

Chapter #3

Analysis & Reconciliation of Data Sources

Local source schema:

The schema of an individual data source.

- ⇒ Local source schema refers to the conceptual schema that describes the data stored in a source without regard to the models and technological solutions implemented.
- ⇒ focus on the representation of the application domain.

- Local source schemas result from the analysis of relations and constraints in relational database.

Global schema: A unified schema that represents the data from all sources.

Definition:

- Refers to the processes used to evaluate, compare and ensure consistency, accuracy, and integration of data from multiple sources before consolidating them in a unified system, such as data warehouse.

(i) Data source analysis:

Data sources can show strong relationships or can be completely independent.

Def: The process of identifying and evaluating the available data sources.

Its Data most designers primary goal is to improve his or her knowledge of data issues in the source analysis phase.

Key question involves:

- What are the data types, formats, and relationships in the data source?
- Are there inconsistencies or missing data that need to be addressed?
- Are there duplicate records or conflict information across sources?
- How does data from one source relate to data from another source?

(ii) Reconciliation process:

- ⇒ Reconciliation process is needed in order to obtain consistent, error free information.
- ⇒ Reconciliation process involves:

- integrating
- cleaning
- transforming

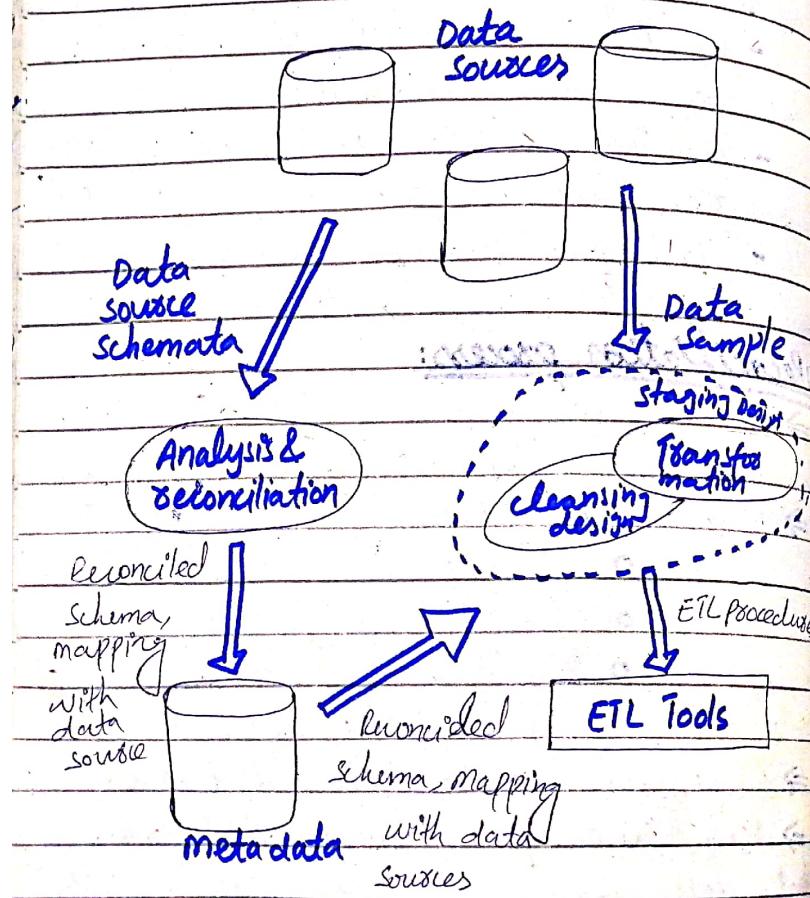
data to create a reconciled database.

⇒ This process requires time and resources.

⇒ Reconciliation process ensures that ensuring that data from different sources matches and correctly integrated.

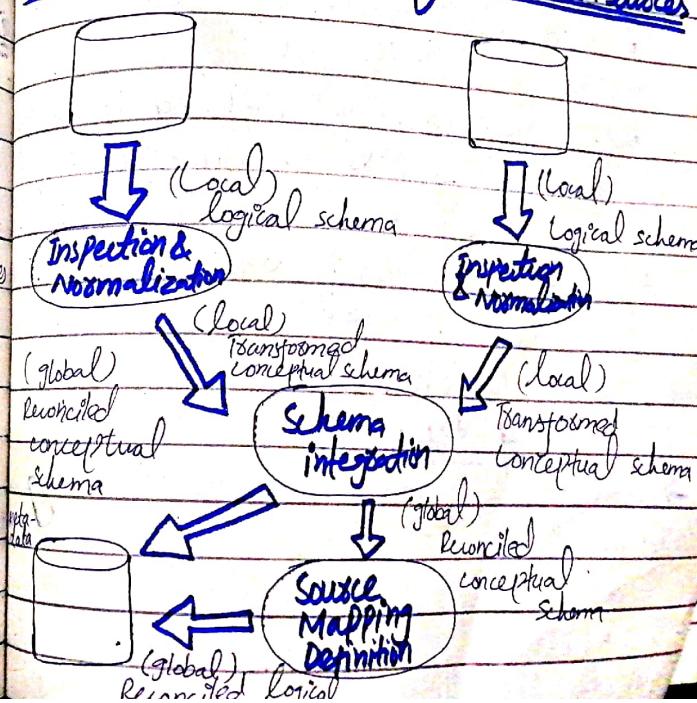
⇒ Integration phase is based on the intensional features of data sources.

