# Database Systems

# CHAPTER 1:  INTRODUCTION

- Purpose of Database Systems

- View of Data

- Data Models

- Data Definition Language

- Data Manipulation Language

- Transaction Management

- Storage Management

- Database Administrator

- Database Users

- Overall System Structure

# DATABASE MANAGEMENT SYSTEM (DBMS)

- Collection of interrelated data
- Set of programs to access the data
- DBMS contains information about a particular enterprise
- DBMS provides an environment that is both *convenient* and *efficient* to use.
- Database Applications:
    - Banking: all transactions
    - Airlines: reservations, schedules
    - Universities:  registration, grades
    - Sales: customers, products, purchases
    - Manufacturing: production, inventory, orders, supply chain
    - Human resources:  employee records, salaries, tax deductions
- Databases touch all aspects of our lives

# PURPOSE OF DATABASE SYSTEM

- In the early days, database applications were built on top of file systems

- Drawbacks of using file systems to store data:

    - Data redundancy and inconsistency

        - Multiple file formats, duplication of information in different files

    - Difficulty in accessing data

        - Need to write a new program to carry out each new task

    - Data isolation — multiple files and formats

    - Integrity problems

        - Integrity constraints  (e.g. account balance > 0) become part of program code

        - Hard to add new constraints or change existing ones

# PURPOSE OF DATABASE SYSTEMS (CONT.)

- Drawbacks of using file systems (cont.)

  - Atomicity of updates

    - Failures may leave database in an inconsistent state with partial updates carried out

    - E.g. transfer of funds from one account to another should either complete or not happen at all

  - Concurrent access by multiple users

    - Concurrent accessed needed for performance

    - Uncontrolled concurrent accesses can lead to inconsistencies

      - E.g. two people reading a balance and updating it at the same time

  - Security problems

- Database systems offer solutions to all the above problems
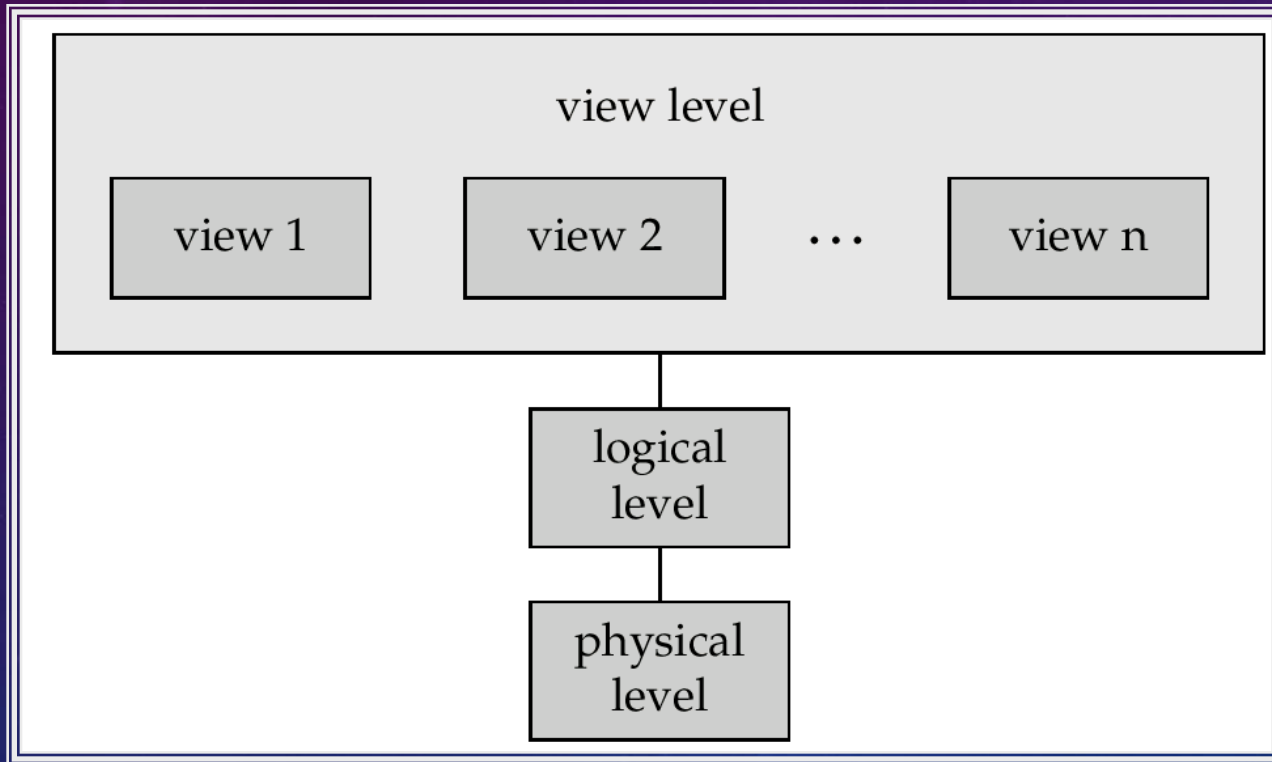
# LEVELS OF ABSTRACTION

- Physical level describes how a record (e.g., customer) is stored.

- Logical level: describes data stored in database, and the relationships among the data.

                        **type** customer = **record**

                                     *name* : string;

                                       *street* : string;

                                       *city* : integer;

                         **end**;

- View level: application programs hide details of data types.  Views can also hide information (e.g., salary) for security purposes.

An architecture for a database system

# VIEW OF DATA

# INSTANCES AND SCHEMAS

- Similar to types and variables in programming languages
- **Schema** – the logical structure of the database
    - e.g., the database consists of information about a set of customers and accounts and the relationship between them)
    - Analogous to type information of a variable in a program
    - **Physical schema**: database design at the physical level
    - **Logical schema**: database design at the logical level
- **Instance** – the actual content of the database at a particular point in time
    - Analogous to the value of a variable
- **Physical Data Independence** – the ability to modify the physical schema without changing the logical schema
    - Applications depend on the logical schema
    - In general, the interfaces between the various levels and components should be well defined so that changes in some parts do not seriously influence others.
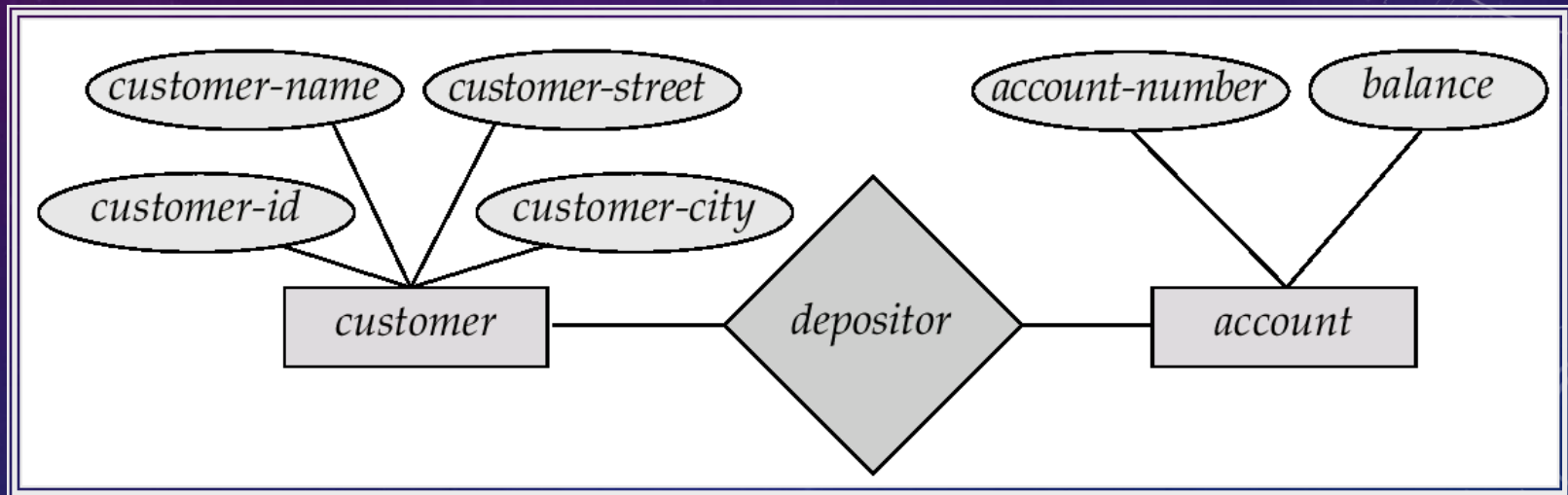
# DATA MODELS

- A collection of tools for describing

  - data

  - data relationships

  - data semantics

  - data constraints

- Entity-Relationship model

- Relational model

- Other models:

  - object-oriented model

  - semi-structured data models

  - Older models: network model and hierarchical model

# ENTITY-RELATIONSHIP MODEL

Example of schema in the entity-relationship model

# ENTITY RELATIONSHIP MODEL (CONT.)

- E-R model of real world
  - Entities (objects)
    - E.g. customers, accounts, bank branch
  - Relationships between entities
    - E.g. Account A-101 is held by customer Johnson
    - Relationship set *depositor* associates customers with accounts
- Widely used for database design
  - Database design in E-R model usually converted to design in the relational model (coming up next) which is used for storage and processing

# RELATIONAL MODEL

- Example of tabular data in the relational model

Attributes

| Customer-id | customer-name | customer-street | customer-city | account-number |
|---|---|---|---|---|
| 192-83-7465 | Johnson | Alma | Palo Alto | A-101 |
| 019-28-3746 | Smith | North | Rye | A-215 |
| 192-83-7465 | Johnson | Alma | Palo Alto | A-201 |
| 321-12-3123 | Jones | Main | Harrison | A-217 |
| 019-28-3746 | Smith | North | Rye | A-201 |

# A SAMPLE RELATIONAL DATABASE

| customer-id | customer-name | customer-street | customer-city |
|---|---|---|---|
| 192-83-7465 | Johnson | 12 Alma St. | Palo Alto |
| 019-28-3746 | Smith | 4 North St. | Rye |
| 677-89-9011 | Hayes | 3 Main St. | Harrison |
| 182-73-6091 | Turner | 123 Putnam Ave. | Stamford |
| 321-12-3123 | Jones | 100 Main St. | Harrison |
| 336-66-9999 | Lindsay | 175 Park Ave. | Pittsfield |
| 019-28-3746 | Smith | 72 North St. | Rye |

(a) The *customer* table

| account-number | balance |
|---|---|
| A-101 | 500 |
| A-215 | 700 |
| A-102 | 400 |
| A-305 | 350 |
| A-201 | 900 |
| A-217 | 750 |
| A-222 | 700 |

(b) The *account* table

| customer-id | account-number |
|---|---|
| 192-83-7465 | A-101 |
| 192-83-7465 | A-201 |
| 019-28-3746 | A-215 |
| 677-89-9011 | A-102 |
| 182-73-6091 | A-305 |
| 321-12-3123 | A-217 |
| 336-66-9999 | A-222 |
| 019-28-3746 | A-201 |

(c) The *depositor* table

# DATA DEFINITION LANGUAGE (DDL)

- Specification notation for defining the database schema

    - E.g.
        **create table** *account* (
        *account-number*   **char**(10),
        *balance*         **integer**)

- DDL compiler generates a set of tables stored in a *data dictionary*

- Data dictionary contains metadata (i.e., data about data)

    - database schema

    - Data *storage and definition* language

        - language in which the storage structure and access methods used by the database system are specified

        - Usually an extension of the data definition language

# DATA MANIPULATION LANGUAGE (DML)

- Language for accessing and manipulating the data organized by the appropriate data model
  - DML also known as query language
- Two classes of languages
  - Procedural – user specifies what data is required and how to get those data
  - Nonprocedural – user specifies what data is required without specifying how to get those data
- SQL is the most widely used query language

# SQL

SQL: widely used non-procedural language

- E.g. find the name of the customer with customer-id 192-83-7465

  **select**   *customer.customer-name*
  **from**     *customer*
  **where**  *customer.customer-id* = '192-83-7465'

- E.g. find the balances of all accounts held by the customer with customer-id 192-83-7465

  **select**   *account.balance*
  **from**     *depositor, account*
  **where**  *depositor.customer-id* = '192-83-7465' **and**
         *depositor.account-number = account.account-number*

- Application programs generally access databases through one of

  - Language extensions to allow embedded SQL

  - Application program interface (e.g. ODBC/JDBC) which allow SQL queries to be sent to a database

# DATABASE USERS

- Users are differentiated by the way they expect to interact with the system

- Application programmers – interact with system through DML calls

- Sophisticated users – form requests in a database query language

- Specialized users – write specialized database applications that do not fit into the traditional data processing framework

- Naïve users – invoke one of the permanent application programs that have been written previously

  - E.g. people accessing database over the web, bank tellers, clerical staff

# DATABASE ADMINISTRATOR

- Coordinates all the activities of the database system; the database administrator has a good understanding of the enterprise's information resources and needs.

- Database administrator's duties include:

  - Schema definition

  - Storage structure and access method definition

  - Schema and physical organization modification

  - Granting user authority to access the database

  - Specifying integrity constraints

  - Acting as liaison with users

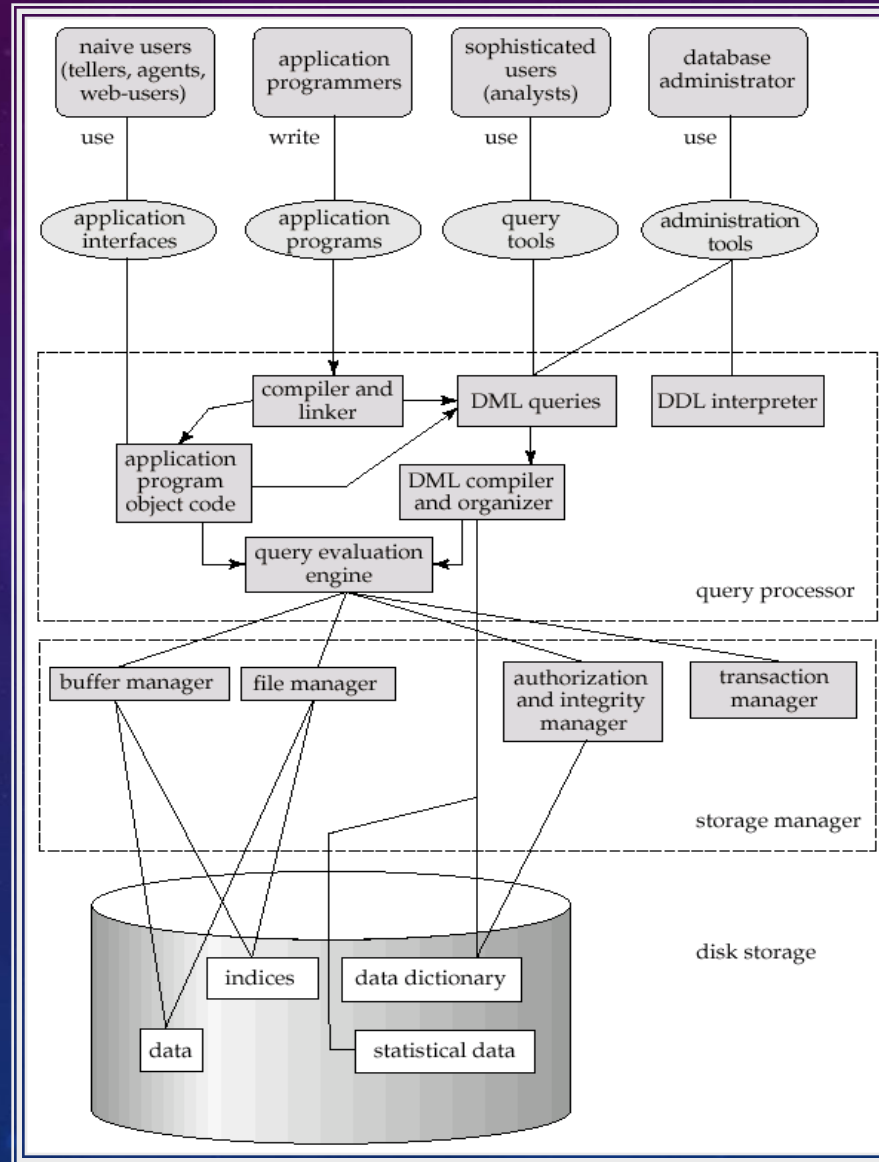  - Monitoring performance and responding to changes in requirements

# TRANSACTION MANAGEMENT

- A *transaction* is a collection of operations that performs a single logical function in a database application

- Transaction-management component ensures that the database remains in a consistent (correct) state despite system failures (e.g., power failures and operating system crashes) and transaction failures.

- Concurrency-control manager controls the interaction among the concurrent transactions, to ensure the consistency of the database.
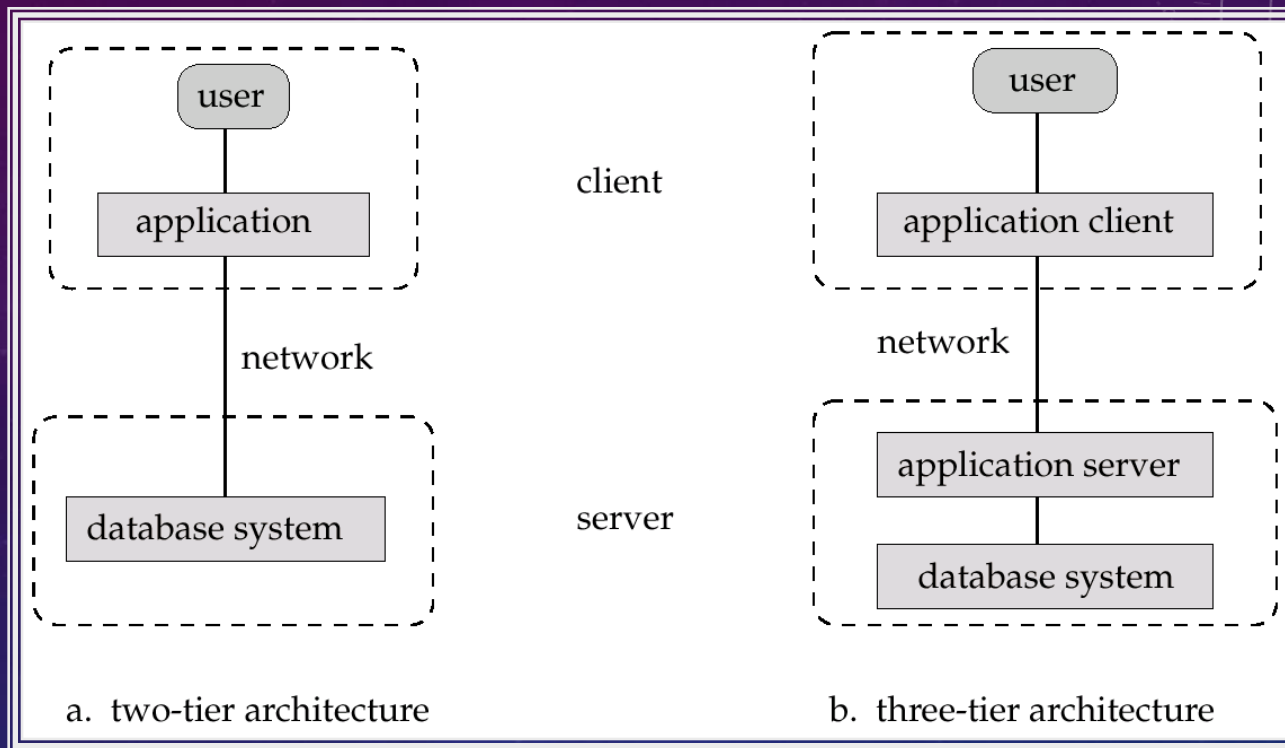
# STORAGE MANAGEMENT

- Storage manager is a program module that provides the interface between the low-level data stored in the database and the application programs and queries submitted to the system.

- The storage manager is responsible to the following tasks:

  - interaction with the file manager

  - efficient storing, retrieving and updating of data

# OVERALL SYSTEM STRUCTURE

# APPLICATION ARCHITECTURES



a. two-tier architecture

b. three-tier architecture

- **Two-tier architecture**:  E.g. client programs using ODBC/JDBC to communicate with a database
- **Three-tier architecture**: E.g. web-based applications, and applications built using "middleware"

# ENTITY SETS

- A *database* can be modeled as:
  - a collection of entities,
  - relationship among entities.
- An *entity* is an object that exists and is distinguishable from other objects.
  - Example:  specific person, company, event, plant
- Entities have *attributes*
  - Example: people have *names* and *addresses*
- An *entity set* is a set of entities of the same type that share the same properties.
  - Example: set of all persons, companies, trees, holidays

# ENTITY SETS *CUSTOMER* AND *LOAN*

| customer-id | customer-name | customer-street | customer-city | | loan-number | amount |
|---|---|---|---|---|---|---|
| 321-12-3123 | Jones | Main | Harrison | | L-17 | 1000 |
| 019-28-3746 | Smith | North | Rye | | L-23 | 2000 |
| 677-89-9011 | Hayes | Main | Harrison | | L-15 | 1500 |
| 555-55-5555 | Jackson | Dupont | Woodside | | L-14 | 1500 |
| 244-66-8800 | Curry | North | Rye | | L-19 | 500 |
| 963-96-3963 | Williams | Nassau | Princeton | | L-11 | 900 |
| 335-57-7991 | Adams | Spring | Pittsfield | | L-16 | 1300 |

*customer*      *loan*

# ATTRIBUTES

- An entity is represented by a set of attributes, that is descriptive properties possessed by all members of an entity set.
  Example:

  *customer = (customer-id, customer-name, customer-street, customer-city)*
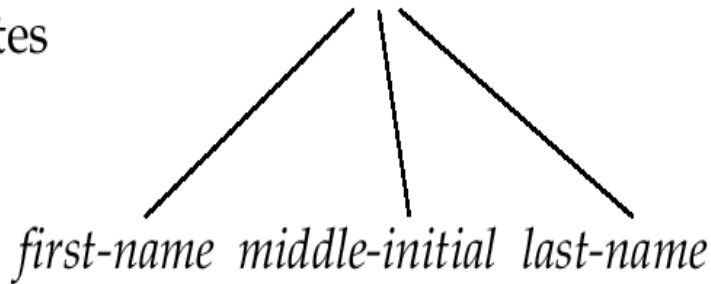  *loan = (loan-number, amount)*

- *Domain* – the set of permitted values for each attribute

- Attribute types:

  - *Simple* and *composite* attributes.

