

Database Management Systems



Syllabus

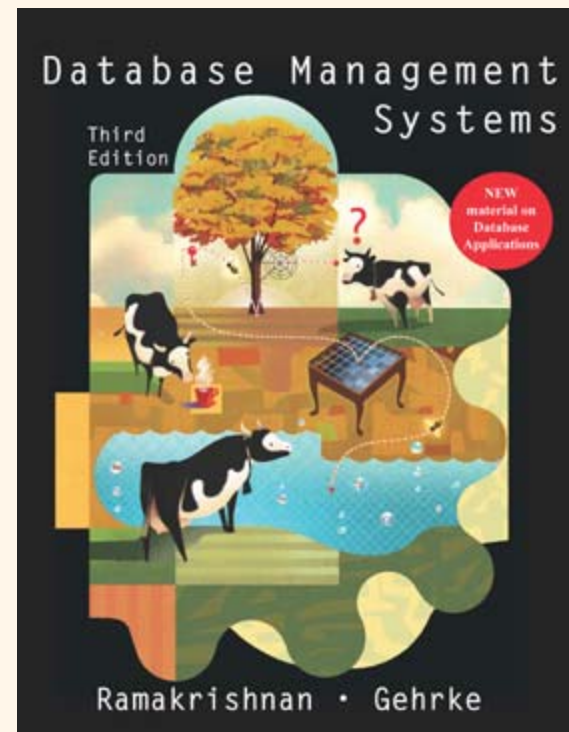
Instructor: Vinnie Costa
vcosta@optonline.net

Course Description

- ❖ This course is designed to provide individuals with an introduction to database concepts and the relational database model. Topics include SQL, normalization, design methodology, DBMS functions, database administration, and other database management approaches, such as client/server databases, object oriented databases, and data warehouses. At the completion of this course, students should be able to understand a user's database requirements and translate those into a valid database design. The emphasis will be on application development rather than system fundamentals.
- ❖ *Prerequisites:* None

Text

- ❖ *Required:*
Raghu Ramakrishnan
and Johannes Gehrke,
*Database Management
Systems, 3/e*, McGraw-
Hill Higher Education,
2003, 1065pp., ISBN 0-
07-246563-8



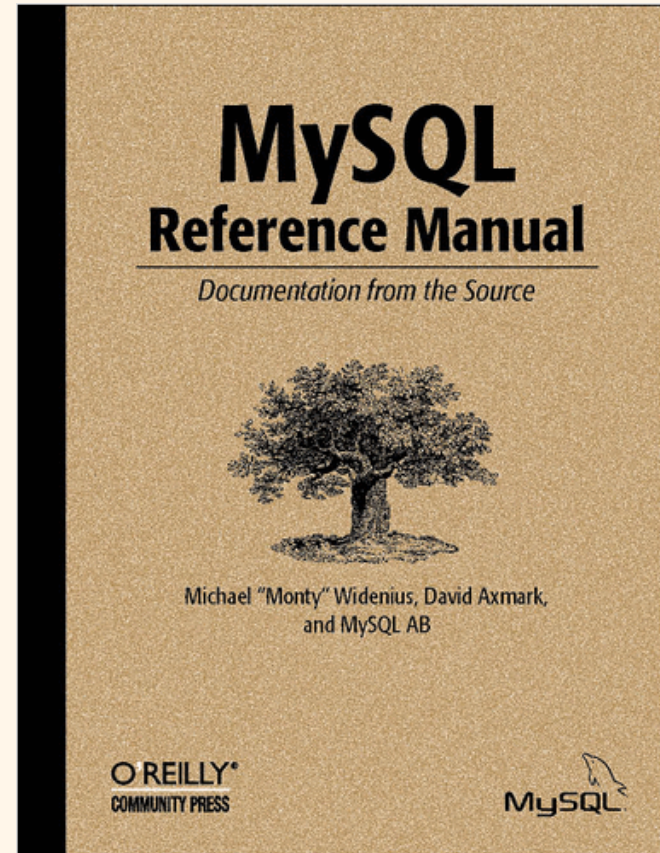
Text

❖ *Reference:*

*Rasmus Lerdorf and Kevin Tatroe,
Programming PHP, O'Reilly & Associates,
Inc., 2002*

*Michael “Monty” Widenius, David Axmark,
and MySQL AB, MySQL Reference Manual,
O'Reilly & Associates, Inc., 2002*

Text



Grading

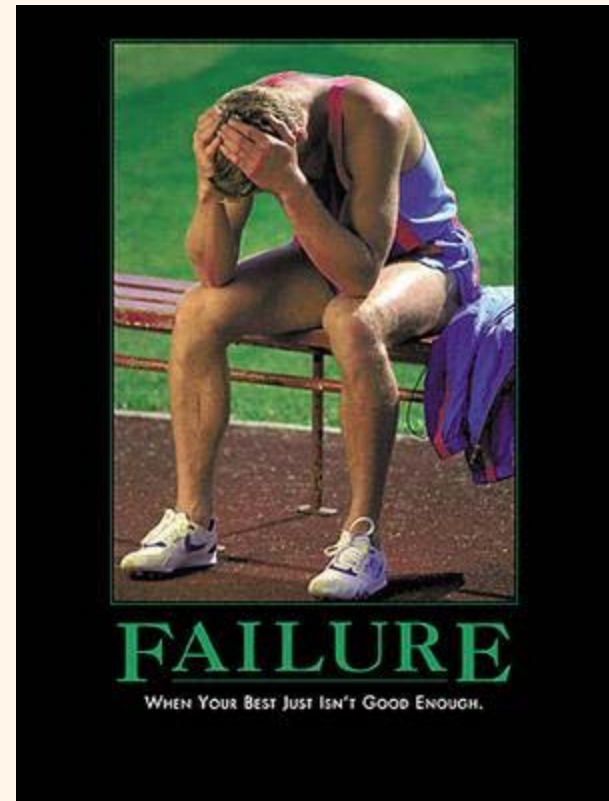
- ❖ Several **assignments**, three count
- ❖ **mid-term** and **end-term**
- ❖ **Class participation**
- ❖ **Final project** or **paper**
- ❖ **No make-up** tests or **extended deadlines**

Point Allocation

- ❖ Assignments 1-3: 5% each
- ❖ Final Project: 30%
- ❖ Mid-Term: 25%
- ❖ End-Term: 25%
- ❖ Participation: 5%

Attendance

- ❖ **Not Mandatory, but...**
- ❖ **...you'll probably fail!**
- ❖ **Participation** is very important
- ❖ **Let me know** if you can't make it



Course Outline

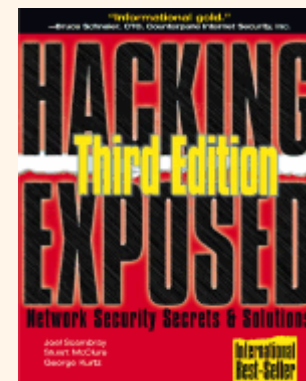
Session	Date	Topic	Comments
1	08/20/05	Overview & Intro to DBMS	
2	08/27/05	Relational Model	Paper Assignment
3	09/10/05	Relational Algebra & Calculus	
4	09/11/05	SQL: Queries, Constraints, Triggers	Mid-term Handout
5	09/17/05	Database Application Development	<i>Mid-term Due</i>
6	09/24/05	Database Internet Applications	
7	10/01/05	Database Internet Applications	
8	10/02/05	Systems Basics: Storage, Transactions	
9	10/08/05	Schema Refinement, Normalization	End-term Handout; <i>Paper Due</i>
10	10/15/05	XML Data Management	<i>End-term Due</i>

