# Introduction to Databases

SARCAR Sayan

Faculty of Library, Information, and Media Science

# Contents

- Purpose of Database System

- View of Data

- Data models

- Data definition language

- Data manipulation language

- SQL

For more detailed information, please visit
http://codex.cs.yale.edu/avi/db-book/db6/slide-dir/index.html

# 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

筑波大学
University of Tsukuba

# Purpose of Database Systems

- 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

筑波大学
*University of Tsukuba*

# Purpose of Database Systems

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

# Level 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.

# 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

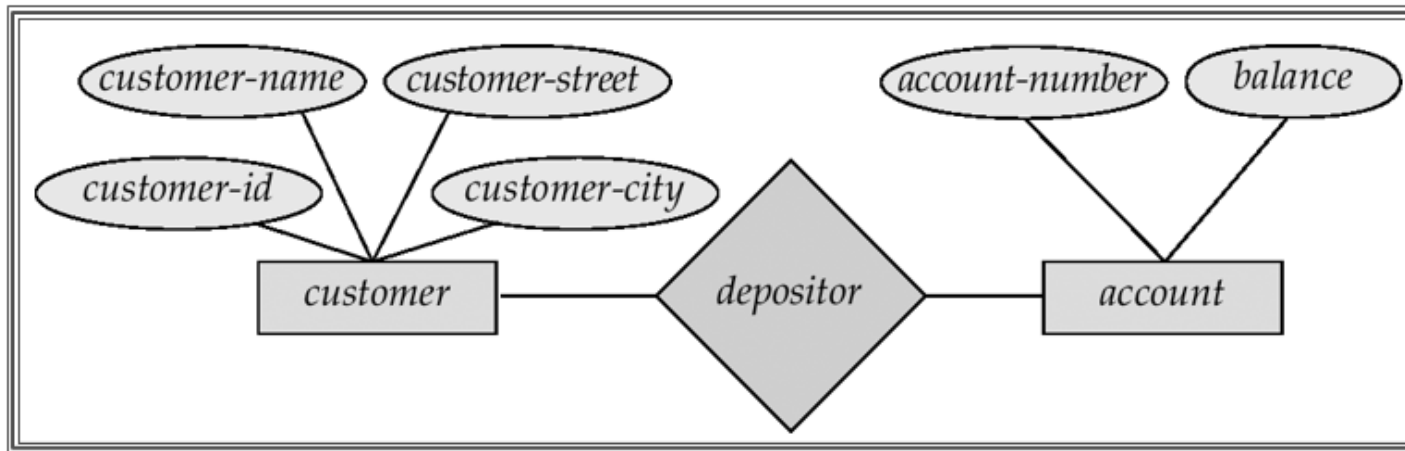Relational model

Other models:
- object-oriented model
- semi-structured data models
- Older models: network model and hierarchical model

筑波大学
University of Tsukuba

# Entity-Relationship Model

Example of schemas in the entity-relationship model

# Entity-relationship Model

- 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

筑波大学
University of Tsukuba

# 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 |

筑波大学
University of Tsukuba

# 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

筑波大学
*University of Tsukuba*

# Data Definition Language

- Specification notation for defining the database schema
  - v 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)
  - v database schema
  - v 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

筑波大学
*University of Tsukuba*

# Data Manipulation Language

- 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

筑波大学
*University of Tsukuba*

# SQL

- SQL: widely used non-procedural language

  - v 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'

  - v 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

  - v Language extensions that allow embedded SQL

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

筑波大学
*University of Tsukuba*

# The Relational Model

- Structure of Relational Databases
- Relational Algebra
- Tuple Relational Calculus
- Domain Relational Calculus
- Extended Relational-Algebra-Operations
- Modification of the Database
- Views

筑波大学
University of Tsukuba

# 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 $D_1$ x $D_2$ x … x$D_n$
  Thus a relation is a set of n-tuples $(a_1, a_2, …, a_n)$ where
  $a_i \in D_i$