  - *Single-valued* and *multi-valued* attributes

    - E.g. multivalued attribute: *phone-numbers*

  - *Derived* attributes

    - Can be computed from other attributes
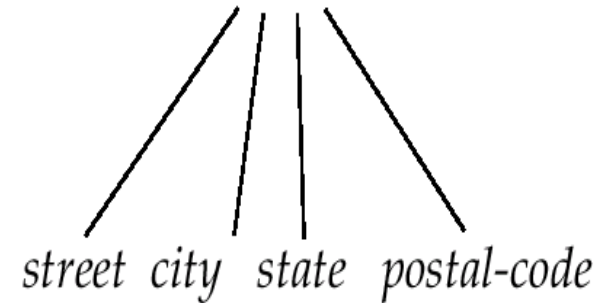
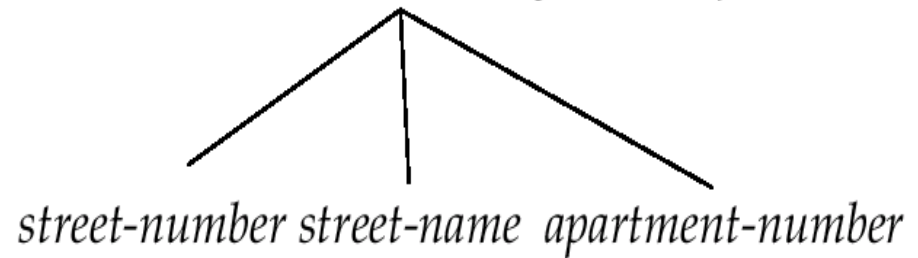    - E.g. *age*, given date of birth

# COMPOSITE ATTRIBUTES

# RELATIONSHIP SETS

- A relationship is an association among several entities

  Example:

  | <u>Hayes</u> | <u>*depositor*</u> | <u>A-102</u> |
  |:---:|:---:|:---:|
  | *customer* entity | relationship set | *account* entity |

- A *relationship* set is a mathematical relation among $n \geq 2$ entities, each taken from entity sets
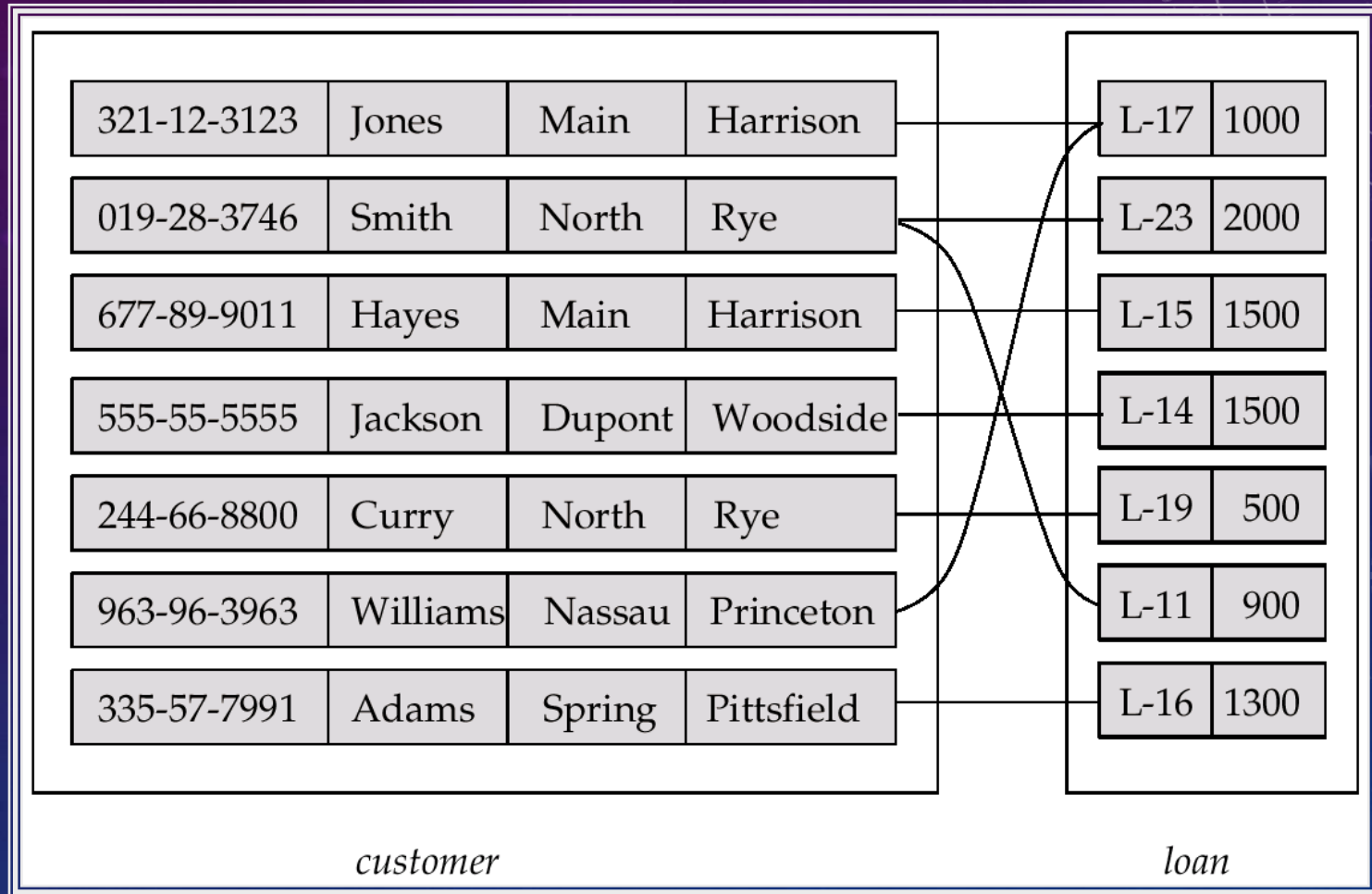
  $$\{(e_1, e_2, \dots e_n) \mid e_1 \in E_1, e_2 \in E_2, \dots, e_n \in E_n\}$$

  where $(e_1, e_2, \dots, e_n)$ is a relationship

  - Example:

    $$(\text{Hayes, A-102}) \in \textit{depositor}$$

# RELATIONSHIP SET *BORROWER*



| | | | | | | |
|---|---|---|---|---|---|---|
| 321-12-3123 | Jones | Main | Harrison | | L-17 | 1000 |
| 019-28-3746 | Smith | North | Rye | | L-23 | 2000 |
| 677-89-9011 | Hayes | Main | Harrison | | L-15 | 1500 |
| 555-55-5555 | Jackson | Dupont | Woodside | | L-14 | 1500 |
| 244-66-8800 | Curry | North | Rye | | L-19 | 500 |
| 963-96-3963 | Williams | Nassau | Princeton | | L-11 | 900 |
| 335-57-7991 | Adams | Spring | Pittsfield | | L-16 | 1300 |

*customer*          *loan*

# RELATIONSHIP SETS (CONT.)

- An *attribute* can also be property of a relationship set.

- For instance, the *depositor* relationship set between entity sets *customer* and *account* may have the attribute *access-date*
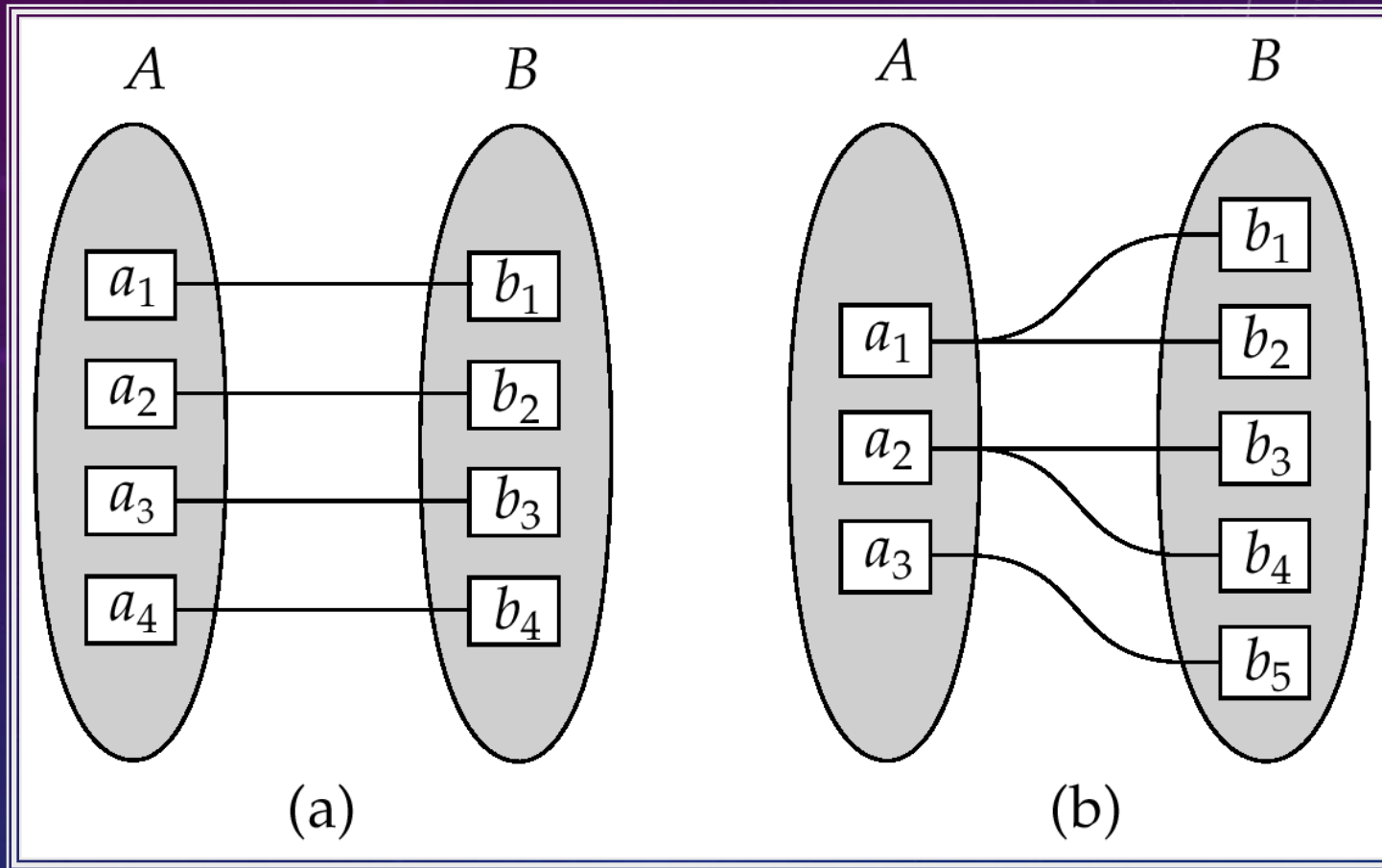
# DEGREE OF A RELATIONSHIP SET

- Refers to number of entity sets that participate in a relationship set.

- Relationship sets that involve two entity sets are *binary* (or degree two). Generally, most relationship sets in a database system are binary.

- Relationship sets may involve more than two entity sets.

  - ☐ E.g. Suppose employees of a bank may have jobs (responsibilities) at multiple branches, with different jobs at different branches. Then there is a ternary relationship set between entity sets *employee, job and branch*

- Relationships between more than two entity sets are rare. Most relationships are binary. (More on this later.)

# MAPPING CARDINALITIES

- Express the number of entities to which another entity can be associated via a relationship set.

- Most useful in describing binary relationship sets.

- For a binary relationship set the mapping cardinality must be one of the following types:

  - One to one

  - One to many

  - Many to one
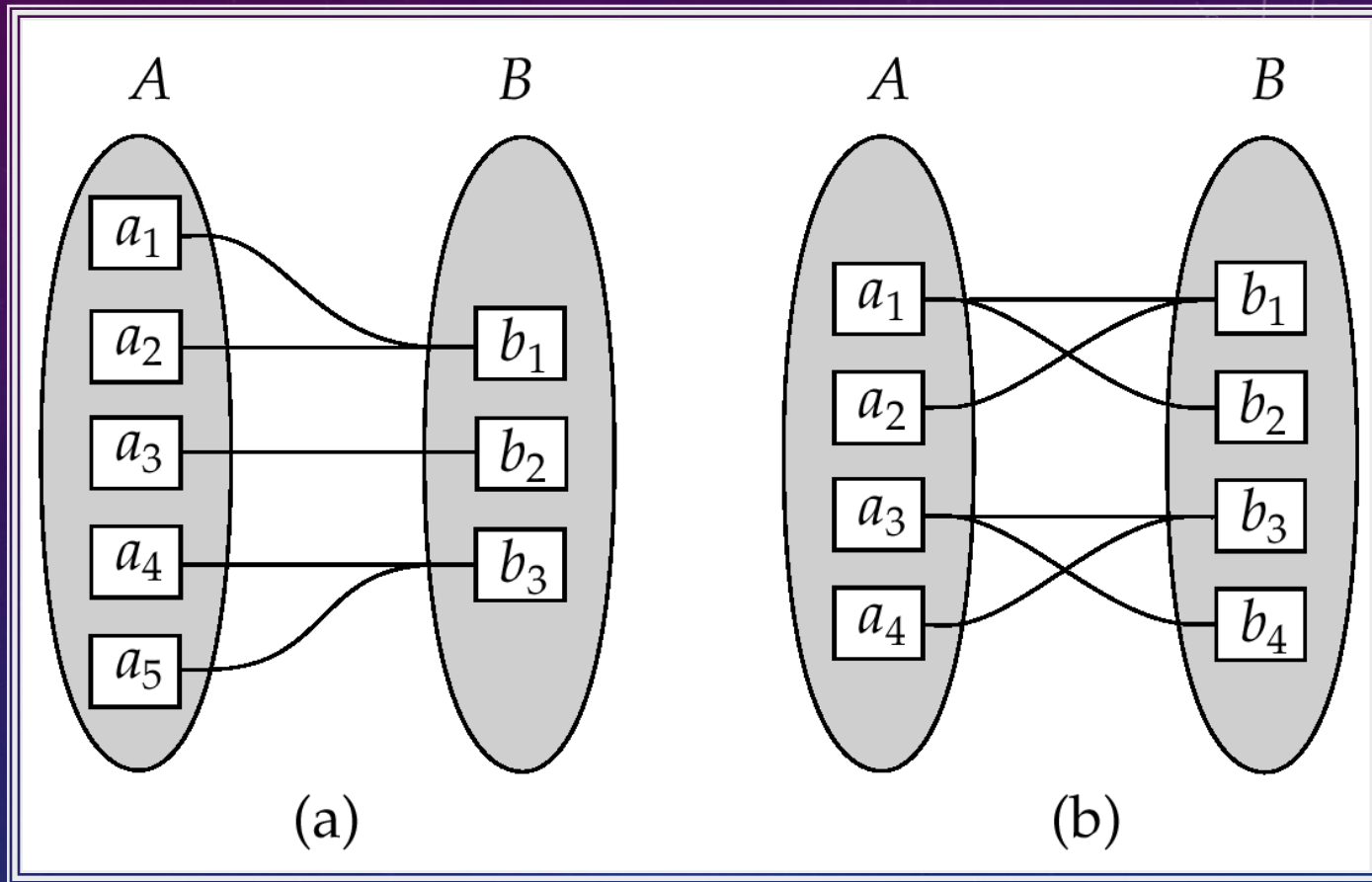
  - Many to many

# MAPPING CARDINALITIES



One to one

One to many

Note: Some elements in A and B may not be mapped to any elements in the other set
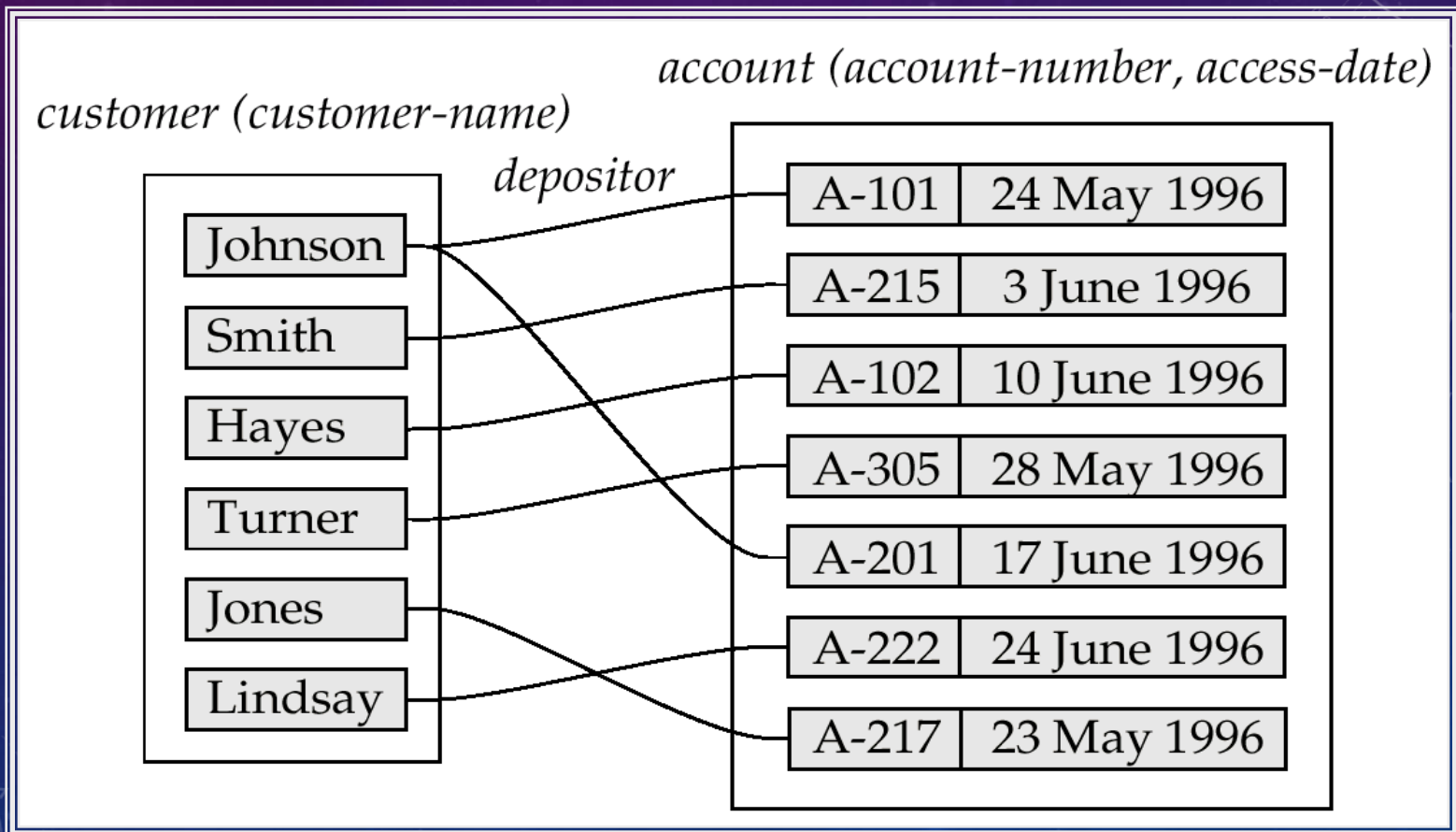
# MAPPING CARDINALITIES



Many to one            Many to many

Note: Some elements in A and B may not be mapped to any elements in the other set

# MAPPING CARDINALITIES AFFECT ER DESIGN

- Can make *access-date* an attribute of account, instead of a relationship attribute, if each account can have only one customer
    - I.e., the relationship from account to customer is many to one, or equivalently, customer to account is one to many



account (account-number, access-date)

customer (customer-name)

depositor

| Johnson | |
| Smith | |
| Hayes | |
| Turner | |
| Jones | |
| Lindsay | |

| A-101 | 24 May 1996 |
| A-215 | 3 June 1996 |
| A-102 | 10 June 1996 |
| A-305 | 28 May 1996 |
| A-201 | 17 June 1996 |
| A-222 | 24 June 1996 |
| A-217 | 23 May 1996 |

# E-R DIAGRAMS



- ☐ **Rectangles** represent entity sets.
- ☐ **Diamonds** represent relationship sets.
- ☐ **Lines** link attributes to entity sets and entity sets to relationship sets.
- ☐ **Ellipses** represent attributes
  - ☐ **Double ellipses** represent multivalued attributes.
  - ☐ **Dashed ellipses** denote derived attributes.
- ☐ **Underline** indicates primary key attributes (will study later)

# E-R DIAGRAM WITH COMPOSITE, MULTIVALUED, AND DERIVED ATTRIBUTES

# RELATIONSHIP SETS WITH ATTRIBUTES

# ROLES

- Entity sets of a relationship need not be distinct

- The labels "manager" and "worker" are called roles; they specify how employee entities interact via the works-for relationship set.

- Roles are indicated in E-R diagrams by labeling the lines that connect diamonds to rectangles.

- Role labels are optional, and are used to clarify semantics of the relationship

# CARDINALITY CONSTRAINTS

- We express cardinality constraints by drawing either a directed line (→), signifying "one," or an undirected line (—), signifying "many," between the relationship set and the entity set.

- E.g.: One-to-one relationship:

  - A customer is associated with at most one loan via the relationship *borrower*

  - A loan is associated with at most one customer via *borrower*

# ONE-TO-MANY RELATIONSHIP

- In the one-to-many relationship a loan is associated with at most one customer via *borrower*, a customer is associated with several (including 0) loans via *borrower*

# MANY-TO-ONE RELATIONSHIPS

- In a many-to-one relationship a loan is associated with several (including 0) customers via *borrower*, a customer is associated with at most one loan via *borrower*

- A customer is associated with several (possibly 0) loans via borrower
- A loan is associated with several (possibly 0) customers via borrower

# PARTICIPATION OF AN ENTITY SET IN A RELATIONSHIP SET

- Total participation (indicated by double line):  every entity in the entity set participates in at least one relationship in the relationship set
  - E.g. participation of *loan* in *borrower* is total
    - every loan must have a customer associated to it via borrower
- Partial participation:  some entities may not participate in any relationship in the relationship set
  - E.g. participation of *customer* in *borrower* is partial

□ Cardinality limits can also express participation constraints

# KEYS

- A *super key* of an entity set is a set of one or more attributes whose values uniquely determine each entity.

- A *candidate key* of an entity set is a minimal super key
  - *Customer-id* is candidate key of *customer*
  - *account-number* is candidate key of *account*

- Although several candidate keys may exist, one of the candidate keys is selected to be the *primary key*.

# KEYS FOR RELATIONSHIP SETS

- The combination of primary keys of the participating entity sets forms a super key of a relationship set.

    - (*customer-id, account-number*) is the super key of *depositor*

    - *NOTE:  this means a pair of entity sets can have at most one relationship in a particular relationship set.*

        - E.g. if we wish to track all access-dates to each account by each customer, we cannot assume a relationship for each access.  We can use a multivalued attribute though

- Must consider the mapping cardinality of the relationship set when deciding the what are the candidate keys

- Need to consider semantics of relationship set in selecting the *primary key*  in case of more than one candidate key

# E-R DIAGRAM WITH A TERNARY RELATIONSHIP

# CARDINALITY CONSTRAINTS ON TERNARY RELATIONSHIP

- We allow at most one arrow out of a ternary (or greater degree) relationship to indicate a cardinality constraint

- E.g. an arrow from *works-on* to *job* indicates each employee works on at most one job at any branch.

- If there is more than one arrow, there are two ways of defining the meaning.

