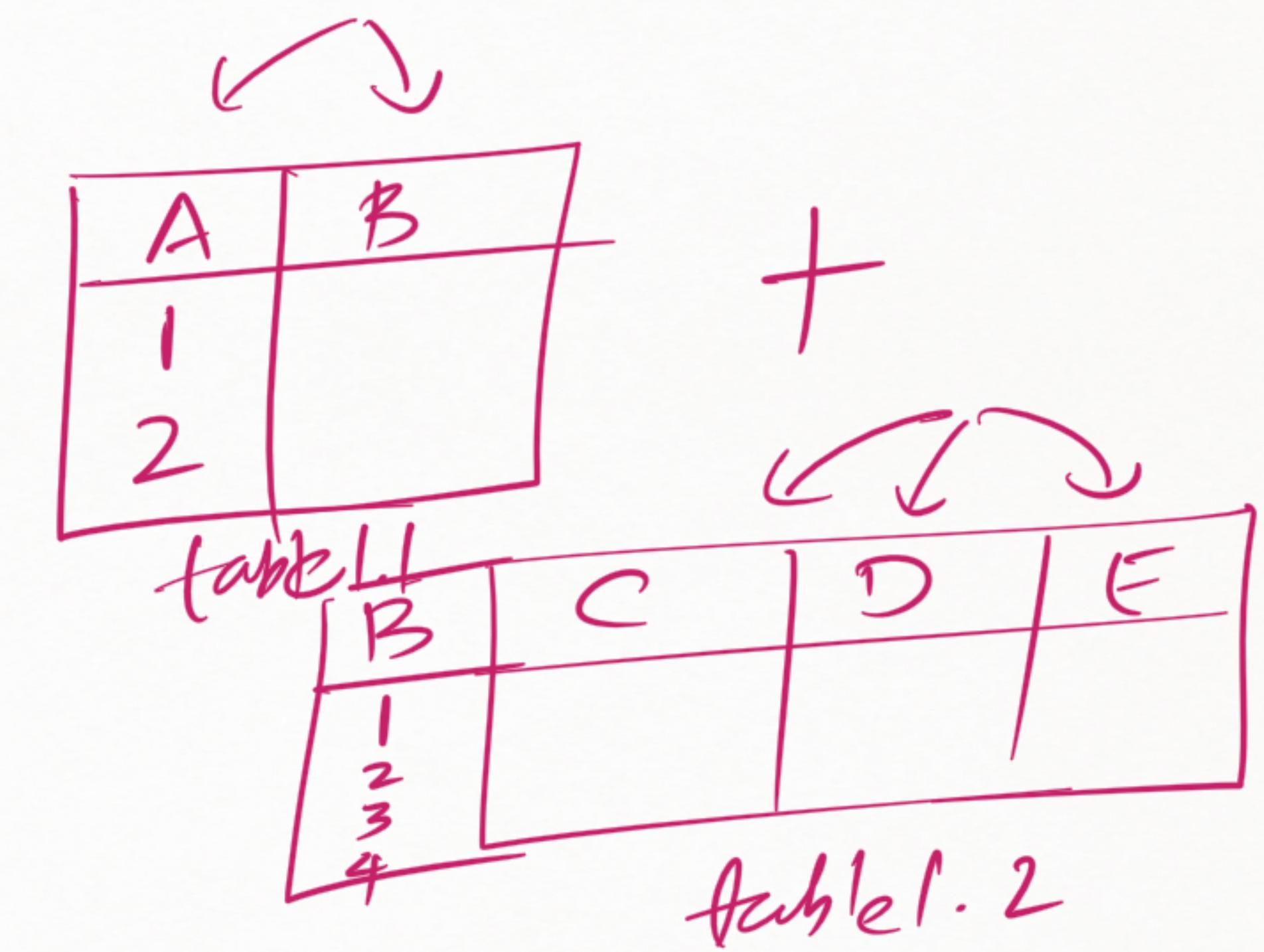
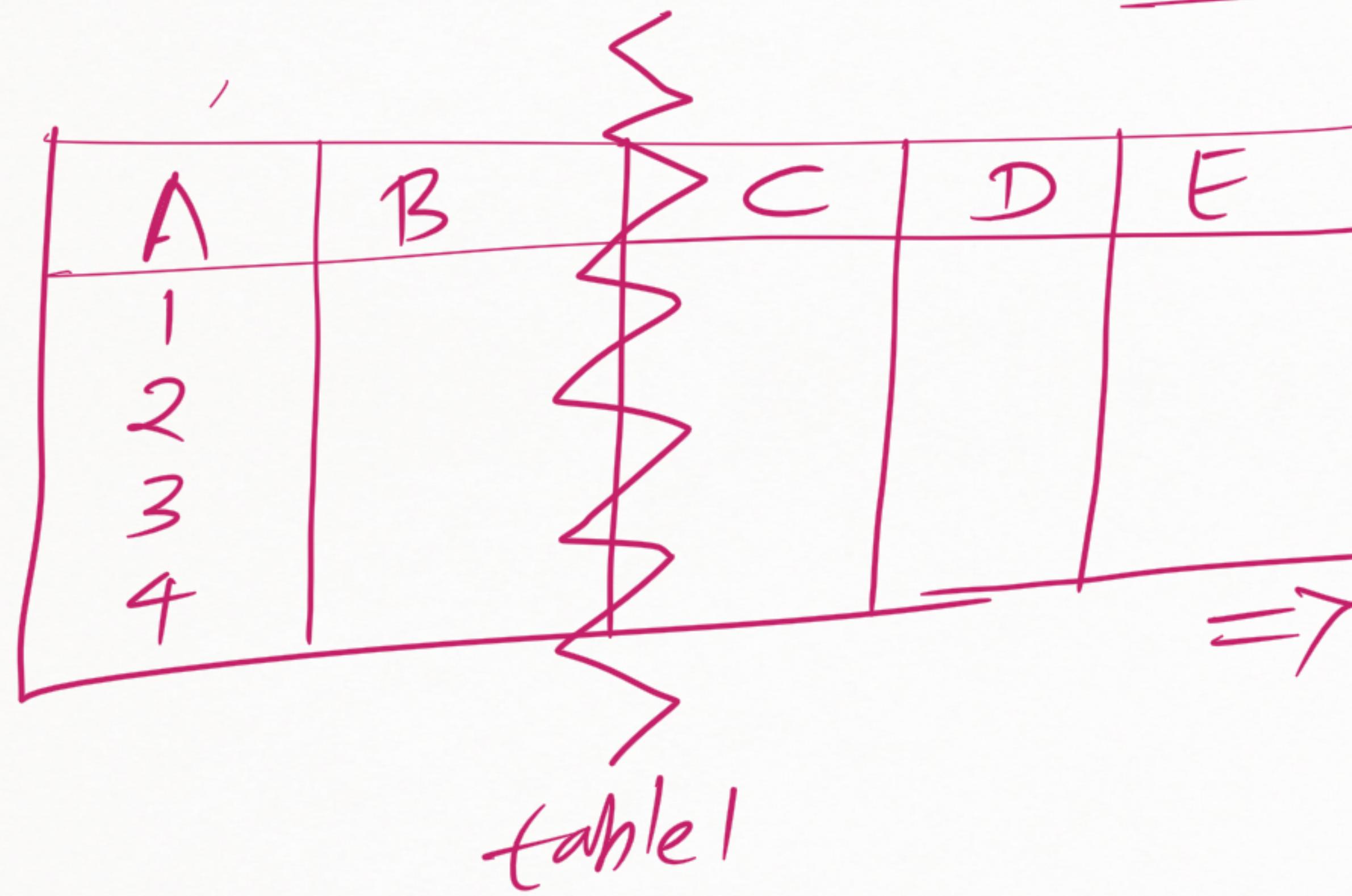


