

Ch: 7

~~Data Modeling~~

XML & JSON

used for
data interchange

- ⇒ Extensible Markup language
- ⇒ Tag is label for a section e.g. `<header>`
- ⇒ Element: Data Section beg. `<tag>` end `</tag>`
- ⇒ Attributes: Represented by `"name = value"`

Inside opening tag \Rightarrow an element

Attributes vs Subelements

- ⇒ Subelement: `<tag1> <tag2> ... </tag2> </tag1>`
- ⇒ Attr. are part of document markup
- ⇒ Subelement are part of basic document
- ⇒ Suggestion: Attr \rightarrow id of elements

sub-ele \rightarrow contents

- ⇒ String that used as tags
- ```
<! [CDATA[<account> .. </account>] >
```

Here `<account>` and `</account>` are treated  
as strings

- ⇒ Document Schema

⇒ DTD (Document Type Definition)

⇒ XML Schema (New)

Date:

Day: M T W T F S

## → JSON

- ⇒ Storing number, object, array, ~~boolean~~, null
- ⇒ Can be used by JSON function.

## Ch: 11

### NO SQL

#### Why

- ⇒ Joins are expensive
- ⇒ Hard to scale
- ⇒ Impedance mismatch
- ⇒ Expensive product cost
- ⇒ Speed
- ⇒ Partition
- ⇒ Availability

#### ⇒ Characteristics

- ⇒ Avoids overhead Acid transactions
- ⇒ Complexity of SQL
- ⇒ Schema design
- ⇒ DBA presence
- ⇒ Provides easy and freq changes to DB
- ⇒ Fast dev
- ⇒ Large data (Big data e.g Google)
- ⇒ Schema less

## When/Why?

- ⇒ RDBMS is restrictive (flexible schema) <sup>needed</sup>
- ⇒ ACID not "really" needed
- ⇒ Logging data from distributed systems
- ⇒ Temporary Data (carts, sessions, favorites)
- ⇒ Polyglot Persistence:

Best Data Store depending on data

## Schema less

### ⇒ In RDB

- ⇒ Can't add record with different datatypes, less attr., or multiple ~~fields~~ data.
- ⇒ Should consider Primary key, JOINS etc

### ⇒ In No-SQL

- ⇒ No Schema, no-wiredcell, no-datatype

## App. Data Models

- ⇒ Key-value (redis) ⇒ Document (Mongo DB)
- ⇒ Col. family (Apache) ⇒ Graph (Neo4j)

### ⇒ Key-Value

- ⇒ Easiest, access data using keys (string, hashed)
- ⇒ Operations: Insert, Fetch, Update, delete

### ⇒ Col-family

Cassandra → Facebook Search

⇒ MySQL > 50GB

⇒ Write: ~300ms, Read: ~350ms

⇒ Cassandra > 50GB

⇒ Write: 0.12ms, read: 15ms

## ⇒ Document

⇒ Pair key with complex data structure

⇒ Indexes with B-trees

⇒ Nested documents

## SQL vs NO-SQL (Diff)

⇒ MySQL vs mongo

⇒ Row tables vs key-value, JSON, XML

⇒ Schemas: Static vs Dynamic

⇒ Scaling: vertical, horizontal vs horizontal

⇒ Data Manipulation: Queries vs API's

## CAP theorem:

⇒ Scalability, Consistency, availability, Fault tolerance  
Impossible at a time.