- 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*
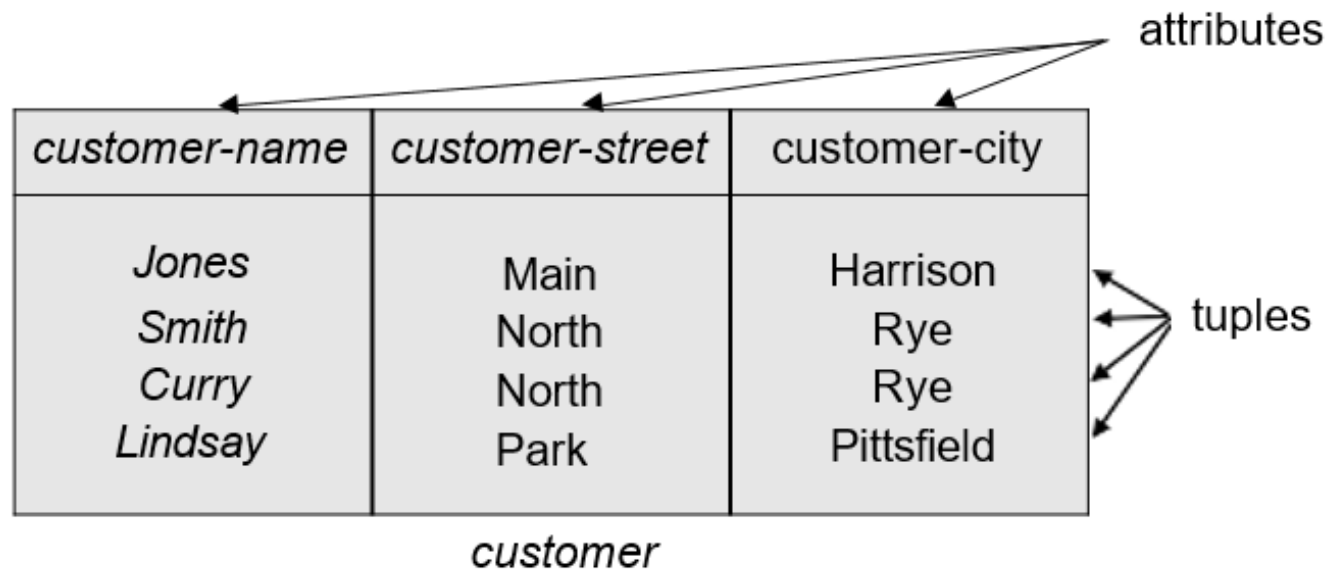
筑波大学
University of Tsukuba

# Attribute Type

- 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

  - v  E.g. multivalued attribute values are not atomic

  - v  E.g. composite attribute values are not atomic

筑波大学
University of Tsukuba

# Relation Schema

- $A_1, A_2, \ldots, A_n$ are *attributes*

- $R = (A_1, A_2, \ldots, A_n)$ 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)*

筑波大学
*University of Tsukuba*

# 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



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

*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

  - v  repetition of information (e.g. two customers own an account)

  - v  the need for null values (e.g. represent a customer without an account)