# Normalization = Breaking down of Table  
Along the columns by following  
some rules



# # Why we Need Normalization

Million

EID	E Name	Dep. ID	Dep. Name	Super ID
101	Syed	A	CS	105 ✓
102	Ankit	B	IT	110
103	Suvanna	A	CS	105 ✓
104	Mohit	=	IT	- ✓

Millions of Records

Repeating

\* All emp of Dep A  
are leaving TOB.

=> Automatically Dep A  
will also be deleted  
because Dep A data  
was stored along  
with employees  
but not separately.

① Redundancy => Repeated data

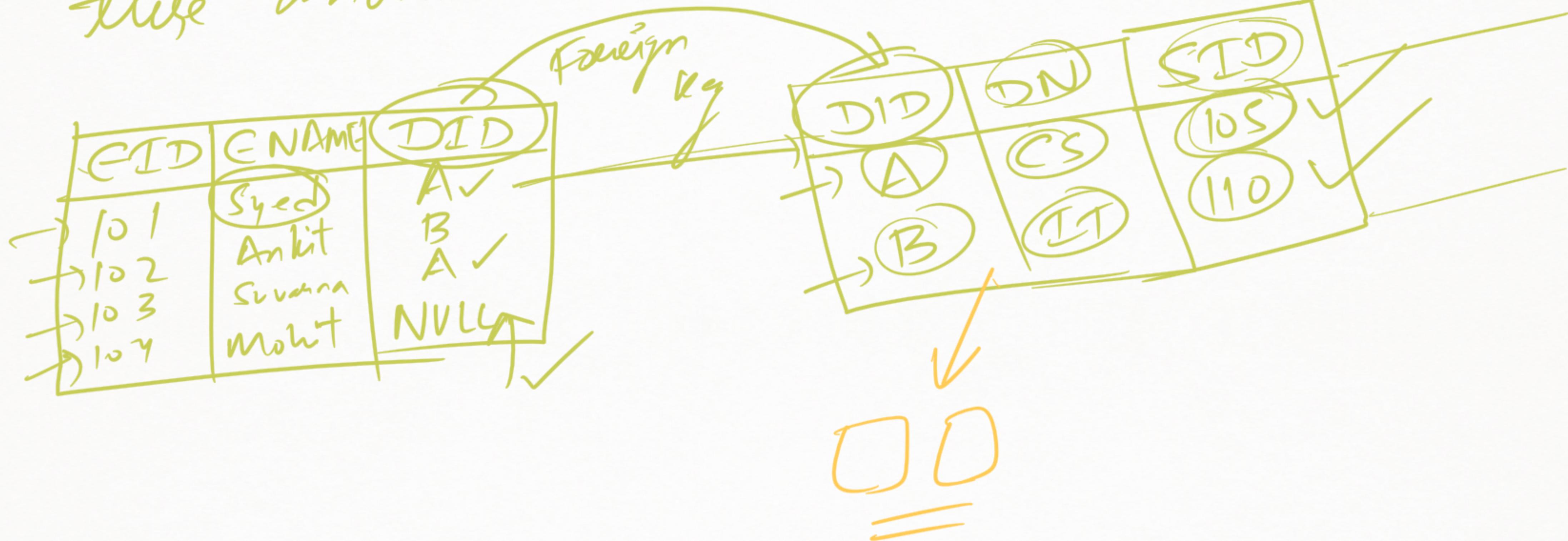
- ② Anomalies =>
- ① Insertion
  - ② Deletion
  - ③ Update each & every one