  - E.g a ternary relationship *R* between *A*, *B* and *C* with arrows to *B* and *C* could mean

  - 1.  each *A* entity is associated with a unique entity from *B* and *C* or

  - 2.  each pair of entities from (*A, B*) is associated with a unique *C* entity,       and each pair (*A, C*) is associated with a unique *B*

  - Each alternative has been used in different formalisms

  - To avoid confusion we outlaw more than one arrow

# BINARY VS. NON-BINARY RELATIONSHIPS

- Some relationships that appear to be non-binary may be better represented using binary relationships

    - E.g.  A ternary relationship *parents*, relating a child to his/her father and mother, is best replaced by two binary relationships,  *father* and *mother*

        - Using two binary relationships allows partial information (e.g. only mother being know)

    - But there are some relationships that are naturally non-binary

        - E.g. *works-on*

# CONVERTING NON-BINARY RELATIONSHIPS TO BINARY FORM

- In general, any non-binary relationship can be represented using binary relationships by creating an artificial entity set.

  - Replace $R$ between entity sets A, B and C by an entity set $E$, and three relationship sets:

  1. $R_A$, relating $E$ and $A$          2. $R_B$, relating $E$ and $B$

  3. $R_C$, relating $E$ and $C$

  - Create a special identifying attribute for $E$

  - Add any attributes of $R$ to $E$

  - For each relationship ($a_i$, $b_i$, $c_i$) in $R$, create

  1. a new entity $e_i$ in the entity set $E$     2. add ($e_i$, $a_i$) to $R_A$

  3. add ($e_i$, $b_i$) to $R_B$          4. add ($e_i$, $c_i$) to $R_C$

# CONVERTING NON-BINARY RELATIONSHIPS (CONT.)

- Also need to translate constraints

    - Translating all constraints may not be possible

    - There may be instances in the translated schema that cannot correspond to any instance of $R$

        - *Exercise: add constraints to the relationships $R_A$, $R_B$ and $R_C$ to ensure that a newly created entity corresponds to exactly one entity in each of entity sets $A$, $B$ and $C$*

    - We can avoid creating an identifying attribute by making E a weak entity set (described shortly) identified by the three relationship sets

# DESIGN ISSUES

- Use of entity sets vs. attributes
  Choice mainly depends on the structure of the enterprise being modeled, and on the semantics associated with the attribute in question.

- Use of entity sets vs. relationship sets
  Possible guideline is to designate a relationship set to describe an action that occurs between entities

- Binary versus $n$-ary relationship sets
  Although it is possible to replace any nonbinary ($n$-ary, for $n > 2$) relationship set by a number of distinct binary relationship sets, a $n$-ary relationship set shows more clearly that several entities participate in a single relationship.

- Placement of relationship attributes

HOW ABOUT DOING AN ER DESIGN INTERACTIVELY ON THE BOARD? SUGGEST AN APPLICATION TO BE MODELED.

# WEAK ENTITY SETS

- An entity set that does not have a primary key is referred to as a *weak entity set*.

- The existence of a weak entity set depends on the existence of a *identifying entity set*

  - it must relate to the identifying entity set via a total, one-to-many relationship set from the identifying to the weak entity set

  - Identifying relationship depicted using a double diamond

- The *discriminator (or partial key)* of a weak entity set is the set of attributes that distinguishes among all the entities of a weak entity set.

- The primary key of a weak entity set is formed by the primary key of the strong entity set on which the weak entity set is existence dependent, plus the weak entity set's discriminator.

# WEAK ENTITY SETS (CONT.)

- We depict a weak entity set by double rectangles.

- We underline the discriminator of a weak entity set with a dashed line.

- *payment-number* – discriminator of the *payment* entity set

- Primary key for *payment* – (*loan-number, payment-number*)

# WEAK ENTITY SETS (CONT.)

- Note: the primary key of the strong entity set is not explicitly stored with the weak entity set, since it is implicit in the identifying relationship.

- If *loan-number* were explicitly stored, *payment* could be made a strong entity, but then the relationship between *payment* and *loan* would be duplicated by an implicit relationship defined by the attribute *loan-number* common to *payment* and *loan*

# MORE WEAK ENTITY SET EXAMPLES

- In a university, a *course* is a strong entity and a *course-offering* can be modeled as a weak entity

- The discriminator of *course-offering* would be *semester* (including year) and *section-number* (if there is more than one section)

- If we model *course-offering* as a strong entity we would model *course-number* as an attribute.

  Then the relationship with *course* would be implicit in the *course-number* attribute

# SPECIALIZATION

- Top-down design process; we designate subgroupings within an entity set that are distinctive from other entities in the set.

- These subgroupings become lower-level entity sets that have attributes or participate in relationships that do not apply to the higher-level entity set.

- Depicted by a *triangle* component labeled ISA (E.g. *customer* "is a" *person*).

- **Attribute inheritance** – a lower-level entity set inherits all the attributes and relationship participation of the higher-level entity set to which it is linked.

# SPECIALIZATION EXAMPLE

# GENERALIZATION

- A bottom-up design process – combine a number of entity sets that share the same features into a higher-level entity set.

- Specialization and generalization are simple inversions of each other; they are represented in an E-R diagram in the same way.

- The terms specialization and generalization are used interchangeably.

# SPECIALIZATION AND GENERALIZATION (CONTD.)

- Can have multiple specializations of an entity set based on different features.

- E.g. *permanent-employee* vs. *temporary-employee*, in addition to *officer* vs. *secretary* vs. *teller*

- Each particular employee would be
    - a member of one of *permanent-employee* or *temporary-employee*,
    - and also a member of one of *officer*, *secretary*, or *teller*

- The ISA relationship also referred to as **superclass - subclass** relationship

# DESIGN CONSTRAINTS ON A SPECIALIZATION/GENERALIZATION

- Constraint on which entities can be members of a given lower-level entity set.
    - condition-defined
        - E.g. all customers over 65 years are members of *senior-citizen* entity set; *senior-citizen* ISA *person*.
    - user-defined
- Constraint on whether or not entities may belong to more than one lower-level entity set within a single generalization.
    - Disjoint
        - an entity can belong to only one lower-level entity set
        - Noted in E-R diagram by writing *disjoint* next to the ISA triangle
    - Overlapping
        - an entity can belong to more than one lower-level entity set

# DESIGN CONSTRAINTS ON A SPECIALIZATION/GENERALIZATION (CONTD.)

- Completeness constraint -- specifies whether or not an entity in the higher-level entity set must belong to at least one of the lower-level entity sets within a generalization.

  - **total** : an entity must belong to one of the lower-level entity sets

  - **partial**: an entity need not belong to one of the lower-level entity sets

# AGGREGATION

 Consider the ternary relationship *works-on*, which we saw earlier

 Suppose we want to record managers for tasks performed by an employee at a branch

# AGGREGATION (CONT.)

- Relationship sets *works-on* and *manages* represent overlapping information

    - Every *manages* relationship corresponds to a *works-on* relationship

    - However, some *works-on* relationships may not correspond to any *manages* relationships

        - So we can't discard the *works-on* relationship

- Eliminate this redundancy via *aggregation*

    - Treat relationship as an abstract entity

    - Allows relationships between relationships

    - Abstraction of relationship into new entity

- Without introducing redundancy, the following diagram represents:

    - An employee works on a particular job at a particular branch

    - An employee, branch, job combination may have an associated manager

# E-R DIAGRAM WITH AGGREGATION

# E-R DESIGN DECISIONS

- The use of an attribute or entity set to represent an object.

- Whether a real-world concept is best expressed by an entity set or a relationship set.

- The use of a ternary relationship versus a pair of binary relationships.

- The use of a strong or weak entity set.

- The use of specialization/generalization – contributes to modularity in the design.

- The use of aggregation – can treat the aggregate entity set as a single unit without concern for the details of its internal structure.

# UML

- UML: Unified Modeling Language

- UML has many components to graphically model different aspects of an entire software system

- UML Class Diagrams correspond to E-R Diagram, but several differences.

# UML CLASS DIAGRAMS (CONTD.)

- Entity sets are shown as boxes, and attributes are shown within the box, rather than as separate ellipses in E-R diagrams.

- Binary relationship sets are represented in UML by just drawing a line connecting the entity sets. The relationship set name is written adjacent to the line.

- The role played by an entity set in a relationship set may also be specified by writing the role name on the line, adjacent to the entity set.

- The relationship set name may alternatively be written in a box, along with attributes of the relationship set, and the box is connected, using a dotted line, to the line depicting the relationship set.

- Non-binary relationships drawn using diamonds, just as in ER diagrams

# UML CLASS DIAGRAM NOTATION (CONT.)



*Note reversal of position in cardinality constraint depiction
*Generalization can use merged or separate arrows independent
  of disjoint/overlapping

# UML CLASS DIAGRAMS (CONTD.)

- Cardinality constraints are specified in the form $l..h$, where $l$ denotes the minimum and $h$ the maximum number of relationships an entity can participate in.

- Beware: the positioning of the constraints is exactly the reverse of the positioning of constraints in E-R diagrams.

- The constraint 0..* on the $E2$ side and 0..1 on the $E1$ side means that each $E2$ entity can participate in at most one relationship, whereas each $E1$ entity can participate in many relationships; in other words, the relationship is many to one from $E2$ to $E1$.

- Single values, such as 1 or * may be written on edges; The single value 1 on an edge is treated as equivalent to 1..1, while * is equivalent to 0..*.

# REDUCTION OF AN E-R SCHEMA TO TABLES

- Primary keys allow entity sets and relationship sets to be expressed uniformly as *tables* which represent the contents of the database.

- A database which conforms to an E-R diagram can be represented by a collection of tables.

- For each entity set and relationship set there is a unique table which is assigned the name of the corresponding entity set or relationship set.

- Each table has a number of columns (generally corresponding to attributes), which have unique names.

- Converting an E-R diagram to a table format is the basis for deriving a relational database design from an E-R diagram.

# REPRESENTING ENTITY SETS AS TABLES

- A strong entity set reduces to a table with the same attributes.

| customer-id | customer-name | customer-street | customer-city |
|---|---|---|---|
| 019-28-3746 | Smith | North | Rye |
| 182-73-6091 | Turner | Putnam | Stamford |
| 192-83-7465 | Johnson | Alma | Palo Alto |
| 244-66-8800 | Curry | North | Rye |
| 321-12-3123 | Jones | Main | Harrison |
| 335-57-7991 | Adams | Spring | Pittsfield |
| 336-66-9999 | Lindsay | Park | Pittsfield |
| 677-89-9011 | Hayes | Main | Harrison |
| 963-96-3963 | Williams | Nassau | Princeton |

# COMPOSITE AND MULTIVALUED ATTRIBUTES

- Composite attributes are flattened out by creating a separate attribute for each component attribute

    - E.g. given entity set *customer* with composite attribute *name* with component attributes *first-name* and *last-name* the table corresponding to the entity set has two attributes *name.first-name* and *name.last-name*

- A multivalued attribute M of an entity E is represented by a separate table EM

    - Table EM has attributes corresponding to the primary key of E and an attribute corresponding to multivalued attribute M

    - E.g. Multivalued attribute *dependent-names* of *employee* is represented by a table *employee-dependent-names*( *employee-id, dname*)

    - Each value of the multivalued attribute maps to a separate row of the table EM

        - E.g., an employee entity with primary key John and dependents Johnson and Johndotir maps to two rows:
        (John, Johnson) and (John, Johndotir)

# REPRESENTING WEAK ENTITY SETS

☐ A weak entity set becomes a table that includes a column for the primary key of the identifying strong entity set

| loan-number | payment-number | payment-date | payment-amount |
|:---:|:---:|:---:|:---:|
| L-11 | 53 | 7 June 2001 | 125 |
| L-14 | 69 | 28 May 2001 | 500 |
| L-15 | 22 | 23 May 2001 | 300 |
| L-16 | 58 | 18 June 2001 | 135 |
| L-17 | 5 | 10 May 2001 | 50 |
| L-17 | 6 | 7 June 2001 | 50 |
| L-17 | 7 | 17 June 2001 | 100 |
| L-23 | 11 | 17 May 2001 | 75 |
| L-93 | 103 | 3 June 2001 | 900 |
| L-93 | 104 | 13 June 2001 | 200 |

# REPRESENTING RELATIONSHIP SETS AS TABLES

- A many-to-many relationship set is represented as a table with columns for the primary keys of the two participating entity sets, and any descriptive attributes of the relationship set.

- E.g.: table for relationship set *borrower*

| customer-id | loan-number |
|---|---|
| 019-28-3746 | L-11 |
| 019-28-3746 | L-23 |
| 244-66-8800 | L-93 |
| 321-12-3123 | L-17 |
| 335-57-7991 | L-16 |
| 555-55-5555 | L-14 |
| 677-89-9011 | L-15 |
| 963-96-3963 | L-17 |

# REDUNDANCY OF TABLES

☐ Many-to-one and one-to-many relationship sets that are total on the many-side can be represented by adding an extra attribute to the many side, containing the primary key of the one side

☐ E.g.: Instead of creating a table for relationship *account-branch*, add an attribute *branch* to the entity set *account*

# REDUNDANCY OF TABLES (CONT.)

- For one-to-one relationship sets, either side can be chosen to act as the "many" side
  - That is, extra attribute can be added to either of the tables corresponding to the two entity sets
- If participation is *partial* on the many side, replacing a table by an extra attribute in the relation corresponding to the "many" side could result in null values
- The table corresponding to a relationship set linking a weak entity set to its identifying strong entity set is redundant.
  - E.g. The *payment* table already contains the information that would appear in the *loan-payment* table (i.e., the columns loan-number and *payment-number*).

# Relational Model

- Structure of Relational Databases

- Relational Algebra

- Tuple Relational Calculus

- Domain Relational Calculus

- Extended Relational-Algebra-Operations

- Modification of the Database

- Views

# Example of a Relation

| account-number | branch-name | balance |
|----------------|-------------|---------|
| A-101 | Downtown | 500 |
| A-102 | Perryridge | 400 |
| A-201 | Brighton | 900 |
| A-215 | Mianus | 700 |
| A-217 | Brighton | 750 |
| A-222 | Redwood | 700 |
| A-305 | Round Hill | 350 |

# Basic Structure

- Formally, given sets $D_1$, $D_2$, .... $D_n$ a **relation** $r$ is a subset of $D_1 \times D_2 \times ... \times D_n$
  Thus a relation is a set of n-tuples $(a_1, a_2, ..., a_n)$ where each $a_i \in D_i$

- Example: if

  *customer-name* = {Jones, Smith, Curry, Lindsay}
  *customer-street* = {Main, North, Park}
  *customer-city*   = {Harrison, Rye, Pittsfield}
  Then $r$ = {  (Jones, Main, Harrison),
          (Smith, North, Rye),
          (Curry, North, Rye),
          (Lindsay, Park, Pittsfield)}
   is a relation over *customer-name x customer-street x customer-city*

# Attribute Types

- Each attribute of a relation has a name

- The set of allowed values for each attribute is called the **domain** of the attribute

- Attribute values are (normally) required to be **atomic**, that is, indivisible
  - E.g. multivalued attribute values are not atomic
  - E.g. composite attribute values are not atomic

- The special value *null* is a member of every domain

- The null value causes complications in the definition of many operations
  - we shall ignore the effect of null values in our main presentation and consider their effect later

# Relation Schema

- $A_1, A_2, ..., A_n$ are *attributes*

- $R = (A_1, A_2, ..., A_n)$ is a *relation schema*

  E.g.  *Customer-schema =*
       *(customer-name, customer-street, customer-city)*

- $r(R)$ is a *relation* on the *relation schema R*

  E.g.  *customer (Customer-schema)*

# Relation Instance

- The current values (*relation instance*) of a relation are specified by a table

- An element *t* of *r* is a *tuple*, represented by a *row* in a table

attributes
(or columns)

| customer-name | customer-street | customer-city |
|---------------|-----------------|---------------|
| Jones<br>Smith<br>Curry<br>Lindsay | Main<br>North<br>North<br>Park | Harrison<br>Rye<br>Rye<br>Pittsfield |

tuples
(or rows)

customer

# Relations are Unordered

☐ Order of tuples is irrelevant (tuples may be stored in an arbitrary order)
☐ E.g. account relation with unordered tuples

| account-number | branch-name | balance |
|:---:|:---:|:---:|
| A-101 | Downtown | 500 |
| A-215 | Mianus | 700 |
| A-102 | Perryridge | 400 |
| A-305 | Round Hill | 350 |
| A-201 | Brighton | 900 |
| A-222 | Redwood | 700 |
| A-217 | Brighton | 750 |

# Database

- A database consists of multiple relations

- Information about an enterprise is broken up into parts, with each relation storing one part of the information

  E.g.: *account* :   stores information about accounts
  *depositor* : stores information about which customer
  owns which account
  *customer* : stores information about customers

- Storing all information as a single relation such as
  *bank*(*account-number, balance, customer-name*, ..)
  results in
  - repetition of information (e.g. two customers own an account)
  - the need for null values  (e.g. represent a customer without an account)

- Normalization theory (Chapter 7) deals with how to design relational schemas

# The *customer* Relation

| customer-name | customer-street | customer-city |
|---------------|-----------------|---------------|
| Adams | Spring | Pittsfield |
| Brooks | Senator | Brooklyn |
| Curry | North | Rye |
| Glenn | Sand Hill | Woodside |
| Green | Walnut | Stamford |
| Hayes | Main | Harrison |
| Johnson | Alma | Palo Alto |
| Jones | Main | Harrison |
| Lindsay | Park | Pittsfield |
| Smith | North | Rye |
| Turner | Putnam | Stamford |
| Williams | Nassau | Princeton |

# The *depositor* Relation

| *customer-name* | *account-number* |
|-----------------|------------------|
| Hayes | A-102 |
| Johnson | A-101 |
| Johnson | A-201 |
| Jones | A-217 |
| Lindsay | A-222 |
| Smith | A-215 |
| Turner | A-305 |

# E-R Diagram for the Banking Enterprise

# Keys

- Let $K \subseteq R$

- *K* is a ***superkey*** of *R* if values for *K* are sufficient to identify a unique tuple of each possible relation *r(R)*
  - by "possible *r*" we mean a relation *r* that could exist in the enterprise we are modeling.
  - Example: {*customer-name, customer-street*} and
    {*customer-name*}
    are both superkeys of *Customer*, if no two customers can possibly have the same name.

- *K* is a ***candidate key*** if *K* is minimal
  Example: {*customer-name*} is a candidate key for *Customer*, since it is a superkey (assuming no two customers can possibly have the same name), and no subset of it is a superkey.

# Determining Keys from E-R Sets

- **Strong entity set**. The primary key of the entity set becomes the primary key of the relation.

- **Weak entity set**. The primary key of the relation consists of the union of the primary key of the strong entity set and the discriminator of the weak entity set.

- **Relationship set**. The union of the primary keys of the related entity sets becomes a super key of the relation.
  - For binary many-to-one relationship sets, the primary key of the "many" entity set becomes the relation's primary key.
  - For one-to-one relationship sets, the relation's primary key can be that of either entity set.
  - For many-to-many relationship sets, the union of the primary keys becomes the relation's primary key

# Schema Diagram for the Banking Enterprise

# Query Languages

- Language in which user requests information from the database.
- Categories of languages
  - procedural
  - non-procedural
- "Pure" languages:
  - Relational Algebra
  - Tuple Relational Calculus
  - Domain Relational Calculus
- Pure languages form underlying basis of query languages that people use.

# Relational Algebra

- Procedural language
- Six basic operators
  - select
  - project
  - union
  - set difference
  - Cartesian product
  - rename
- The operators take two or more relations as inputs and give a new relation as a result.

# Select Operation – Example

- Relation r

| A | B | C | D |
|---|---|---|---|
| ⍰ | ⍰ | 1 | 7 |
| ⍰ | ⍰ | 5 | 7 |
| ⍰ | ⍰ | 12 | 3 |
| ⍰ | ⍰ | 23 | 10 |

- $\sigma_{A=B \wedge D > 5}(r)$

| A | B | C | D |
|---|---|---|---|
| ⍰ | ⍰ | 1 | 7 |
| ⍰ | ⍰ | 23 | 10 |

# Select Operation

- Notation: $\sigma_p(r)$

- $p$ is called the selection predicate

- Defined as:

$$\sigma_p(r) = \{t \mid t \in r \text{ and } p(t)\}$$

Where $p$ is a formula in propositional calculus consisting of terms connected by : $\wedge$ (**and**), $\vee$ (**or**), $\neg$ (**not**)
Each term is one of:

<attribute> *op* <attribute> or <constant>

where *op* is one of: $=, \neq, >, \geq. <. \leq$

- Example of selection:

$\sigma_{branch\text{-}name="Perryridge"}(account)$

# Project Operation – Example

- Relation *r*:

| A | B | C |
|---|---|---|
| ? | 10 | 1 |
| ? | 20 | 1 |
| ? | 30 | 1 |
| ? | 40 | 2 |

n   $\Pi_{A,C}$ (r)

| A | C |
|---|---|
| ? | 1 |
| ? | 1 |
| ? | 1 |
| ? | 2 |

=

| A | C |
|---|---|
| ? | 1 |
| ? | 1 |
| ? | 2 |

# Project Operation

- Notation:

$$\prod_{A_1, A_2, \ldots, Ak} (r)$$

where $A_1$, $A_2$ are attribute names and $r$ is a relation name.

- The result is defined as the relation of $k$ columns obtained by erasing the columns that are not listed

- Duplicate rows removed from result, since relations are sets

- E.g. To eliminate the *branch-name* attribute of *account*
  $$\prod_{account\text{-}number,\ balance} (account)$$

# Union Operation – Example

- Relations *r, s:*

| A | B |
|---|---|
| ? | 1 |
| ? | 2 |
| ? | 1 |

r

| A | B |
|---|---|
| ? | 2 |
| ? | 3 |

s

r ⬚ s:

| A | B |
|---|---|
| ? | 1 |
| ? | 2 |
| ? | 1 |
| ? | 3 |

# Union Operation

- Notation: $r \cup s$

- Defined as:

$$r \cup s = \{t \mid t \in r \text{ or } t \in s\}$$

- For $r \cup s$ to be valid.