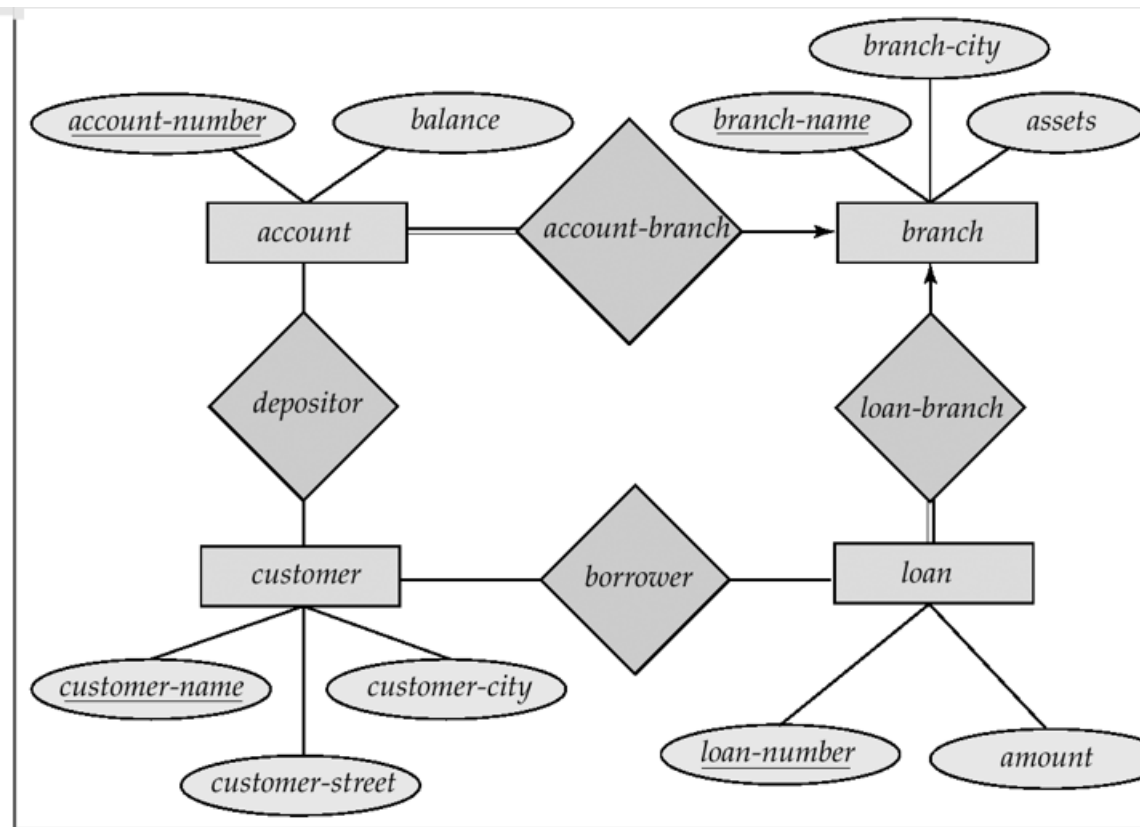- Normalization theory deals with how to design relational schemas