New Employee joins

Training

IT → Information Technology

# Normalization if done correctly can avoid these anomalies & redundancy.



Employee

emp_id	first_name	last_name	birth_date	sex	salary	super_id	branch_id
100	David	Wallace	1967-11-17	M	250,000	NULL	1
101	Jan	Levinson	1961-05-11	F	110,000	100	1
102	Michael	Scott	1964-03-15	M	75,000	100	2
103	Angela	Martin	1971-06-25	F	63,000	102	2
104	Kelly	Kapoor	1980-02-05	F	55,000	102	2
105	Stanley	Hudson	1958-02-19	M	69,000	102	2
106	Josh	Porter	1969-09-05	M	78,000	100	3
107	Andy	Bernard	1973-07-22	M	65,000	106	3
108	Jim	Halpert	1978-10-01	M	71,000	106	3

Branch

branch_id	branch_name	mgr_id	mgr_start_date
1	Corporate	106	2006-02-09
2	Scranton	102	1992-04-06
3	Stamford	106	1996-02-13

Works\_With

emp_id	client_id	total_sales
105	400	55,000
102	401	267,000
108	402	22,500
107	403	5,000
108	403	12,000
105	404	33,000
107	405	26,000
102	406	15,000
105	406	130,000

Client

client_id	client_name	branch_id
400	Dunmore Highschool	2
401	Lackawana Country	2
402	FedEx	3
403	John Daly Law, LLC	3
404	Scranton Whitepages	2
405	Times Newspaper	3
406	FedEx	2