Three-Phase Reconciled layers design from data source layers:



In short:
Reconciliation and analysis phase
do the following:
o Inspect and normalize all the local schemata to generate a set of comprehensive, locally consistent, conceptual schemata.
o Then do integration, resulting in a globally consistent conceptual schema.

Design steps forming the analysis & reconciliation phase of two data sources.



Inspecting & Normalizing Schemata

⇒ Detailed knowledge of data sources is necessary before data mart. conceptual design phase.

⇒ For doing so, need to perform the following tasks:

- Inspection
- Normalization

Inspection:

- Close investigation of local schemata.
- Refers to examining and understanding the structure, format, and quality of raw data before performing transformation.
- Goal is to identify any issues or patterns that might impact data integration into the data warehouse.

Key steps:

⇒ Data profiling - null values or missing data

- Inconsistent data formats
- outliers or anomalies
- Duplicate Records

- Inconsistent naming convention
- Data type verification of fields.
- Data integrity check
- Data consistency
- Error detection

Normalization:

- An attempt to rectify local schemata to model the database house domain as accurately as possible.
- Process of organizing data into structured format, before put into the data warehouses, where vast amount of data from different sources must be harmonized.
- Denormalization might be necessary in data warehousing to optimize query performance to speed up read operations.

⇒ In analysis phase, designers collaborate with the experts, such as data processing center staff, managers etc.

⇒ Also they need to check or find any relationships that was evident missed or left.

For example,

Explicitly representation functional dependencies skipped and new relationships between entities.

⇒ Both processes can brought changes in local schema for adjusting the specific data warehousing project needs.

⇒ All changes applied to local schemata, clarify explicitly all the concepts extracted from the data stored in data warehouse.

⇒ Normalization and inspection should be even single data source available.

⇒ For more than one data sources, both processes repeated for every local schema.

⇒ Data designer should also check for which portion of local schemata need.

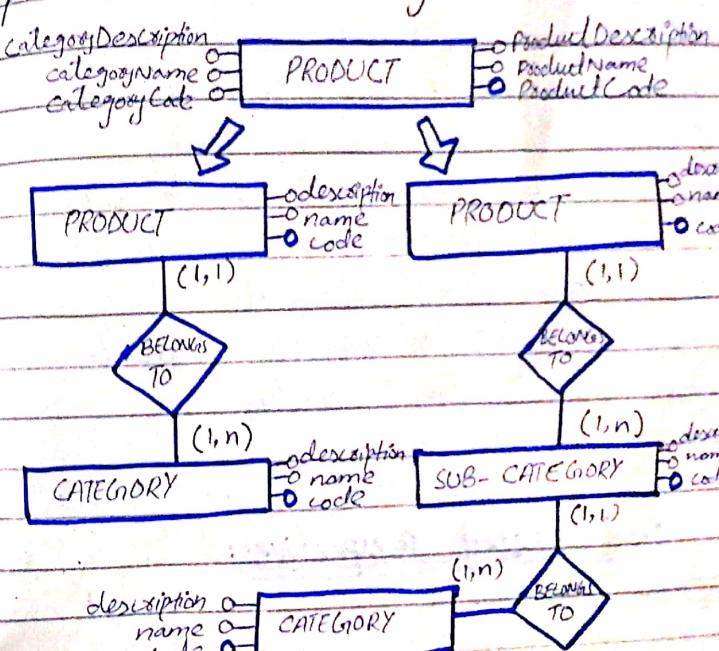
Example

⇒ Suppose consider the ER-schema modeling of specific company production.

- ⇒ First change:
- ProductCode → CategoryCode can be apply bcz functional dependency is true.
- ⇒ Second change:

- BELONGS TO association

Cardinality cannot be applied because the operation data does not contain the information required to represent explicitly any products into subcategories.