  1. $r$, $s$ must have the *same arity* (same number of attributes)

  2. The attribute domains must be *compatible* (e.g., 2nd column of $r$ deals with the same type of values as does the 2nd column of $s$)

- E.g. to find all customers with either an account or a loan
  $\prod_{customer\text{-}name} (depositor) \cup \prod_{customer\text{-}name} (borrower)$

# Set Difference Operation – Example

- Relations *r, s:*

| A | B |
|---|---|
| ⍰ | 1 |
| ⍰ | 2 |
| ⍰ | 1 |

r

| A | B |
|---|---|
| ⍰ | 2 |
| ⍰ | 3 |

s

r – s:

| A | B |
|---|---|
| ⍰ | 1 |
| ⍰ | 1 |

# Set Difference Operation

- Notation $r - s$

- Defined as:

$$r - s = \{t \mid t \in r \textbf{ and } t \notin s\}$$

- Set differences must be taken between *compatible* relations.
  - $r$ and $s$ must have the *same arity*
  - attribute domains of $r$ and $s$ must be compatible

# Cartesian-Product Operation-Example

Relations r, s:

| A | B |
|---|---|
| ⍰ | 1 |
| ⍰ | 2 |

r

| C | D | E |
|---|---|---|
| ⍰ | 10 | a |
| ⍰ | 10 | a |
| ⍰ | 20 | b |
| ⍰ | 10 | b |

s

r x s:

| A | B | C | D | E |
|---|---|---|---|---|
| ⍰ | 1 | ⍰ | 10 | a |
| ⍰ | 1 | ⍰ | 10 | a |
| ⍰ | 1 | ⍰ | 20 | b |
| ⍰ | 1 | ⍰ | 10 | b |
| ⍰ | 2 | ⍰ | 10 | a |
| ⍰ | 2 | ⍰ | 10 | a |
| ⍰ | 2 | ⍰ | 20 | b |
| ⍰ | 2 | ⍰ | 10 | b |

# Cartesian-Product Operation

- Notation *r* x *s*

- Defined as:

$$r \times s = \{t\ q \mid t \in r \textbf{ and } q \in s\}$$

- Assume that attributes of r(R) and s(S) are disjoint. (That is, $R \cap S = \varnothing$).

- If attributes of *r(R)* and *s(S)* are not disjoint, then renaming must be used.

# Composition of Operations

- Can build expressions using multiple operations

- Example: $\sigma_{A=C}(r \times s)$

- $r \times s$

| A | B | C | D | E |
|---|---|---|---|---|
| ⍰ | 1 | ⍰ | 10 | a |
| ⍰ | 1 | ⍰ | 10 | a |
| ⍰ | 1 | ⍰ | 20 | b |
| ⍰ | 1 | ⍰ | 10 | b |
| ⍰ | 2 | ⍰ | 10 | a |
| ⍰ | 2 | ⍰ | 10 | a |
| ⍰ | 2 | ⍰ | 20 | b |
| ⍰ | 2 | ⍰ | 10 | b |

- $\sigma_{A=C}(r \times s)$

| A | B | C | D | E |
|---|---|---|---|---|
| ⍰ | 1 | ⍰ | 10 | a |
| ⍰ | 2 | ⍰ | 20 | a |
| ⍰ | 2 | ⍰ | 20 | b |

# Rename Operation

- Allows us to name, and therefore to refer to, the results of relational-algebra expressions.

- Allows us to refer to a relation by more than one name.

Example:

$$\rho_X(E)$$

returns the expression $E$ under the name $X$

If a relational-algebra expression $E$ has arity $n$, then

$$\rho_{X(A_1, A_2, \ldots, A_n)}(E)$$

returns the result of expression $E$ under the name $X$, and with the attributes renamed to $A_1, A_2, \ldots, A_n$.

# Banking Example

*branch (branch-name, branch-city, assets)*

*customer (customer-name, customer-street, customer-only)*

*account (account-number, branch-name, balance)*

*loan (loan-number, branch-name, amount)*

*depositor (customer-name, account-number)*

*borrower (customer-name, loan-number)*

# Example Queries

- Find all loans of over $1200

$$\sigma_{amount > 1200} (loan)$$

Find the loan number for each loan of an amount greater than $1200

$$\Pi_{loan-number} (\sigma_{amount > 1200} (loan))$$

# Example Queries

- Find the names of all customers who have a loan, an account, or both, from the bank

$$\Pi_{customer\text{-}name} \ (borrower) \cup \Pi_{customer\text{-}name} \ (depositor)$$

Find the names of all customers who have a loan and an account at bank.

$$\Pi_{customer\text{-}name} \ (borrower) \cap \Pi_{customer\text{-}name} \ (depositor)$$

# Example Queries

- Find the names of all customers who have a loan at the Perryridge branch.

$$\Pi_{customer\text{-}name} \, (\sigma_{branch\text{-}name="Perryridge"}$$
$$(\sigma_{borrower.loan\text{-}number = loan.loan\text{-}number}(borrower \times loan)))$$

☐ Find the names of all customers who have a loan at the Perryridge branch but do not have an account at any branch of the bank.

$$\Pi_{customer\text{-}name} \, (\sigma_{branch\text{-}name = "Perryridge"}$$
$$(\sigma_{borrower.loan\text{-}number = loan.loan\text{-}number}(borrower \times loan))) -$$
$$\Pi_{customer\text{-}name}(depositor)$$

# Example Queries

- Find the names of all customers who have a loan at the Perryridge branch.

 – Query 1

 $\Pi_{\text{customer-name}}(\sigma_{\text{branch-name = "Perryridge"}} ($
 $\sigma_{\text{borrower.loan-number = loan.loan-number}}(\text{borrower x loan})))$

 $\sigma$     Query 2

 $\Pi_{\text{customer-name}}(\sigma_{\text{loan.loan-number = borrower.loan-number}} ($
 $(\sigma_{\text{branch-name = "Perryridge"}}(\text{loan})) \text{ x } \text{borrower}))$

# Example Queries

Find the largest account balance

- Rename *account* relation as *d*

- The query is:

$$\Pi_{\text{balance}}(\text{account}) - \Pi_{\text{account.balance}}$$
$$(\sigma_{\text{account.balance < d.balance}} (\text{account} \times \rho_d (\text{account})))$$

# Formal Definition

- A basic expression in the relational algebra consists of either one of the following:
  - A relation in the database
  - A constant relation

- Let $E_1$ and $E_2$ be relational-algebra expressions; the following are all relational-algebra expressions:
  - $E_1 \cup E_2$
  - $E_1 - E_2$
  - $E_1 \times E_2$
  - $\sigma_p (E_1)$, $P$ is a predicate on attributes in $E_1$
  - $\prod_S (E_1)$, $S$ is a list consisting of some of the attributes in $E_1$
  - $\rho_x (E_1)$, x is the new name for the result of $E_1$

# Additional Operations

We define additional operations that do not add any power to the relational algebra, but that simplify common queries.

- Set intersection
- Natural join
- Division
- Assignment

# Set-Intersection Operation

- Notation: $r \cap s$
- Defined as:
- $r \cap s = \{\, t \mid t \in r \textbf{ and } t \in s \,\}$
- Assume:
  - $r$, $s$ have the *same arity*
  - attributes of r and s are compatible
- Note: $r \cap s = r - (r - s)$

# Set-Intersection Operation - Example

- Relation r, s:

| A | B |
|---|---|
| ? | 1 |
| ? | 2 |
| ? | 1 |

r

| A | B |
|---|---|
| ? | 2 |
| ? | 3 |

s

- r ∩ s

| A | B |
|---|---|
| ? | 2 |

# Natural-Join Operation

- Notation: $r \bowtie s$

- Let $r$ and $s$ be relations on schemas $R$ and $S$ respectively. Then, $r \bowtie s$ is a relation on schema $R \cup S$ obtained as follows:
  - Consider each pair of tuples $t_r$ from $r$ and $t_s$ from $s$.
  - If $t_r$ and $t_s$ have the same value on each of the attributes in $R \cap S$, add a tuple $t$ to the result, where
    - $t$ has the same value as $t_r$ on $r$
    - $t$ has the same value as $t_s$ on $s$

- Example:
  $R = (A, B, C, D)$
  $S = (E, B, D)$
  - Result schema = $(A, B, C, D, E)$
  - $r \bowtie s$ is defined as:

    $$\Pi_{r.A,\ r.B,\ r.C,\ r.D,\ s.E}\ (\sigma_{r.B\ =\ s.B\ \wedge\ r.D\ =\ s.D}\ (r \times s))$$

    $\bowtie$

# Natural Join Operation – Example

- Relations r, s:

| A | B | C | D |
|---|---|---|---|
| ⍰ | 1 | ⍰ | a |
| ⍰ | 2 | ⍰ | a |
| ⍰ | 4 | ⍰ | b |
| ⍰ | 1 | ⍰ | a |
| ⍰ | 2 | ⍰ | b |

r

| B | D | E |
|---|---|---|
| 1 | a | ⍰ |
| 3 | a | ⍰ |
| 1 | a | ⍰ |
| 2 | b | ⍰ |
| 3 | b | ⍰ |

s

r ⋈ s

| A | B | C | D | E |
|---|---|---|---|---|
| ⍰ | 1 | ⍰ | a | ⍰ |
| ⍰ | 1 | ⍰ | a | ⍰ |
| ⍰ | 1 | ⍰ | a | ⍰ |
| ⍰ | 1 | ⍰ | a | ⍰ |
| ⍰ | 2 | ⍰ | b | ⍰ |

# Division Operation

$$r \div s$$

- Suited to queries that include the phrase "for all".

- Let $r$ and $s$ be relations on schemas R and S respectively where
  - $R = (A_1, ..., A_m, B_1, ..., B_n)$
  - $S = (B_1, ..., B_n)$
  
  The result of $r \div s$ is a relation on schema
  
  $R - S = (A_1, ..., A_m)$

$$r \div s = \{\, t \mid t \in \prod_{R\text{-}S}(r) \land \forall\, u \in s\,(\, tu \in r\,)\,\}$$

# Division Operation – Example

Relations r, s:

| A | B |
|---|---|
| ⍰ | 1 |
| ⍰ | 2 |
| ⍰ | 3 |
| ⍰ | 1 |
| ⍰ | 1 |
| ⍰ | 1 |
| ⍰ | 3 |
| ⍰ | 4 |
| ⍰ | 6 |
| ⍰ | 1 |
| ⍰ | 2 |

r

| B |
|---|
| 1 |
| 2 |

s

r ÷ s:

| A |
|---|
| ⍰ |
| ⍰ |

# Another Division Example

Relations r, s:

| A | B | C | D | E |
|---|---|---|---|---|
| ⍰ | a | ⍰ | a | 1 |
| ⍰ | a | ⍰ | a | 1 |
| ⍰ | a | ⍰ | b | 1 |
| ⍰ | a | ⍰ | a | 1 |
| ⍰ | a | ⍰ | b | 3 |
| ⍰ | a | ⍰ | a | 1 |
| ⍰ | a | ⍰ | b | 1 |
| ⍰ | a | ⍰ | b | 1 |

r

| D | E |
|---|---|
| a | 1 |
| b | 1 |

s

r ⍰ s:

| A | B | C |
|---|---|---|
| ⍰ | a | ⍰ |
| ⍰ | a | ⍰ |

# Division Operation (Cont.)

- Property
  - Let $q - r \div s$
  - Then $q$ is the largest relation satisfying $q \times s \subseteq r$

- Definition in terms of the basic algebra operation
  Let $r(R)$ and $s(S)$ be relations, and let $S \subseteq R$

$$r \div s = \prod_{R\text{-}S} (r) - \prod_{R\text{-}S} ( (\prod_{R\text{-}S} (r) \times s) - \prod_{R\text{-}S,S}(r))$$

To see why
  - $\prod_{R\text{-}S,S}(r)$ simply reorders attributes of $r$

  - $\prod_{R\text{-}S}(\prod_{R\text{-}S} (r) \times s) - \prod_{R\text{-}S,S}(r))$ gives those tuples t in

    $\prod_{R\text{-}S} (r)$ such that for some tuple $u \in s$, $tu \notin r$.

# Assignment Operation

- The assignment operation ($\leftarrow$) provides a convenient way to express complex queries.
    - Write query as a sequential program consisting of
    - a series of assignments
    - followed by an expression whose value is displayed as a result of the query.
    - Assignment must always be made to a temporary relation variable.

- Example:  Write $r \div s$ as

$$temp1 \leftarrow \prod_{R\text{-}S} (r)$$
$$temp2 \leftarrow \prod_{R\text{-}S} ((temp1 \times s) - \prod_{R\text{-}S,S} (r))$$
$$result = temp1 - temp2$$

- The result to the right of the $\leftarrow$ is assigned to the relation variable on the left of the $\leftarrow$.

- May use variable in subsequent expressions.

# Example Queries

- Find all customers who have an account from at least the "Downtown" and the Uptown" branches.

Query 1

$$\Pi_{CN}(\sigma_{BN=\text{"Downtown"}}(\text{depositor} \bowtie \text{account})) \cap$$

$$\Pi_{CN}(\sigma_{BN=\text{"Uptown"}}(\text{depositor} \bowtie \text{account}))$$

where CN denotes customer-name and BN denotes branch-name.

Query 2

$$\Pi_{\text{customer-name, branch-name}}(\text{depositor} \bowtie \text{account})$$

$$\div \rho_{\text{temp(branch-name)}}(\{(\text{"Downtown"}), (\text{"Uptown"})\})$$

# Example Queries

- Find all customers who have an account at all branches located in Brooklyn city.

$\Pi_{\text{customer-name, branch-name}}$ (depositor ⋈ account)

$\div$ $\Pi_{\text{branch-name}}$ ($\sigma_{\text{branch-city = "Brooklyn"}}$ (branch))

# Extended Relational-Algebra-Operations

- Generalized Projection

- Outer Join

- Aggregate Functions

# Generalized Projection

- Extends the projection operation by allowing arithmetic functions to be used in the projection list.

$$\prod_{F_1, F_2, ..., Fn}(E)$$

- $E$ is any relational-algebra expression

- Each of $F_1$, $F_2$, ..., $F_n$ are are arithmetic expressions involving constants and attributes in the schema of $E$.

- Given relation *credit-info(customer-name, limit, credit-balance)*, find how much more each person can spend:

$$\prod_{customer\text{-}name,\ limit\ -\ credit\text{-}balance} (credit\text{-}info)$$

# Aggregate Functions and Operations

- **Aggregation function** takes a collection of values and returns a single value as a result.

    **avg**:  average value
    **min**:  minimum value
    **max**:  maximum value
    **sum**:  sum of values
    **count**:  number of values

- **Aggregate operation** in relational algebra

$$_{G1, G2, ..., Gn}\, g\, _{F1(\,A1\,),\, F2(\,A2\,),...,\, Fn(\,An\,)} (E)$$

  - $E$ is any relational-algebra expression
  - $G_1, G_2 ..., G_n$ is a list of attributes on which to group (can be empty)
  - Each $F_i$ is an aggregate function
  - Each $A_i$ is an attribute name

# Aggregate Operation – Example

• Relation *r*:

| A | B | C |
|---|---|---|
| ? | ? | 7 |
| ? | ? | 7 |
| ? | ? | 3 |
| ? | ? | 10 |

$g_{\text{sum(c)}}{}^{(r)}$

| sum-C |
|-------|
| 27 |

# Aggregate Operation – Example

- Relation *account* grouped by *branch-name*:

| branch-name | account-number | balance |
|---|---|---|
| Perryridge | A-102 | 400 |
| Perryridge | A-201 | 900 |
| Brighton | A-217 | 750 |
| Brighton | A-215 | 750 |
| Redwood | A-222 | 700 |

$$_{\text{branch-name}}\ g\ _{\text{sum(balance)}}\ (account)$$

| branch-name | balance |
|---|---|
| Perryridge | 1300 |
| Brighton | 1500 |
| Redwood | 700 |

# Aggregate Functions (Cont.)

- Result of aggregation does not have a name
  - Can use rename operation to give it a name
  - For convenience, we permit renaming as part of aggregate operation

$$\text{branch-name } g \text{ } \mathbf{sum}(\text{balance}) \text{ } \mathbf{as} \text{ sum-balance } (\text{account})$$

# Outer Join

- An extension of the join operation that avoids loss of information.

- Computes the join and then adds tuples form one relation that does not match tuples in the other relation to the result of the join.

- Uses *null* values:
  - *null* signifies that the value is unknown or does not exist
  - All comparisons involving *null* are (roughly speaking) **false** by definition.
    - Will study precise meaning of comparisons with nulls later

# Outer Join – Example

- Relation *loan*

| loan-number | branch-name | amount |
|-------------|-------------|--------|
| L-170 | Downtown | 3000 |
| L-230 | Redwood | 4000 |
| L-260 | Perryridge | 1700 |

 Relation borrower

| customer-name | loan-number |
|---------------|-------------|
| Jones | L-170 |
| Smith | L-230 |
| Hayes | L-155 |

# Outer Join – Example

- **Inner Join**

*loan ⋈ Borrower*

| loan-number | branch-name | amount | customer-name |
|-------------|-------------|--------|---------------|
| L-170<br>L-230 | Downtown<br>Redwood | 3000<br>4000 | Jones<br>Smith |

## ⬜ Left Outer Join

loan ⟕ Borrower

| loan-number | branch-name | amount | customer-name |
|-------------|-------------|--------|---------------|
| L-170<br>L-230<br>L-260 | Downtown<br>Redwood<br>Perryridge | 3000<br>4000<br>1700 | Jones<br>Smith<br>null |

# Outer Join – Example

- **Right Outer Join**

  *loan* ⟖ *borrower*

| loan-number | branch-name | amount | customer-name |
|---|---|---|---|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |
| L-155 | null | null | Hayes |

**⬚ Full Outer Join**

loan ⟗ borrower

| loan-number | branch-name | amount | customer-name |
|---|---|---|---|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |
| L-260 | Perryridge | 1700 | null |
| L-155 | null | null | Hayes |

# Null Values

- It is possible for tuples to have a null value, denoted by *null*, for some of their attributes

- *null* signifies an unknown value or that a value does not exist.

- The result of any arithmetic expression involving *null* is *null.*

- Aggregate functions simply ignore null values
  - Is an arbitrary decision.  Could have returned null as result instead.
  - We follow the semantics of SQL in its handling of null values

- For duplicate elimination and grouping, null is treated like any other value, and two nulls are assumed to be  the same
  - Alternative: assume each null is different from each other
  - Both are arbitrary decisions,  so we simply follow SQL

# Null Values

- Comparisons with null values return the special truth value *unknown*
  - If *false* was used instead of *unknown*, then   *not (A < 5)*
                    would not be equivalent to        *A >= 5*

- Three-valued logic using the truth value *unknown*:
  - OR: (*unknown* **or** *true*)        = *true*,
     (*unknown* **or** *false*)       = *unknown*
     (*unknown* **or** *unknown*) = *unknown*
  - AND:   *(true* **and** *unknown)*        = *unknown*,
        *(false* **and** *unknown)*        = *false*,
        *(unknown* **and** *unknown) = unknown*
  - NOT*:  (***not*** *unknown) = unknown*
  - In SQL "*P* **is unknown**" evaluates to true if predicate *P* evaluates to *unknown*

- Result of select  predicate is treated as *false* if it evaluates to *unknown*

# Modification of the Database

- The content of the database may be modified using the following operations:
  - Deletion
  - Insertion
  - Updating

- All these operations are expressed using the assignment operator.

# Deletion

- A delete request is expressed similarly to a query, except instead of displaying tuples to the user, the selected tuples are removed from the database.

- Can delete only whole tuples; cannot delete values on only particular attributes

- A deletion is expressed in relational algebra by:

$$r \leftarrow r - E$$

where $r$ is a relation and $E$ is a relational algebra query.

# Deletion Examples

- Delete all account records in the Perryridge branch.

$$account \leftarrow account - \sigma_{branch\text{-}name = \text{"Perryridge"}} (account)$$

☐ Delete all loan records with amount in the range of 0 to 50

$$loan \leftarrow loan - \sigma_{amount \geq 0\ and\ amount \leq 50} (loan)$$

☐ Delete all accounts at branches located in Needham.

$r_1 \leftarrow \sigma_{branch\text{-}city = \text{"Needham"}} (account \bowtie branch)$

$r_2 \leftarrow \Pi_{branch\text{-}name,\ account\text{-}number,\ balance} (r_1)$

$r_3 \leftarrow \Pi_{customer\text{-}name,\ account\text{-}number} (r_2 \bowtie depositor)$

$account \leftarrow account - r_2$

$depositor \leftarrow depositor - r_3$

# Insertion

- To insert data into a relation, we either:
  - specify a tuple to be inserted
  - write a query whose result is a set of tuples to be inserted

- in relational algebra, an insertion is expressed by:

$$r \leftarrow r \cup E$$

  where $r$ is a relation and $E$ is a relational algebra expression.

- The insertion of a single tuple is expressed by letting $E$ be a constant relation containing one tuple.

# Insertion Examples

- Insert information in the database specifying that Smith has $1200 in account A-973 at the Perryridge branch.

$$account \leftarrow account \cup \{(\text{"Perryridge"}, A\text{-}973, 1200)\}$$
$$depositor \leftarrow depositor \cup \{(\text{"Smith"}, A\text{-}973)\}$$

☐ Provide as a gift for all loan customers in the Perryridge branch, a $200 savings account.  Let the loan number serve as the account number for the new savings account.

$$r_1 \leftarrow (\sigma_{branch\text{-}name = \text{"Perryridge"}} (borrower \bowtie loan))$$
$$account \leftarrow account \cup \Pi_{branch\text{-}name, account\text{-}number, 200} (r_1)$$
$$depositor \leftarrow depositor \cup \Pi_{customer\text{-}name, loan\text{-}number}(r_1)$$

# Updating

- A mechanism to change a value in a tuple without charging *all* values in the tuple

- Use the generalized projection operator to do this task

$$r \leftarrow \prod_{F_1, F_2, \ldots, Fl,} (r)$$

- Each $F_i$ is either
  - the $i$th attribute of $r$, if the $i$th attribute is not updated, or,
  - if the attribute is to be updated $F_i$ is an expression, involving only constants and the attributes of $r$, which gives the new value for the attribute

# Update Examples

- Make interest payments by increasing all balances by 5 percent.

$$\text{account} \leftarrow \Pi_{AN, BN, BAL * 1.05} (\text{account})$$

where AN, BN and BAL stand for account-number, branch-name and balance, respectively.

☐ Pay all accounts with balances over $10,000 6 percent interest and pay all others 5 percent

$$\text{account} \leftarrow \quad \Pi_{AN, BN, BAL * 1.06} (\sigma_{BAL > 10000} (\text{account}))$$

$$\cup \; \Pi_{AN, BN, BAL * 1.05} (\sigma_{BAL \leq 10000} (\text{account}))$$

# Views

- In some cases, it is not desirable for all users to see the entire logical model (i.e., all the actual relations stored in the database.)

- Consider a person who needs to know a customer's loan number but has no need to see the loan amount.  This person should see a relation described, in the relational algebra, by

$$\prod_{customer\text{-}name,\ loan\text{-}number} (borrower \bowtie loan)$$

- Any relation that is not of the conceptual model but is made visible to a user as a "virtual relation" is called a **view**.

# View Definition

- A view is defined using the **create view** statement which has the form

    **create view** *v* **as** <query expression

    where <query expression> is any legal relational algebra query expression. The view name is represented by *v*.

- Once a view is defined, the view name can be used to refer to the virtual relation that the view generates.

- View definition is not the same as creating a new relation by evaluating the query expression
    - Rather, a view definition causes the saving of an expression; the expression is substituted into queries using the view.

# View Examples

- Consider the view (named *all-customer*) consisting of branches and their customers.

**create view** all-customer **as**

$\Pi_{\text{branch-name, customer-name}}$ (depositor $\bowtie$ account)

$\cup \ \Pi_{\text{branch-name, customer-name}}$ (borrower $\bowtie$ loan)

□ We can find all customers of the Perryridge branch by writing:

$\Pi_{\text{branch-name}}$
($\sigma_{\text{branch-name = "Perryridge"}}$ (all-customer))

# Updates Through View

- Database modifications expressed as views must be translated to modifications of the actual relations in the database.