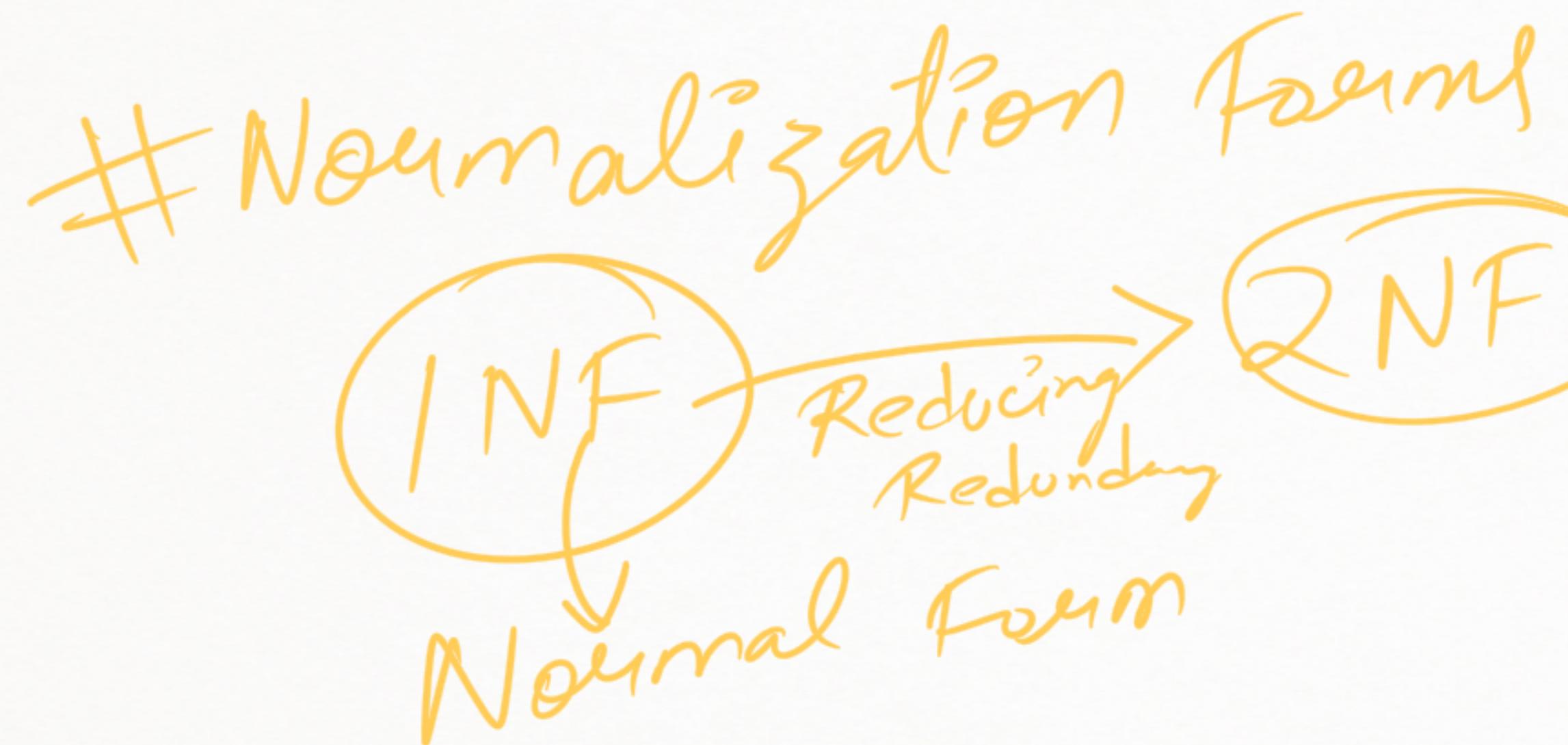
Branch Supplier

branch_id	supplier_name	supply_type
2	Hammer Mill	Paper
2	Uni-ball	Writing Utensils
3	Patriot Paper	Paper
2	J.T. Forms & Labels	Custom Forms