- Show applicable & inapplicable changes in the source inspection & normalization phase

The Integration problems

→ Integrating a set of homogeneous data sources, such as relational databases, data files, and legacy sources, means detecting any relationships between the concepts represented in local schemata and solving the conflicts found to create one single global schema, whose elements can be mapped onto local schema elements.

→ If each data source modeled independently, then there will be no problem while integration process. However, this is not practically doing mostly.

→ The problems that need to be resolved during the integration phase are:

- Different Perspectives
- Equivalent modeling constructs
- Incompatible Specifications
- Common Concepts
- Interrelated Concepts

1. Different perspectives: (Different views may have different data)

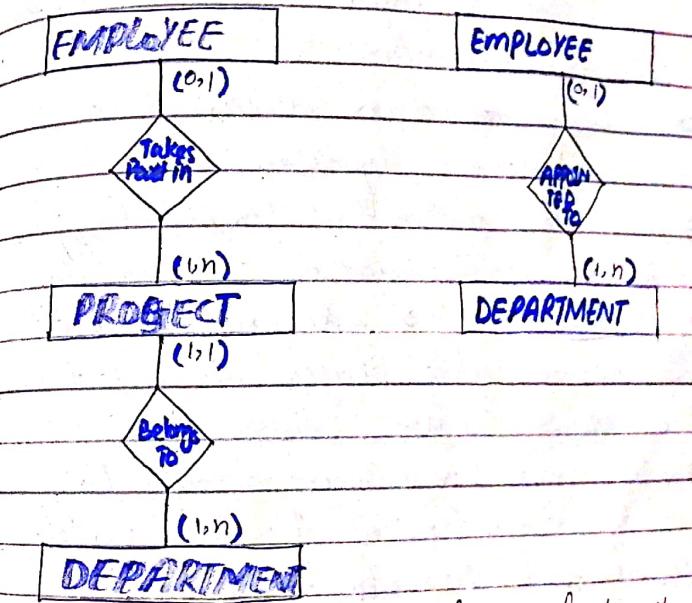
→ User groups may have different expectations for how data should be presented and analyzed.

→ For example, executives might need high-level summaries, while analysts might need detailed data.

→ Depending on the points relevant for tasks users have to perform need to be considered.

Example:

Considered to model the appointment of an employee to a department there two perspectives:



→ The left hand modeling includes the project in which the employee is involved.

→ This modeling could be suitable for

a database used for business organization chart management.

- The right hand modeling does not include this concept because this is not relevant.

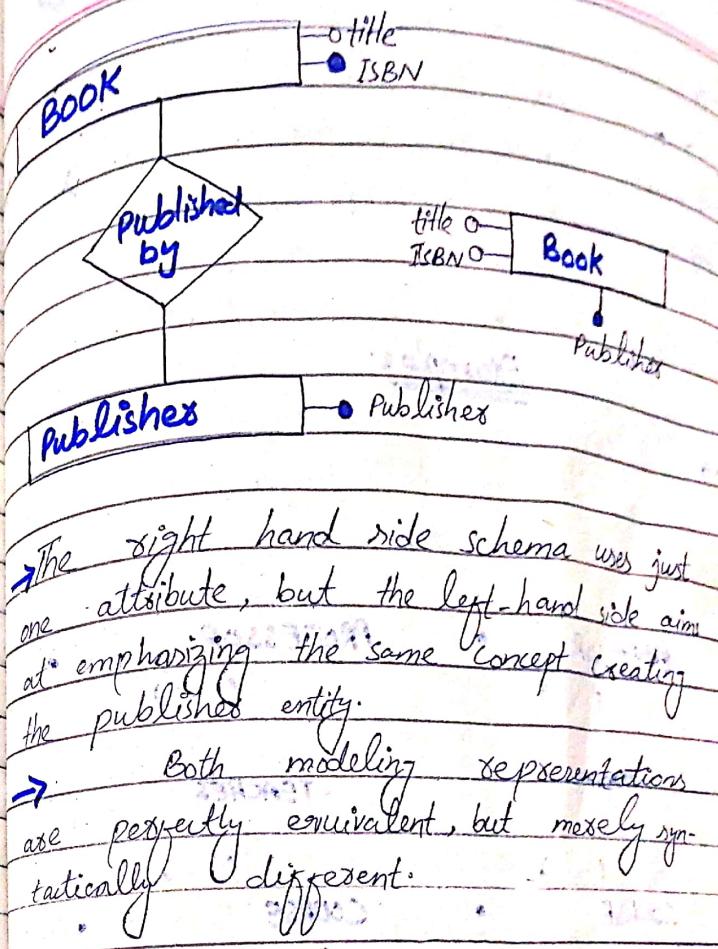
2-Employee+

↑ Different modeling techniques can represent the same information in different ways

2. Equivalent Modeling Constructs

- To represent individual concept, every modeling formalism can use different construct combinations.
- They may ~~be~~ not use identical data structures, terms or definitions.
- This include variations in terminology, data organization, definitions and terminology.
- For example, one system might use "Customer-ID" while another use "client-no" for the same data.
- To solve these issues, need to map and standardize the data models and transform data into a compatible format for cohesive integration.

