

CPT-S 415

Big Data

Yinghui Wu

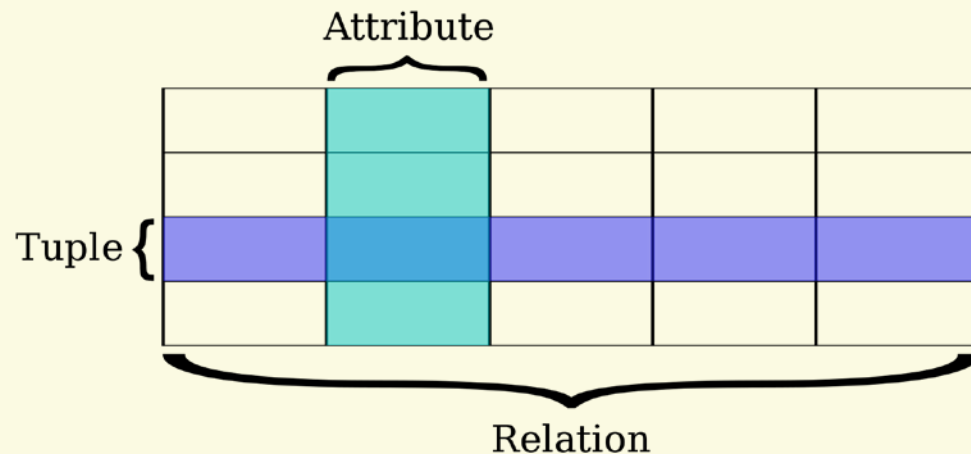
EME B45

CPT-S 415

Big Data

Relational data model and DBMS

- ✓ Relational Model Concepts
- ✓ Relational Model Constraints and Schemas
- ✓ Update Operations and Dealing with Constraint Violations



Data Models

- ✓ A collection of tools for describing
 - Data
 - Data relationships
 - Data semantics
 - Data constraints
 - ✓ **Relational model**
 - ✓ Entity-Relationship data model (mainly for database design)
 - ✓ Object-based data models (Object-oriented and Object-relational)
 - ✓ **Semistructured data model (XML and graphs)**
 - ✓ Other models:
 - *Network model*
 - *Hierarchical model*
- “What goes around comes around”, by Michael Stonebraker

Relational Model Concepts

- ✓ A Relation is a mathematical concept based on the ideas of sets
- ✓ The model was first proposed by Dr. E.F. Codd of IBM Research in 1970:
 - "A Relational Model for Large Shared Data Banks," Communications of the ACM, June 1970
- ✓ The above paper caused a major revolution in the field of database management and earned Dr. Codd the coveted ACM Turing Award
- ✓ *Turing Award winners who are databases researchers*
 - 1973, Charles W. Bachman
 - 1981, Edgar F. Codd
 - 1998, Jim Gray
 - 2014, Michael Stonebraker

Informal Definitions

- ✓ Informally, a **relation** looks like a **table** of values.
- ✓ A relation typically contains a **set of rows**.
- ✓ The data elements in each **row** represent certain facts that correspond to a real-world **entity** or **relationship**
 - In the formal model, rows are called **tuples**
- ✓ Each **column** has a column header that gives an indication of the meaning of the data items in that column
 - In the formal model, the column header is called an **attribute name** (or just **attribute**)