3	Uni-ball	Writing Utensils
3	Hammer Mill	Paper
3	Stamford Lables	Custom Forms

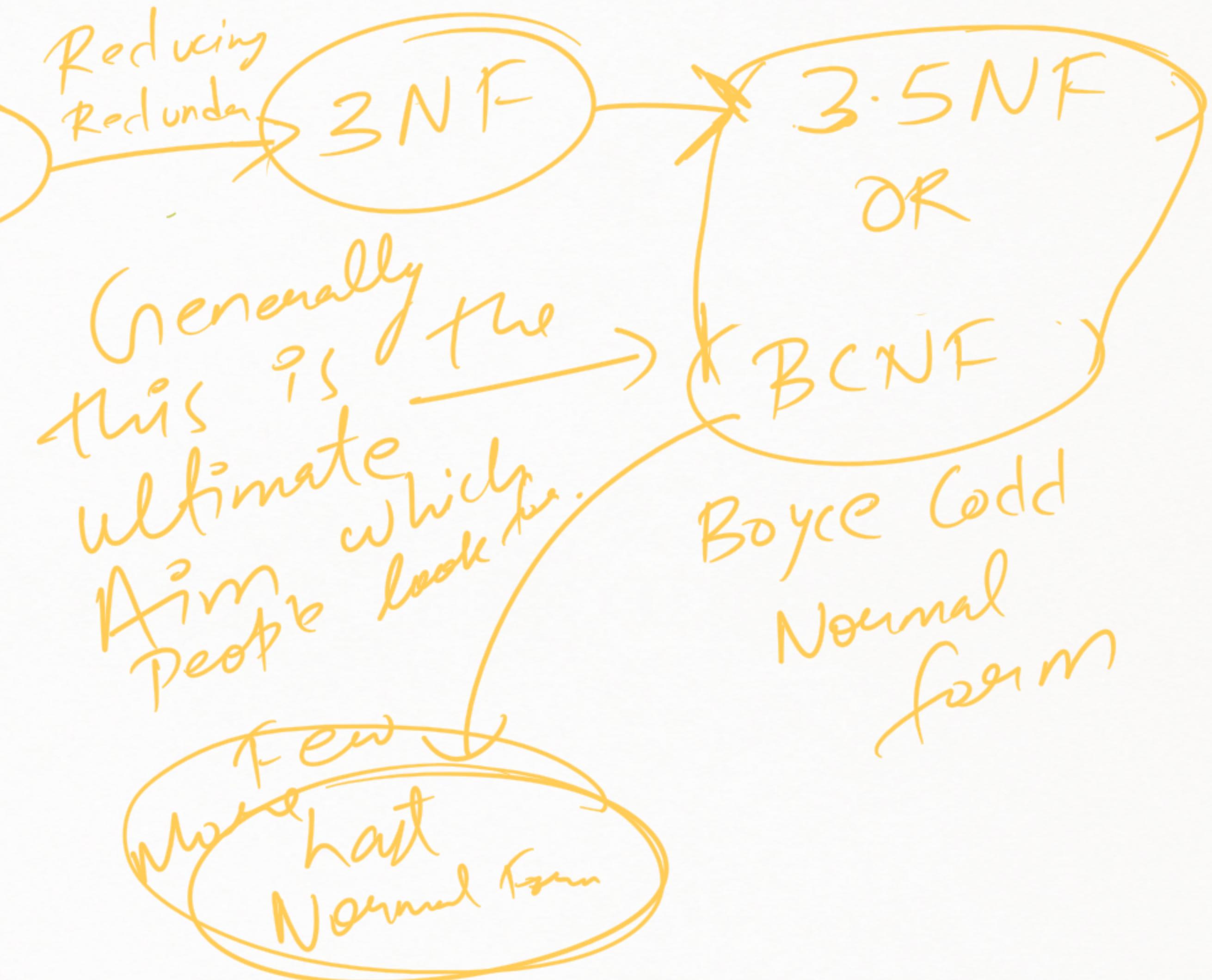
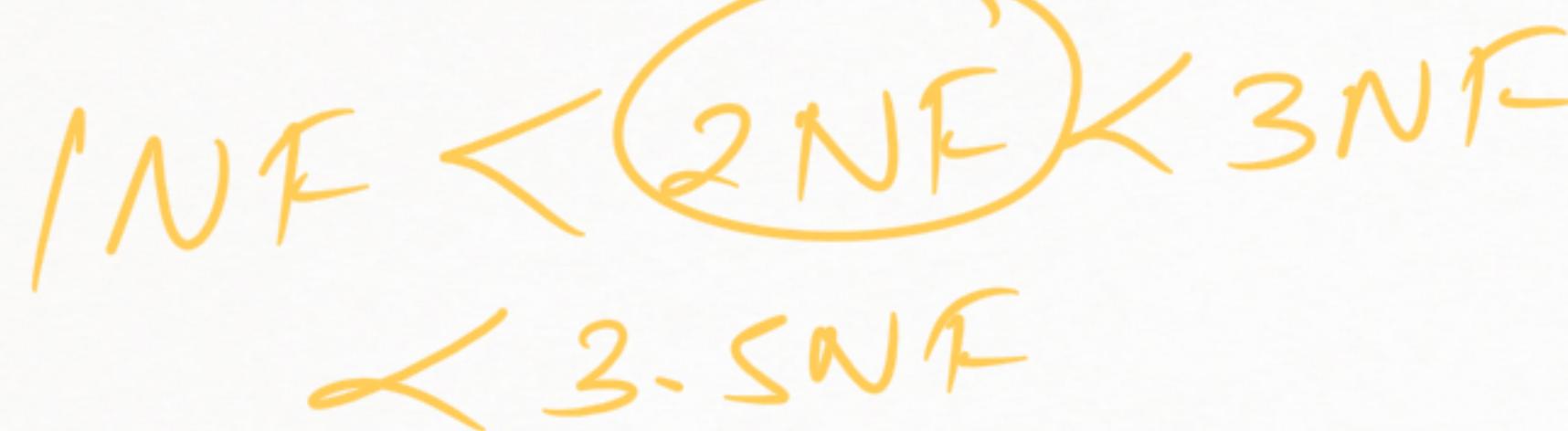
# What if we create just single Table for this huge data?