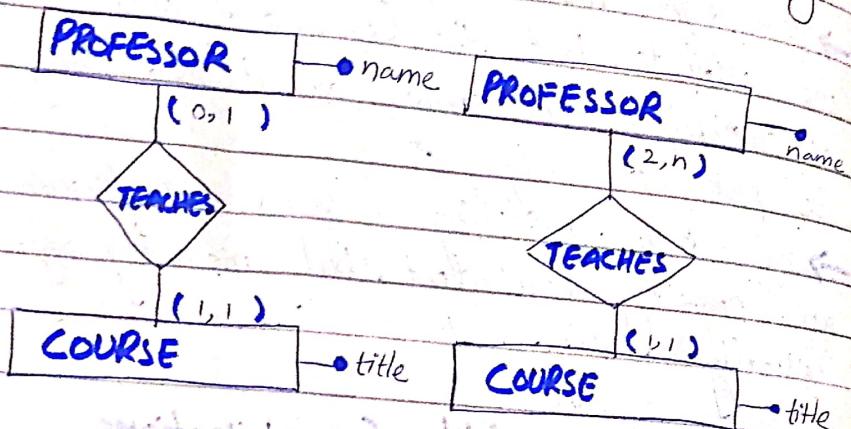
Example:



3. Incompatible Specifications

- Issues arises when integrating data from various sources.
- Incompatible specifications means different schemata, modeling different, contrasting concepts.
- Differences can be due to:

- Incorrect design decisions (different naming conventions, data types)
- Integrity constraints, discrepancies
- If data is timestamped in different time zones, they may reflect different versions and aligning this can be difficult.
- Two unlikely modeling for the association between university professors and courses.



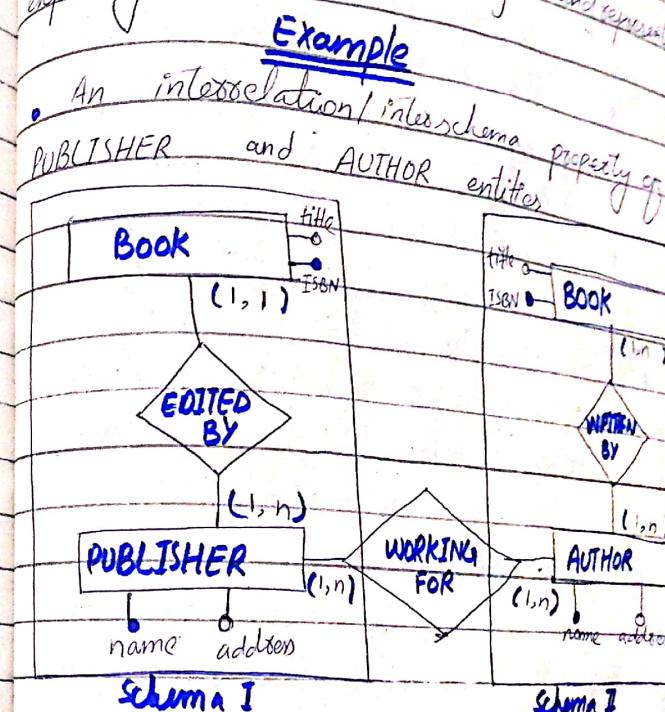
- Both modeling representation could be result of error because:
- Unlikely that professor is not allowed to teach more than one course at the same time (on left side).
 - He is required to teach more than one course (on right side).

4. Interrelated concepts:

Once integration phase is completed, various interrelated concepts will be part of the integrated schema.

New relationships are created without losing existing previous relationships.

These relationships are known as inter-schema properties and identified and expressed explicitly.



- Association between Author and publisher could not be represented in same schema.

5. Common Concepts:

- Semantic relationships exist between common concepts that are modeled in different ways in separate schemata.
- Four types of relationships exist b/w separate representation for the same concepts

identity:

- Two or more data records from different sources represent the same identity.
 - R_1 and R_2 are exactly coincide.
- Equivalence: (Two schemas are equivalent if they represent the same info in different ways)
- Occurs when two or more data elements from different sources contain identical information or values, even if represented in slightly different formats.