Example of a Relation

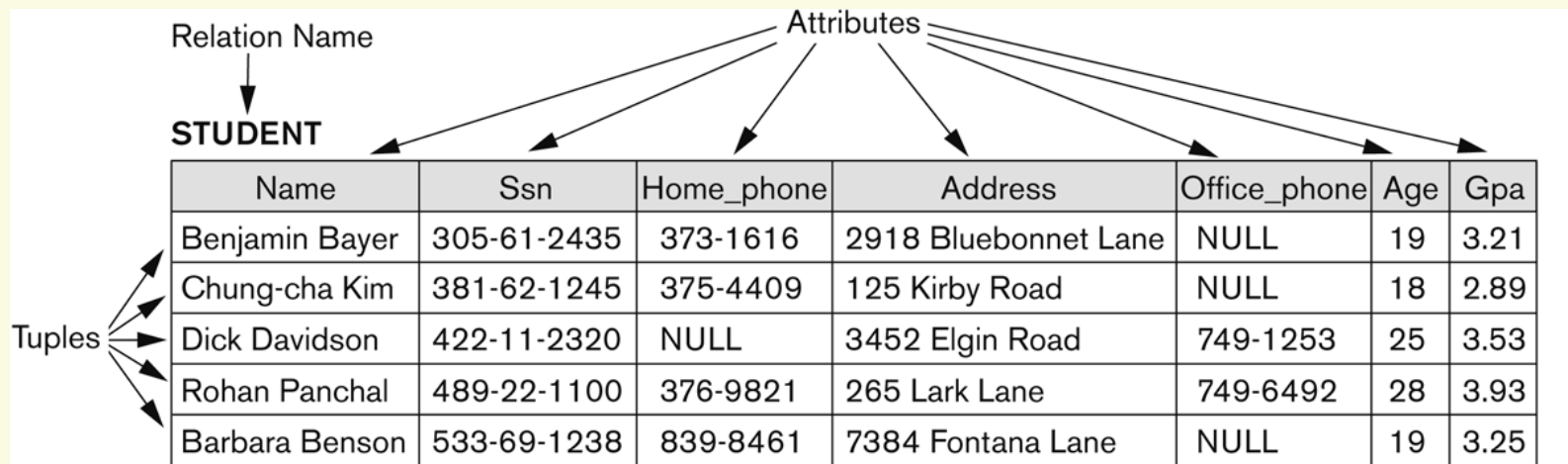


Figure 5.1

The attributes and tuples of a relation STUDENT.

Informal Definitions

✓ Key of a Relation:

- Each row has a value of a data item (or set of items) that uniquely identifies that row in the table
 - Called the *key*
- In the STUDENT table, SSN is the key
- Sometimes row-ids or sequential numbers are assigned as keys to identify the rows in a table
 - Called *artificial key* or *surrogate key*

Formal Definitions - Schema

✓ The **Schema** (or description) of a Relation:

- Denoted by $R(A_1, A_2, \dots, A_n)$
- R is the **name** of the relation
- The **attributes** of the relation are A_1, A_2, \dots, A_n

✓ Example:

CUSTOMER (Cust-id, Cust-name, Address, Phone#)

- CUSTOMER is the relation name
 - Defined over the four attributes: Cust-id, Cust-name, Address, Phone#
- ✓ Each attribute has a **domain** or a set of valid values.
- For example, the domain of Cust-id is 6 digit numbers.

Formal Definitions - Tuple

- ✓ A **tuple** is an ordered n-tuple of values (enclosed in angled brackets '< ... >')
- ✓ Each value is derived from an appropriate *domain*.
- ✓ A row in the CUSTOMER relation is a 4-tuple and would consist of four values, for example:
 - <632895, "John Smith", "101 Main St. Atlanta, GA 30332", "(404) 894-2000">
 - This is called a 4-tuple as it has 4 values
 - A tuple (row) in the CUSTOMER relation.
- ✓ A relation is a **set** of such tuples (rows)

Formal Definitions - Domain

- ✓ A **domain** has a logical definition:
 - Example: “USA_phone_numbers” are the set of 10 digit phone numbers valid in the U.S.
- ✓ A domain also has a data-type or a format.
 - The USA_phone_numbers may have a format: (ddd)ddd-dddd where each d is a decimal digit.
 - Dates have various formats such as year, month, date formatted as yyyy-mm-dd, or as dd mm,yyyy etc.
- ✓ The attribute name designates the role played by a domain in a relation:
 - Used to interpret the meaning of the data elements corresponding to that attribute
 - Example: The domain Date may be used to define two attributes named “Invoice-date” and “Payment-date” with different meanings

Formal Definitions - State

- ✓ The **relation state** is a subset of the Cartesian product of the domains of its attributes
 - each domain contains the set of all possible values the attribute can take.
- ✓ Example: attribute Cust-name is defined over the domain of character strings of maximum length 25
 - $\text{dom}(\text{Cust-name})$ is `varchar(25)`
- ✓ The role these strings play in the CUSTOMER relation is that of the *name of a customer*.