# What if we combine Employee along with the Branch

A LOT OF CHAOS



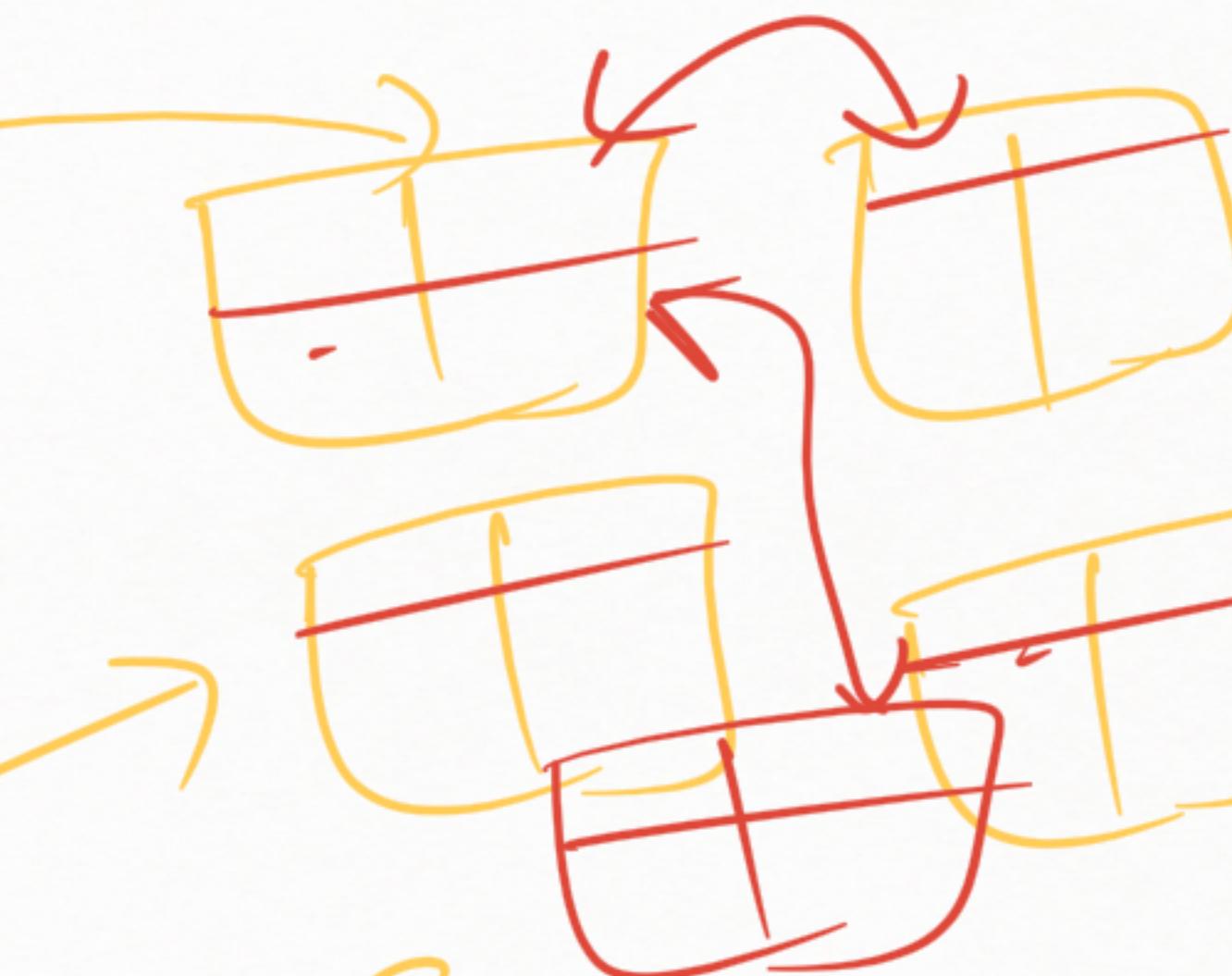
Normalization Criteria



$3\text{-}5\text{NF} = 3\text{NF} + \text{Something extra}$

$3\text{NF} = 2\text{NF} + \text{"}$

$2\text{NF} = \text{1NF} + \text{"}$



This method can  
reduce redundancy  
completely & falls  
in highest NF.

So many tables  
But Query will increase  
so much

When we  
have to  
gain so  
many tables  
for retrieval  
of Data.

## How to Do Normalization

- We have to follow some rules
- To understand those rules, we will have to understand few concepts like Partial Dependency & Functional Dependency.

# Functional Dependency

A	B	C
1	a	b
2	a	c
3	d	e
4	f	g

$A \rightarrow BC$

A determines BC

Given the value of A, you will be able to identify B & C

$$x^2 - y = x^2 - 2$$

## # Formal Definition of FD

$X \rightarrow Y$

Collection of Attributes      Collection of Attributes

	A	B
$a_1$	1 ✓	a ✓
$a_2$	1 ✓	a ✓
$a_3$	2 ✓	b ✓
$a_4$	3 ✓	c ✓



We don't have to worry

Ex:  $AB \rightarrow CD$   
 $A \rightarrow BC$

$X \rightarrow Y$

CASE	X	Y	✓
1	If $a_1, a_2$ are same here ✓	They must be same here also	
2	If $a_1, a_2$ are different here	They may or not be diff here ✓	

## # How to rule out the wrong FD

A	B	C
a	b	c
a	b	e
a	b	f
a	b	g

↓  
↓  
there can be  
so many  
rows, only  
few are  
shown to  
you !!

$$A \rightarrow B$$

Q. If it is a wrong  
FD

No

Case 1  $\Rightarrow$

Whenever A  
repeats, B  
should also  
repeat.