For example:

A product price stored as \$100.00 in one system & 100 in another system is equivalent.

⇒ Two schemata R_1 and R_2 are equivalent if their instances can have one-to-one association.

Compatibility: 1) both schemas can be integrated without modification
 When R_1 and R_2 are neither identical nor equivalent but the constraints used and designer's point of views are comparable.

Incompatibility:

When R_1 and R_2 are incompatible because of inconsistent specifications, structural mismatches, semantic mismatches or conflicting values

Two instances for the equivalent schemata:

Book	Publisher	Publisher
ISBN	Title	Publisher

Book	ISBN	Title	Publisher

Conflict:

→ A conflict b/w two representations R_1 and R_2 for the same concept every time both representations are not identical

INTEGRATION Phases

- After integrating local schemata, new features emerged and their representation requires complex set of tasks.
- Different methodological approaches are used for handling those tasks:

1. Pre-integration
2. Schema Comparison
3. Schema Alignment
4. Merging & Resolving Schemata

Preintegration

(This phase involves analyzing the data and identifying parts of the schema that can be integrated & discarded)

- During this Phase:
 - Data sources thoroughly analyzed
 - Define the general integration standards.

→ Decisions taken during this stage are:

- what parts of schemata to be integrated? — bcz all operational data is not useful, therefore, some parts of the schemata are scrapped.
- which integration strategy should be used?
- The integration techniques can be:

- i. Binary Techniques
- ii. N-Ary Techniques

N-ary techniques

more than two schemata integrated at the same time.

- Pros: Reduce the total amount of comparisons b/w concepts because each concept is only processed only once.
- N-Ary technique can be
 1. Single-step
 2. Iterative

1. Single-step:

- All N-schemata integrated at once.
- Complexity is very high.
- Used when have a clear understanding of all schemas and their relationship.



2. Iterative:

- All N-schemata are integrated incrementally not at once.
- And also map schemata gradually

- Lowers complexity
- Testing is done after each iteration.

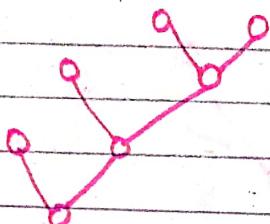
Binary Techniques:

- Always involves pairs of schemata.
- Simpler than N-ary. Proper technique.
- Binary techniques can be:

1. Laddex
2. Balanced

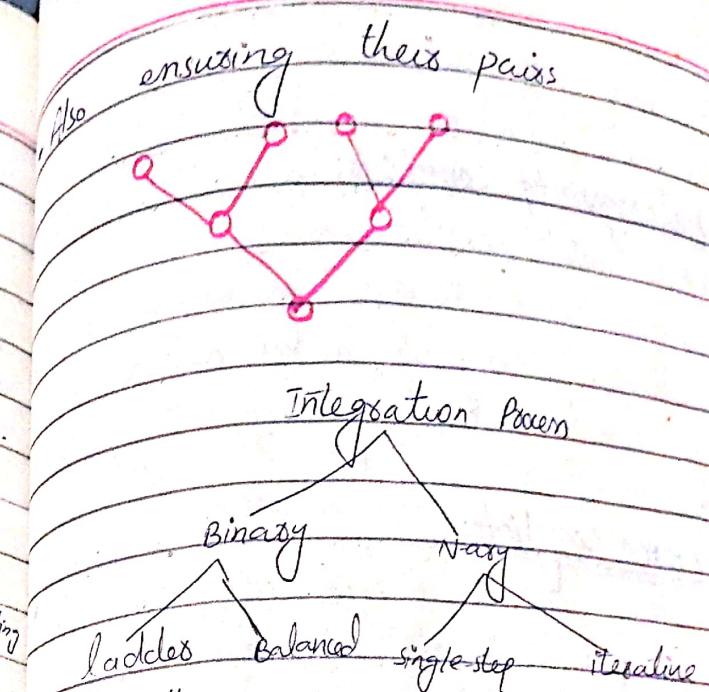
1. Laddex: (Prioritizing data sources and integrating them one-by-one).

- New schemata are integrated into a current schema. Prioritize source schema & order in which they are analyzed.
- Sequential and step-by-step strategy
- New or current schemata add the functionality of previous schemata and testing at each level before moving to the next.



2. Balanced:

multiple schemata integrated at the same time in parallel manner.



Schema Comparison

(compares the schemas of different data sources to identify similarities, differences & conflicts).

→ This Phase consists of analysis of schema. Relationships and conflicts discovered in schemata.

→ Effectiveness of this phase depends on the knowledge of designer that should do comparison b/w different schemata. For this purpose they should collaborate.