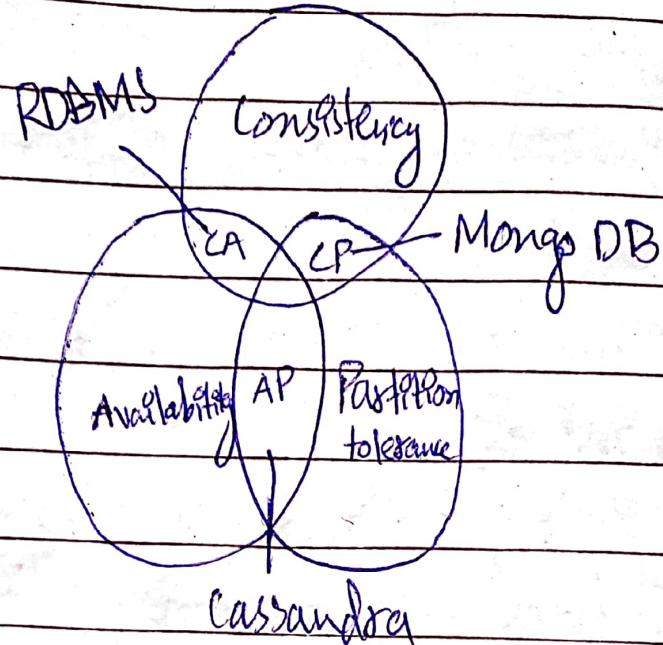
⇒ Impossible for any shared data-system

to guarantee simultaneously all three

⇒ consistency → Availability → Partition tolerance

Date: \_\_\_\_\_

Day: MTWTFSS



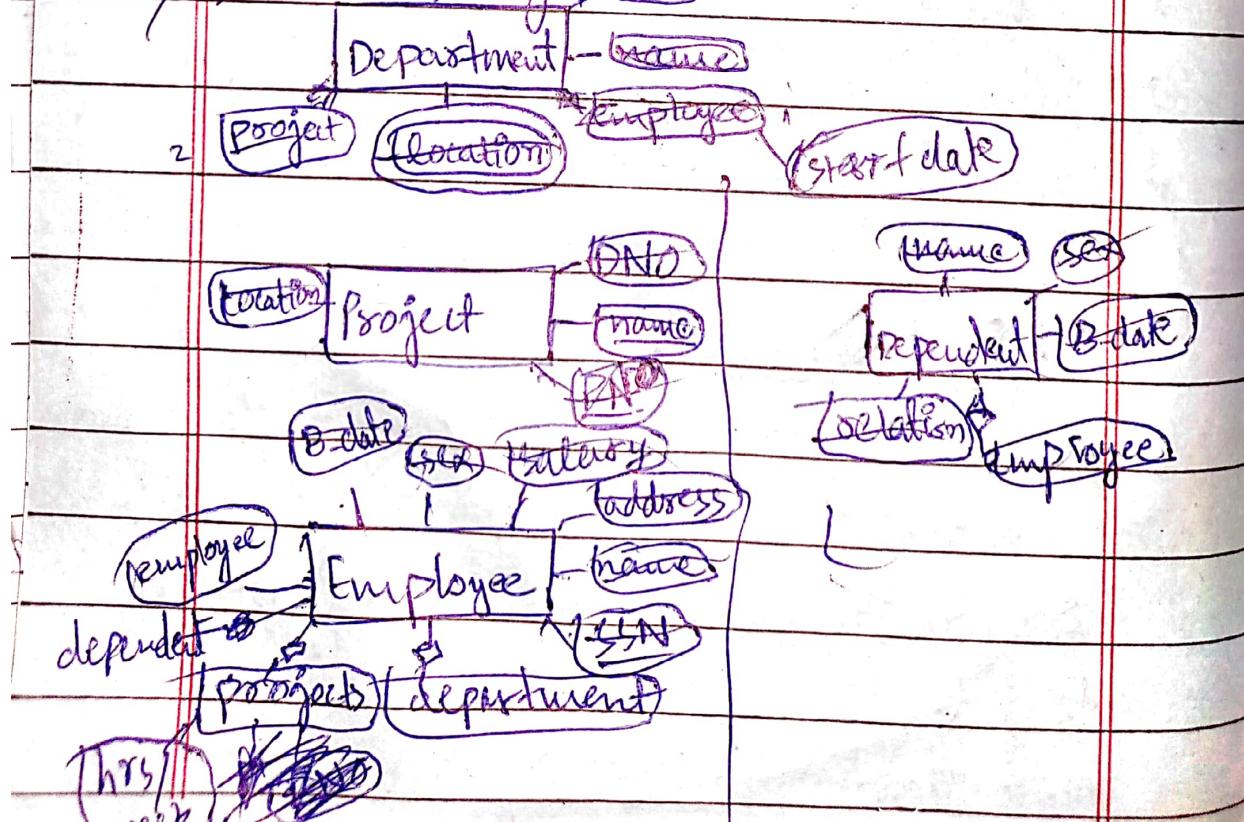
NoSQL : AP

RDBMS : CP

## Ch:7

### Company(ER)

#### 3) Entity Design (EER)



Date: \_\_\_\_\_

