

Database Systems

Instructor: Hina Iqbal

Data Normalization

- Primarily a tool to validate and improve a logical design so that it satisfies certain constraints that ***avoid unnecessary duplication of data***
- The process of decomposing relations with anomalies to produce smaller, ***well-structured*** relations

Well-Structured Relations

- A relation that contains minimal data redundancy and allows users to insert, delete, and update rows without causing data inconsistencies
- Goal is to avoid anomalies
 - **Insertion Anomaly**—adding new rows forces user to create duplicate data
 - **Deletion Anomaly**—deleting rows may cause a loss of data that would be needed for other future rows
 - **Modification Anomaly**—changing data in a row forces changes to other rows because of duplication

General rule of thumb: A table should not pertain to more than one entity type

Example

EMPLOYEE2

<u>Emp_ID</u>	Name	Dept_Name	Salary	Course_Title	Date_Completed
100	Margaret Simpson	Marketing	48,000	SPSS	6/19/200X
100	Margaret Simpson	Marketing	48,000	Surveys	10/7/200X
140	Alan Beeton	Accounting	52,000	Tax Acc	12/8/200X
110	Chris Lucero	Info Systems	43,000	Visual Basic	1/12/200X
110	Chris Lucero	Info Systems	43,000	C++	4/22/200X
190	Lorenzo Davis	Finance	55,000		
150	Susan Martin	Marketing	42,000	SPSS	6/19/200X
150	Susan Martin	Marketing	42,000	Java	8/12/200X

Question—Is this a relation?

Answer—Yes: Unique rows and no multivalued attributes

Question—What's the primary key?

Answer—Composite: Emp_ID, Course_Title

Anomalies in this Table

- **Insertion**—can't enter a new employee without having the employee take a class
- **Deletion**—if we remove employee 140, we lose information about the existence of a Tax Acc class
- **Modification**—giving a salary increase to employee 100 forces us to update multiple records

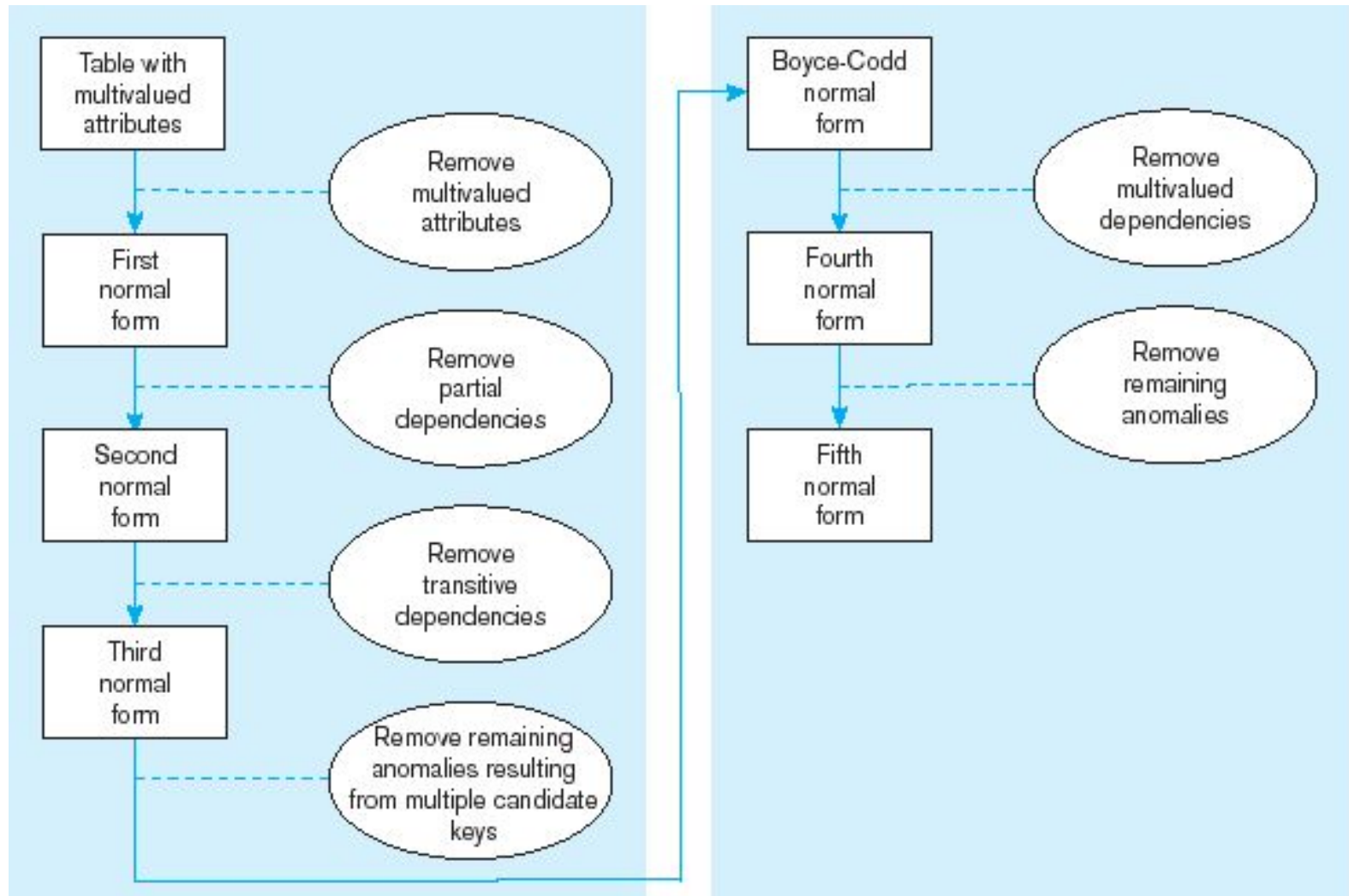
Why do these anomalies exist?