# 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 |

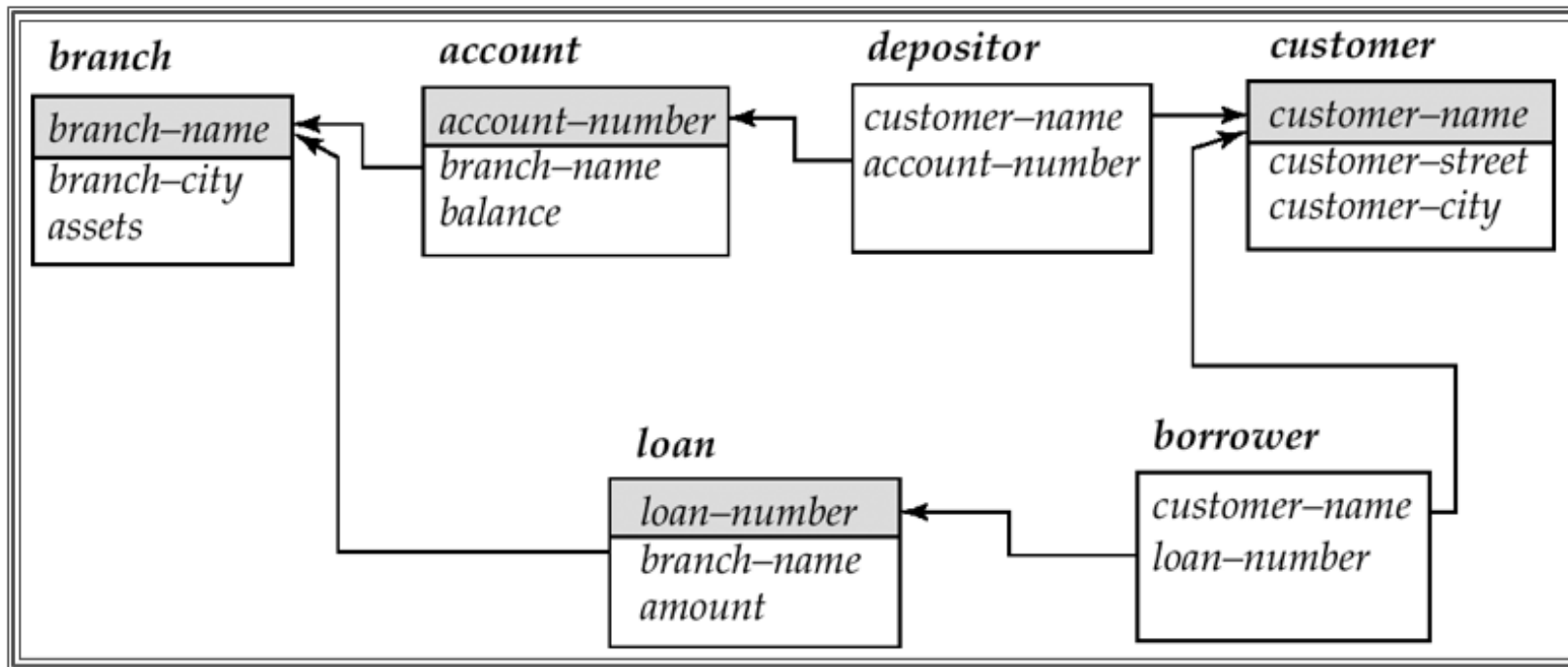# ER Diagram for the Banking Enterprise

# 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
  \qquad\qquad {*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.

筑波大学
*University of Tsukuba*

# Determining Keys from the ER 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.

  - v   For binary many-to-one relationship sets, the primary key of the "many" entity set becomes the relation's primary key.

  - v   For one-to-one relationship sets, the relation's primary key can be that of either entity set.

  - v   For many-to-many relationship sets, the union of the primary keys becomes the relation's primary key

筑波大学
University of Tsukuba

# Schema Diagram for the Banking Enterprise

# Query language

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 |
|---|---|---|---|
| α | α | 1 | 7 |
| α | β | 5 | 7 |
| β | β | 12 | 3 |
| β | β | 23 | 10 |

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

| A | B | C | D |
|---|---|---|---|
| α | α | 1 | 7 |
| β | β | 23 | 10 |

University of Tsukuba

# Select operation

- Notation: $\sigma_p(r)$

- $p$ is called the selection predicate

- Defined as:

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

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 |
|---|---|---|
| α | 10 | 1 |
| α | 20 | 1 |
| β | 30 | 1 |
| β | 40 | 2 |

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

| A | C |
|---|---|
| α | 1 |
| α | 1 |
| β | 1 |
| β | 2 |

=

| A | C |
|---|---|
| α | 1 |
| β | 1 |
| β | 2 |

# Project operation

- Notation:

$$\Pi_{A1, A2, ..., Ak}(r)$$

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*

$$\Pi_{account\text{-}number, balance}(account)$$

# Union operation - Example

- Relations *r, s*:

| A | B |
|---|---|
| α | 1 |
| α | 2 |
| β | 1 |

*r*

| A | B |
|---|---|
| α | 2 |
| β | 3 |

*s*

r ∪ s:

| A | B |
|---|---|
| α | 1 |
| α | 2 |
| β | 1 |
| β | 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
$$\Pi_{customer\text{-}name} (depositor) \cup \Pi_{customer\text{-}name} (borrower)$$

# Set Difference - Example

- Relations *r, s*:

| A | B |
|---|---|
| α | 1 |
| α | 2 |
| β | 1 |

*r*

| A | B |
|---|---|
| α | 2 |
| β | 3 |

*s*

*r − s:*

| A | B |
|---|---|
| α | 1 |
| β | 1 |

# Set Difference

- Notation $r - s$

- Defined as:

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

- 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 - Example

Relations $r$, $s$:

| A | B |
|---|---|
| α | 1 |
| β | 2 |

$r$

| C | D | E |
|---|---|---|
| α | 10 | a |
| β | 10 | a |
| β | 20 | b |
| γ | 10 | b |

$s$

$r$ x $s$:

| A | B | C | D | E |
|---|---|---|---|---|
| α | 1 | α | 10 | a |
| α | 1 | β | 19 | a |
| α | 1 | β | 20 | b |
| α | 1 | γ | 10 | b |
| β | 2 | α | 10 | a |
| β | 2 | β | 10 | a |
| β | 2 | β | 20 | b |
| β | 2 | γ | 10 | b |

# Cartesian Product Operation

- Notation $r \times s$

- 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.