→ Conflicts that are detected are:

1. Heterogeneity conflicts
2. Name conflicts

3. Semantic conflicts

→ Structural conflict

1. Heterogeneity conflicts:

- Due to differences in data structures & formats
- Point out inconsistencies due to the differences in the structure of source schemata.
- Heterogeneity conflicts can be reduced by using formal method for modeling the schemata.

2. Name conflicts: → Some data element is represented by different names in different schemas by different designers in various data sources.

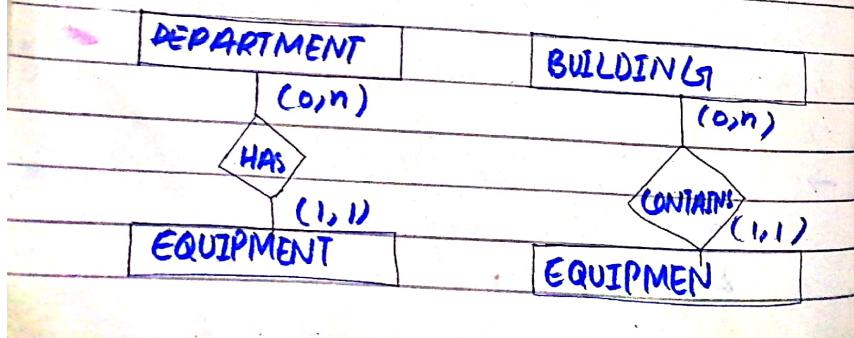
→ Two main name conflicts occurs are:

1. Homonyms
2. Synonyms

→ Same term used for two different concepts

→ Different terms used for same concept

e.g.

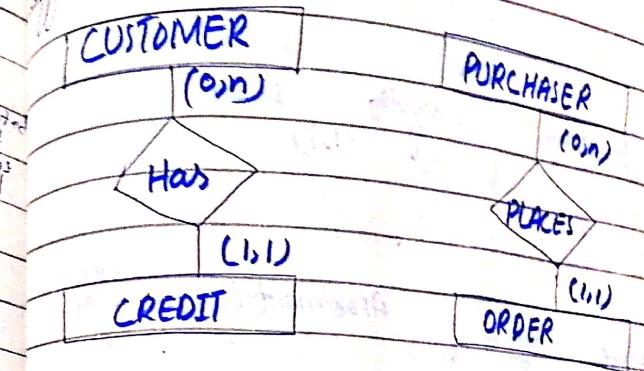


3. Semantic Conflicts

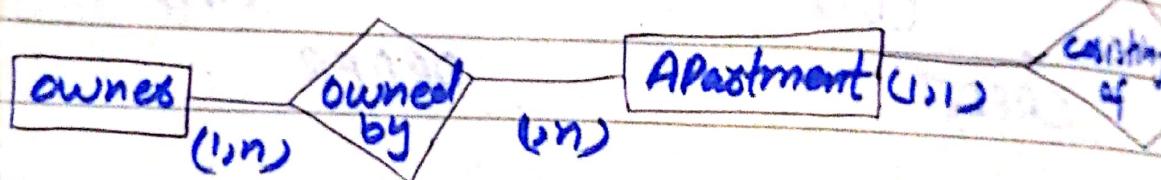
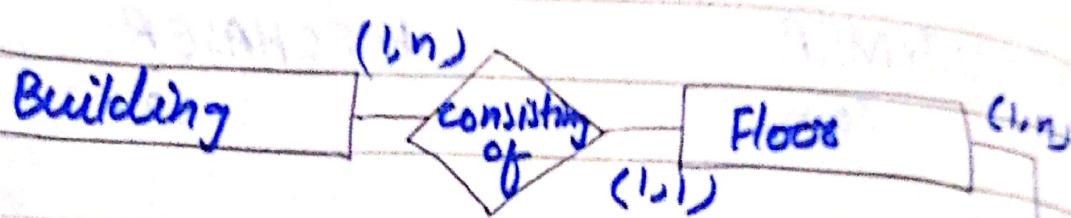
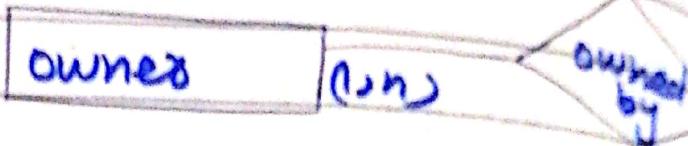
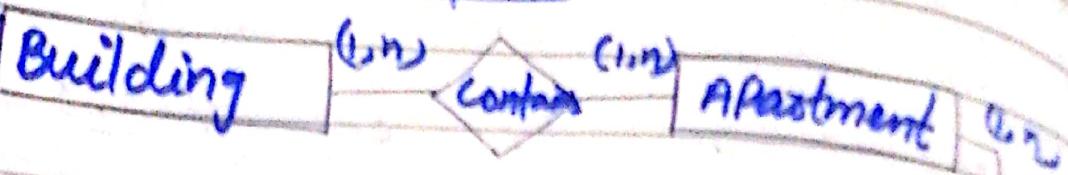
→ Occurs when two source schemata represent similar concepts or entities in different ways & leading to mismatching when integrating. This due to

