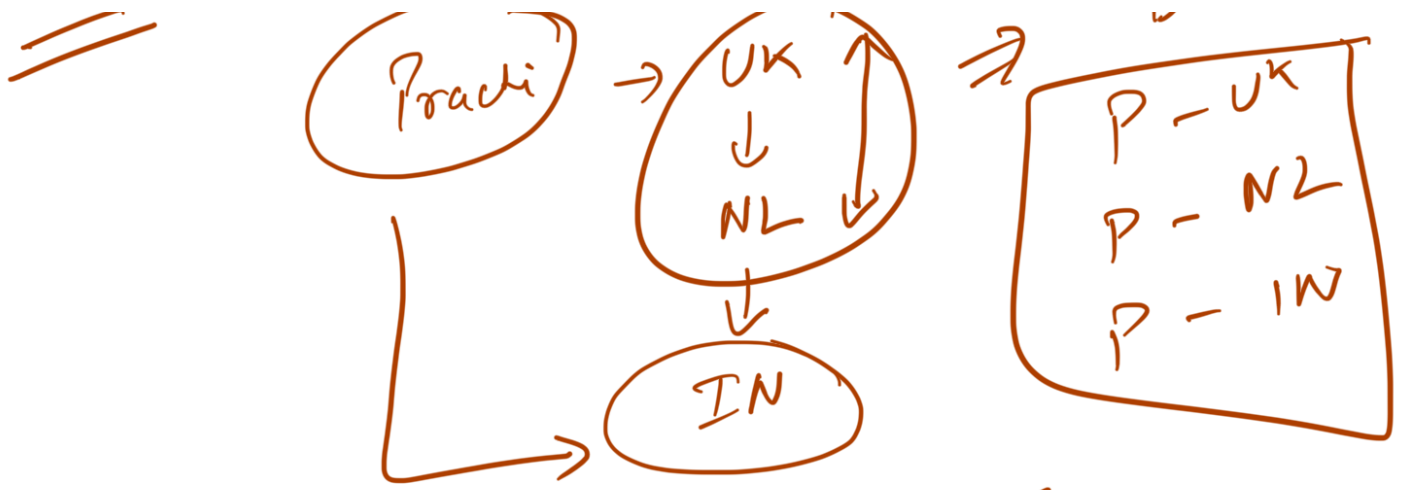
- ⑤ Higher no. of joins
- ⑥ When the dimension table is relatively big in size, Snowflaking is a better choice as it reduces data redundancy
- ⑦ Dim tables are in Normalized form but fact table is still in Denormalized

⑤ same

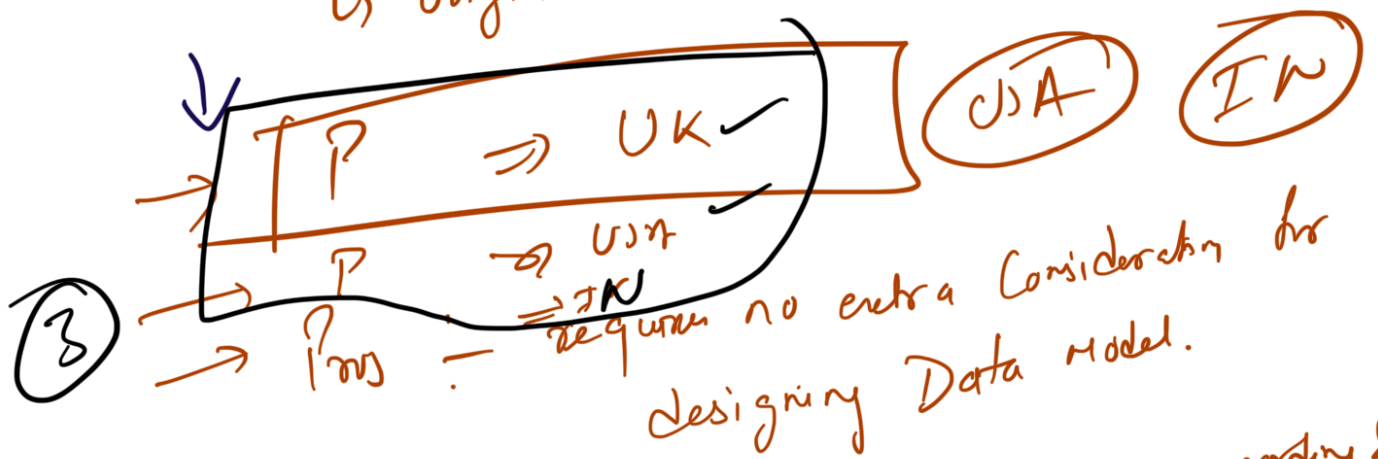
- ⑥ When the Dim table contains less no. of rows, Go with star schema