筑波大学
University of Tsukuba

# Rename Operation - Example

- Allows us to refer to a relation, (say E) by more than one name.

$$\rho_x (E)$$

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

- Relations $r$

| A | B |
|---|---|
| $\alpha$ | 1 |
| $\beta$ | 2 |

$r$

- $r \times \rho_s (r)$

| $r.A$ | $r.B$ | $s.A$ | $s.B$ |
|-------|-------|-------|-------|
| $\alpha$ | 1 | $\alpha$ | 1 |
| $\alpha$ | 1 | $\beta$ | 2 |
| $\beta$ | 2 | $\alpha$ | 1 |
| $\beta$ | 2 | $\beta$ | 2 |

筑波大学
University of Tsukuba

# 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

$$\rho_{X(A1, A2, ..., An)}(E)$$

returns the result of expression $E$ under the name $X$, and with the attributes renamed to $A1, A2, ...., An$.

# 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

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

# 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.

$$\Pi_{customer\text{-}name} \left( \sigma_{branch\text{-}name=\text{"Perryridge"}} \right.$$

$$\left. \left( \sigma_{borrower.loan\text{-}number = loan.loan\text{-}number}(borrower \times loan) \right) \right)$$

- 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.

$$\Pi_{customer\text{-}name} \left( \sigma_{branch\text{-}name = \text{"Perryridge"}} \right.$$

$$\left. \left( \sigma_{borrower.loan\text{-}number = loan.loan\text{-}number}(borrower \times loan) \right) \right)$$

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

# Example Queries

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

    - Query 1

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

    - Query 2

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

筑波大学
University of Tsukuba

# Example Queries

Find the largest account balance

    v  Rename *account* relation as *d*

    v  The query then is:

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

筑波大学
*University of Tsukuba*

# Formal Definitions

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

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

- $r \cap s = \{\, t \mid t \in r \text{ and } t \in s \,\}$

- Assume:

  - v  $r, s$ have the *same arity*

  - v  attributes of r and s are compatible

- Note: $r \cap s = r - (r - s)$

# Set-Intersection Operation - Example

- Relation r, s:

| A | B |
|---|---|
| α | 1 |
| α | 2 |
| β | 1 |

r

| A | B |
|---|---|
| α | 2 |
| β | 3 |

s

- r ∩ s

| A | B |
|---|---|
| α | 2 |

# Natural-Join Operation

- Notation: $r \bowtie s$

- Let $r$ and $s$ be relations on schemas $R$ and $S$ respectively. The result is a relation on schema $R \cup S$ which is obtained by considering 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 $R \cap S$, a tuple $t$ is added to the result, where

    v   $t$ has the same value as $t_r$ on $r$

    v   $t$ 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:

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

# Natural-Join Operation - Example

- Relations r, s:

| A | B | C | D |
|---|---|---|---|
| α | 1 | α | a |
| β | 2 | γ | a |
| γ | 4 | β | b |
| α | 1 | γ | a |
| δ | 2 | β | b |

r

| B | D | E |
|---|---|---|
| 1 | a | α |
| 3 | a | β |
| 1 | a | γ |
| 2 | b | δ |
| 3 | b | ∈ |

s

r ⋈

| A | B | C | D | E |
|---|---|---|---|---|
| α | 1 | α | a | α |
| α | 1 | α | a | γ |
| α | 1 | γ | a | α |
| α | 1 | γ | a | γ |
| δ | 2 | β | b | δ |

# 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

  - $R = (A_1, \ldots, A_m, B_1, \ldots, B_n)$
  - $S = (B_1, \ldots, B_n)$

  The result of r ÷ s is a relation on schema

  $R - S = (A_1, \ldots, A_m)$

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

# Division Operation - Example

Relations $r$, $s$:

| A | B |
|---|---|
| α | 1 |
| α | 2 |
| α | 3 |
| β | 1 |
| γ | 1 |
| δ | 1 |
| δ | 3 |
| δ | 4 |
| ∈ | 6 |
| ∈ | 1 |
| β | 2 |