- Consider the person who needs to see all loan data in the *loan* relation except *amount*. The view given to the person, *branch-loan*, is defined as:

  **create view** *branch-loan* **as**

  $$\prod_{branch\text{-}name,\ loan\text{-}number} (loan)$$

- Since we allow a view name to appear wherever a relation name is allowed, the person may write:

  *branch-loan* $\leftarrow$ *branch-loan* $\cup$ {("Perryridge", L-37)}

# Updates Through Views (Cont.)

- The previous insertion must be represented by an insertion into the actual relation *loan* from which the view *branch-loan* is constructed.

- An insertion into *loan* requires a value for *amount*. The insertion can be dealt with by either.
  - rejecting the insertion and returning an error message to the user.
  - inserting a tuple ("L-37", "Perryridge", *null*) into the *loan* relation

- Some updates through views are impossible to translate into database relation updates
  - create view v as $\sigma_{branch\text{-}name\ =\ "Perryridge"}(account))$
    v ← v ∪ (L-99, Downtown, 23)

- Others cannot be translated uniquely
  - *all-customer* ← *all-customer* ∪ {("Perryridge", "John")}
    - Have to choose loan or account, and create a new loan/account number!

# Views Defined Using Other Views

- One view may be used in the expression defining another view

- A view relation $v_1$ is said to *depend directly* on a view relation $v_2$ if $v_2$ is used in the expression defining $v_1$

- A view relation $v_1$ is said to *depend on* view relation $v_2$ if either $v_1$ depends directly to $v_2$ or there is a path of dependencies from $v_1$ to $v_2$

- A view relation $v$ is said to be *recursive* if it depends on itself.

# View Expansion

- A way to define the meaning of views defined in terms of other views.

- Let view $v_1$ be defined by an expression $e_1$ that may itself contain uses of view relations.

- View expansion of an expression repeats the following replacement step:

  **repeat**
  Find any view relation $v_i$ in $e_1$
  Replace the view relation $v_i$ by the expression defining $v_i$
  **until** no more view relations are present in $e_1$

- As long as the view definitions are not recursive, this loop will terminate

# Tuple Relational Calculus

- A nonprocedural query language, where each query is of the form

$$\{t \mid P(t)\}$$

- It is the set of all tuples $t$ such that predicate $P$ is true for $t$

- $t$ is a *tuple variable*, $t[A]$ denotes the value of tuple $t$ on attribute $A$

- $t \in r$ denotes that tuple $t$ is in relation $r$

- $P$ is a *formula* similar to that of the predicate calculus

# Predicate Calculus Formula

1. Set of attributes and constants

2. Set of comparison operators: (e.g., $<, \leq, =, \neq, >, \geq$)

3. Set of connectives: and ($\wedge$), or ($\vee$), not ($\neg$)

4. Implication ($\Rightarrow$): $x \Rightarrow y$, if x if true, then y is true

$$x \Rightarrow y \equiv \neg x \vee y$$

5. Set of quantifiers:
   - $\exists\, t \in r\,(Q(t)) \equiv$ "there exists" a tuple in $t$ in relation $r$ such that predicate $Q(t)$ is true
   - $\forall t \in r\,(Q(t)) \equiv Q$ is true "for all" tuples $t$ in relation $r$

# Banking Example

- *branch (branch-name, branch-city, assets)*
- *customer (customer-name, customer-street, customer-city)*
- *account (account-number, branch-name, balance)*
- *loan (loan-number, branch-name, amount)*
- *depositor (customer-name, account-number)*
- *borrower (customer-name, loan-number)*

# Example Queries

- Find the *loan-number, branch-name,* and *amount* for loans of over $1200

$$\{t \mid t \in loan \land t\,[amount] > 1200\}$$

☐Find the loan number for each loan of an amount greater than $1200

$$\{t \mid \exists\, s \in loan\, (t[loan\text{-}number] = s[loan\text{-}number] \land s\,[amount] > 1200)\}$$

Notice that a relation on schema [loan-number] is implicitly defined by the query

# Example Queries

- Find the names of all customers having a loan, an account, or both at the bank

$$\{t \mid \exists s \in borrower(\ t[customer\text{-}name] = s[customer\text{-}name])$$
$$\lor \exists u \in depositor(\ t[customer\text{-}name] = u[customer\text{-}name])$$

□ Find the names of all customers who have a loan and an account at the bank

$$\{t \mid \exists s \in borrower(\ t[customer\text{-}name] = s[customer\text{-}name])$$
$$\land \exists u \in depositor(\ t[customer\text{-}name] = u[customer\text{-}name])$$

# Example Queries

- Find the names of all customers having a loan at the Perryridge branch

$$\{t \mid \exists s \in \text{borrower}(t[\text{customer-name}] = s[\text{customer-name}]$$
$$\wedge \ \exists u \in \text{loan}(u[\text{branch-name}] = \text{"Perryridge"}$$
$$\wedge \ u[\text{loan-number}] = s[\text{loan-number}]))\}$$

◻ Find the names of all customers who have a loan at the Perryridge branch, but no account at any branch of the bank

$$\{t \mid \exists s \in \text{borrower}( t[\text{customer-name}] = s[\text{customer-name}]$$
$$\wedge \ \exists u \in \text{loan}(u[\text{branch-name}] = \text{"Perryridge"}$$
$$\wedge \ u[\text{loan-number}] = s[\text{loan-number}]))$$
$$\wedge \ \textbf{not} \ \exists v \in \text{depositor} (v[\text{customer-name}] = $$
$$t[\text{customer-name}]) \}$$

# Example Queries

- Find the names of all customers having a loan from the Perryridge branch, and the cities they live in

$$\{t \mid \exists s \in loan(s[\text{branch-name}] = \text{``Perryridge''}$$
$$\land \exists u \in borrower \ (u[\text{loan-number}] = s[\text{loan-number}]$$
$$\land \ t \ [\text{customer-name}] = u[\text{customer-name}])$$
$$\land \ \exists v \in customer \ (u[\text{customer-name}] = v[\text{customer-name}]$$
$$\land \ t[\text{customer-city}] = v[\text{customer-city}]))\}$$

# Example Queries

- Find the names of all customers who have an account at all branches located in Brooklyn:

$$\{t \mid \exists\, c \in \text{customer} \ (t[\text{customer.name}] = c[\text{customer-name}]) \ \land$$
$$\forall\, s \in \text{branch}(s[\text{branch-city}] = \text{“Brooklyn”} \Rightarrow$$
$$\exists\, u \in \text{account} \ ( s[\text{branch-name}] = u[\text{branch-name}]$$
$$\land\ \exists\, s \in \text{depositor} \ ( t[\text{customer-name}] = s[\text{customer-name}]$$
$$\land\ s[\text{account-number}] = u[\text{account-number}] \ )) \ )\}$$

# Safety of Expressions

- It is possible to write tuple calculus expressions that generate infinite relations.

- For example, $\{t \mid \neg\, t \in r\}$ results in an infinite relation if the domain of any attribute of relation $r$ is infinite

- To guard against the problem, we restrict the set of allowable expressions to safe expressions.

- An expression $\{t \mid P(t)\}$ in the tuple relational calculus is *safe* if every component of $t$ appears in one of the relations, tuples, or constants that appear in $P$
  - NOTE: this is more than just a syntax condition.
    - E.g. $\{\, t \mid t[A]{=}5 \vee \mathbf{true}\, \}$ is not safe --- it defines an infinite set with attribute values that do not appear in any relation or tuples or constants in $P$.

# Domain Relational Calculus

- A nonprocedural query language equivalent in power to the tuple relational calculus

- Each query is an expression of the form:

$$\{<x_1, x_2, ..., x_n> \mid P(x_1, x_2, ..., x_n)\}$$

  - $x_1, x_2, ..., x_n$ represent domain variables
  - $P$ represents a formula similar to that of the predicate calculus

# Safety of Expressions

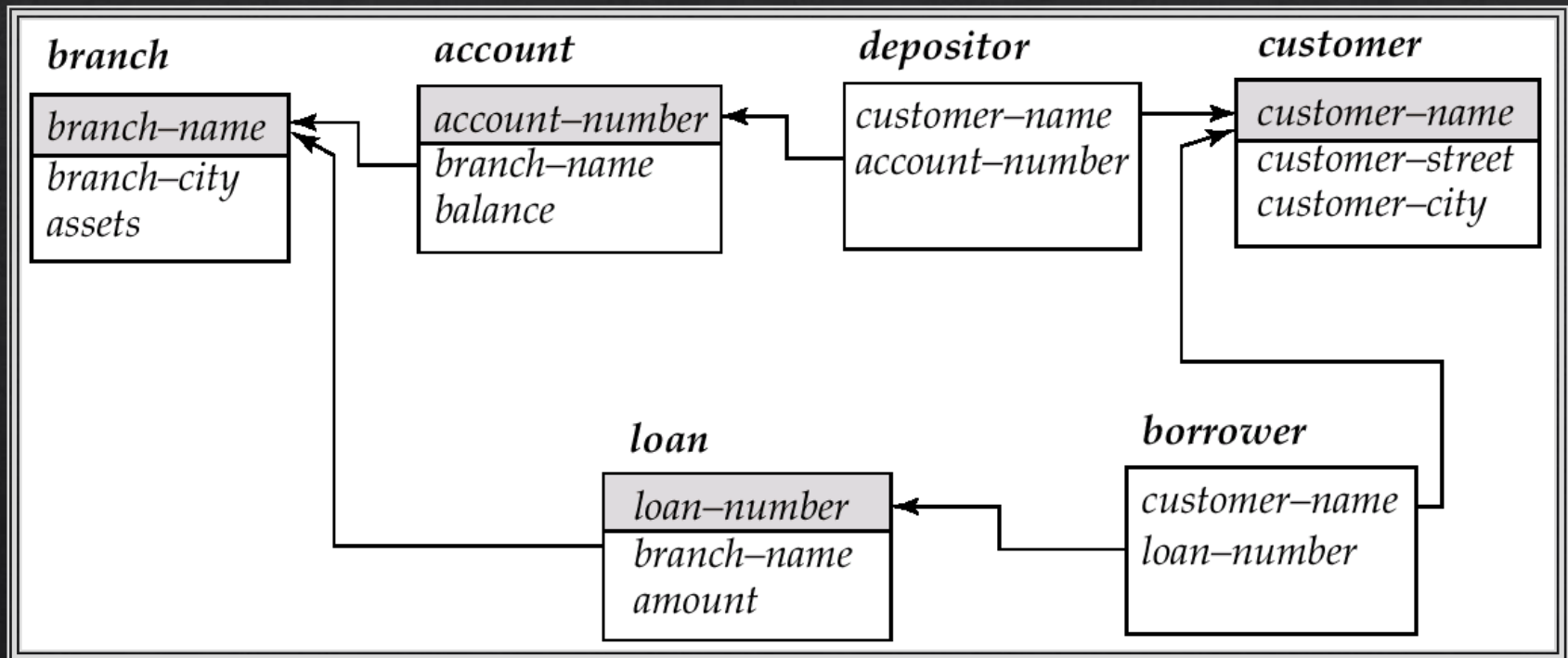$$\{<x_1, x_2, ..., x_n>\,|\,P(x_1, x_2, ..., x_n)\}$$

is safe if all of the following hold:

1.  All values that appear in tuples of the expression are values from $dom(P)$ (that is, the values appear either in $P$ or in a tuple of a relation mentioned in $P$).

2.  For every "there exists" subformula of the form $\exists\,x\,(P_1(x))$, the subformula is true if and only if there is a value of $x$ in $dom(P_1)$ such that $P_1(x)$ is true.

3. For every "for all" subformula of the form $\forall_x\,(P_1\,(x))$, the subformula is true if and only if $P_1(x)$ is true for all values $x$ from $dom\,(P_1)$.

# SQL

- ◈ Basic Structure
- ◈ Set Operations
- ◈ Aggregate Functions
- ◈ Null Values
- ◈ Nested Subqueries
- ◈ Derived Relations
- ◈ Views
- ◈ Modification of the Database
- ◈ Joined Relations
- ◈ Data Definition Language

# Schema Used in Examples

# Basic Structure

⬥ SQL is based on set and relational operations with certain modifications and enhancements

⬥ A typical SQL query has the form:

**select** $A_1, A_2, …, A_n$
**from** $r_1, r_2, …, r_m$
**where** $P$

  ⬥ $A_i s$ represent attributes

  ⬥ $r_i s$ represent relations

  ⬥ $P$ is a predicate.

⬥ This query is equivalent to the relational algebra expression.

$$\Pi_{A1,\ A2,\ …,\ An}(\sigma_P(r_1 \ \text{x} \ r_2 \ \text{x} \ … \ \text{x} \ r_m))$$

⬥ The result of an SQL query is a relation.

# The select Clause

- The **select** clause list the attributes desired in the result of a query

  - corresponds to the projection operation of the relational algebra

- E.g. find the names of all branches in the *loan* relation
  **select** *branch-name*
  **from** *loan*

- In the "pure" relational algebra syntax, the query would be:

$$\Pi_{\text{branch-name}}(loan)$$

- NOTE:  SQL does not permit the '-' character in names,

  - Use, e.g., *branch_name* instead of *branch-name* in a real implementation.

  - We use '-' since it looks nicer!

- NOTE:  SQL names are case insensitive, i.e. you can use capital or small letters.

  - You may wish to use upper case where-ever we use bold font.

# The select Clause (Cont.)

✧ SQL allows duplicates in relations as well as in query results.

✧ To force the elimination of duplicates, insert the keyword **distinct** after **select.**

✧ Find the names of all branches in the *loan* relations, and remove duplicates

> **select distinct** *branch-name*
> **from** *loan*

✧ The keyword **all** specifies that duplicates not be removed.

> **select all** *branch-name*
> **from** *loan*

# The select Clause (Cont.)

◆ An asterisk in the select clause denotes "all attributes"

> **select** *
> **from** *loan*

◆ The **select** clause can contain arithmetic expressions involving the operation, +, −, *, and /, and operating on constants or attributes of tuples.

◆ The query:

> **select** *loan-number, branch-name, amount* * 100
> **from** *loan*

would return a relation which is the same as the *loan* relations, except that the attribute *amount* is multiplied by 100.

# The where Clause

◆ The **where** clause specifies conditions that the result must satisfy

  ◇ corresponds to the selection predicate of the relational algebra.

◆ To find all loan number for loans made at the Perryridge branch with loan amounts greater than $1200.

> **select** *loan-number*
> **from** *loan*
> **where** *branch-name* = 'Perryridge' **and** *amount* > 1200

◆ Comparison results can be combined using the logical connectives **and, or,** and **not.**

◆ Comparisons can be applied to results of arithmetic expressions.

# The where Clause (Cont.)

◈ SQL includes a **between** comparison operator

◈ E.g. Find the loan number of those loans with loan amounts between $90,000 and $100,000 (that is, ≥$90,000 and ≤$100,000)

```
select loan-number
    from loan
    where amount between 90000 and 100000
```

# The from Clause

◈ The **from** clause lists the relations involved in the query

　◇ corresponds to the Cartesian product operation of the relational algebra.

◈ Find the Cartesian product *borrower x loan*        **select** *

        **from** *borrower, loan*

☐ Find the name, loan number and loan amount of all customers having a loan at the Perryridge branch.

      **select** *customer-name, borrower.loan-number, amount*
        **from** *borrower, loan*
        **where**   *borrower.loan-number = loan.loan-number* **and**
           *branch-name =* 'Perryridge'

# The Rename Operation

◇ The SQL allows renaming relations and attributes using the **as** clause:
*old-name* **as** *new-name*

◇ Find the name, loan number and loan amount of all customers; rename the column name *loan-number* as *loan-id.*

**select** *customer-name, borrower.loan-number* **as** *loan-id, amount*
**from** *borrower, loan*
**where** *borrower.loan-number = loan.loan-number*

# Tuple Variables

◈ Tuple variables are defined in the **from** clause via the use of the **as** clause.

◈ Find the customer names and their loan numbers for all customers having a loan at some branch.

> **select** *customer-name, T.loan-number, S.amount*
>     **from** *borrower* **as** *T, loan* **as** *S*
>     **where** *T.loan-number = S.loan-number*

▯ Find the names of all branches that have greater assets than some branch located in Brooklyn.

> **select distinct** *T.branch-name*
>     **from** *branch* **as** *T, branch* **as** *S*
>     **where** *T.assets* > *S.assets* **and** *S.branch-city = 'Brooklyn'*

# String Operations

- SQL includes a string-matching operator for comparisons on character strings. Patterns are described using two special characters:

  - percent (%). The % character matches any substring.

  - underscore (_). The _ character matches any character.

- Find the names of all customers whose street includes the substring "Main".

    **select** *customer-name*
    **from** *customer*
    **where** *customer-street* **like** '%Main%'

- Match the name "Main%"

    **like** 'Main\%' **escape** '\'

- SQL supports a variety of string operations such as

  - concatenation (using "||")

  - converting from upper to lower case (and vice versa)

  - finding string length, extracting substrings, etc.

# Ordering the Display of Tuples

◈ List in alphabetic order the names of all customers having a loan in Perryridge branch

> **select distinct** *customer-name*
> **from**    *borrower, loan*
> **where** *borrower loan-number - loan.loan-number* **and**
>          *branch-name* = 'Perryridge'
> **order by** *customer-name*

◈ We may specify **desc** for descending order or **asc** for ascending order, for each attribute; ascending order is the default.

   ◈ E.g.  **order by** *customer-name* **desc**

# Duplicates

◇ In relations with duplicates, SQL can define how many copies of tuples appear in the result.

◇ *Multiset* versions of some of the relational algebra operators – given multiset relations $r_1$ and $r_2$:

1. $\sigma_\theta(r_1)$: If there are $c_1$ copies of tuple $t_1$ in $r_1$, and $t_1$ satisfies selections $\sigma_\theta$, then there are $c_1$ copies of $t_1$ in $\sigma_\theta(r_1)$.

2. $\Pi_A(r)$: For each copy of tuple $t_1$ in $r_1$, there is a copy of tuple $\Pi_A(t_1)$ in $\Pi_A(r_1)$ where $\Pi_A(t_1)$ denotes the projection of the single tuple $t_1$.

3. $r_1 \times r_2$ : If there are $c_1$ copies of tuple $t_1$ in $r_1$ and $c_2$ copies of tuple $t_2$ in $r_2$, there are $c_1 \times c_2$ copies of the tuple $t_1. t_2$ in $r_1 \times r_2$

# Duplicates (Cont.)

⬥ Example: Suppose multiset relations $r_1$ $(A, B)$ and $r_2$ $(C)$ are as follows:

$$r_1 = \{(1, a)\ (2, a)\} \qquad r_2 = \{(2), (3), (3)\}$$

⬥ Then $\Pi_B(r_1)$ would be $\{(a), (a)\}$, while $\Pi_B(r_1) \times r_2$ would be

$$\{(a,2), (a,2), (a,3), (a,3), (a,3), (a,3)\}$$

⬥ SQL duplicate semantics:

> **select** $A_1,, A_2, \ldots, A_n$
> **from** $r_1, r_2, \ldots, r_m$
> **where** $P$

is equivalent to the *multiset* version of the expression:

$$\Pi_{A1,, A2, \ldots, An}(\sigma_P (r_1 \times r_2 \times \ldots \times r_m))$$

# Set Operations

◇ The set operations **union, intersect,** and **except** operate on relations and correspond to the relational algebra operations $\cup$, $\cap$, $-$.

◇ Each of the above operations automatically eliminates duplicates; to retain all duplicates use the corresponding multiset versions **union all, intersect all** and **except all.**

Suppose a tuple occurs $m$ times in $r$ and $n$ times in $s$, then, it occurs:

  ◇ $m + n$ times in $r$ **union all** $s$

  ◇ $\min(m,n)$ times in $r$ **intersect all** $s$

  ◇ $\max(0, m - n)$ times in $r$ **except all** $s$

# Set Operations

⬦ Find all customers who have a loan, an account, or both:

> (**select** *customer-name* **from** *depositor*)
> **union**
> (**select** *customer-name* **from** *borrower)*

□ Find all customers who have both a loan and an account.

> (**select** *customer-name* **from** *depositor*)
> **intersect**
> (**select** *customer-name* **from** *borrower)*

□ Find all customers who have an account but no loan.

> (**select** *customer-name* **from** *depositor*)
> **except**
> (**select** *customer-name* **from** *borrower)*

# Aggregate Functions

◆ These functions operate on the multiset of values of a column of a relation, and return a value

**avg:** average value
**min:**  minimum value
**max:**  maximum value
**sum:**  sum of values
**count:**  number of values

# Aggregate Functions (Cont.)

◈ Find the average account balance at the Perryridge branch.

> **select avg** *(balance)*
> **from** *account*
> **where** *branch-name* = 'Perryridge'

▯ Find the number of tuples in the *customer* relation.

> **select count** (*)
> **from** *customer*

▯ Find the number of depositors in the bank.

> **select count (distinct** *customer-name)*
> **from** *depositor*

# Aggregate Functions – Group By

◈ Find the number of depositors for each branch.

> **select** *branch-name,* **count (distinct** *customer-name)*
> **from** *depositor, account*
> **where** *depositor.account-number = account.account-number*
> **group by** *branch-name*

Note:  Attributes in **select** clause outside of aggregate functions must
   appear in **group by** list

# Aggregate Functions – Having Clause

◈ Find the names of all branches where the average account balance is more than $1,200.

> **select** *branch-name,* **avg** *(balance)*
> **from** *account*
> **group by** *branch-name*
> **having avg** *(balance)* > 1200

Note:  predicates in the **having** clause are applied after the formation of groups whereas predicates in the **where** clause are applied before forming groups

# Null Values

◈ It is possible for tuples to have a null value, denoted by *null*, for some of their attributes

◈ *null* signifies an unknown value or that a value does not exist.

◈ The predicate **is null** can be used to check for null values.

    ◇ E.g. Find all loan number which appear in the *loan* relation with null values for *amount*.

    **select** *loan-number*
    **from** *loan*
    **where** *amount* **is null**

◈ The result of any arithmetic expression involving *null* is *null*

    ◇ E.g.  5 + null  returns null

◈ However, aggregate functions simply ignore nulls

    ◇ more on this shortly

# Null Values and Three Valued Logic

◇ Any comparison with *null* returns *unknown*

  ◇ *E.g.  5 < null   or   null <> null    or    null = null*

◇ Three-valued logic using the truth value *unknown*:

  ◇ OR: (*unknown* **or** *true*) = *true*, (*unknown* **or** *false*) = *unknown*
      (*unknown* **or** *unknown*) = *unknown*

  ◇ AND: *(true* **and** *unknown) = unknown,    (false* **and** *unknown) = false,*
       *(unknown* **and** *unknown) = unknown*

  ◇ NOT*: (***not** *unknown) = unknown*

  ◇ "*P* **is unknown**" evaluates to true if predicate *P* evaluates to *unknown*

◇ Result of **where** clause predicate is treated as *false* if it evaluates to *unknown*

# Null Values and Aggregates

◇ Total all loan amounts

> **select sum** (*amount*)
> **from** *loan*

   ◇ Above statement ignores null amounts

   ◇ result is null if there is no non-null amount, that is the

◇ All aggregate operations except **count(\*)** ignore tuples with null values on the aggregated attributes.