→ Similar data have different interpretations meanings across different sources

↳ Problem occurs bcz EQUIPMENT hard schema represent scientific equipment and right hand side for building purpose. Reconciled schema should suitable make sense b/w them.



examples



→ Two schemata in semantic conflict. They model a common part of the application domain at different level of abstraction.

4. Structural Conflicts:

- Arise when different options for modeling the same concept or categories into just types:
 1. Type
 2. Dependency
 3. Key
 4. Behavior (on registers)

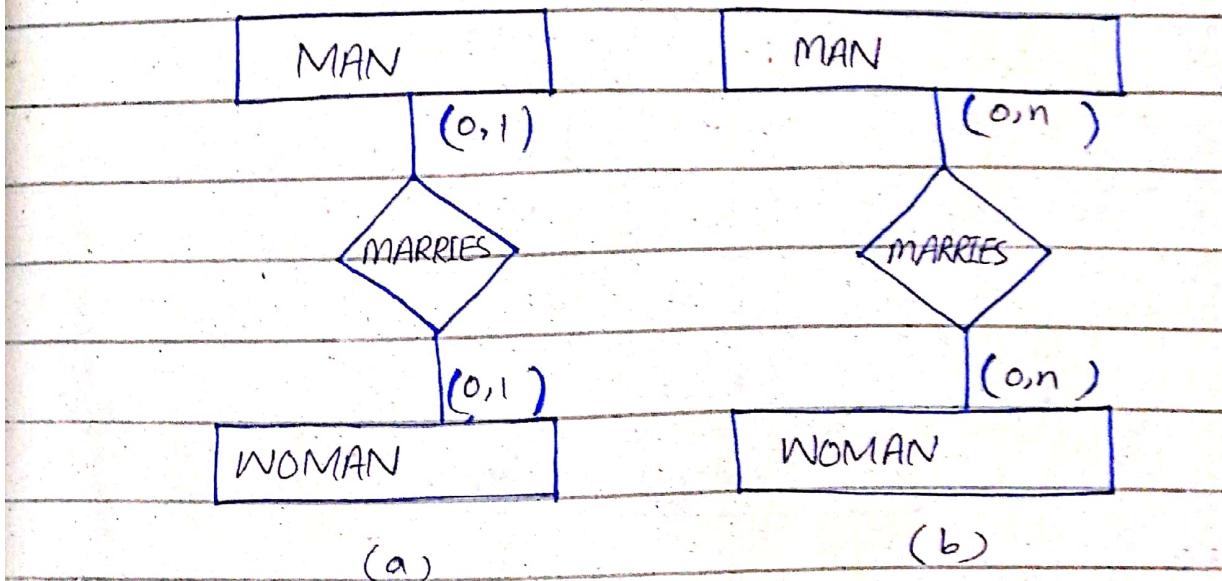
ex: one schema may use VARCHAR for
field while others use INT.

Type Conflict

- o Occurs when two different constructs are used to model the same concept.

Dependency Conflict: → there are contradiction in the s/s b/w data elements in different schemas.

- o Occurs when two or more constructs are interrelated with different dependencies in separate schemata.



(a) Define wedding between man and woman by one-to-one association.

(b) define wedding between man and woman by one-to-many association.

Key conflicts

- Occurs when different attributes in separate schemata are used to identify the same concept.
For example,

⇒ People could be identified on the basis of their Social Security numbers or progressive numbering.

Behavior conflict:

- Occurs when different actions take place to delete/edit data on the same concept in separate schemata.

For example:

Information on a customer is deleted from schemata if the customer has no pending orders, but another schema still stores the customer's information for future use.

3. Schema Alignment

⇒ The purpose of this phase:

- Solve any conflicts that occur during previous step.

⇒ Can be done by applying transformation primitives to source schemata.