$r$

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

$s$

$r \div s$:

| A |
|---|
| α |
| β |

筑波大学
University of Tsukuba

# Another Division Example

Relations $r$, $s$:

| A | B | C | D | E |
|---|---|---|---|---|
| α | a | α | a | 1 |
| α | a | γ | a | 1 |
| α | a | γ | b | 1 |
| β | a | γ | a | 1 |
| β | a | γ | b | 3 |
| γ | a | γ | a | 1 |
| γ | a | γ | b | 1 |
| γ | a | β | b | 1 |

$r$

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

$s$

$r \div s$:

| A | B | C |
|---|---|---|
| α | a | γ |
| γ | a | γ |

# Division Operation

- Property
  - v   Let $q - r \div s$
  - v   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 $S \subseteq R$

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

To see why

- v   $\Pi_{R-S,S}(r)$ simply reorders attributes of $r$

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

  $\Pi_{R-S}(r)$ 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

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

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

  v  May use variable in subsequent expressions.

筑波大学
University of Tsukuba

# Assignment Operation - Example

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

  - v 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.

  - v 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.

$$\Pi_{customer\text{-}name,\ branch\text{-}name}\ (depositor \bowtie account)$$
$$\div\ \Pi_{branch\text{-}name}\ (\sigma_{branch\text{-}city\ =\ \text{``Brooklyn''}}\ (branch))$$

# Extended Relational Algebra Operations

- Generalized Projection
- Aggregate Functions

# Generalized Projections

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

$$\Pi_{F1, F2, ..., Fn}(E)$$

- $E$ 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:

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

筑波大学
University of Tsukuba

# 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

$$G1, G2, ..., Gn \, @ \, F1(A1), F2(A2),..., Fn(An)(E)$$

v $E$ is any relational-algebra expression

v $G_1, G_2 ..., G_n$ is a list of attributes on which to group (can be empty)

v Each $F_i$ is an aggregate function

v Each $A_i$ is an attribute name

# Aggregate Operation - Example

- Relation $r$:

| A | B | C |
|---|---|---|
| α | α | 7 |
| α | β | 7 |
| β | β | 3 |
| β | β | 10 |

@ $\mathbf{sum(c)}^{(r)}$

| 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 |

*branch-name sum(balance)* $(account)$

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

# Aggregate Function

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

$$_{branch\text{-}name} \; \mathbf{sum}(balance) \; \mathbf{as} \; sum\text{-}balance \, (account)$$

筑波大学
University of Tsukuba

# Modification of Database

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

- All these operations are expressed using the assignment operator.

筑波大学
University of Tsukuba

# 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

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

- 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:
  - v specify a tuple to be inserted
  - v 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.

筑波大学
*University of Tsukuba*

# Insertion Example

- 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.

$$r_1 \leftarrow (\sigma_{branch\text{-}name = \text{"Perryridge"}} (borrower \bowtie loan)) \quad account$$

$$\leftarrow account \cup \Pi_{branch\text{-}name, \ account\text{-}number, 200} (r_1)$$

$$depositor \leftarrow depositor \cup \Pi_{customer\text{-}name, \ loan\text{-}number,} (r_1)$$

# Update

- 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

$$r \leftarrow \Pi_{F1,\ F2,\ \dots,\ Fl,}(r)$$

- Each $F$, 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 Example

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

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

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

$$account \leftarrow \Pi_{AN,\ BN,\ BAL\ *\ 1.06}(\sigma_{BAL > 10000}(account))$$
$$\cup\ \Pi_{AN,\ BN,\ BAL\ *\ 1.05}(\sigma_{BAL \leq 10000}(account))$$

筑波大学
*University of Tsukuba*

# SQL

- Basic Structure
- Set Operations
- Aggregate Functions
- Nested Subqueries
- Derived Relations
- Modification of the Database
- Data Definition Language

筑波大学
University of Tsukuba

# 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.

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

- The result of an SQL query is a relation.