Slides, Links & News

- ❖ <http://www.cs.hofstra.edu/~cscvjc/Fall05>
- ❖ Check frequently!!!
- ❖ E-Mail
- ❖ Jabber – vcosta@jabber.org

Class Rules

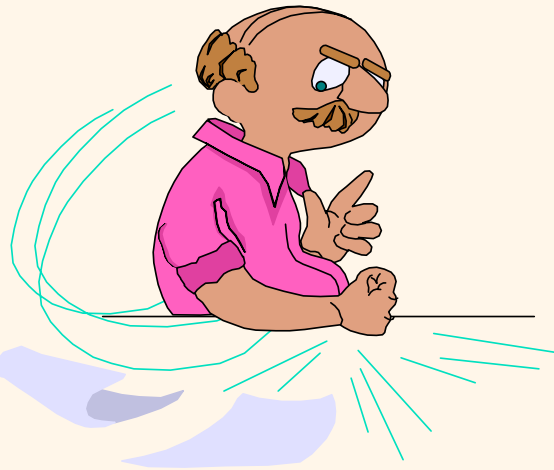
- ❖ Assignments are to be completed individually
- ❖ Academic **honesty** taken seriously
- ❖ *Any attempt to gain unauthorized access to any system will be dealt with harshly*



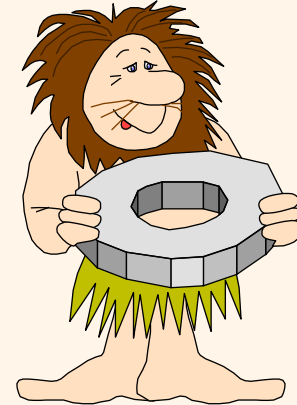
Database Management Systems

Chapter 1

Instructor: Vinnie Costa
vcosta@optonline.net



What Is a DBMS?



- ❖ A very large, integrated collection of data.
- ❖ Models real-world enterprise.
 - Entities (e.g., students, courses)
 - Relationships (e.g., Madonna is taking CS564)
- ❖ A Database Management System (DBMS) is a software package designed to store and manage databases.

Files vs. DBMS

- ❖ Application must stage large datasets between main memory and secondary storage (e.g., buffering, page-oriented access, 32-bit addressing, etc.)
- ❖ Special code for different queries
- ❖ Must protect data from inconsistency due to multiple concurrent users
- ❖ Crash recovery
- ❖ Security and access control

Why Use a DBMS?



- ❖ Data independence and efficient access.
- ❖ Reduced application development time.
- ❖ Data integrity and security.
- ❖ Uniform data administration.
- ❖ Concurrent access, recovery from crashes.

Why Study Databases??



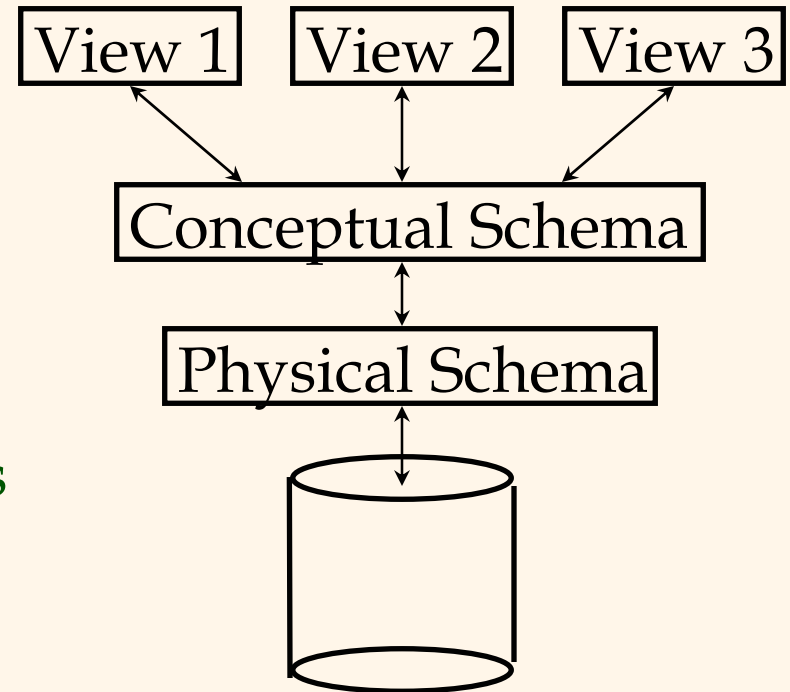
- ❖ Shift from computation to information
 - at the “low end”: scramble to webspace (a mess!)
 - at the “high end”: scientific applications
- ❖ Datasets increasing in diversity and volume.
 - Digital libraries, interactive video, Human Genome project, EOS project
 - ... need for DBMS exploding
- ❖ DBMS encompasses most of CS
 - OS, languages, theory, AI, multimedia, logic

Data Models

- ❖ A data model is a collection of concepts for describing data.
- ❖ A schema is a description of a particular collection of data, using the a given data model.
- ❖ The relational model of data is the most widely used model today.
 - Main concept: relation, basically a table with rows and columns.
 - Every relation has a schema, which describes the columns, or fields.

Levels of Abstraction

- ❖ Many views, single conceptual (logical) schema and physical schema.
 - Views describe how users see the data.
 - Conceptual schema defines logical structure
 - Physical schema describes the files and indexes used.



* *Schemas are defined using DDL; data is modified/queried using DML.*

Example: University Database

❖ Conceptual schema:

- *Students(sid: string, name: string, login: string, age: integer, gpa:real)*
- *Courses(cid: string, cname:string, credits:integer)*
- *Enrolled(sid:string, cid:string, grade:string)*

❖ Physical schema:

- Relations stored as unordered files.
- Index on first column of Students.

❖ External Schema (View):

- *Course_info(cid:string,enrollment:integer)*

*Data Independence **

- ❖ Applications insulated from how data is structured and stored.
- ❖ *Logical data independence*: Protection from changes in *logical* structure of data.
- ❖ *Physical data independence*: Protection from changes in *physical* structure of data.

** One of the most important benefits of using a DBMS!*

Concurrency Control

- ❖ Concurrent execution of user programs is essential for good DBMS performance.
 - Because disk accesses are frequent, and relatively slow, it is important to keep the cpu humming by working on several user programs concurrently.
- ❖ Interleaving actions of different user programs can lead to inconsistency: e.g., check is cleared while account balance is being computed.
- ❖ DBMS ensures such problems don't arise: users can pretend they are using a single-user system.

Transaction: An Execution of a DB Program

- ❖ Key concept is transaction, which is an *atomic* sequence of database actions (reads/writes).
- ❖ Each transaction, executed completely, must leave the DB in a consistent state if DB is consistent when the transaction begins.
 - Users can specify some simple integrity constraints on the data, and the DBMS will enforce these constraints.
 - Beyond this, the DBMS does not really understand the semantics of the data. (e.g., it does not understand how the interest on a bank account is computed).
 - Thus, ensuring that a transaction (run alone) preserves consistency is ultimately the *user's* responsibility!