# Nested Subqueries

- SQL provides a mechanism for the nesting of subqueries.

- A subquery is a **select-from-where** expression that is nested within another query.

- A common use of subqueries is to perform tests for set membership, set comparisons, and set cardinality.

# Example Query

◇ Find all customers who have both an account and a loan at the bank.

**select distinct** *customer-name*
    **from** *borrower*
    **where** *customer-name* **in (select** *customer-name*
                      **from** depositor)

Find all customers who have a loan at the bank but do not have an account at the bank

**select distinct** *customer-name*
    **from** *borrower*
    **where** *customer-name* **not in (select** *customer-name*
                      **from** *depositor)*

# Example Query

◈ Find all customers who have both an account and a loan at the Perryridge branch

**select distinct** *customer-name*
    **from** *borrower, loan*
    **where** *borrower.loan-number = loan.loan-number* **and**
        *branch-name = "Perryridge"* **and**
        *(branch-name, customer-name)* **in**
          **(select** *branch-name, customer-name*
          **from** *depositor, account*
          **where** *depositor.account-number =*
                *account.account-number)*

☐ Note: Above query can be written in a much simpler manner. The formulation above is simply to illustrate SQL features.

(Schema used in this example)

# Set Comparison

◈ Find all branches that have greater assets than some branch located in Brooklyn.

> **select distinct**  *T.branch-name*
>     **from** *branch* **as** *T, branch* **as** *S*
>     **where**  *T.assets* > *S.assets* **and**
>             *S.branch-city* = 'Brooklyn'

▯ Same query using > **some** clause

> **select** *branch-name*
>     **from** *branch*
>     **where** *assets* > **some**
>         **(select** *assets*
>          **from** *branch*
>             **where** *branch-city* = 'Brooklyn'**)**

# Definition of Some Clause

◈ F <comp> **some** $r \Leftrightarrow \exists\, t \in r$ *s.t.* (F <comp> t)
Where <comp> can be: $<, \leq, >, =, \neq$

(5< **some** [0 / 5 / 6] ) = true
         (read:  5 < some tuple in the relation)

(5< **some** [0 / 5] ) = false

(5 = **some** [0 / 5] ) = true

(5 $\neq$ **some** [0 / 5] ) = true (since $0 \neq 5$)

(= **some**) ≡ **in**
However, ($\neq$ **some**) ≡ **not in**

# Definition of all Clause

◇ F \<comp\> **all** $r \Leftrightarrow \forall\, t \in r$ (F \<comp\> $t$)

$(5 < \textbf{all}\;\boxed{\begin{matrix}0\\5\\6\end{matrix}}\;) = \text{false}$

$(5 < \textbf{all}\;\boxed{\begin{matrix}6\\10\end{matrix}}\;) = \text{true}$

$(5 = \textbf{all}\;\boxed{\begin{matrix}4\\5\end{matrix}}\;) = \text{false}$

$(5 \neq \textbf{all}\;\boxed{\begin{matrix}4\\6\end{matrix}}\;) = \text{true (since } 5 \neq 4 \text{ and } 5 \neq 6)$

$(\neq \textbf{all}) \equiv \textbf{not in}$
However, $(= \textbf{all}) \equiv \textbf{in}$

# Example Query

◈ Find the names of all branches that have greater assets than all branches located in Brooklyn.

> **select** *branch-name*
>     **from** *branch*
>     **where** *assets* **> all**
>         **(select** *assets*
>         **from** *branch*
>         **where** *branch-city* = 'Brooklyn')

# Test for Empty Relations

- The **exists** construct returns the value **true** if the argument subquery is nonempty.

- **exists** $r \Leftrightarrow r \neq \varnothing$

- **not exists** $r \Leftrightarrow r = \varnothing$

# Example Query

◈ Find all customers who have an account at all branches located in Brooklyn.

**select distinct** *S.customer-name*
    **from** *depositor* **as** *S*
    **where not exists (**
        **(select** *branch-name*
        **from** *branch*
        **where** *branch-city* = 'Brooklyn')
      **except**
        **(select** *R.branch-name*
        **from** *depositor* **as** *T, account* **as** *R*
        **where** *T.account-number* = *R.account-number* **and**
            *S.customer-name* = *T.customer-name))*

☐ (Schema used in this example)

☐ Note that $X - Y = \emptyset \iff X \subseteq Y$

☐ *Note:* Cannot write this query using = **all** and its variants

# Test for Absence of Duplicate Tuples

- The **unique** construct tests whether a subquery has any duplicate tuples in its result.

- Find all customers who have at most one account at the Perryridge branch.

  **select** *T.customer-name*
  **from** *depositor* **as** *T*
  **where unique** (

      **select** *R.customer-name*
      **from** *account, depositor* **as** *R*
      **where** *T.customer-name = R.customer-name* **and**
            *R.account-number = account.account-number* **and**
            *account.branch-name =* 'Perryridge')

- (Schema used in this example)

# Example Query

◈ Find all customers who have at least two accounts at the Perryridge branch.

> **select distinct** *T.customer-name*
> **from** *depositor T*
> **where not unique** (
> > **select** *R.customer-name*
> > **from** *account, depositor* **as** *R*
> > **where** *T.customer-name = R.customer-name*
>
> **and**
> > *R.account-number = account.account-number*
>
> **and**
> > *account.branch-name* = 'Perryridge')

☐ (Schema used in this example)

# Views

◈ Provide a mechanism to hide certain data from the view of certain users.  To create a view we use the command:

**create view** *v* **as** <query expression>

where:

☐ <query expression> is any legal expression

☐ The view name is represented by *v*

# Example Queries

◈ A view consisting of branches and their customers

> **create view** *all-customer* **as**
>   (**select** *branch-name, customer-name*
>    **from** *depositor, account*
>    **where** *depositor.account-number = account.account-number)*
>     **union**
>   (**select** *branch-name, customer-name*
>    **from** *borrower, loan*
>    **where** *borrower.loan-number = loan.loan-number)*

☐ Find all customers of the Perryridge branch

> **select** *customer-name*
>       **from** *all-customer*
>       **where** *branch-name* = 'Perryridge'

# Derived Relations

◇ Find the average account balance of those branches where the average account balance is greater than $1200.

> **select** *branch-name, avg-balance*
> **from (select** *branch-name,* **avg** *(balance)*
>         **from** *account*
>         **group by** *branch-name)*
>         **as** *result (branch-name, avg-balance)*
> **where** *avg-balance* > 1200

Note that we do not need to use the **having** clause, since we compute the temporary (view) relation *result* in the **from** clause, and the attributes of *result* can be used directly in the **where** clause.

# With Clause

◈ With clause allows views to be defined locally to a query, rather than globally. Analogous to procedures in a programming language.

◈ Find all accounts with the maximum balance

> **with** *max-balance*(*value*) **as**
>     **select** max (*balance*)
>     **from** *account*
>  **select** *account-number*
> **from** *account, max-balance*
> **where** *account.balance = max-balance.value*

# Complex Query using With Clause

◈ Find all branches where the total account deposit is greater than the average of the total account deposits at all branches.

      **with** *branch-total* (*branch-name*, *value*) **as**
        **select** *branch-name*, **sum** (*balance*)
        **from** *account*
        **group by** *branch-name*
      **with** *branch-total-avg*(*value*) **as**
        **select avg** (*value*)
        **from** *branch-total*
      **select** *branch-name*
      **from** *branch-total*, *branch-total-avg*
      **where** *branch-total.value >= branch-total-avg.value*

# Modification of the Database – Deletion

◈ Delete all account records at the Perryridge branch

> **delete from** *account*
> **where** *branch-name* = 'Perryridge'

◈ Delete all accounts at every branch located in Needham city.

**delete from** *account*
**where** *branch-name* **in** (**select** *branch-name*
　　　　　　　　　　　　　　**from** *branch*
　　　　　　　　　　　　　　**where** *branch-city* = 'Needham')

*delete from depositor*
**where** *account-number* **in**
　　　　　　(**select** *account-number*
　　　　　　 **from** *branch, account*
　　　　　　 **where** *branch-city* = 'Needham'
　　　　　　　 **and** *branch.branch-name* = *account.branch-name)*

◈ (Schema used in this example)

# Example Query

◈ Delete the record of all accounts with balances below the average at the bank.

**delete from** *account*
    **where** *balance* < (**select avg** *(balance)*
    **from** *account)*

☐  Problem:  as we delete tuples from *deposit,* the average balance changes

☐  Solution used in SQL:

1.  First, compute **avg** balance and find all tuples to delete

2.  Next, delete all tuples found above (without recomputing **avg** or retesting the tuples)

# Modification of the Database – Insertion

◈ Add a new tuple to *account*

> **insert into** *account*
> > **values** ('A-9732', 'Perryridge',1200)

or equivalently

**insert into** *account (branch-name, balance, account-number)*
> **values** ('Perryridge', 1200, 'A-9732')

◈ Add a new tuple to *account* with *balance* set to null

> **insert into** *account*
> > **values** ('A-777','Perryridge', *null*)

# Modification of the Database – Insertion

◈ Provide as a gift for all loan customers of the Perryridge branch, a $200 savings account.  Let the loan number serve as the account number for the new savings account

> **insert into** *account*
>    **select** *loan-number, branch-name,*  200
>    **from** *loan*
>    **where** *branch-name* = 'Perryridge'
> **insert into** *depositor*
>    **select** *customer-name, loan-number*
>    **from** *loan, borrower*
>    **where** branch-name =  'Perryridge'
>          **and** *loan.account-number = borrower.account-number*

◈ The select from where statement is fully evaluated before any of its results are inserted into the relation (otherwise queries like
>    **insert into** *table*1 **select** * **from** *table*1
would cause problems

# Modification of the Database – Updates

◈ Increase all accounts with balances over $10,000 by 6%, all other accounts receive 5%.

 ◇ Write two **update** statements:

> **update** *account*
> **set** *balance = balance* ∗ 1.06
> **where** *balance* > 10000

> **update** *account*
> **set** *balance = balance* ∗ 1.05
> **where** *balance* ≤ 10000

 ◇ The order is important

 ◇ Can be done better using the **case** statement (next slide)

# Case Statement for Conditional Updates

◇ Same query as before: Increase all accounts with balances over $10,000 by 6%, all other accounts receive 5%.

```
update account
set balance =  case
                    when balance <= 10000 then balance *1.05
                    else   balance * 1.06
                end
```

# Update of a View

◈ Create a view of all loan data in *loan* relation, hiding the *amount* attribute

> **create view** *branch-loan* **as**
> **select** *branch-name, loan-number*
> **from** *loan*

◈ Add a new tuple to *branch-loan*

> **insert into** *branch-loan*
> **values** ('Perryridge', 'L-307')

This insertion must be represented by the insertion of the tuple

('L-307', 'Perryridge', *null*)

into the *loan* relation

◈ Updates on more complex views are difficult or impossible to translate, and hence are disallowed.

◈ Most SQL implementations allow updates only on simple views (without aggregates) defined on a single relation

# Transactions

- A transaction is a sequence of queries and update statements executed as a single unit
  - Transactions are started implicitly and terminated by one of
    - **commit work:** makes all updates of the transaction permanent in the database
    - **rollback work:** undoes all updates performed by the transaction.
- Motivating example
  - Transfer of money from one account to another involves two steps:
    - deduct from one account and credit to another
  - If one steps succeeds and the other fails, database is in an inconsistent state
  - Therefore, either both steps should succeed or neither should
- If any step of a transaction fails, all work done by the transaction can be undone by **rollback work.**
- Rollback of incomplete transactions is done automatically, in case of system failures

# Transactions (Cont.)

◈ In most database systems, each SQL statement that executes successfully is automatically committed.

   ◇ Each transaction would then consist of only a single statement

   ◇ Automatic commit can usually be turned off, allowing multi-statement transactions, but how to do so depends on the database system

   ◇ Another option in SQL:1999: enclose statements within
   **begin atomic**
      **…**
      **end**

# Joined Relations

◈ Join operations take two relations and return as a result another relation.

◈ These additional operations are typically used as subquery expressions in the **from** clause

◈ Join condition – defines which tuples in the two relations match, and what attributes are present in the result of the join.

◈ Join type – defines how tuples in each relation that do not match any tuple in the other relation (based on the join condition) are treated.

| Join Types |
|---|
| **inner join** |
| **left outer join** |
| **right outer join** |
| **full outer join** |

| Join Conditions |
|---|
| **natural** |
| **on** <predicate> |
| **using** ($A_1$, $A_2$, ..., $A_n$) |

# Joined Relations – Datasets for Examples

◇ Relation *loan*

| loan-number | branch-name | amount |
|-------------|-------------|--------|
| L-170 | Downtown | 3000 |
| L-230 | Redwood | 4000 |
| L-260 | Perryridge | 1700 |

☐ Relation *borrower*

| customer-name | loan-number |
|---------------|-------------|
| Jones | L-170 |
| Smith | L-230 |
| Hayes | L-155 |

☐ Note: borrower information missing for L-260 and loan information missing for L-155

# Joined Relations – Examples

◈ *loan* **inner join** *borrower* **on**
*loan.loan-number = borrower.loan-number*

| loan-number | branch-name | amount | customer-name | loan-number |
|---|---|---|---|---|
| L-170 | Downtown | 3000 | Jones | L-170 |
| L-230 | Redwood | 4000 | Smith | L-230 |

▢ *loan* **left outer join** *borrower* **on**
*loan.loan-number = borrower.loan-number*

| loan-number | branch-name | amount | customer-name | loan-number |
|---|---|---|---|---|
| L-170 | Downtown | 3000 | Jones | L-170 |
| L-230 | Redwood | 4000 | Smith | L-230 |
| L-260 | Perryridge | 1700 | null | null |

# Joined Relations – Examples

◈ *loan* **natural inner join** *borrower*

| loan-number | branch-name | amount | customer-name |
|:-----------:|-------------|:------:|---------------|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |

☐ loan **natural right outer join** *borrower*

| loan-number | branch-name | amount | customer-name |
|:-----------:|-------------|:------:|---------------|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |
| L-155 | null | null | Hayes |

# Joined Relations – Examples

◇ *loan* **full outer join** *borrower* **using** *(loan-number)*

| loan-number | branch-name | amount | customer-name |
|---|---|---|---|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |
| L-260 | Perryridge | 1700 | *null* |
| L-155 | null | null | Hayes |

☐ Find all customers who have either an account or a loan (but not both) at the bank.

    **select** *customer-name*
        **from** (*depositor* **natural full outer join** *borrower*)
        **where** *account-number* **is** *null* **or** *loan-number* **is** *null*

# Data Definition Language (DDL)

Allows the specification of not only a set of relations but also information about each relation, including:

⬦ The schema for each relation.

⬦ The domain of values associated with each attribute.

⬦ Integrity constraints

⬦ The set of indices to be maintained for each relations.

⬦ Security and authorization information for each relation.

⬦ The physical storage structure of each relation on disk.

# Domain Types in SQL

- ◈ **char(n).** Fixed length character string, with user-specified length $n$.

- ◈ **varchar(n).** Variable length character strings, with user-specified maximum length $n$.

- ◈ **int.** Integer (a finite subset of the integers that is machine-dependent).

- ◈ **smallint.** Small integer (a machine-dependent subset of the integer domain type).

- ◈ **numeric(p,d).** Fixed point number, with user-specified precision of $p$ digits, with $n$ digits to the right of decimal point.

- ◈ **real, double precision.** Floating point and double-precision floating point numbers, with machine-dependent precision.

- ◈ **float(n).** Floating point number, with user-specified precision of at least $n$ digits.

- ◈ Null values are allowed in all the domain types. Declaring an attribute to be **not null** prohibits null values for that attribute.

- ◈ **create domain** construct in SQL-92 creates user-defined domain types

    **create domain** *person-name* **char**(20) **not null**

# Date/Time Types in SQL (Cont.)

◇ **date.** Dates, containing a (4 digit) year, month and date
  - ◇ E.g. **date** '2001-7-27'
◇ **time.** Time of day, in hours, minutes and seconds.
  - ◇ E.g. **time** '09:00:30'   **time** '09:00:30.75'
◇ **timestamp**: date plus time of day
  - ◇ E.g. **timestamp** '2001-7-27 09:00:30.75'
◇ **Interval**: period of time
  - ◇ E.g. Interval '1' day
  - ◇ Subtracting a date/time/timestamp value from another gives an interval value
  - ◇ Interval values can be added to date/time/timestamp values
◇ Can extract values of individual fields from date/time/timestamp
  - ◇ E.g. **extract (year from** r.starttime)
◇ Can cast string types to date/time/timestamp
  - ◇ E.g. **cast** <string-valued-expression> **as date**

# Storage and File Structure

▶ Overview of Physical Storage Media

▶ Magnetic Disks

▶ RAID

▶ Tertiary Storage

▶ Storage Access

▶ File Organization

▶ Organization of Records in Files

▶ Data-Dictionary Storage

# Classification of Physical Storage Media

- Speed with which data can be accessed
- Cost per unit of data
- Reliability
  - data loss on power failure or system crash
  - physical failure of the storage device
- Can differentiate storage into:
  - **volatile storage:** loses contents when power is switched off
  - **non-volatile storage**:
    - Contents persist even when power is switched off.
    - Includes secondary and tertiary storage, as well as batter-backed up main-memory.

# Physical Storage Media

- **Cache** – fastest and most costly form of storage; volatile; managed by the computer system hardware.

- **Main memory**:
  - fast access (10s to 100s of nanoseconds; 1 nanosecond = $10^{-9}$ seconds)
  - generally too small (or too expensive) to store the entire database
    - capacities of up to a few Gigabytes widely used currently
    - Capacities have gone up and per-byte costs have decreased steadily and rapidly  (roughly factor of 2 every 2 to 3 years)
  - **Volatile** — contents of main memory are usually lost if a power failure or system crash occurs.

# Physical Storage Media (Cont.)

- **Flash memory**
  - Data survives power failure
  - Data can be written at a location only once, but location can be erased and written to again
    - Can support only a limited number of write/erase cycles.
    - Erasing of memory has to be done to an entire bank of memory
  - Reads are roughly as fast as main memory
  - But writes are slow (few microseconds), erase is slower
  - Cost per unit of storage roughly similar to main memory
  - Widely used in embedded devices such as digital cameras
  - also known as EEPROM (Electrically Erasable Programmable Read-Only Memory)

# Physical Storage Media (Cont.)

- **Magnetic-disk**
  - Data is stored on spinning disk, and read/written magnetically
  - Primary medium for the long-term storage of data; typically stores entire database.
  - Data must be moved from disk to main memory for access, and written back for storage
    - Much slower access than main memory (more on this later)
  - **direct-access** – possible to read data on disk in any order, unlike magnetic tape
  - Hard disks  vs  floppy disks
  - Capacities range up to roughly 100 GB currently
    - Much larger capacity and cost/byte than main memory/flash memory
    - Growing constantly and rapidly with technology improvements (factor of 2 to 3 every 2 years)
  - Survives power failures and system crashes
    - disk failure can destroy data, but is very rare

# Physical Storage Media (Cont.)

- **Optical storage**
  - non-volatile, data is read optically from a spinning disk using a laser
  - CD-ROM (640 MB) and DVD (4.7 to 17 GB) most popular forms
  - Write-one, read-many (WORM) optical disks used for archival storage (CD-R and DVD-R)
  - Multiple write versions also available (CD-RW, DVD-RW, and DVD-RAM)
  - Reads and writes are slower than with magnetic disk
  - **Juke-box** systems, with large numbers of removable disks, a few drives, and a mechanism for automatic loading/unloading of disks available for storing large volumes of data

# Physical Storage Media (Cont.)

- **Tape storage**
  - non-volatile, used primarily for backup (to recover from disk failure), and for archival data
  - **sequential-access** – much slower than disk
  - very high capacity (40 to 300 GB tapes available)
  - tape can be removed from drive $\Rightarrow$ storage costs much cheaper than disk, but drives are expensive
  - Tape jukeboxes available for storing massive amounts of data
    - hundreds of terabytes (1 terabyte = $10^9$ bytes) to even a petabyte (1 petabyte = $10^{12}$ bytes)

# Storage Hierarchy

# Storage Hierarchy (Cont.)

- **primary storage:** Fastest media but volatile (cache, main memory).
- **secondary storage:** next level in hierarchy, non-volatile, moderately fast access time
  - also called **on-line storage**
  - E.g. flash memory, magnetic disks
- **tertiary storage:** lowest level in hierarchy, non-volatile, slow access time
  - also called **off-line storage**
  - E.g. magnetic tape, optical storage

# Magnetic Hard Disk Mechanism



**NOTE: Diagram is schematic, and simplifies the structure of actual disk drives**

# Magnetic Disks

- **Read-write head**
  - Positioned very close to the platter surface (almost touching it)
  - Reads or writes magnetically encoded information.
- Surface of platter divided into circular **tracks**
  - Over 16,000 tracks per platter on typical hard disks
- Each track is divided into **sectors.**
  - A sector is the smallest unit of data that can be read or written.
  - Sector size typically 512 bytes
  - Typical sectors per track: 200 (on inner tracks) to 400 (on outer tracks)
- To read/write a sector
  - disk arm swings to position head on right track
  - platter spins continually; data is read/written as sector passes under head
- Head-disk assemblies
  - multiple disk platters on a single spindle (typically 2 to 4)
  - one head per platter, mounted on a common arm.
- **Cylinder** $i$ consists of $i^{th}$ track of all the platters

# Magnetic Disks (Cont.)

- Earlier generation disks were susceptible to head-crashes
  - Surface of earlier generation disks had metal-oxide coatings which would disintegrate on head crash and damage all data on disk
  - Current generation disks are less susceptible to such disastrous failures, although individual sectors may get corrupted
- **Disk controller** – interfaces between the computer system and the disk drive hardware.
  - accepts high-level commands to read or write a sector
  - initiates actions such as moving the disk arm to the right track and actually reading or writing the data
  - Computes and attaches **checksums** to each sector to verify that data is read back correctly
    - If data is corrupted, with very high probability stored checksum won't match recomputed checksum
  - Ensures successful writing by reading back sector after writing it
  - Performs remapping of bad sectors

# Disk Subsystem



- ▶ Multiple disks connected to a computer system through a controller
    - ▶ Controllers functionality (checksum, bad sector remapping) often carried out by individual disks; reduces load on controller
- ▶ Disk interface standards families
    - ▶ ATA (AT adaptor) range of standards
    - ▶ SCSI (Small Computer System Interconnect) range of standards
    - ▶ Several variants of each standard (different speeds and capabilities)