筑波大学
University of Tsukuba

# The Select Clause

- The **select** clause corresponds to the projection operation of the relational algebra. It is used to list the attributes desired in the result of a query.

- 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:

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

- An asterisk in the select clause denotes "all attributes"
  **select** *
  **from** *loan*

NOTES:

- ᵥ  SQL does not permit the '-' character in names, so you would use, for example, *branch_name* instead of *branch-name* in a real implementation. We use '-' since it looks nicer!

- ᵥ  SQL names are case insensitive.

筑波大学
*University of Tsukuba*

# 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*

筑波大学
*University of Tsukuba*

# The Select Clause (Cont.)

- 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.

筑波大学
*University of Tsukuba*

# The Where Clause

- The **where** clause corresponds to the selection predicate of the relational algebra. If consists of a predicate involving attributes of the relations that appear in the **from** clause.

- The find all loan number for loans made a 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 in order to simplify **where** clauses that specify that a value be less than or equal to some value and greater than or equal to some other value.

- Find the loan number of those loans with loan amounts between $90,000 and $100,000 (that is, $\geq$ $90,000 and $\leq$ $100,000)

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

# The From Clause

- The **from** clause corresponds to the Cartesian product operation of the relational algebra. It lists the relations to be scanned in the evaluation of the expression.

- 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:

$$old\text{-}name \textbf{ as } new\text{-}name$$

. 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'*

筑波大学
*University of Tsukuba*

# String Operations

- SQL includes a string-matching operator for comparisons on character strings. Patterns are described using two special characters:
    - v  percent (%). The % character matches any substring.
    - v  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** ˋ%Main%ˊ

- Match the name "Main%"

> **like** ˋMain\%ˊ **escape** ˋ\ˊ

- SQL supports a variety of string operations such as
    - v  concatenation (using "||")
    - v  converting from upper to lower case (and vice versa)
    - v  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.

  v    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. 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. 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. 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:

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

- 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** $A_1,\ A_2,\ ...,\ A_n$
  **from** $r_1,\ r_2,\ ...,\ r_m$
  **where** $P$

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

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

# 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** $s$

筑波大学
*University of Tsukuba*

# 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*)

筑波大学
*University of Tsukuba*

# 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

筑波大学
*University of Tsukuba*

# 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*

筑波大学
*University of Tsukuba*

# 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

# 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.

筑波大学
*University of Tsukuba*

# 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.

筑波大学
*University of Tsukuba*

# Set Comparison

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

$$\textbf{select distinct } \textit{T.branch-name}$$
$$\textbf{from } \textit{branch } \textbf{as } \textit{T, branch } \textbf{as } \textit{S}$$
$$\textbf{where } \textit{T.assets} > \textit{S.assets} \textbf{ and}$$
$$\textit{S.branch-city} = \text{`Brooklyn'}$$

- Same query using > **some** clause

$$\textbf{select } \textit{branch-name}$$
$$\textbf{from } \textit{branch}$$
$$\textbf{where } \textit{assets} > \textbf{some}$$
$$\quad (\textbf{select } \textit{assets}$$
$$\quad \ \textbf{from } \textit{branch}$$
$$\quad \ \ \textbf{where } \textit{branch-city} = \text{`Brooklyn'})$$

# 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** $\begin{array}{|c|} \hline 0 \\ \hline 5 \\ \hline 6 \\ \hline \end{array}$ ) = true

(read: 5 < some tuple in the relation)

(5< **some** $\begin{array}{|c|} \hline 0 \\ \hline 5 \\ \hline \end{array}$ ) = false

(5 = **some** $\begin{array}{|c|} \hline 0 \\ \hline 5 \\ \hline \end{array}$ ) = true

(5 $\neq$ **some** $\begin{array}{|c|} \hline 0 \\ \hline 5 \\ \hline \end{array}$ ) = true (since $0 \neq 5$)