Because there are two themes (entity types) in this one relation. This results in data duplication and an unnecessary dependency between the entities

Functional Dependencies and Keys

- Functional Dependency: The value of one attribute (the ***determinant***) determines the value of another attribute
- Candidate Key:
 - A unique identifier. One of the candidate keys will become the primary key
 - E.g. perhaps there is both credit card number and SS# in a table...in this case both are candidate keys
 - Each non-key field is functionally dependent on every candidate key

Steps in normalization



First Normal Form

- No multivalued attributes
- Every attribute value is atomic
- Fig. 5-25 *is not* in 1st Normal Form (multivalued attributes) □ it is not a relation
- Fig. 5-26 *is* in 1st Normal form
- ***All relations are in 1st Normal Form***

Table with multivalued attributes, not in 1st normal form

Figure 5-25
INVOICE data (Pine Valley Furniture Company)

<u>Order ID</u>	Order_ Date	Customer_ ID	Customer_ Name	Customer_ Address	<u>Product ID</u>	Product_ Description	Product_ Finish	Unit_ Price	Ordered_ Quantity
1006	10/24/2006	2	Value Furniture	Plano, TX	7	Dining Table	Natural Ash	800.00	2
					5	Writer's Desk	Cherry	325.00	2
					4	Entertainment Center	Natural Maple	650.00	1
1007	10/25/2006	6	Furniture Gallery	Boulder, CO	11	4-Dr Dresser	Oak	500.00	4
					4	Entertainment Center	Natural Maple	650.00	3

Note: this is NOT a relation

Table with no multivalued attributes and unique rows, in 1st normal form

<u>Order_ID</u>	Order_ Date	Customer_ ID	Customer_ Name	Customer_ Address	<u>Product_ID</u>	Product_ Description	Product_ Finish	Unit_ Price	Ordered_ Quantity
1006	10/24/2006	2	Value Furniture	Plano, TX	7	Dining Table	Natural Ash	800.00	2
1006	10/24/2006	2	Value Furniture	Plano, TX	5	Writer's Desk	Cherry	325.00	2
1006	10/24/2006	2	Value Furniture	Plano, TX	4	Entertainment Center	Natural Maple	650.00	1
1007	10/25/2006	6	Furniture Gallery	Boulder, CO	11	4-Dr Dresser	Oak	500.00	4
1007	10/25/2006	6	Furniture Gallery	Boulder, CO	4	Entertainment Center	Natural Maple	650.00	3