# Performance Measures of Disks

- **Access time** – the time it takes from when a read or write request is issued to when data transfer begins.  Consists of:
  - **Seek time** – time it takes to reposition the arm over the correct track.
    - Average seek time is 1/2 the worst case seek time.
      - Would be 1/3 if all tracks had the same number of sectors, and we ignore the time to start and stop arm movement
    - 4 to 10 milliseconds on typical disks
  - **Rotational latency** – time it takes for the sector to be accessed to appear under the head.
    - Average latency is 1/2 of the worst case latency.
    - 4 to 11 milliseconds on typical disks (5400 to 15000 r.p.m.)
- **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.
  - 4 to 8 MB per second is typical
  - Multiple disks may share a controller, so rate that controller can handle is also important
    - E.g. ATA-5: 66 MB/second,  SCSI-3: 40 MB/s
    - Fiber Channel: 256 MB/s

# Performance Measures (Cont.)

- **Mean time to failure (MTTF)** – the average time the disk is expected to run continuously without any failure.
  - Typically 3 to 5 years
  - Probability of failure of new disks is quite low, corresponding to a "theoretical MTTF" of 30,000 to 1,200,000 hours for a new disk
    - E.g., an MTTF of 1,200,000 hours for a new disk means that given 1000 relatively new disks, on an average one will fail every 1200 hours
  - MTTF decreases as disk ages

# Optimization of Disk-Block Access

- **Block** – a contiguous sequence of sectors from a single track
  - data is transferred between disk and main memory in blocks
  - sizes range from 512 bytes to several kilobytes
    - Smaller blocks: more transfers from disk
    - Larger blocks:  more space wasted due to partially filled blocks
    - Typical block sizes today range from 4 to 16 kilobytes
- **Disk-arm-scheduling** algorithms order pending accesses to tracks so that disk arm movement is minimized
  - **elevator algorithm** : move disk arm in one direction (from outer to inner tracks or vice versa), processing next request in that direction, till no more requests in that direction, then reverse direction and repeat

# Optimization of Disk Block Access (Cont.)

- **File organization** – optimize block access time by organizing the blocks to correspond to how data will be accessed
    - E.g. Store related information on the same or nearby cylinders.
    - Files may get **fragmented** over time
        - E.g. if data is inserted to/deleted from the file
        - Or free blocks on disk are scattered, and newly created file has its blocks scattered over the disk
        - Sequential access to a fragmented file results in increased disk arm movement
    - Some systems have utilities to defragment the file system, in order to speed up file access

# Optimization of Disk Block Access (Cont.)

- **Nonvolatile write buffers** speed up disk writes by writing blocks to a non-volatile RAM buffer immediately
  - Non-volatile RAM:  battery backed up RAM or flash memory
    - Even if power fails, the data is safe and will be written to disk when power returns
  - Controller then writes to disk whenever the disk has no other requests or request has been pending for some time
  - Database operations that require data to be safely stored before continuing can continue without waiting for data to be written to disk
  - *Writes can be reordered to minimize disk arm movement*
- **Log disk** – a disk devoted to writing a sequential log of block updates
  - Used exactly like nonvolatile RAM
    - Write to log disk is very fast since no seeks are required
    - No need for special hardware (NV-RAM)
- File systems typically reorder writes to disk to improve performance
  - **Journaling file systems** write data in safe order to NV-RAM or log disk
  - Reordering without journaling: risk of corruption of file system data

# RAID

- **RAID: Redundant Arrays of Independent Disks**
  - disk organization techniques that manage a large numbers of disks, providing a view of a single disk of
    - high capacity and high speed  by using multiple disks in parallel, and
    - high reliability by storing data redundantly, so that data can be recovered even if  a disk fails
- The chance that some disk out of a set of $N$ disks will fail is much higher than the chance that a specific single disk will fail.
  - E.g., a system with 100 disks, each with MTTF of 100,000 hours (approx. 11 years), will have a system MTTF of 1000 hours (approx. 41 days)
  - Techniques for using redundancy to avoid data loss are critical with large numbers of disks
- Originally a cost-effective alternative to large, expensive disks
  - I in RAID originally stood for ``inexpensive''
  - Today RAIDs are used for their higher reliability and bandwidth.
    - The "I" is interpreted as independent

# Improvement of Reliability via Redundancy

- **Redundancy** – store extra information that can be used to rebuild information lost in a disk failure

- E.g., **Mirroring** (or **shadowing**)
  - Duplicate every disk.  Logical disk consists of two physical disks.
  - Every write is carried out on both disks
    - Reads can take place from either disk
  - If one disk in a pair fails, data still available in the other
    - Data loss would occur only if a disk fails, and its mirror disk also fails before the system is repaired
      - Probability of combined event is very small
        - Except for dependent failure modes such as fire or building collapse or electrical power surges

- Mean time to data loss depends on mean time to failure, and mean time to repair
  - E.g. MTTF of 100,000 hours, mean time to repair of 10 hours gives mean time to data loss of $500*10^6$ hours (or 57,000 years) for a mirrored pair of disks (ignoring dependent failure modes)

# Improvement in Performance via Parallelism

- Two main goals of parallelism in a disk system:
    1. Load balance multiple small accesses to increase throughput
    2. Parallelize large accesses to reduce response time.
- Improve transfer rate by striping data across multiple disks.
- **Bit-level striping** – split the bits of each byte across multiple disks
    - In an array of eight disks, write bit $i$ of each byte to disk $i$.
    - Each access can read data at eight times the rate of a single disk.
    - But seek/access time worse than for a single disk
        - Bit level striping is not used much any more
- **Block-level striping** – with $n$ disks, block $i$ of a file goes to disk $(i \bmod n) + 1$
    - Requests for different blocks can run in parallel if the blocks reside on different disks
    - A request for a long sequence of blocks can utilize all disks in parallel

# RAID Levels

▶ Schemes to provide redundancy at lower cost by using disk striping combined with parity bits

▶ Different RAID organizations, or RAID levels, have differing cost, performance and reliability characteristics

**RAID Level 0**:  Block striping; non-redundant.

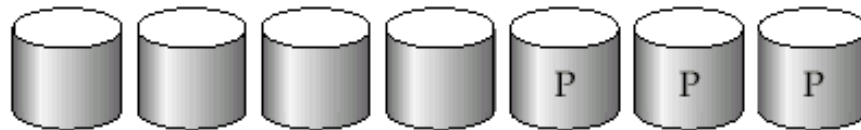☐ Used in high-performance applications where data lost is not critical.

**RAID Level 1**:  Mirrored disks with block striping

☐ Offers best write performance.

☐ Popular for applications such as storing log files in a database system.



(a) RAID 0: nonredundant striping

(b) RAID 1: mirrored disks

# RAID Levels (Cont.)

▶ **RAID Level 2**: Memory-Style Error-Correcting-Codes (ECC) with bit striping.

▶ **RAID Level 3**: Bit-Interleaved Parity

  ▶ a single parity bit is enough for error correction, not just detection, since we know which disk has failed

    ▶ When writing data, corresponding parity bits must also be computed and written to a parity bit disk

    ▶ To recover data in a damaged disk, compute XOR of bits from other disks (including parity bit disk)



(c) RAID 2: memory-style error-correcting codes

(d) RAID 3: bit-interleaved parity

# RAID Levels (Cont.)

- ► RAID Level 3 (Cont.)
  - ► Faster data transfer than with a single disk, but fewer I/Os per second since every disk has to participate in every I/O.
  - ► Subsumes Level 2 (provides all its benefits, at lower cost).
- ► **RAID Level 4:** Block-Interleaved Parity; uses block-level striping, and keeps a parity block on a separate disk for corresponding blocks from *N* other disks.
  - ► When writing data block, corresponding block of parity bits must also be computed and written to parity disk
  - ► To find value of a damaged block, compute XOR of bits from corresponding blocks (including parity block) from other disks.

(e) RAID 4: block-interleaved parity

# RAID Levels (Cont.)

▶ RAID Level 4 (Cont.)
  ▶ Provides higher I/O rates for independent block reads than Level 3
    ▶ block read goes to a single disk, so blocks stored on different disks can be read in parallel
  ▶ Provides high transfer rates for reads of multiple blocks than no-striping
  ▶ Before writing a block, parity data must be computed
    ▶ Can be done by using old parity block, old value of current block and new value of current block (2 block reads + 2 block writes)
    ▶ Or by recomputing the parity value using the new values of blocks corresponding to the parity block
      ▶ More efficient for writing large amounts of data sequentially
  ▶ Parity block becomes a bottleneck for independent block writes since every block write also writes to parity disk

# RAID Levels (Cont.)

▶ **RAID Level 5:** Block-Interleaved Distributed Parity; partitions data and parity among all *N* + 1 disks, rather than storing data in *N* disks and parity in 1 disk.

   ▶ E.g., with 5 disks, parity block for *n*th set of blocks is stored on disk (*n mod* 5) + 1, with the data blocks stored on the other 4 disks.

(f) RAID 5: block-interleaved distributed parity

| | | | | |
|---|---|---|---|---|
| P0 | 0 | 1 | 2 | 3 |
| 4 | P1 | 5 | 6 | 7 |
| 8 | 9 | P2 | 10 | 11 |
| 12 | 13 | 14 | P3 | 15 |
| 16 | 17 | 18 | 19 | P4 |

# RAID Levels (Cont.)



(g) RAID 6: P + Q redundancy

- **RAID Level 5 (Cont.)**

  - Higher I/O rates than Level 4.

    - Block writes occur in parallel if the blocks and their parity blocks are on different disks.

  - Subsumes Level 4: provides same benefits, but avoids bottleneck of parity disk.

- **RAID Level 6**: P+Q Redundancy scheme; similar to Level 5, but stores extra redundant information to guard against multiple disk failures.

  - Better reliability than Level 5 at a higher cost; not used as widely.

# Choice of RAID Level

- Factors in choosing RAID level
  - Monetary cost
  - Performance: Number of I/O operations per second, and bandwidth during normal operation
  - Performance during failure
  - Performance during rebuild of failed disk
    - Including time taken to rebuild failed disk
- RAID 0 is used only when data safety is not important
  - E.g. data can be recovered quickly from other sources
- Level 2 and 4 never used since they are subsumed by 3 and 5
- Level 3 is not used anymore since bit-striping forces single block reads to access all disks, wasting disk arm movement, which block striping (level 5) avoids
- Level 6 is rarely used since levels 1 and 5 offer adequate safety for almost all applications
- So competition is between 1 and 5 only

# Choice of RAID Level (Cont.)

- ▶ Level 1 provides much better write performance than level 5
  - ▶ Level 5 requires at least 2 block reads and 2 block writes to write a single block, whereas Level 1 only requires 2 block writes
  - ▶ Level 1 preferred for high update environments such as log disks
- ▶ Level 1 had higher storage cost than level 5
  - ▶ disk drive capacities increasing rapidly (50%/year) whereas disk access times have decreased much less (x 3 in 10 years)
  - ▶ I/O requirements have increased greatly, e.g. for Web servers
  - ▶ When enough disks have been bought to satisfy required rate of I/O, they often have spare storage capacity
    - ▶ so there is often no extra monetary cost for Level 1!
- ▶ Level 5 is preferred for applications with low update rate, and large amounts of data
- ▶ Level 1 is preferred for all other applications

# Hardware Issues

- **Software RAID**:  RAID implementations done entirely in software, with no special hardware support

- **Hardware RAID**:  RAID implementations with special hardware
  - Use non-volatile RAM to record writes that are being executed
  - Beware:  power failure during write can result in corrupted disk
    - E.g. failure after writing one block but before writing the second in a mirrored system
    - Such corrupted data must be detected when power is restored
      - Recovery from corruption is similar to recovery from failed disk
      - NV-RAM helps to efficiently detected potentially corrupted blocks
        - Otherwise all blocks of disk must be read and compared with mirror/parity block

# Hardware Issues (Cont.)

- **Hot swapping**: replacement of disk while system is running, without power down
  - Supported by some hardware RAID systems,
  - reduces time to recovery, and improves availability greatly
- Many systems maintain spare disks which are kept online, and used as replacements for failed disks immediately on detection of failure
  - Reduces time to recovery greatly
- Many hardware RAID systems ensure that a single point of failure will not stop the functioning of the system by using
  - Redundant power supplies with battery backup
  - Multiple controllers and multiple interconnections to guard against controller/interconnection failures

# Optical Disks

- Compact disk-read only memory (CD-ROM)
  - Disks can be loaded into or removed from a drive
  - High storage capacity (640 MB per disk)
  - High seek times or about 100 msec (optical read head is heavier and slower)
  - Higher latency (3000 RPM) and lower data-transfer rates (3-6 MB/s) compared to magnetic disks
- Digital Video Disk (DVD)
  - DVD-5  holds 4.7 GB , and DVD-9 holds 8.5 GB
  - DVD-10 and DVD-18 are double sided formats with capacities of 9.4 GB and 17 GB
  - Other characteristics similar to CD-ROM
- Record once versions (CD-R and DVD-R) are becoming popular
  - data can only be written once, and cannot be erased.
  - high capacity and long lifetime; used for archival storage
  - Multi-write versions (CD-RW, DVD-RW and DVD-RAM) also available

# Magnetic Tapes

- Hold large volumes of data and provide high transfer rates
  - Few GB for DAT (Digital Audio Tape) format, 10-40 GB with DLT (Digital Linear Tape) format, 100 GB+ with Ultrium format, and 330 GB with Ampex helical scan format
  - Transfer rates from few to 10s of MB/s
- Currently the cheapest storage medium
  - Tapes are cheap, but cost of drives is very high
- Very slow access time in comparison to magnetic disks and optical disks
  - limited to sequential access.
  - Some formats (Accelis) provide faster seek (10s of seconds) at cost of lower capacity
- Used mainly for backup, for storage of infrequently used information, and as an off-line medium for transferring information from one system to another.
- Tape jukeboxes used for very large capacity storage
  - (terabyte ($10^{12}$ bytes) to petabye ($10^{15}$ bytes)

# Storage Access

- A database file is partitioned into fixed-length storage units called **blocks**. Blocks are units of both storage allocation and data transfer.

- Database system seeks to minimize the number of block transfers between the disk and memory. We can reduce the number of disk accesses by keeping as many blocks as possible in main memory.

- **Buffer** – portion of main memory available to store copies of disk blocks.

- **Buffer manager** – subsystem responsible for allocating buffer space in main memory.

# Buffer Manager

▶ Programs call on the buffer manager when they need a block from disk.

1. If the block is already in the buffer, the requesting program is given the address of the block in main memory

2. If the block is not in the buffer,

   1. the buffer manager allocates space in the buffer for the block, replacing (throwing out) some other block, if required, to make space for the new block.

   2. The block that is thrown out is written back to disk only if it was modified since the most recent time that it was written to/fetched from the disk.

   3. Once space is allocated in the buffer, the buffer manager reads the block from the disk to the buffer, and passes the address of the block in main memory to requester.

# File Organization

▶ The database is stored as a collection of *files*. Each file is a sequence of *records*. A record is a sequence of fields.

▶ One approach:

   ▶ assume record size is fixed

   ▶ each file has records of one particular type only

   ▶ different files are used for different relations

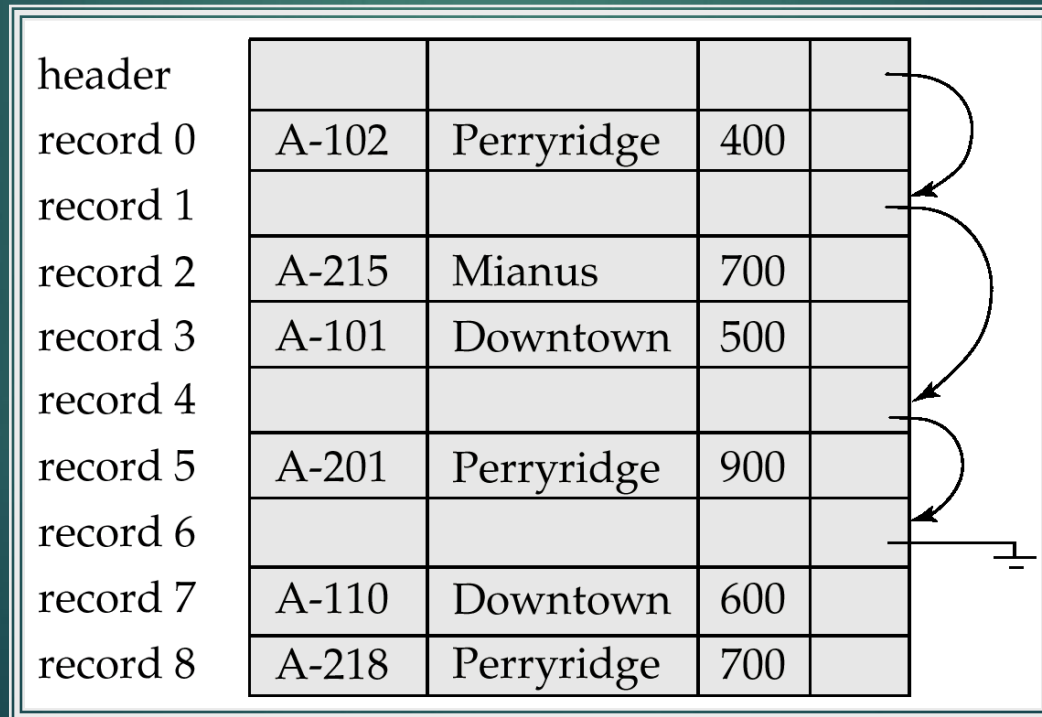   This case is easiest to implement; will consider variable length records

   later.

# Fixed-Length Records

▶ Simple approach:

  ▶ Store record $i$ starting from byte $n * (i - 1)$, where $n$ is the size of each record.

  ▶ Record access is simple but records may cross blocks

    ▶ Modification: do not allow records to cross block boundaries

▶ Deletion of record $l$: alternatives:

  ▶ move records $i + 1, \ldots, n$ to $i, \ldots, n - 1$

  ▶ move record $n$ to $i$

  ▶ do not move records, but link all free records on a *free list*

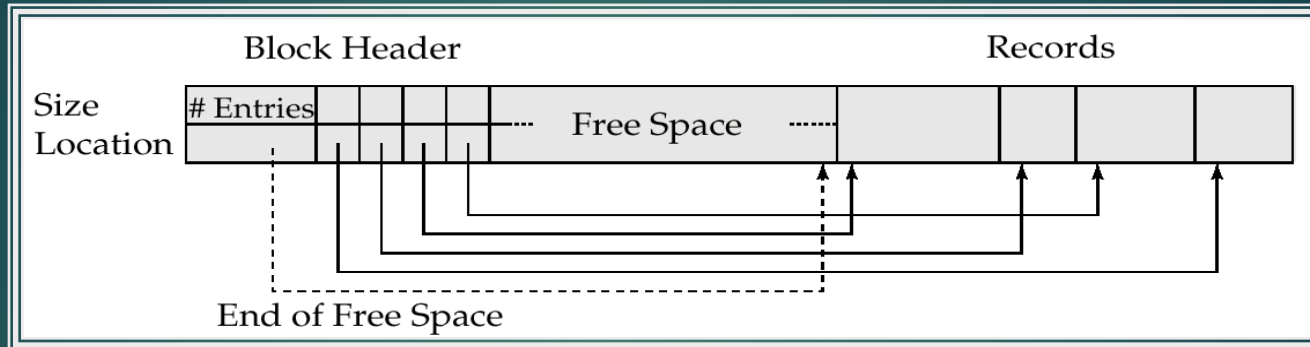| | | | |
|---|---|---|---|
| record 0 | A-102 | Perryridge | 400 |
| record 1 | A-305 | Round Hill | 350 |
| record 2 | A-215 | Mianus | 700 |
| record 3 | A-101 | Downtown | 500 |
| record 4 | A-222 | Redwood | 700 |
| record 5 | A-201 | Perryridge | 900 |
| record 6 | A-217 | Brighton | 750 |
| record 7 | A-110 | Downtown | 600 |
| record 8 | A-218 | Perryridge | 700 |

# Free Lists

▶ Store the address of the first deleted record in the file header.

▶ Use this first record to store the address of the second deleted record, and so on

▶ Can think of these stored addresses as pointers since they "point" to the location of a record.

▶ More space efficient representation: reuse space for normal attributes of free records to store pointers. (No pointers stored in in-use records.)

| | | | | |
|---|---|---|---|---|
| header | | | | |
| record 0 | A-102 | Perryridge | 400 | |
| record 1 | | | | |
| record 2 | A-215 | Mianus | 700 | |
| record 3 | A-101 | Downtown | 500 | |
| record 4 | | | | |
| record 5 | A-201 | Perryridge | 900 | |
| record 6 | | | | |
| record 7 | A-110 | Downtown | 600 | |
| record 8 | A-218 | Perryridge | 700 | |

# Variable-Length Records

▶ Variable-length records arise in database systems in several ways:

  ▶ Storage of multiple record types in a file.

  ▶ Record types that allow variable lengths for one or more fields.

  ▶ Record types that allow repeating fields (used in some older data models).

▶ Byte string representation

  ▶ Attach an *end-of-record* (⊥) control character to the end of each record

  ▶ Difficulty with deletion

  ▶ Difficulty with growth

# Variable-Length Records: Slotted Page Structure



- Slotted page header contains:
  - number of record entries
  - end of free space in the block
  - location and size of each record
- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated.
- Pointers should not point directly to record — instead they should point to the entry for the record in header.

# Variable-Length Records (Cont.)

▶ Fixed-length representation:

  ▶ reserved space

  ▶ pointers

▶ Reserved space – can use fixed-length records of a known maximum length; unused space in shorter records filled with a null or end-of-record symbol.

| 0 | Perryridge | A-102 | 400 | A-201 | 900 | A-218 | 700 |
|---|------------|-------|-----|-------|-----|-------|-----|
| 1 | Round Hill | A-305 | 350 | ⊥ | ⊥ | ⊥ | ⊥ |
| 2 | Mianus | A-215 | 700 | ⊥ | ⊥ | ⊥ | ⊥ |
| 3 | Downtown | A-101 | 500 | A-110 | 600 | ⊥ | ⊥ |
| 4 | Redwood | A-222 | 700 | ⊥ | ⊥ | ⊥ | ⊥ |
| 5 | Brighton | A-217 | 750 | ⊥ | ⊥ | ⊥ | ⊥ |

# Pointer Method

| 0 | Perryridge | A-102 | 400 | |
|---|---|---|---|---|
| 1 | Round Hill | A-305 | 350 | |
| 2 | Mianus | A-215 | 700 | |
| 3 | Downtown | A-101 | 500 | |
| 4 | Redwood | A-222 | 700 | |
| 5 | | A-201 | 900 | |
| 6 | Brighton | A-217 | 750 | |
| 7 | | A-110 | 600 | |
| 8 | | A-218 | 700 | |

▶ Pointer method

   ▶ A variable-length record is represented by a list of fixed-length records, chained together via pointers.

   ▶ Can be used even if the maximum record length is not known

# Pointer Method (Cont.)

▶ Disadvantage to pointer structure; space is wasted in all records except the first in a a chain.