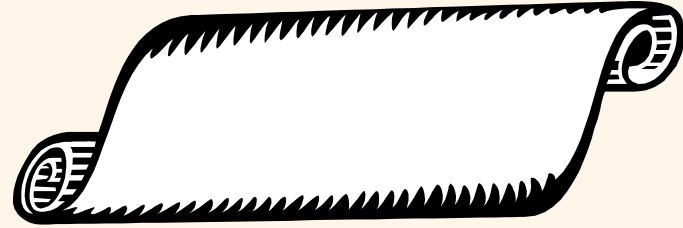
Scheduling Concurrent Transactions

- ❖ DBMS ensures that execution of $\{T_1, \dots, T_n\}$ is equivalent to some serial execution $T_1' \dots T_n'$.
 - Before reading/writing an object, a transaction requests a lock on the object, and waits till the DBMS gives it the lock. All locks are released at the end of the transaction. (Strict 2PL locking protocol.)
 - **Idea:** If an action of T_i (say, writing X) affects T_j (which perhaps reads X), one of them, say T_i , will obtain the lock on X first and T_j is forced to wait until T_i completes; this effectively orders the transactions.
 - What if T_j already has a lock on Y and T_i later requests a lock on Y ? (Deadlock!) T_i or T_j is aborted and restarted!

Ensuring Atomicity

- ❖ DBMS ensures *atomicity* (all-or-nothing property) even if system crashes in the middle of a Xact.
- ❖ **Idea:** Keep a log (history) of all actions carried out by the DBMS while executing a set of Xacts:
 - **Before** a change is made to the database, the corresponding log entry is forced to a safe location. (WAL protocol; OS support for this is often inadequate.)
 - After a crash, the effects of partially executed transactions are undone using the log. (Thanks to WAL, if log entry wasn't saved before the crash, corresponding change was not applied to database!)

The Log

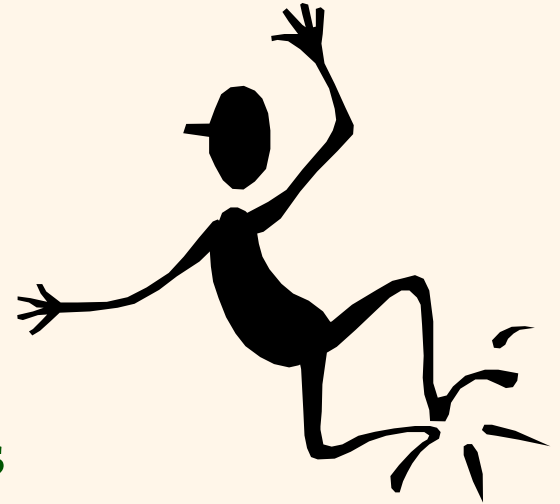


- ❖ The following actions are recorded in the log:
 - *Ti writes an object*: The old value and the new value.
 - Log record must go to disk before the changed page!
 - *Ti commits/aborts*: A log record indicating this action.
- ❖ Log records chained together by Xact id, so it's easy to undo a specific Xact (e.g., to resolve a deadlock).
- ❖ Log is often *duplexed* and *archived* on “stable” storage.
- ❖ All log related activities (and in fact, all CC related activities such as lock/unlock, dealing with deadlocks etc.) are handled transparently by the DBMS.

Databases make these folks happy ...

- ❖ End users and DBMS vendors
- ❖ DB application programmers
 - E.g., smart webmasters
- ❖ *Database administrator (DBA)*
 - Designs logical / physical schemas
 - Handles security and authorization
 - Data availability, crash recovery
 - Database tuning as needs evolve

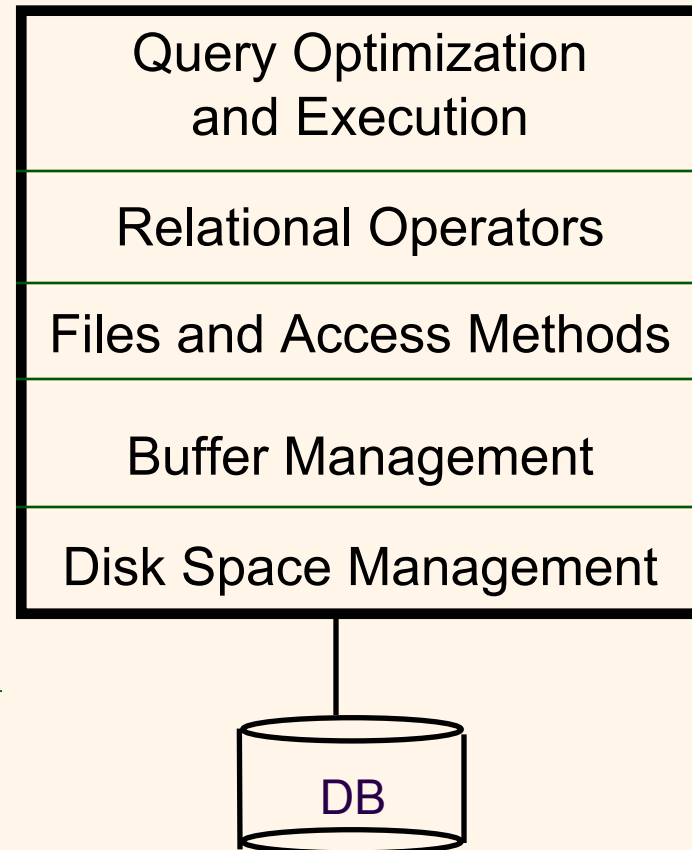
Must understand how a DBMS works!



Structure of a DBMS

These layers must consider concurrency control and recovery

- ❖ A typical DBMS has a layered architecture.
- ❖ The figure does not show the concurrency control and recovery components.
- ❖ This is one of several possible architectures; each system has its own variations.



Structure of a DBMS

- ❖ p20, Figure 1.3 – detailed diagram
- ❖ n-tiered architecture
- ❖ Virtualization, GRIDs
- ❖ Proprietary vs Open
- ❖ Licensing Costs

Summary

- ❖ DBMS used to maintain, query large datasets.
- ❖ Benefits include recovery from system crashes, concurrent access, quick application development, data integrity and security.
- ❖ Levels of abstraction give data independence.
- ❖ A DBMS typically has a layered architecture.
- ❖ DBAs hold responsible jobs and are **well-paid!** 😊
- ❖ DBMS R&D is one of the broadest, most exciting areas in CS.