Transformation primitives:

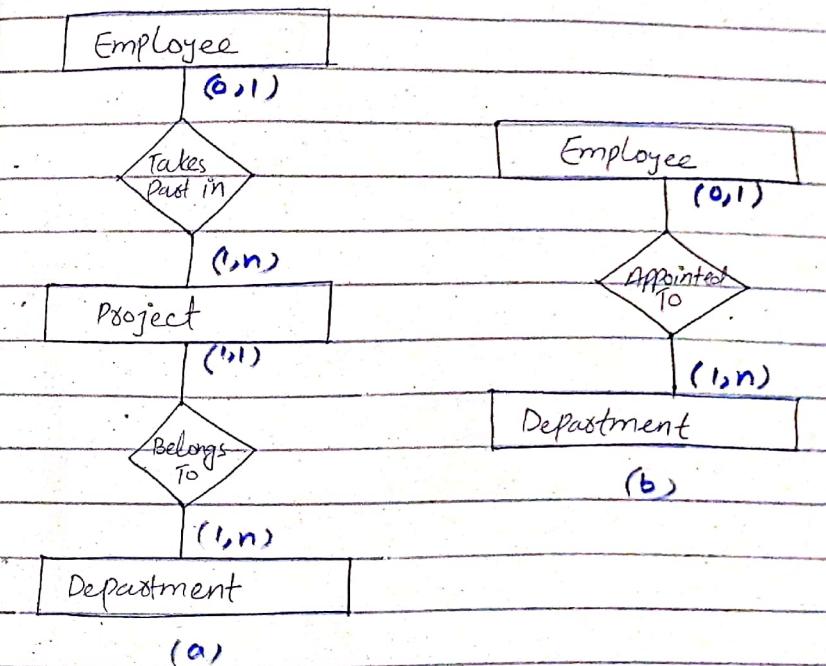
⇒ Typically deal with:

- Changes in names
- Attribute types
- Functional dependencies

- Existing constraints applied to schemata.

→ Conflicts cannot always be solved, because they stem from the basic inconsistencies of the source database.

Example:



→ If semantic conflict occurs then, project concept is introduced on right hand side ~~in~~ schema.

Combining the compatible parts of the schema into unified schema.

4. Merging & Restructuring Schemata

→ During this phase:

- All the aligned schemata merges into a single reconciled schema.
- After that, transformation processes apply to improve the reconciled schema structure.
- For doing this, following properties need to check:

1. Completeness:

- The merged schema should include all relevant attributes and entities from the source systems.
- After overlapping source schemata, additional interschema properties may emerge.
- A complete schema ensures that no data is lost.
- Designers should determine new property that can be added by new associations or generalization hierarchies.

2. Minimality:

(there are no redundant element).

{ Reconciling
to eliminate
redundancy
& improve
minimality)

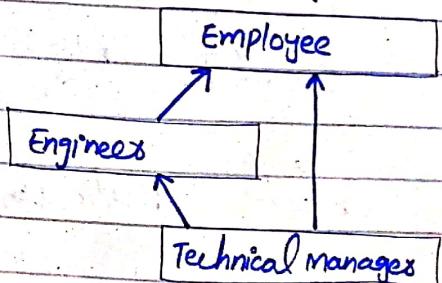
→ The merged schema should avoid redundancy and duplication of data when overlap many schemata.

→ May be mutually derived in reconciled schemata.

for example:

→ Generalization between TECHNICAL MANAGER and EMPLOYEE is redundant and can be deleted.

→ Cyclic relationship between them also cause redundancy.



3. Readability:

(refers to the clarity & understandability of the schema).

- The merged schema should be easy to understand and interpret.
- Use clear and meaningful attribute and entity names.
- Typically, it is difficult to quantify the readability differences b/w two schemas.

Defining Mappings

→ The result of source schema analysis targeted to integration consists of two elements:

1. Reconciled schema where the conflicts between local schemata are solved

2. The set of relationships between the source schema and reconciled schema elements.

→ This mapping process ensures that data is accurately transformed and loaded into data warehouse.

Two approaches that define relationship b/w source system and the target data warehouse:

GAV

- Global as view defines the target data warehouse as a view over the source systems.

- This means that the target schema is defined in terms of SQL queries that extract from the source systems.

Pros:

1. Simplicity

2. Flexibility

↑ where the global schema is defined
as a set of views of the local schema.

Local as View (LAV)

→ LAV defines the target data warehouse house as a collection of materialized views that are populated from the source systems.

→ These materialized views are physical tables that store the transformed and aggregated data.

→ This approach requires very complex transformation processes, generally called ETL.

Example of Mapping

Source: A customer table with fields:
customer-id, first-name, last-name,
address.

Target: A customer dimension table with fields:

customer-key, customer-name,
address-line1, address-line2.

Mapping Rules:

1. Customer-id from source maps to customer-key
2. Concatenate first-name & last-name from source to create customer-name in target.
3. Split address from source into address-line1 & address-line2.