What's the difference between a state of a relation and an instance of a relation?

Formal Definitions - Summary

- ✓ Formally,
 - Given $R(A_1, A_2, \dots, A_n)$
 - $r(R) \subset \text{dom}(A_1) \times \text{dom}(A_2) \times \dots \times \text{dom}(A_n)$
- ✓ $R(A_1, A_2, \dots, A_n)$ is the **schema** of the relation
- ✓ R is the **name** of the relation
- ✓ A_1, A_2, \dots, A_n are the **attributes** of the relation
- ✓ $r(R)$: a specific **state** (or "value" or "population") of relation R – this is a *set of tuples* (rows)
 - $r(R) = \{t_1, t_2, \dots, t_n\}$ where each t_i is an n -tuple
 - $t_i = \langle v_1, v_2, \dots, v_n \rangle$ where each v_j *element-of* $\text{dom}(A_j)$

Formal Definitions - Example

- ✓ Let $R(A1, A2)$ be a relation schema:
 - Let $\text{dom}(A1) = \{0,1\}$
 - Let $\text{dom}(A2) = \{a,b,c\}$
- ✓ Then: $\text{dom}(A1) \times \text{dom}(A2)$ is all possible combinations:
 $\{ \langle 0,a \rangle , \langle 0,b \rangle , \langle 0,c \rangle , \langle 1,a \rangle , \langle 1,b \rangle , \langle 1,c \rangle \}$
- ✓ The relation state $r(R) \subset \text{dom}(A1) \times \text{dom}(A2)$
- ✓ For example: $r(R)$ could be $\{ \langle 0,a \rangle , \langle 0,b \rangle , \langle 1,c \rangle \}$
 - this is one possible state (or “population” or “extension”) r of the relation R , defined over $A1$ and $A2$.
 - It has three 2-tuples: $\langle 0,a \rangle , \langle 0,b \rangle , \langle 1,c \rangle$

Definition Summary

<u>Informal Terms</u>		<u>Formal Terms</u>
Table		Relation
Column Header		Attribute
All possible Column Values		Domain
Row		Tuple
Table Definition		Schema of a Relation
Populated Table		State of the Relation

Example – A relation STUDENT

Relation Name

STUDENT

Attributes

Tuples

Name	Ssn	Home_phone	Address	Office_phone	Age	Gpa
Benjamin Bayer	305-61-2435	373-1616	2918 Bluebonnet Lane	NULL	19	3.21
Chung-cha Kim	381-62-1245	375-4409	125 Kirby Road	NULL	18	2.89
Dick Davidson	422-11-2320	NULL	3452 Elgin Road	749-1253	25	3.53
Rohan Panchal	489-22-1100	376-9821	265 Lark Lane	749-6492	28	3.93
Barbara Benson	533-69-1238	839-8461	7384 Fontana Lane	NULL	19	3.25

Figure 5.1

The attributes and tuples of a relation STUDENT.

Characteristics Of Relations

- ✓ Ordering of tuples in a relation $r(R)$:
 - The tuples are *not considered to be ordered*, even though they appear to be in the tabular form.
- ✓ Ordering of attributes in a relation schema R (and of values within each tuple):
 - We consider the attributes in $R(A_1, A_2, \dots, A_n)$ and the values in $t = \langle v_1, v_2, \dots, v_n \rangle$ to be ordered.