Useful Websites



<http://www.oracle.com/>



<http://www.mysql.com/>

❖ <http://www.cs.wisc.edu/~dbbook/>

Beyond Relational Databases

- ❖ <http://www.acmqueue.org/modules.php?name=Content&pa=showpage&pid=299>
- ❖ *Margo Seltzer, SleepyCat*
- ❖ *ACM Queue vol. 3, no. 3 - April 2005*

Homework

- ❖ Read Chapter One
- ❖ Exercises pp.23-24: 1.1, 1.4, 1.6, 1.9
- ❖ Read Beyond Relational Databases

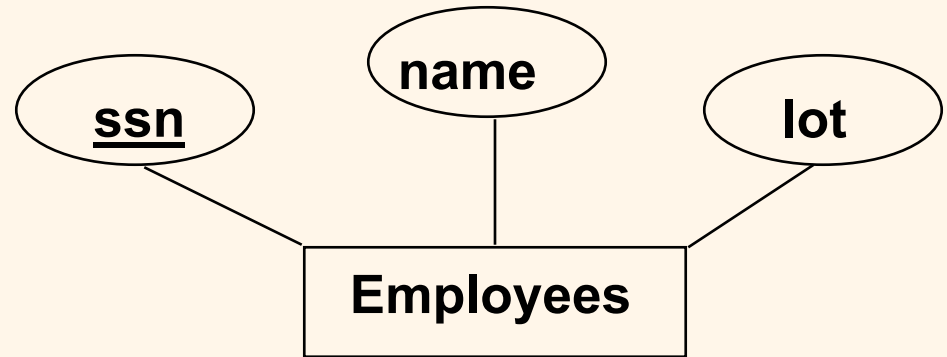
The Entity-Relationship Model

Chapter 2

Overview of Database Design

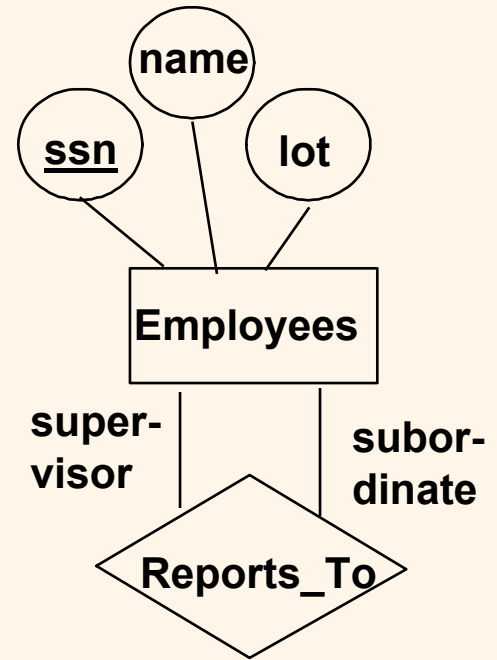
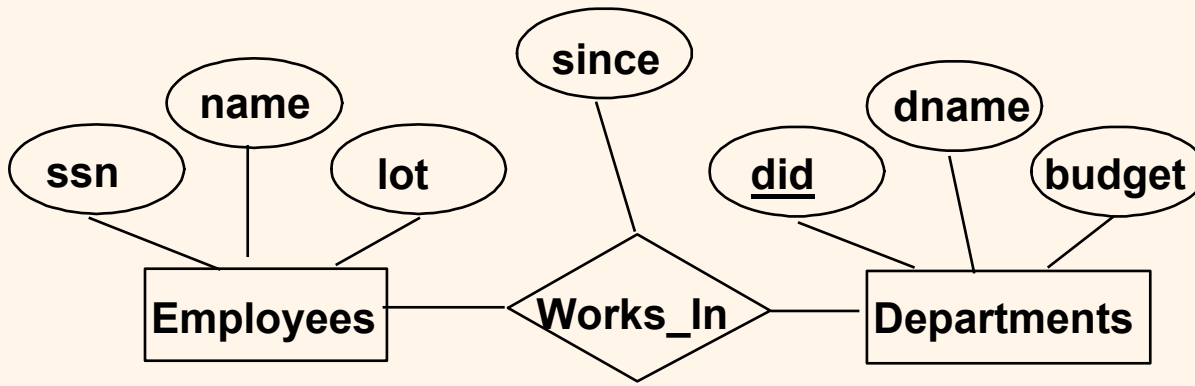
- ❖ Conceptual design: (*ER Model is used at this stage.*)
 - What are the *entities* and *relationships* in the enterprise?
 - What information about these entities and relationships should we store in the database?
 - What are the *integrity constraints* or *business rules* that hold?
 - A database 'schema' in the ER Model can be represented pictorially (*ER diagrams*).
 - Can map an ER diagram into a relational schema.

ER Model Basics



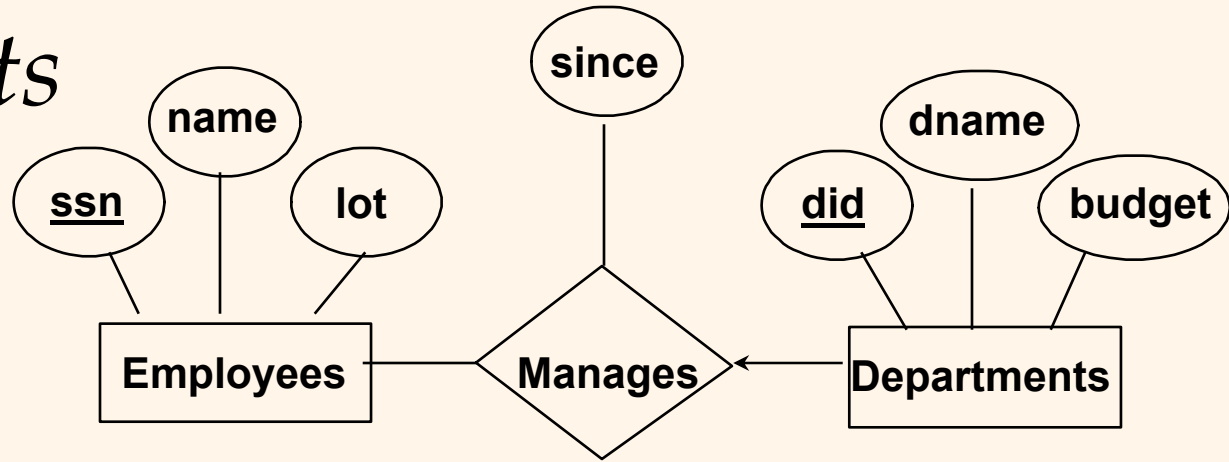
- ❖ Entity: Real-world object distinguishable from other objects. An entity is described (in DB) using a set of attributes.
- ❖ Entity Set: A collection of similar entities. E.g., all employees.
 - All entities in an entity set have the same set of attributes. (Until we consider ISA hierarchies, anyway!)
 - Each entity set has a *key*.
 - Each attribute has a *domain*.

ER Model Basics (Contd.)



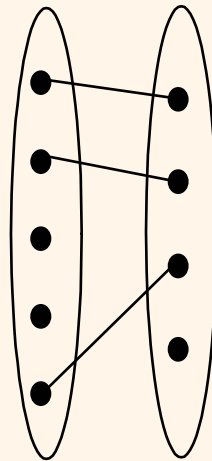
- ❖ **Relationship**: Association among two or more entities. E.g., Attishoo works in Pharmacy department.
- ❖ **Relationship Set**: Collection of similar relationships.
 - An n-ary relationship set R relates n entity sets $E_1 \dots E_n$; each relationship in R involves entities $e_1 \in E_1, \dots, e_n \in E_n$
 - Same entity set could participate in different relationship sets, or in different “roles” in same set.