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

# Definition of All Clause

. $\quad$ F \<comp\> **all** $r \Leftrightarrow \forall\ t \in r$ (F \<comp\> $t$)

$$(5 < \textbf{all}\ \begin{array}{|c|} \hline 0 \\ \hline 5 \\ \hline 6 \\ \hline \end{array}\ ) = \text{false}$$

$$(5 < \textbf{all}\ \begin{array}{|c|} \hline 6 \\ \hline 10 \\ \hline \end{array}\ ) = \text{true}$$

$$(5 = \textbf{all}\ \begin{array}{|c|} \hline 4 \\ \hline 5 \\ \hline \end{array}\ ) = \text{false}$$

$$(5 \neq \textbf{all}\ \begin{array}{|c|} \hline 4 \\ \hline 6 \\ \hline \end{array}\ ) = \text{true (since } 5 \neq 4 \text{ and } 5 \neq 6)$$

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

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

# 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** $r \Leftrightarrow r \neq \emptyset$

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

筑波大学
University of Tsukuba

# 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))*

- Note that $X - Y = \emptyset \Leftrightarrow X \subseteq Y$

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

筑波大学
*University of Tsukuba*

# 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')

筑波大学
*University of Tsukuba*

# 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‚)

# 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 relation *result* in the **from** clause, and the attributes of *result* can be used directly in the **where** clause.

# 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*)

筑波大学
University of Tsukuba

# 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)*

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

  - v  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

筑波大学
University of Tsukuba

# Modification of the Database - Updates

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

  - v  Write two **update** statements:

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

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

  - v  The order is important!

# 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.

筑波大学
University of Tsukuba

# 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.

# Create Table Construct

- An SQL relation is defined using the **create table** command:

$$\textbf{create table } r\ (A_1\ D_1,\ A_2\ D_2,\ ...,\ A_n\ D_n,$$
$$(\text{integrity-constraint}_1),$$
$$...,$$
$$(\text{integrity-constraint}_k))$$

  - v    $r$ is the name of the relation

  - v    each $A_i$ is an attribute name in the schema of relation $r$

  - v    $D_i$ is the data type of values in the domain of attribute $A_i$

- Example:

$$\textbf{create table } branch$$
$$(branch\text{-}name\ \text{char(15)}\ \textbf{not null,}$$
$$branch\text{-}city\qquad \text{char(30)},$$
$$assets\qquad\qquad \text{integer})$$

筑波大学
*University of Tsukuba*

# Integrity Constraints in Create Table

- **not null**

- **primary key** $(A_1, ..., A_n)$

- **check** *(P)*, where *P* is a predicate

Example: Declare *branch-name* as the primary key for *branch* and ensure that the values of *assets* are non-negative.

> **create table** *branch*
>     *(branch-name* char(15),
>     *branch-city* char(30)
>     *assets*              integer,
>     **primary key** *(branch-name),*
>     **check** *(assets >= 0))*

**primary key** declaration on an attribute automatically ensures **not null** in SQL-92 onwards, needs to be explicitly stated in SQL-89

# Drop and Alter Table Constructs

- The **drop table** command deletes all information about the dropped relation from the database.

- The **after table** command is used to add attributes to an existing relation. All tuples in the relation are assigned *null* as the value for the new attribute. The form of the **alter table** command is

**alter table** *r* **add** *A D*

where *A* is the name of the attribute to be added to relation *r* and *D* is the domain of *A*.

- The **alter table** command can also be used to drop attributes of a relation

**alter table** *r* **drop** *A*

where *A* is the name of an attribute of relation *r*

  v Dropping of attributes not supported by many databases

筑波大学
*University of Tsukuba*

# SQL Data Definition for Part of the Bank Database

```
create table customer
    (customer-name      char(20),
     customer-street    char(30),
     customer-city      char(30),
     primary key (customer-name))

create table branch
    (branch-name        char(15),
     branch-city        char(30),
     assets             integer,
     primary key (branch-name),
     check (assets >= 0))

create table account
    (account-number     char(10),
     branch-name        char(15),
     balance            integer,
     primary key (account-number),
     check (balance >= 0))

create table depositor
    (customer-name      char(20),
     account-number     char(10),
     primary key (customer-name, account-number))
```

筑波大学
University of Tsukuba

# Q & A
**Please write any feedback regarding class to
sayans@slis.tsukuba.ac.jp
Sub: Informatics class feedback**