- ⑦ Both Dim/Fact are in Denormalized form.





SCD Type 0: The Passive Method ✓
 ↳ No change is Captured.
 ↳ original Data will be retained.



Cons: Provide no Control over recording & analyzing historical data & state change.

② SCD Type 1: overwriting the old values

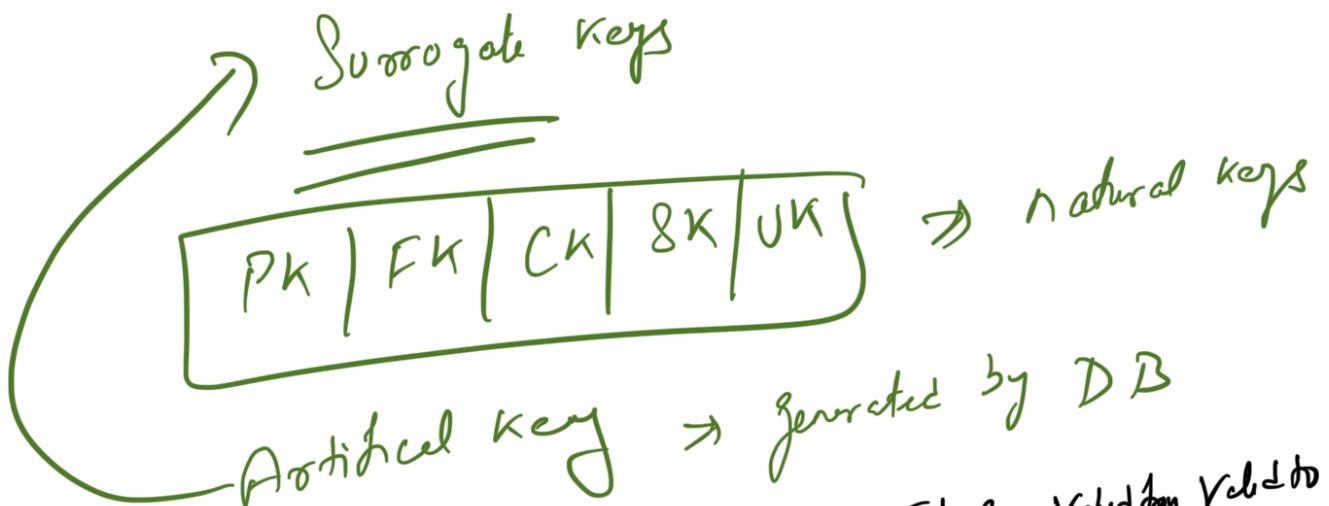
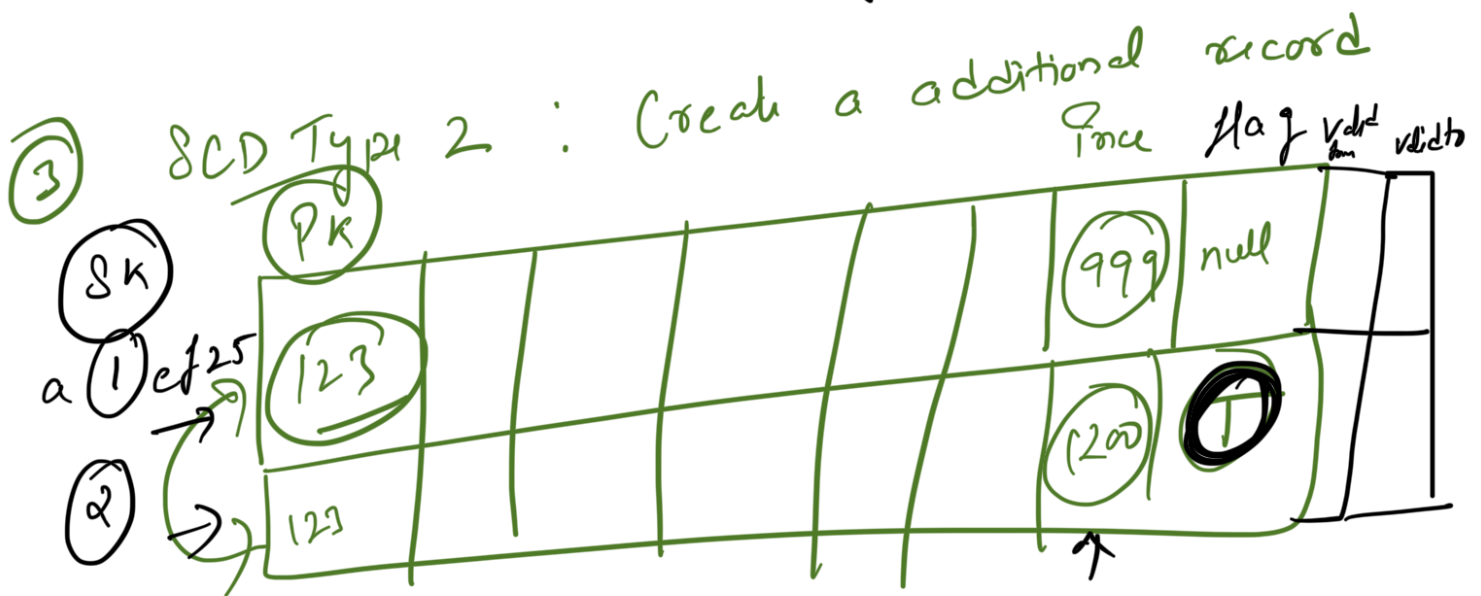
UP	IN
123	1001 1788
Crazy Rich Asian	999 X INR
1200	

→

Pros: Straightforward to implement previous values.

- ↳ we ignore ...
- ↳ Uses less disk space only 1 current record exists.

Cons :- stores no historical data for analysis
 ↳ Hard to find previous states before reaching the current one.



SK	ASEN		Flag	Valid from	Valid to
1	1a2b3c	123	0	01/01/23	01/01/24
1	1a2b3c	123	0	02/01/24	09/09/24
1	1a2b3d	123	1	09/09/24	10/04/24
1	123	-	800		

1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 | 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |

Pros :- Unlimited legacy
 ↳ store more info as to the data changes.
 ↳ Original table schema doesn't change.

Cons :- More disk space.

↳ Schema design must be decided at the time of Data modelling.

↳ changes in SCD-2 could be very expensive ops.

SCD Type 3 = Adding a new column

ASus	BN		CUP	POP	JNL
101			1200	999	

Pros :- Suitable for tracking only the most recent change.

Cons :- Design limited to the no of columns designated for storing historical data.

SCD Type 4 :- Using historical table
 ↳ clear distinction b/w

Pros :- Provides a current
historical or Current data
→ audit table works

Cons :- Two tables every time for every
Dimension or fact table
→ Disk Space is highest

* SCD Type 6 :-

SCD (1)+(2)+(3)