Key Constraints

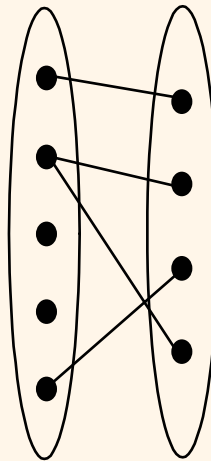


❖ Consider Works_In:
An employee can work in many departments; a dept can have many employees.

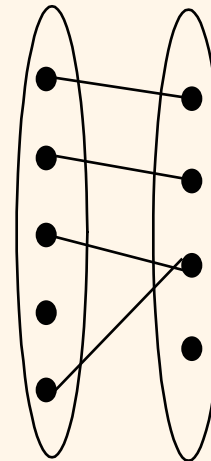
❖ In contrast, each dept has at most one manager, according to the key constraint on Manages.



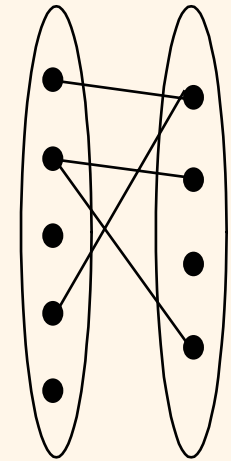
1-to-1



1-to Many



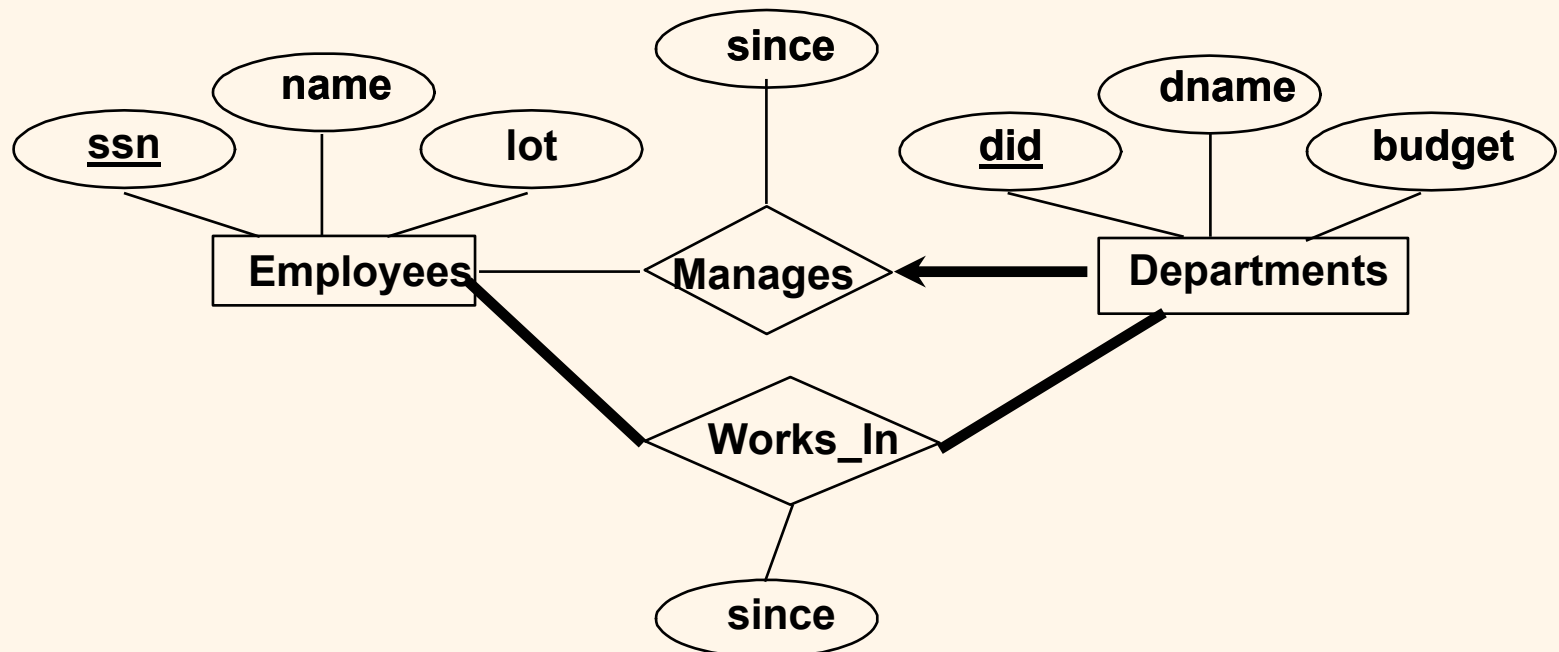
Many-to-1



Many-to-Many

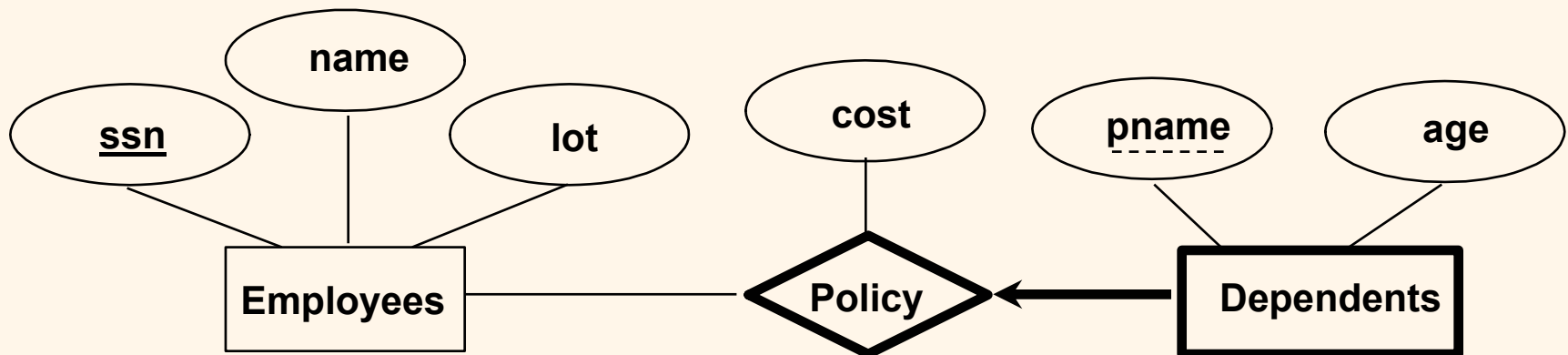
Participation Constraints

- ❖ Does every department have a manager?
 - If so, this is a *participation constraint*: the participation of Departments in Manages is said to be *total* (vs. *partial*).
 - Every Departments entity must appear in an instance of the Manages relationship.