Figure 5-26

INVOICE relation (1NF) (Fine Valley
Furniture Company)

Product_ID → Product_Description, Product_Finish, Unit_Price
Order_ID, Product_ID → Ordered_Quantity

Note: this is relation, but not a well-structured one

Anomalies in this Table

- **Insertion**—if new product is ordered for order 1007 of existing customer, customer data must be re-entered, causing duplication
- **Deletion**—if we delete the Dining Table from Order 1006, we lose information concerning this item's finish and price
- **Update**—changing the price of product ID 4 requires update in several records

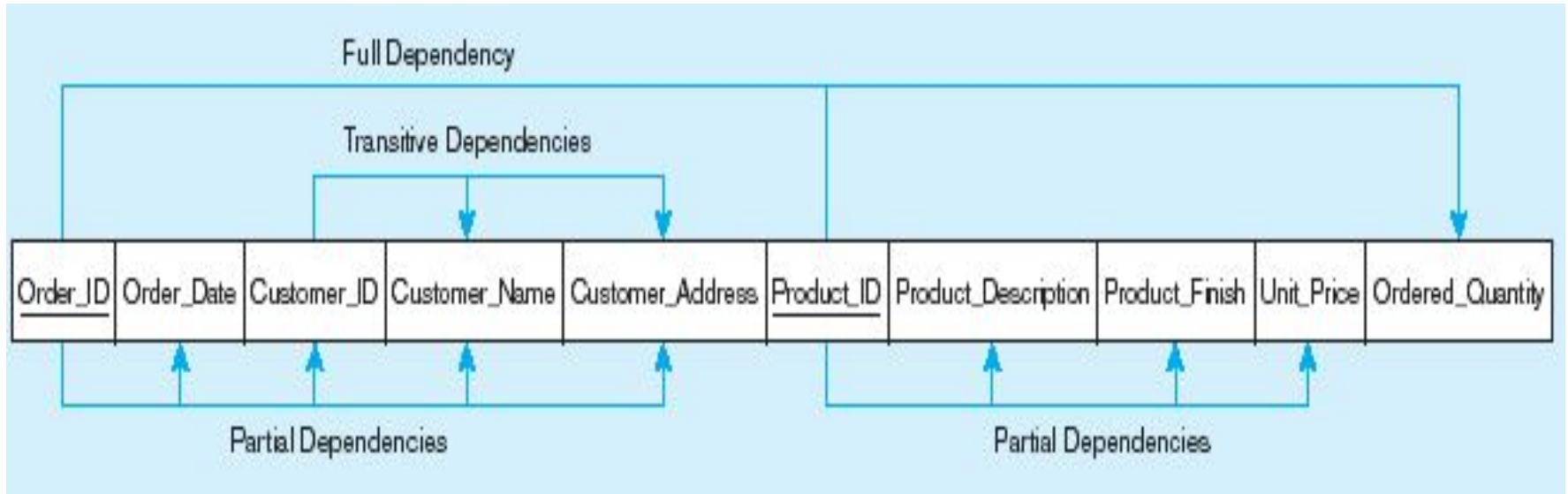
Why do these anomalies exist?

Because there are multiple themes (entity types) in one relation. This results in duplication and an unnecessary dependency between the entities

Second Normal Form

- 1NF PLUS ***every non-key attribute is fully functionally dependent on the ENTIRE primary key***
 - Every non-key attribute must be defined by the entire key, not by only part of the key
 - No partial functional dependencies

