



CP mr 7 data lake

Left ← notebooks → Right

Udf spark SQL

PySpark

Delta Table + Notebooks

Tables Logic

Adf - Fabric

## Set Type 2

Src ID	Data	
1	A	
2	B	
3	D	
4	E	
5	F	
1	B	

needs to update or already exists

exists already

new values

changed value

⊗ wrong

1 → 0

1 → 0

1 → 0

will be inserted 3 times if ID = 10 & Egt. status = 0

① Dedup to get the latest data row with latest changes

Trg

ID	Data
2	B
3	D
4	E
1	B

ID	Data	Status	End date
1	A	1	→ 0
2	B	1	
1	C	0	

ID = ID & trg.status = 0

or

src.data < trg.data

this is wrong multiple insert for multiple true value.

② update → Enddate = SYS DATE TIME ()

trg.status = 0

if (src.ID = trg.ID & trg.status = 1 & src.data < trg.data) if true

③ Insert Left join src.ID = trg.ID & trg.status = 1

if trg.ID IS NULL

src ID	tgt ID	tgt Status
2	2	
3	Null	
4	Null	
5	Null	
1	Null	