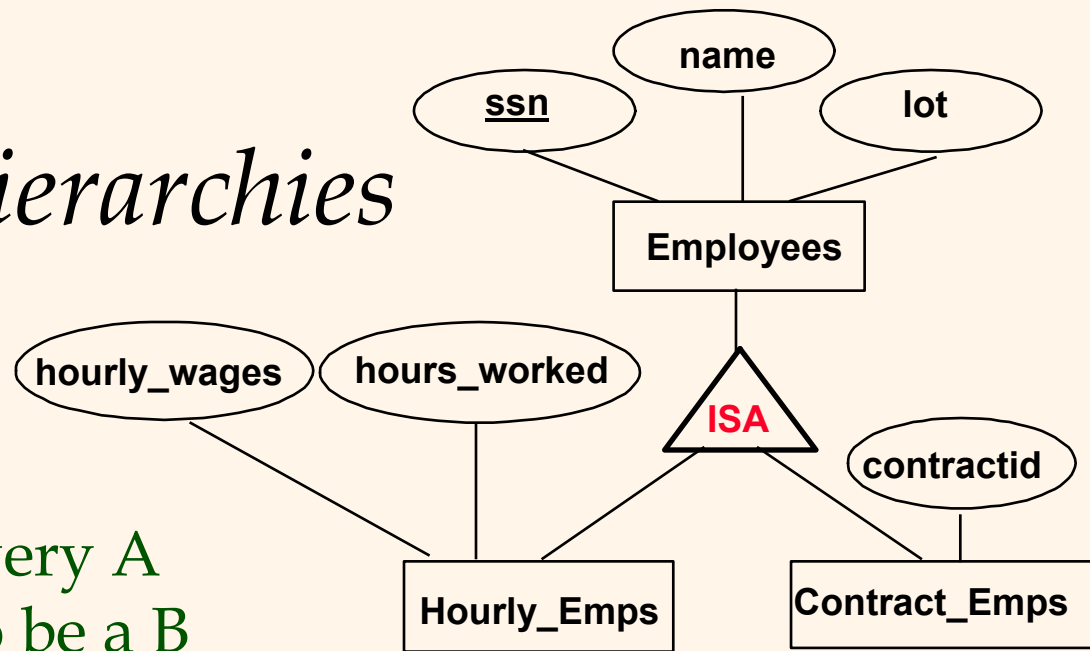
Weak Entities

- ❖ A *weak entity* can be identified uniquely only by considering the primary key of another (*owner*) entity.
 - Owner entity set and weak entity set must participate in a one-to-many relationship set (one owner, many weak entities).
 - Weak entity set must have total participation in this *identifying* relationship set.



ISA ('is a') Hierarchies

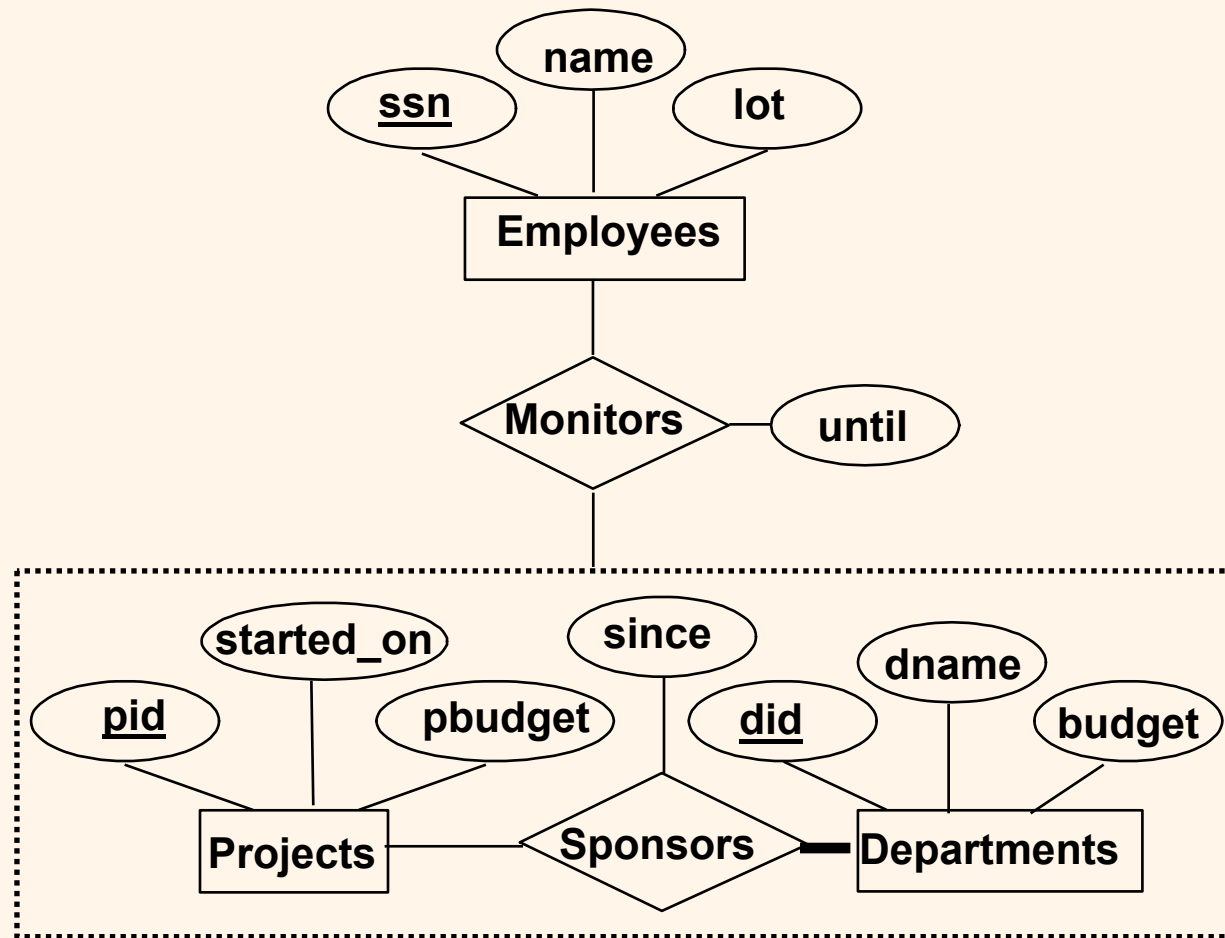
- ❖ As in C++, or other PLs, attributes are inherited.
- ❖ If we declare A **ISA** B, every A entity is also considered to be a B entity.



- ❖ *Overlap constraints*: Can Joe be an Hourly_Emps as well as a Contract_Emps entity? (*Allowed/disallowed*)
- ❖ *Covering constraints*: Does every Employees entity also have to be an Hourly_Emps or a Contract_Emps entity? (*Yes/no*)
- ❖ Reasons for using ISA:
 - To add descriptive attributes specific to a subclass.
 - To identify entities that participate in a relationship.

Aggregation

- ❖ Used when we have to model a relationship involving (entity sets and) a *relationship set*.
 - Aggregation allows us to treat a relationship set as an entity set for purposes of participation in (other) relationships.



- * *Aggregation vs. ternary relationship:*
 - ❖ **Monitors** is a distinct relationship, with a descriptive attribute.
 - ❖ Also, can say that each sponsorship is monitored by at most one employee.

Conceptual Design Using the ER Model

❖ Design choices:

- Should a concept be modeled as an entity or an attribute?
- Should a concept be modeled as an entity or a relationship?
- Identifying relationships: Binary or ternary?
Aggregation?