▶ Solution is to allow two kinds of block in file:

   ▶ Anchor block – contains the first records of chain

   ▶ Overflow block – contains records other than those that are the first records of chairs.

| anchor block | | | |
|---|---|---|---|
| Perryridge | A-102 | 400 | |
| Round Hill | A-305 | 350 | |
| Mianus | A-215 | 700 | |
| Downtown | A-101 | 500 | |
| Redwood | A-222 | 700 | |
| Brighton | A-217 | 750 | |

| overflow block | | |
|---|---|---|
| A-201 | 900 | |
| A-218 | 700 | |
| A-110 | 600 | |

# Organization of Records in Files

▶ **Heap** – a record can be placed anywhere in the file where there is space

▶ **Sequential** – store records in sequential order, based on the value of the search key of each record

▶ **Hashing** – a hash function computed on some attribute of each record; the result specifies in which block of the file the record should be placed

▶ Records of each relation may be stored in a separate file. In a **clustering file organization** records of several different relations can be stored in the same file

    ▶ Motivation: store related records on the same block to minimize I/O

# Sequential File Organization

▶ Suitable for applications that require sequential processing of the entire file

▶ The records in the file are ordered by a search-key

| A-217 | Brighton   | 750 |
|-------|------------|-----|
| A-101 | Downtown   | 500 |
| A-110 | Downtown   | 600 |
| A-215 | Mianus     | 700 |
| A-102 | Perryridge | 400 |
| A-201 | Perryridge | 900 |
| A-218 | Perryridge | 700 |
| A-222 | Redwood    | 700 |
| A-305 | Round Hill | 350 |

# Sequential File Organization (Cont.)

▶ Deletion – use pointer chains

▶ Insertion –locate the position where the record is to be inserted

  ▶ if there is free space insert there

  ▶ if no free space, insert the record in an overflow block

  ▶ In either case, pointer chain must be updated

▶ Need to reorganize the file from time to time to restore sequential order

| A-217 | Brighton | 750 | |
|-------|----------|-----|--|
| A-101 | Downtown | 500 | |
| A-110 | Downtown | 600 | |
| A-215 | Mianus | 700 | |
| A-102 | Perryridge | 400 | |
| A-201 | Perryridge | 900 | |
| A-218 | Perryridge | 700 | |
| A-222 | Redwood | 700 | |
| A-305 | Round Hill | 350 | |

| A-888 | North Town | 800 | |
|-------|------------|-----|--|

# Clustering File Organization

▶ Simple file structure stores each relation in a separate file

▶ Can instead store several relations in one file using a **clustering** file organization

▶ E.g., clustering organization of *customer* and *depositor:*

| Hayes | Main | Brooklyn |
|-------|-------|----------|
| Hayes | A-102 | |
| Hayes | A-220 | |
| Hayes | A-503 | |
| Turner | Putnam | Stamford |
| Turner | A-305 | |

☐ good for queries involving depositor ⋈ customer, and for queries involving one single customer and his accounts

☐ bad for queries involving only customer

☐ results in variable size records

# Data Dictionary Storage

Data dictionary (also called system catalog) stores metadata:  that is, data about data, such as

- Information about relations
    - names of relations
    - names and types of attributes of each relation
    - names and definitions of views
    - integrity constraints
- User and accounting information, including passwords
- Statistical and descriptive data
    - number of tuples in each relation
- Physical file organization information
    - How relation is stored (sequential/hash/…)
    - Physical location of relation
        - operating system file name or
        - disk addresses of blocks containing records of the relation
- Information about indices (Chapter 12)

# Data Dictionary Storage (Cont.)

▶ Catalog structure:  can use either

 ▶ specialized data structures designed for efficient access

 ▶ a set of relations, with existing system features used to ensure efficient access

 The latter alternative is usually preferred

▶ A possible catalog representation:


*Relation-metadata = (relation-name, number-of-attributes,*
*storage-organization, location)*
*Attribute-metadata = (attribute-name, relation-name, domain-type,*
*position, length)*
*User-metadata = (user-name, encrypted-password, group)*
*Index-metadata = (index-name, relation-name, index-type,*
*index-attributes)*
*View-metadata = (view-name, definition)*

# DATABASE SYSTEM ARCHITECTURES

- Centralized Systems

- Client--Server Systems

- Parallel Systems

- Distributed Systems

- Network Types

# CENTRALIZED SYSTEMS

- Run on a single computer system and do not interact with other computer systems.

- General-purpose computer system: one to a few CPUs and a number of device controllers that are connected through a common bus that provides access to shared memory.

- Single-user system (e.g., personal computer or workstation): desk-top unit, single user, usually has only one CPU  and one or two hard disks; the OS may support only one user.

- Multi-user system: more disks, more memory, multiple CPUs, and a multi-user OS. Serve a large number of users who are connected to the system vie terminals. Often called *server* systems.

# A CENTRALIZED COMPUTER SYSTEM

# CLIENT-SERVER SYSTEMS

- Server systems satisfy requests generated at $m$ client systems, whose general structure is shown below:

# CLIENT-SERVER SYSTEMS (CONT.)

- Database functionality can be divided into:

  - **Back-end**: manages access structures, query evaluation and optimization, concurrency control and recovery.

  - **Front-end**: consists of tools such as *forms*, *report-writers*, and graphical user interface facilities.

- The interface between the front-end and the back-end is through SQL or through an application program interface.

| SQL user-interface | forms interface | report writer | graphical interface | front-end |
|---|---|---|---|---|
| | | | | interface (SQL + API) |
| | SQL engine | | | back-end |

# CLIENT-SERVER SYSTEMS (CONT.)

- Advantages of replacing mainframes with networks of workstations or personal computers connected to back-end server machines:
  - better functionality for the cost
  - flexibility in locating resources and expanding facilities
  - better user interfaces
  - easier maintenance
- Server systems can be broadly categorized into two kinds:
  - **transaction servers** which are widely used in relational database systems, and
  - **data servers**, used in object-oriented database systems

# TRANSACTION SERVERS

- Also called **query server** systems or SQL *server* systems; clients send requests to the server system where the transactions are executed, and results are shipped back to the client.

- Requests specified in SQL, and communicated to the server through a *remote procedure call (RPC)* mechanism.

- Transactional RPC allows many RPC calls to collectively form a transaction.

- *Open Database Connectivity (ODBC)* is a C language application program interface standard from Microsoft for connecting to a server, sending SQL requests, and receiving results.

- JDBC standard similar to ODBC, for Java

# TRANSACTION SERVER PROCESS STRUCTURE

- A typical transaction server consists of multiple processes accessing data in shared memory.

- Server processes
  - These receive user queries (transactions), execute them and send results back
  - Processes may be **multithreaded**, allowing a single process to execute several user queries concurrently
  - Typically multiple multithreaded server processes

- Lock manager process
  - More on this later

- Database writer process
  - Output modified buffer blocks to disks continually

# TRANSACTION SERVER PROCESSES (CONT.)

- Log writer process

    - Server processes simply add log records to log record buffer

    - Log writer process outputs log records to stable storage.

- Checkpoint process

    - Performs periodic checkpoints

- Process monitor process

    - Monitors other processes, and takes recovery actions if any of the other processes fail

        - E.g. aborting any transactions being executed by a server process and restarting it

# TRANSACTION SYSTEM PROCESSES (CONT.)

- Shared memory contains shared data
  - Buffer pool
  - Lock table
  - Log buffer
  - Cached query plans (reused if same query submitted again)
- All database processes can access shared memory
- To ensure that no two processes are accessing the same data structure at the same time, databases systems implement **mutual exclusion** using either
  - Operating system semaphores
  - Atomic instructions such as test-and-set

# TRANSACTION SYSTEM PROCESSES (CONT.)

- To avoid overhead of interprocess communication for lock request/grant, each database process operates directly on the lock table data structure (Section 16.1.4) instead of sending requests to lock manager process
    - Mutual exclusion ensured on the lock table using semaphores, or more commonly, atomic instructions
    - If a lock can be obtained, the lock table is updated directly in shared memory
    - If a lock cannot be immediately obtained, a lock request is noted in the lock table and the process (or thread) then waits for lock to be granted
    - When a lock is released, releasing process updates lock table to record release of lock, as well as grant of lock to waiting requests (if any)
    - Process/thread waiting for lock may either:
        - Continually scan lock table to check for lock grant, or
        - Use operating system semaphore mechanism to wait on a semaphore.
            - Semaphore identifier is recorded in the lock table
            - When a lock is granted, the releasing process signals the semaphore to tell the waiting process/thread to proceed
- Lock manager process still used for deadlock detection

# DATA SERVERS

- Used in LANs, where there is a very high speed connection between the clients and the server, the client machines are comparable in processing power to the server machine, and the tasks to be executed are compute intensive.

- Ship data to client machines where processing is performed, and then ship results back to the server machine.

- This architecture requires full back-end functionality at the clients.

- Used in many object-oriented database systems

- Issues:

  - Page-Shipping versus Item-Shipping

  - Locking

  - Data Caching

  - Lock Caching

# DATA SERVERS (CONT.)

- **Page-Shipping** versus **Item-Shipping**
  - Smaller unit of shipping $\Rightarrow$ more messages
  - Worth **prefetching** related items along with requested item
  - Page shipping can be thought of as a form of prefetching
- Locking
  - Overhead of requesting and getting locks from server is high due to message delays
  - Can grant locks on requested and prefetched items; with page shipping, transaction is granted lock on whole page.
  - Locks on a prefetched item can be P{called back} by the server, and returned by client transaction if the prefetched item has not been used.
  - Locks on the page can be **deescalated** to locks on items in the page when there are lock conflicts. Locks on unused items can then be returned to server.

# DATA SERVERS (CONT.)

- **Data Caching**
  - Data can be cached at client even in between transactions
  - But check that data is up-to-date before it is used (**cache coherency**)
  - Check can be done when requesting lock on data item

- **Lock Caching**
  - Locks can be retained by client system even in between transactions
  - Transactions can acquire cached locks locally, without contacting server
  - Server **calls back** locks from clients when it receives conflicting lock request.  Client returns lock once no local transaction is using it.
  - Similar to deescalation, but across transactions.

# PARALLEL SYSTEMS

- Parallel database systems consist of multiple processors and multiple disks connected by a fast interconnection network.

- A **coarse-grain parallel** machine consists of a small number of powerful processors

- A **massively parallel** or **fine grain parallel** machine utilizes thousands of smaller processors.

- Two main performance measures:
  - **throughput** --- the number of tasks that can be completed in a given time interval
  - **response time** --- the amount of time it takes to complete a single task from the time it is submitted

# SPEED-UP AND SCALE-UP

- **Speedup**: a fixed-sized problem executing on a small system is given to a system which is *N*-times larger.

    - Measured by:

    *speedup = small system elapsed time*

    *large system elapsed time*

    - Speedup is **linear** if equation equals N.

- **Scaleup**: increase the size of both the problem and the system

    - *N*-times larger system used to perform *N*-times larger job

    - Measured by:

    *scaleup = small system small problem elapsed time*

    *big system big problem elapsed time*

    - Scale up is **linear** if equation equals 1.

# SPEED-UP AND SCALE-UP SPEEDUP



Speedup

# SPEED-UP AND SCALE-UP
# SCALEUP



Scaleup

# BATCH AND TRANSACTION SCALEUP

- **Batch scaleup**:
  - A single large job; typical of most database queries and scientific simulation.
  - Use an $N$-times larger computer on $N$-times larger problem.

- **Transaction scaleup**:
  - Numerous small queries submitted by independent users to a shared database; typical transaction processing and timesharing systems.
  - $N$-times as many users submitting requests (hence, $N$-times as many requests) to an $N$-times larger database, on an $N$-times larger computer.
  - Well-suited to parallel execution.

# FACTORS LIMITING SPEEDUP AND SCALEUP

Speedup and scaleup are often sublinear due to:

- **Startup costs**: Cost of starting up multiple processes may dominate computation time, if the degree of parallelism is high.

- **Interference**: Processes accessing shared resources (e.g.,system bus, disks, or locks) compete with each other, thus spending time waiting on other processes, rather than performing useful work.

- **Skew**: Increasing the degree of parallelism increases the variance in service times of parallely executing tasks. Overall execution time determined by **slowest** of parallely executing tasks.

# INTERCONNECTION NETWORK ARCHITECTURES

- **Bus**. System components send data on and receive data from a single communication bus;
  - Does not scale well with increasing parallelism.

- **Mesh**. Components are arranged as nodes in a grid, and each component is connected to all adjacent components
  - Communication links grow with growing number of components, and so scales better.
  - But may require $2\sqrt{n}$ hops to send message to a node (or $\sqrt{n}$ with wraparound connections at edge of grid).

- **Hypercube**. Components are numbered in binary; components are connected to one another if their binary representations differ in exactly one bit.
  - $n$ components are connected to $log(n)$ other components and can reach each other via at most $log(n)$ links; reduces communication delays.

# INTERCONNECTION ARCHITECTURES
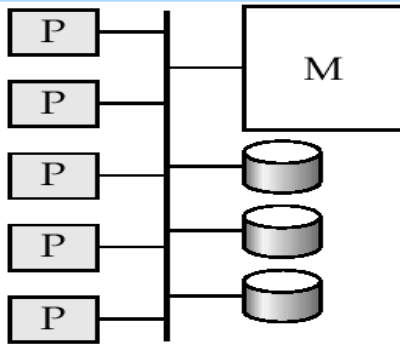


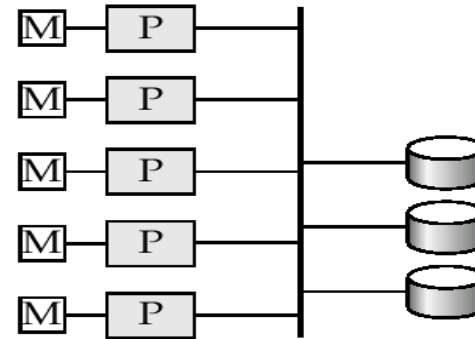(a) bus       (b) mesh       (c) hypercube

# PARALLEL DATABASE ARCHITECTURES

- **Shared memory** -- processors share a common memory

- **Shared disk** -- processors share a common disk

- **Shared nothing** -- processors share neither a common memory nor common disk

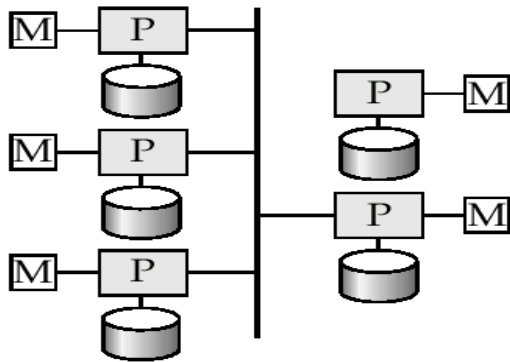- **Hierarchical** -- hybrid of the above architectures

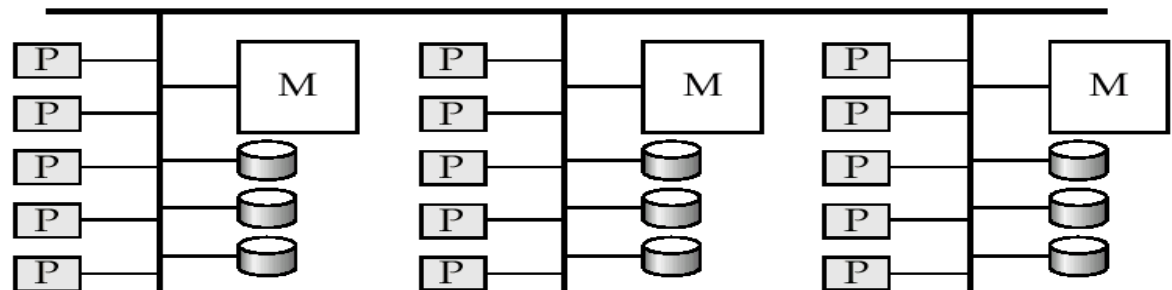# PARALLEL DATABASE ARCHITECTURES



(a) shared memory

(b) shared disk

(c) shared nothing

(d) hierarchical

# SHARED MEMORY

- Processors and disks have access to a common memory, typically via a bus or through an interconnection network.

- Extremely efficient communication between processors — data in shared memory can be accessed by any processor without having to move it using software.

- Downside – architecture is not scalable beyond 32 or 64 processors since the bus or the interconnection network becomes a bottleneck

- Widely used for lower degrees of parallelism (4 to 8).

# SHARED DISK

- All processors can directly access all disks via an interconnection network, but the processors have private memories.
    - The memory bus is not a bottleneck
    - Architecture provides a degree of **fault-tolerance** — if a processor fails, the other processors can take over its tasks since the database is resident on disks that are accessible from all processors.
- Examples: IBM Sysplex and DEC clusters (now part of Compaq) running Rdb (now Oracle Rdb) were early commercial users
- Downside: bottleneck now occurs at interconnection to the disk subsystem.
- Shared-disk systems can scale to a somewhat larger number of processors, but communication between processors is slower.
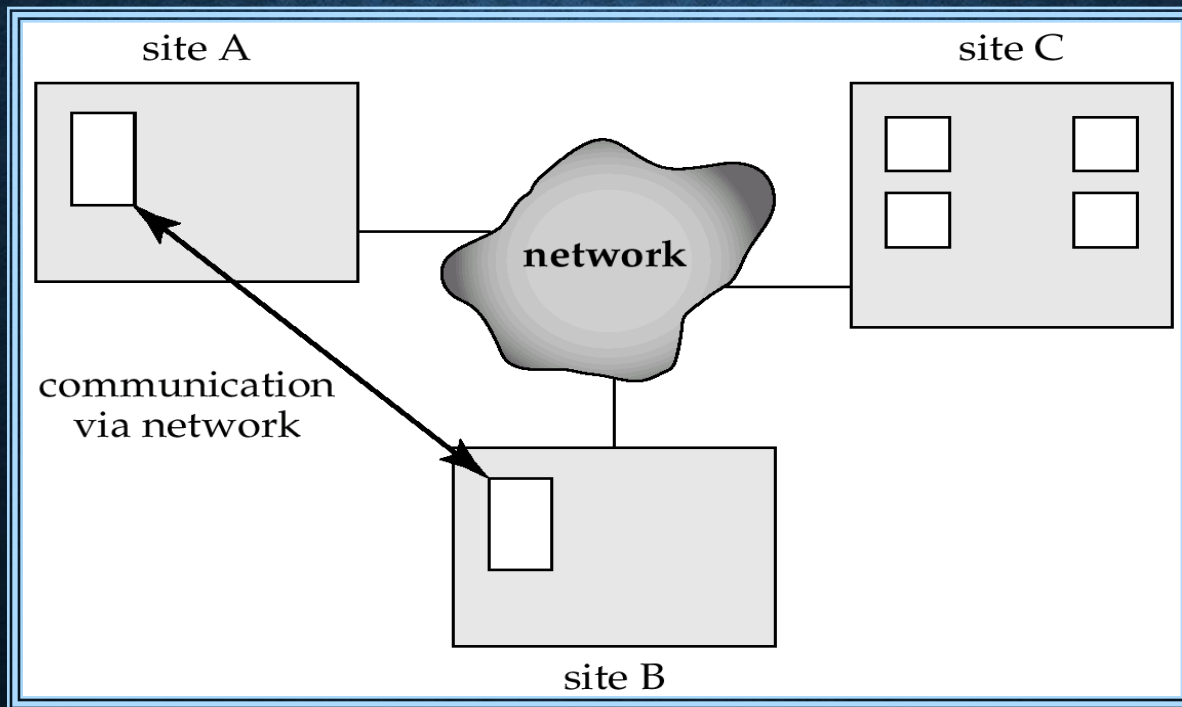
# SHARED NOTHING

- Node consists of a processor, memory, and one or more disks. Processors at one node communicate with another processor at another node using an interconnection network. A node functions as the server for the data on the disk or disks the node owns.

- Examples: Teradata, Tandem, Oracle-n CUBE

- Data accessed from local disks (and local memory accesses) do not pass through interconnection network, thereby minimizing the interference of resource sharing.

- Shared-nothing multiprocessors can be scaled up to thousands of processors without interference.

- Main drawback: cost of communication and non-local disk access; sending data involves software interaction at both ends.

# HIERARCHICAL

- Combines characteristics of shared-memory, shared-disk, and shared-nothing architectures.

- Top level is a shared-nothing architecture – nodes connected by an interconnection network, and do not share disks or memory with each other.

- Each node of the system could be a shared-memory system with a few processors.

- Alternatively, each node could be a shared-disk system, and each of the systems sharing a set of disks could be a shared-memory system.

- Reduce the complexity of programming such systems by **distributed virtual-memory** architectures
  - Also called **non-uniform memory architecture (NUMA)**

# DISTRIBUTED SYSTEMS

- Data spread over multiple machines (also referred to as **sites** or **nodes**.

- Network interconnects the machines

- Data shared by users on multiple machines

# DISTRIBUTED DATABASES

- Homogeneous distributed databases
  - Same software/schema on all sites, data may be partitioned among sites
  - Goal: provide a view of a single database, hiding details of distribution
- Heterogeneous distributed databases
  - Different software/schema on different sites
  - Goal: integrate existing databases to provide useful functionality
- Differentiate between *local* and *global* transactions
  - A local transaction accesses data in the *single* site at which the transaction was initiated.
  - A global transaction either accesses data in a site different from the one at which the transaction was initiated or accesses data in several different sites.

# TRADE-OFFS IN DISTRIBUTED SYSTEMS

- Sharing data – users at one site able to access the data residing at some other sites.

- Autonomy – each site is able to retain a degree of control over data stored locally.

- Higher system availability through redundancy — data can be replicated at remote sites, and system can function even if a site fails.

- Disadvantage: added complexity required to ensure proper coordination among sites.
  - Software development cost.
  - Greater potential for bugs.
  - Increased processing overhead.

# IMPLEMENTATION ISSUES FOR DISTRIBUTED DATABASES

- Atomicity needed even for transactions that update data at multiple site
  - Transaction cannot be committed at one site and aborted at another
- The two-phase commit protocol (2PC) used to ensure atomicity
  - Basic idea: each site executes transaction till just before commit, and the leaves final decision to a coordinator
  - Each site must follow decision of coordinator: even if there is a failure while waiting for coordinators decision
    - To do so, updates of transaction are logged to stable storage and transaction is recorded as "waiting"
  - More details in Sectin 19.4.1
- 2PC is not always appropriate: other transaction models based on persistent messaging, and workflows, are also used
- Distributed concurrency control (and deadlock detection) required
- Replication of data items required for improving data availability
- Details of above in Chapter 19

# NETWORK TYPES

- **Local-area networks (**LANs) – composed of processors that are distributed over small geographical areas, such as a single building or a few adjacent buildings.

- **Wide-area networks (**WANs) – composed of processors distributed over a large geographical area.

- *Discontinuous connection* – WANs, such as those based on periodic dial-up (using, e.g., UUCP), that are connected only for part of the time.

- *Continuous connection* – WANs, such as the Internet, where hosts are connected to the network at all times.

# NETWORKS TYPES (CONT.)

- WANs with continuous connection are needed for implementing distributed database systems

- Groupware applications such as Lotus notes can work on WANs with discontinuous connection:
  - Data is replicated.
  - Updates are propagated to replicas periodically.
  - No global locking is possible, and copies of data may be independently updated.
  - Non-serializable executions can thus result. Conflicting updates may have to be detected, and resolved in an application dependent manner.

# END OF DATABASE SYSTEMS