Functional dependency diagram for INVOICE



Order_ID \square **Order_Date, Customer_ID, Customer_Name, Customer_Address**

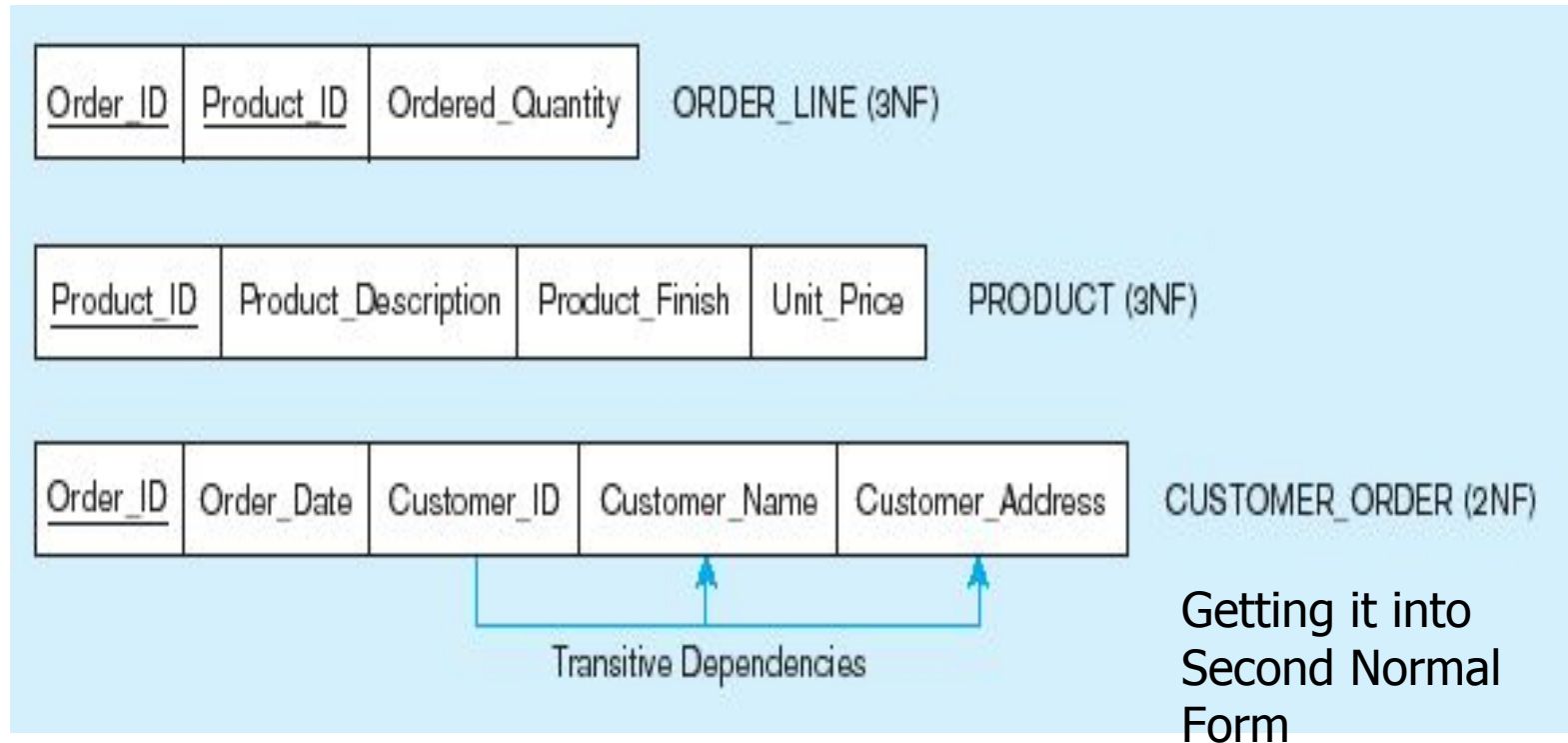
Customer_ID \square **Customer_Name, Customer_Address**

Product_ID \square **Product_Description, Product_Finish, Unit_Price**

Order_ID, Product_ID \square **Order_Quantity**

Therefore, NOT in 2nd Normal Form

Removing partial dependencies



Partial dependencies are removed, but there are still transitive dependencies

Third Normal Form

- 2NF PLUS ***no transitive dependencies*** (functional dependencies on non-primary-key attributes)
- Note: This is called transitive, because the primary key is a determinant for another attribute, which in turn is a determinant for a third
- Solution: Non-key determinant with transitive dependencies go into a new table; non-key determinant becomes primary key in the new table and stays as foreign key in the old table

Removing partial dependencies

<u>Order_ID</u>	Order_Date	----- Customer_ID
-----------------	------------	----------------------

ORDER (3NF)

Getting it into
Third Normal
Form

<u>Customer_ID</u>	Customer_Name	Customer_Address
--------------------	---------------	------------------

CUSTOMER (3NF)

Transitive dependencies are removed