❖ Constraints in the ER Model:

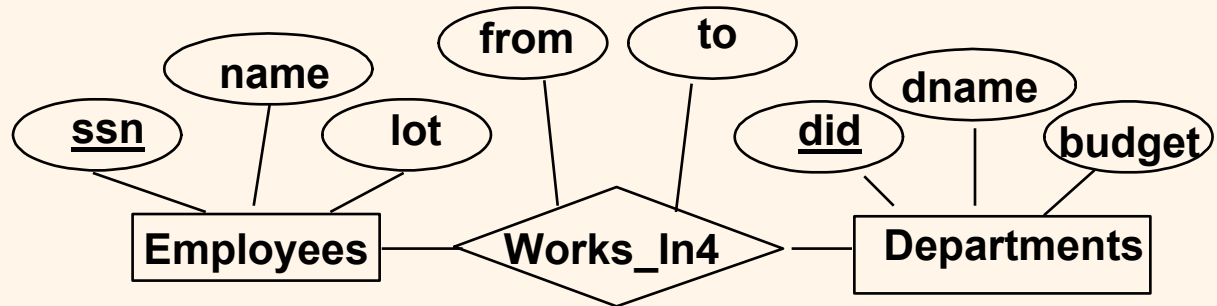
- A lot of data semantics can (and should) be captured.
- But some constraints cannot be captured in ER diagrams.

Entity vs. Attribute

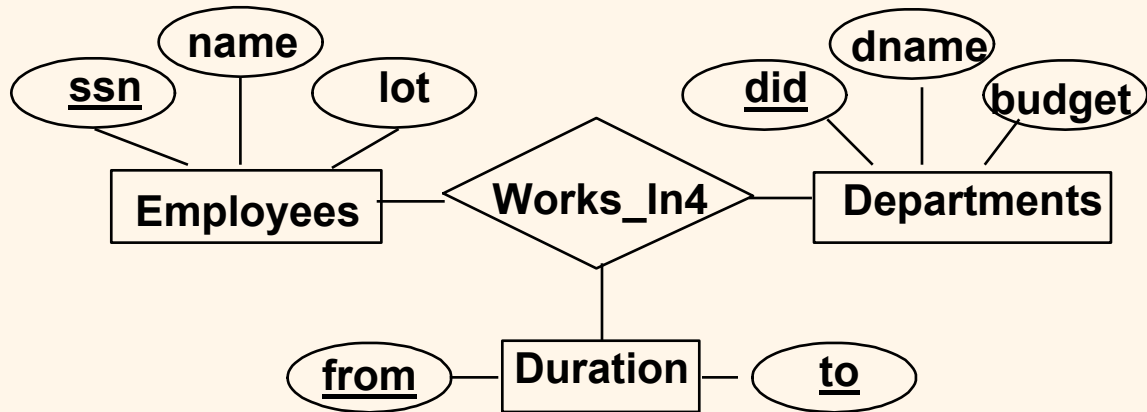
- ❖ Should *address* be an attribute of Employees or an entity (connected to Employees by a relationship)?
- ❖ Depends upon the use we want to make of address information, and the semantics of the data:
 - If we have several addresses per employee, *address* must be an entity (since attributes cannot be set-valued).
 - If the structure (city, street, etc.) is important, e.g., we want to retrieve employees in a given city, *address* must be modeled as an entity (since attribute values are atomic).

Entity vs. Attribute (Contd.)

❖ Works_In4 does not allow an employee to work in a department for two or more periods.

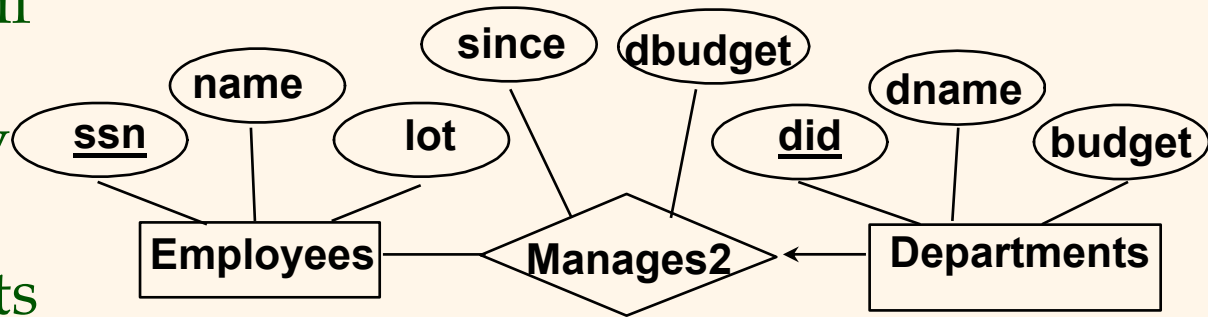


❖ Similar to the problem of wanting to record several addresses for an employee: We want to record *several values of the descriptive attributes for each instance of this relationship*. Accomplished by introducing new entity set, Duration.



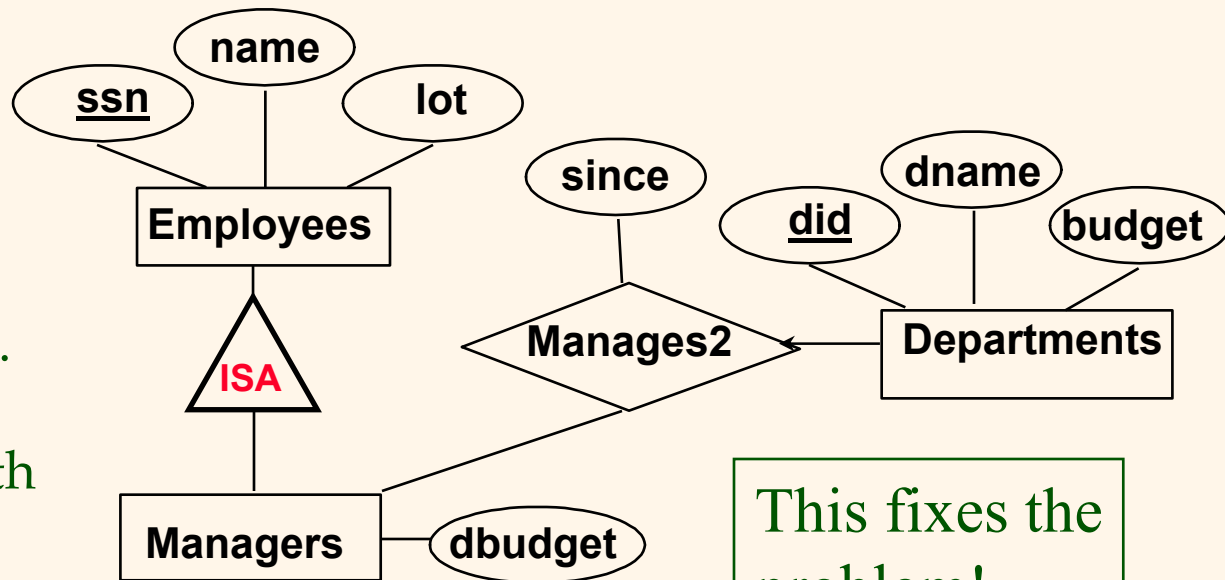
Entity vs. Relationship

❖ First ER diagram OK if a manager gets a separate discretionary budget for each dept.