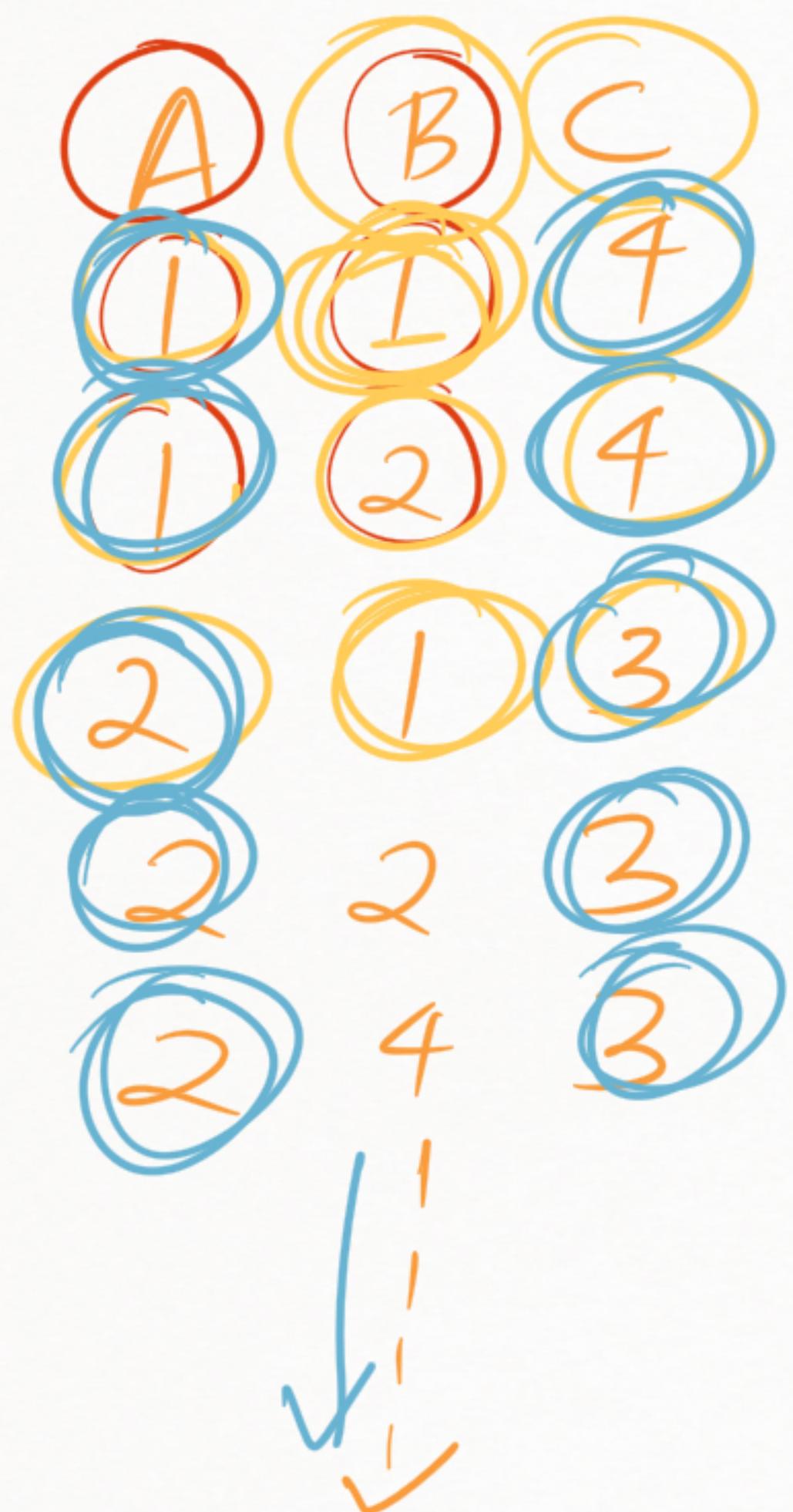
$$A \rightarrow C$$

Yes

Whenever

A  
subject  
C

does  
not.



~~FDS~~  
 A → B  
 B → C  
 B → A  
 C → B  
 C → A <sup>DK</sup>  
 I don't know  
 A → C <sup>DK</sup>

Which are  
 wrong F.D's  
 ⇒ 1, 2, 3, 4

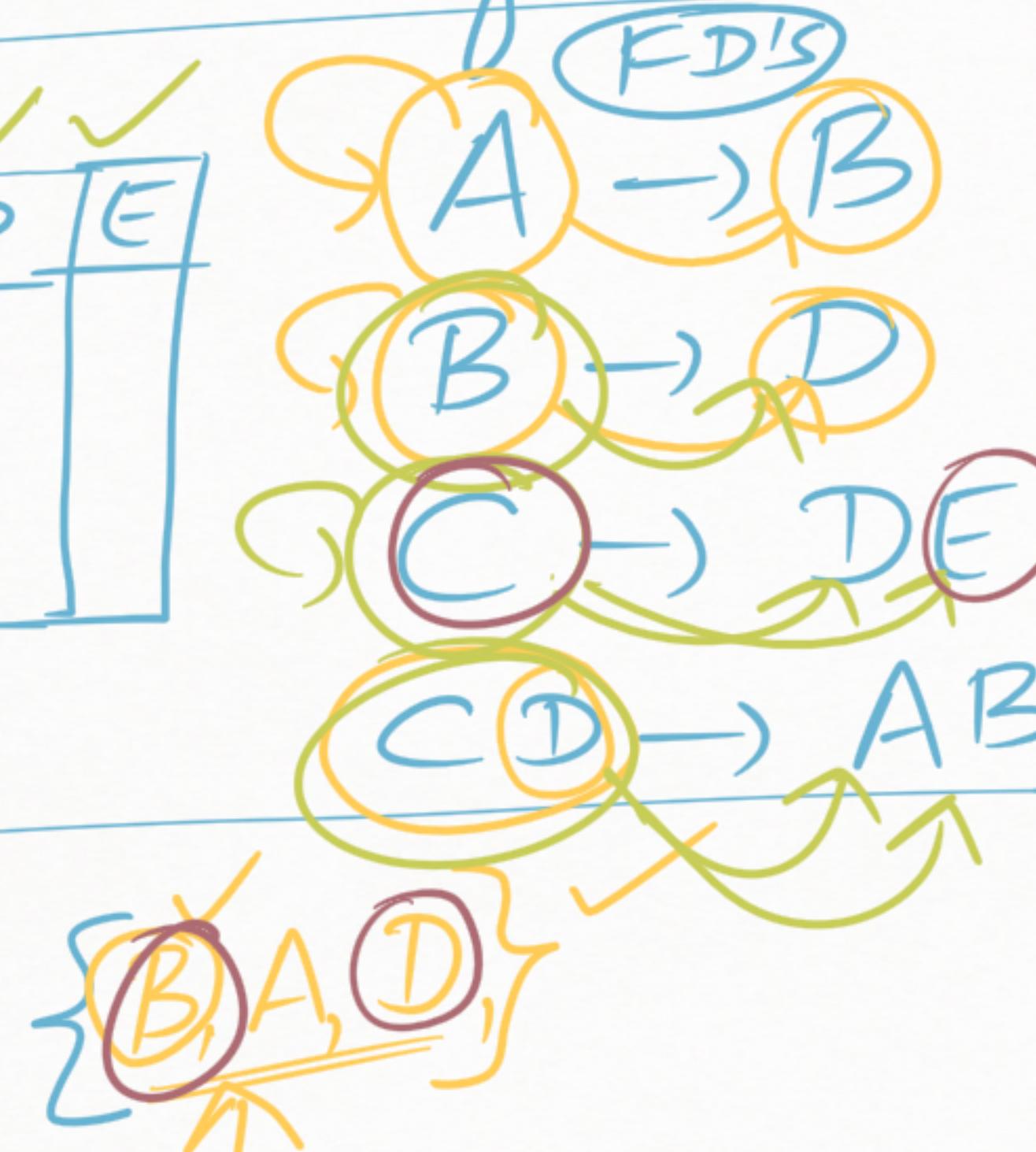
## # Closure Set of Attributes

A	B	C	D	E

$$CA^+ = \{all\checkmark\}$$

$$A^+$$

(closure set of A)