Same state as previous Figure (but with different order of tuples)

Figure 5.2

The relation STUDENT from Figure 5.1 with a different order of tuples.

STUDENT

Name	Ssn	Home_phone	Address	Office_phone	Age	Gpa
Dick Davidson	422-11-2320	NULL	3452 Elgin Road	749-1253	25	3.53
Barbara Benson	533-69-1238	839-8461	7384 Fontana Lane	NULL	19	3.25
Rohan Panchal	489-22-1100	376-9821	265 Lark Lane	749-6492	28	3.93
Chung-cha Kim	381-62-1245	375-4409	125 Kirby Road	NULL	18	2.89
Benjamin Bayer	305-61-2435	373-1616	2918 Bluebonnet Lane	NULL	19	3.21

Characteristics Of Relations

- ✓ Values in a tuple:
 - All values are considered atomic (indivisible).
 - Each value in a tuple must be from the domain of the attribute for that column
 - If tuple $t = \langle v_1, v_2, \dots, v_n \rangle$ is a tuple (row) in the relation state r of $R(A_1, A_2, \dots, A_n)$
 - Then each v_i must be a value from $dom(A_i)$
 - **component values** of a tuple t : $t[A_i]$ or $t.A_i$
 - Similarly, $t[A_u, A_v, \dots, A_w]$ refers to the subtuple of t containing the values of attributes A_u, A_v, \dots, A_w , respectively in t
 - A special **null** value is used to represent values that are unknown or inapplicable to certain tuples.

Relational Integrity Constraints

- ✓ Constraints are **conditions** that must hold on **all** valid relation states.
- ✓ There are three *main types* of constraints in the relational model:
 - **Key** constraints
 - **Entity integrity** constraints
 - **Referential integrity** constraints
- ✓ Another implicit constraint is the **domain** constraint
 - Every value in a tuple must be from the *domain of its attribute* (or it could be **null**, if allowed for that attribute)

Key Constraints

✓ Superkey of R:

- Is a set of attributes SK of R with the following condition:
 - No two tuples in any valid relation state $r(R)$ will have the same value for SK
 - That is, for any distinct tuples t_1 and t_2 in $r(R)$, $t_1[SK] \neq t_2[SK]$
 - This condition must hold in *any valid state* $r(R)$

✓ Key of R:

- A "minimal" superkey
- That is, a key is a superkey K such that removal of any attribute from K results in a set of attributes that is not a superkey (does not possess the superkey uniqueness property)

Key Constraints (continued)

- ✓ Example: Consider the CAR relation schema:
 - CAR(State, Reg#, SerialNo, Make, Model, Year)
 - CAR has two keys:
 - Key1 = {State, Reg#}
 - Key2 = {SerialNo}
 - Both are also superkeys of CAR
 - {SerialNo, Make} is a superkey but *not* a key.
- ✓ In general:
 - Any *key* is a *superkey* (but not vice versa)
 - Any set of attributes that *includes a key* is a *superkey*
 - A *minimal* superkey is a key

Key Constraints (continued)

- ✓ If a relation has several **candidate keys**, one is chosen to be the **primary key**.
 - The primary key attributes are underlined.
- ✓ Example: Consider the CAR relation schema:
 - CAR(State, Reg#, SerialNo, Make, Model, Year)
 - We chose SerialNo as the primary key
- ✓ The primary key value is used to *uniquely identify* each tuple in a relation
 - Provides the tuple identity
- ✓ Also used to *reference* the tuple from another tuple
 - General rule: Choose as primary key the smallest of the candidate keys (in terms of size)
 - Not always applicable – choice is sometimes subjective

CAR table with two candidate keys – LicenseNumber chosen as Primary Key

Figure 5.4

The CAR relation, with
two candidate keys:
License_number and
Engine_serial_number.