❖ What if a manager gets a discretionary budget that covers *all* managed depts?

- **Redundancy:** *dbudget* stored for each dept managed by manager.
- **Misleading:** Suggests *dbudget* associated with department-mgr combination.

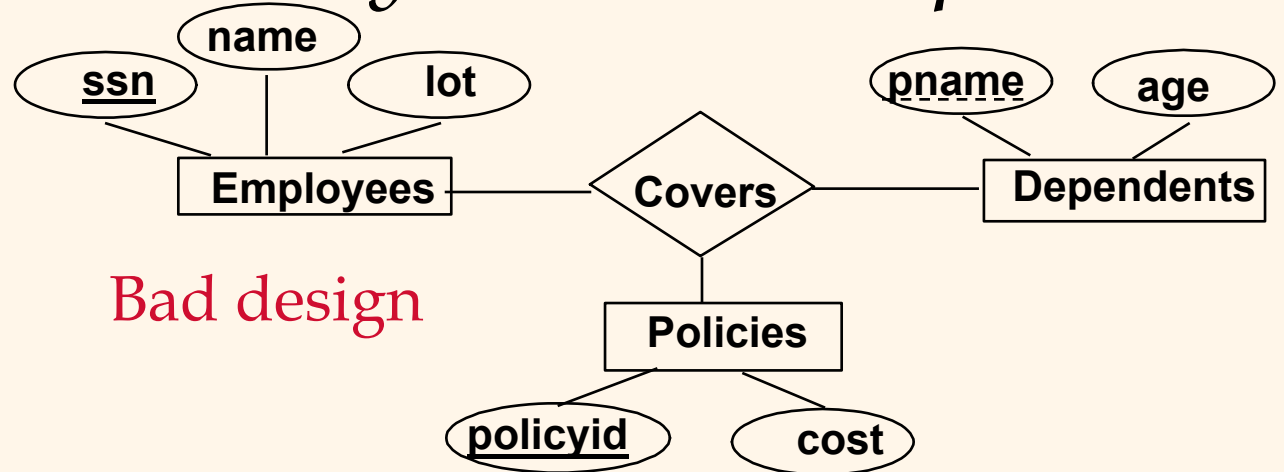


This fixes the problem!

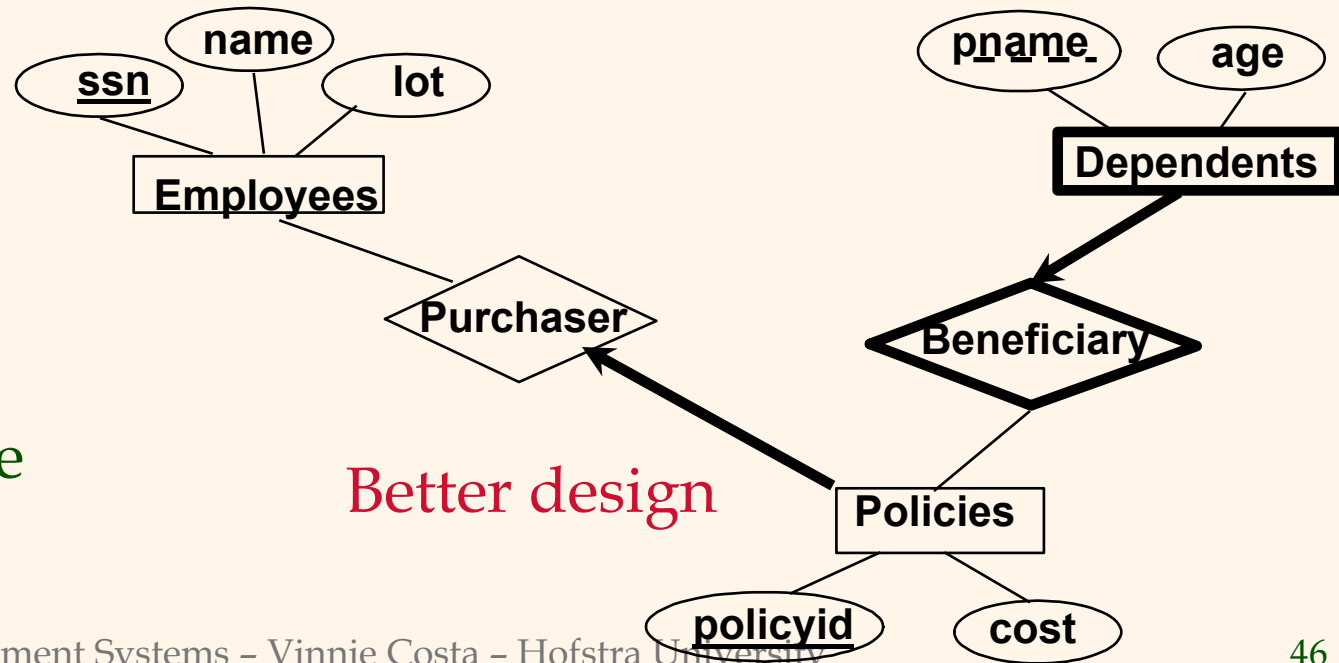
Binary vs. Ternary Relationships

❖ If each policy is owned by just 1 employee, and each dependent is tied to the covering policy, first diagram is inaccurate.

❖ What are the additional constraints in the 2nd diagram?



Bad design



Better design

Binary vs. Ternary Relationships (Contd.)

- ❖ Previous example illustrated a case when two binary relationships were better than one ternary relationship.
- ❖ An example in the other direction: a ternary relation **Contracts** relates entity sets **Parts**, **Departments** and **Suppliers**, and has descriptive attribute *qty*. No combination of binary relationships is an adequate substitute:
 - S “can-supply” P, D “needs” P, and D “deals-with” S does not imply that D has agreed to buy P from S.
 - How do we record *qty*?

Summary of Conceptual Design

- ❖ *Conceptual design follows requirements analysis,*
 - Yields a high-level description of data to be stored
- ❖ ER model popular for conceptual design
 - Constructs are expressive, close to the way people think about their applications.
- ❖ Basic constructs: *entities, relationships, and attributes* (of entities and relationships).
- ❖ Some additional constructs: *weak entities, ISA hierarchies, and aggregation.*
- ❖ Note: There are many variations on ER model.

Summary of ER (Contd.)

- ❖ Several kinds of integrity constraints can be expressed in the ER model: *key constraints*, *participation constraints*, and *overlap/covering constraints* for ISA hierarchies. Some *foreign key constraints* are also implicit in the definition of a relationship set.
 - Some constraints (notably, *functional dependencies*) cannot be expressed in the ER model.
 - Constraints play an important role in determining the best database design for an enterprise.

Summary of ER (Contd.)

- ❖ ER design is *subjective*. There are often many ways to model a given scenario! Analyzing alternatives can be tricky, especially for a large enterprise. Common choices include:
 - Entity vs. attribute, entity vs. relationship, binary or n-ary relationship, whether or not to use ISA hierarchies, and whether or not to use aggregation.
- ❖ Ensuring good database design: resulting relational schema should be analyzed and refined further. FD information and normalization techniques are especially useful.

Useful Websites

❖ <http://www.omg.org/>

Homework

- ❖ Read Chapter Two
- ❖ Exercises p.52: 2.1, 2.2