Day: M T W T F S

### → Relationship types

✓ → manages emp → dep

✓ → works on emp → proj

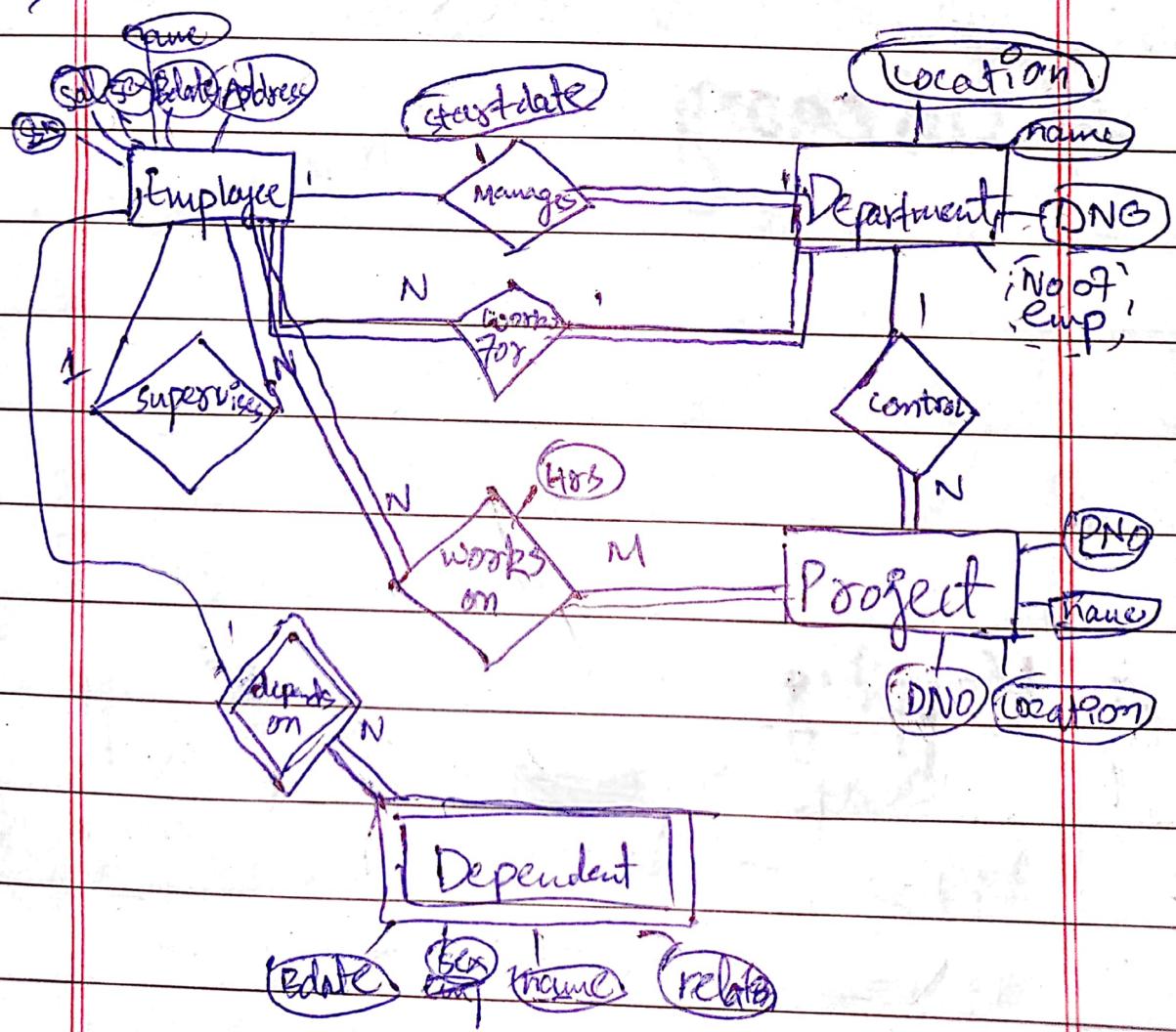
✓ → depends on emp → depend

✓ → works for emp → dep

✓ → supervises emp → emp

✓ → controls dep → proj

### 2) ER



# XML to JSON

Date:

July 2023

Day: M T W T F S

- ⇒ Cardinality (max 1:M)
- ⇒ Modality (zero, one or more) min

## High Degree Relations

⇒ Uni we kon ga

## ER Class

- ⇒ d must use
- ⇒ no or o means optional
- ⇒ refinement ⇒ spec synthesis ⇒ guess

## Category (Union type)

## Mapping

Weak ⇒ ~~1 PK~~ include ~~1:1~~ ~~1:N~~ of strong

1:1 ⇒ -PK → 1:1

1:N ⇒ 1 PK → NFK

M:N ⇒ Make new table and  
include both PR

Multivalued ⇒ New table include PK  
and attr.

## Ch 15

### Normalization

#### ⇒ Informal Design

##### ⇒ Semantics

⇒ Bottom up ⇒ adding attr to form relations.

⇒ Topdown ⇒ starts grouping attr to form relations

##### ⇒ Guidelines

- semantics
- Redundancy
- null
- Spurious

#### ⇒ 1<sup>st</sup> Guide:

⇒ each tuple should represent one entity or relationship

⇒ Attr. of diff relation shouldn't mix.

⇒ Only FK for referencing

⇒ Entity and relationship should be kept apart.

#### ⇒ 2<sup>nd</sup> Guide

⇒ Mixing attr. that depends of others

⇒ Anomalies (Insert, delete, update)

⇒ Use views and update them by triggers.

## 3rd Guide:

- ⇒ Separate table with PK
- ⇒ Reason (not applicable, invalid, unknown, unavailable)

## 4th Guide

⇒ Bad Design may result in erroneous results for JOIN, use lossless JOIN

⇒ Relations should be designed for spurious lossless JOIN and doesn't spurious tuples

⇒ There are 2 imp decompr. prop.

- lossless or non-additive
- Preservation of FD

⇒ FD: constraint b/w 2 attrs interrel.

⇒ Used to define NF

⇒ A set of attrs X functionally determines a set of Y if X determines unique Y

⇒  $X \rightarrow Y$ , 2 tuples same X must have same Y

## 5) Inference Rules (Armstrong)

⇒ IR<sub>1</sub> (reflexive): Y subset of X :  $X \rightarrow Y$

⇒ IR<sub>2</sub> (Augmentation): If  $X \rightarrow Y$  then  $XZ \rightarrow YZ$

⇒ IR<sub>3</sub> (Transitive): If  $X \rightarrow Y$  :  $X \rightarrow Z$ ,  $Y \rightarrow Z$  then

## Additional

⇒ Decomposition:  $X \rightarrow YZ$ :  $X \rightarrow Y, X \rightarrow Z$

⇒ Union:  $X \rightarrow Y$  and  $X \rightarrow Z$  then  $X \rightarrow YZ$

⇒ Pseudotransitivity:  $X \rightarrow Y, Y \rightarrow Z$  then  $X \rightarrow Z$

## ⇒ Normal Forms

⇒ Normalization: process of decomp-  
bad relations by breaking into multiple.

⇒ 1, 2, 3 FDs, 4 MVD, 5, JD

⇒ Denormalization (Ultra Krd)

⇒ Prime attr are part of <sup>candidate</sup> composite key.

## ⇒ 1<sup>st</sup> NF:

⇒ Atomic values

⇒ Composite, multivalued, nested rel

⇒ Make multiple tuples

## ⇒ 2<sup>nd</sup> NF

⇒ Nonprimes are fully FD on PK <sup>A in R</sup>