Whenever something is here we can use that

$$CB = \{B, all\checkmark\}$$

Closure of B

$$B^+ = \{B, D\}$$

$$\begin{aligned} AB^+ &= \{A, B, D\} \\ AC^+ &= \{A \cup B, D\} \\ &\quad E \} \end{aligned}$$

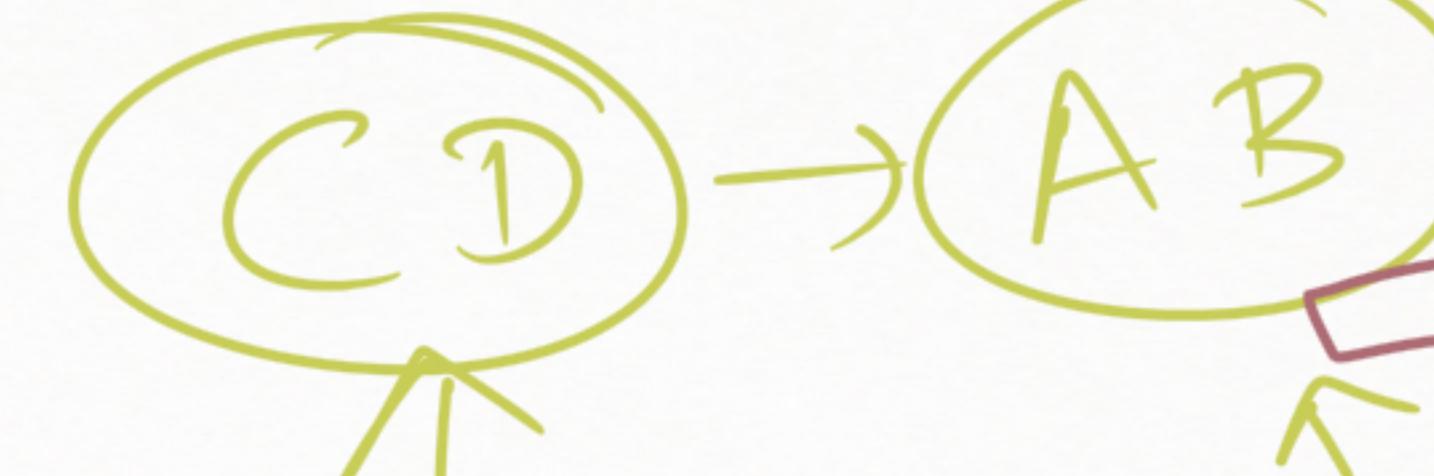
all ✓

Closure of C

$$C^+ = \{C, D, E, A, B\}$$

$$D^+ = \{D\}$$

$$\begin{aligned} AD^+ &= \\ BD^+ &= \end{aligned}$$



If we have complete  $CD$  only then

we can determine  $AB$  ✓

## # Candidate Key & Super Key

① Candidate Key  $\Rightarrow$  Given an attribute or a collection of attribute, we will be able to uniquely identify a row.

Note  $\Rightarrow$  If collection of Attributes, than those attributes have to be minimal.

$\Rightarrow$  In the previous case  
possible to be minimised  
 $C_A = \{ \text{all } \checkmark \}$   
 $C = \{ \text{all } \checkmark \}$

But it will violate the Note  $C_k$

Not a  $C_k$

ID	ID2	Name	Clan
1	A ✓	Mohan	II
2	B ✓	Rohan	II
3	C ✓	Mohan	I

What is the Name for  $ID1 = 2$   
 Can you ans it without any doubt?  $\Rightarrow$  YES  
 Rohan ✓

For Clan II, what is Name  
 Can you ans it without Doubt?  $\Rightarrow$  NO  
 Mohan OR ???  
 Rohan ?!

$ID1^+ \neq \{ \text{Name}, ID1, ID2, \text{Clan} \}$   
 Can we minimize?  
 $ID2 \subset_k \{ \text{Name}, ID1, ID2, \text{Clan} \}$

\* Primary Key  $\Rightarrow$  It is just one of the Candidate Key that we wish to choose.

$\Rightarrow$  There is a difference when we say Primary key is a collection of attributes/columns vs a ~~collection~~ wrong it does not happen. There are a ~~collection~~ of Primary keys

There is only 1 and 1 Primary key

# Super Key = Candidate Key +  
example

# example

D1	D2	Nnn	Cn
Z	A S	M R K	T E F T
3			GK

IDI

DD2

Diagram illustrating the relationship between three data fields:

- ID1 + Name
- ID1 + Email + Name
- Upper Key

Relationships shown:

- ID1 + Name is connected to ID1 + Email + Name.
- ID1 + Name is connected to Upper Key.

Sub	Teacher Name	Dep
A	Syed	IT
B	Ankit	CS
A	Mohit	CS
C	Syed	Maths

$C_k$  could be a collection of attributes

Minimal

What are  $C_k$ 's here?

What is the dep.

of Syed?

Can you ans  
it without  
doubt? NO

Any other  $C_k$ ?  
~~Teacher Name + Dept~~  
~~Sub + Dept~~

Sub + Teacher Name	Dep
A Syed	IT
B Ankit	CS
A Mohit	CS
C Syed	Maths

Can we determine  
dep from  
Sub + Teacher Name  
without doubt?

Sub + Teacher Name  
 $C_k$