CAR

<u>License_number</u>	Engine_serial_number	Make	Model	Year
Texas ABC-739	A69352	Ford	Mustang	02
Florida TVP-347	B43696	Oldsmobile	Cutlass	05
New York MPO-22	X83554	Oldsmobile	Delta	01
California 432-TFY	C43742	Mercedes	190-D	99
California RSK-629	Y82935	Toyota	Camry	04
Texas RSK-629	U028365	Jaguar	XJS	04

COMPANY Database Schema

EMPLOYEE

Fname	Minit	Lname	<u>Ssn</u>	Bdate	Address	Sex	Salary	Super_ssn	Dno
-------	-------	-------	------------	-------	---------	-----	--------	-----------	-----

DEPARTMENT

Dname	<u>Dnumber</u>	Mgr_ssn	Mgr_start_date
-------	----------------	---------	----------------

DEPT_LOCATIONS

<u>Dnumber</u>	<u>Dlocation</u>
----------------	------------------

PROJECT

Pname	<u>Pnumber</u>	Plocation	Dnum
-------	----------------	-----------	------

WORKS_ON

<u>Essn</u>	<u>Pno</u>	Hours
-------------	------------	-------

DEPENDENT

<u>Essn</u>	<u>Dependent_name</u>	Sex	Bdate	Relationship
-------------	-----------------------	-----	-------	--------------

Figure 5.5
Schema diagram for
the COMPANY
relational database
schema.

Entity Integrity

✓ Entity Integrity:

- The *primary key attributes* PK of each relation schema R in S cannot have null values in any tuple of $r(R)$.
 - This is because primary key values are used to *identify* the individual tuples.
 - $t[PK] \neq \text{null}$ for any tuple t in $r(R)$
 - If PK has several attributes, null is not allowed in any of these attributes
- Note: Other attributes of R may be constrained to disallow null values, even though they are not members of the primary key.

Referential Integrity

- ✓ A constraint involving **two** relations
 - The previous constraints involve a single relation.
- ✓ Used to specify a **relationship** among tuples in two relations:
 - The **referencing relation** and the **referenced relation**.
- ✓ Tuples in the **referencing relation** R1 have attributes FK (called **foreign key** attributes) that reference the primary key attributes PK of the **referenced relation** R2.
 - A tuple t1 in R1 is said to **reference** a tuple t2 in R2 if $t1[FK] = t2[PK]$.
- ✓ A referential integrity constraint can be displayed in a relational database schema as a directed arc from R1.FK to R2.

Referential Integrity (or foreign key) Constraint

- ✓ Statement of the constraint
 - The value in the foreign key column (or columns) FK of the the **referencing relation** R1 can be **either**:
 - (1) a value of an existing primary key value of a corresponding primary key PK in the **referenced relation** R2, or
 - (2) a **null**.
- ✓ In case (2), the FK in R1 should **not** be a part of its own primary key.

Displaying a relational database schema and its constraints

- ✓ Each relation schema can be displayed as a row of attribute names
- ✓ The name of the relation is written above the attribute names
- ✓ The primary key attribute (or attributes) will be underlined
- ✓ A foreign key (referential integrity) constraints is displayed as a directed arc (arrow) from the foreign key attributes to the referenced table
 - Can also point the primary key of the referenced relation for clarity
- ✓ Next slide shows the COMPANY **relational schema diagram**

Referential Integrity Constraints for COMPANY database

Figure 5.7

Referential integrity constraints displayed on the COMPANY relational database schema.

EMPLOYEE

Fname	Minit	Lname	<u>Ssn</u>	Bdate	Address	Sex	Salary	Super_ssn	Dno
-------	-------	-------	------------	-------	---------	-----	--------	-----------	-----

DEPARTMENT

Dname	<u>Dnumber</u>	Mgr_ssn	Mgr_start_date
-------	----------------	---------	----------------

DEPT_LOCATIONS

<u>Dnumber</u>	<u>Dlocation</u>
----------------	------------------

PROJECT

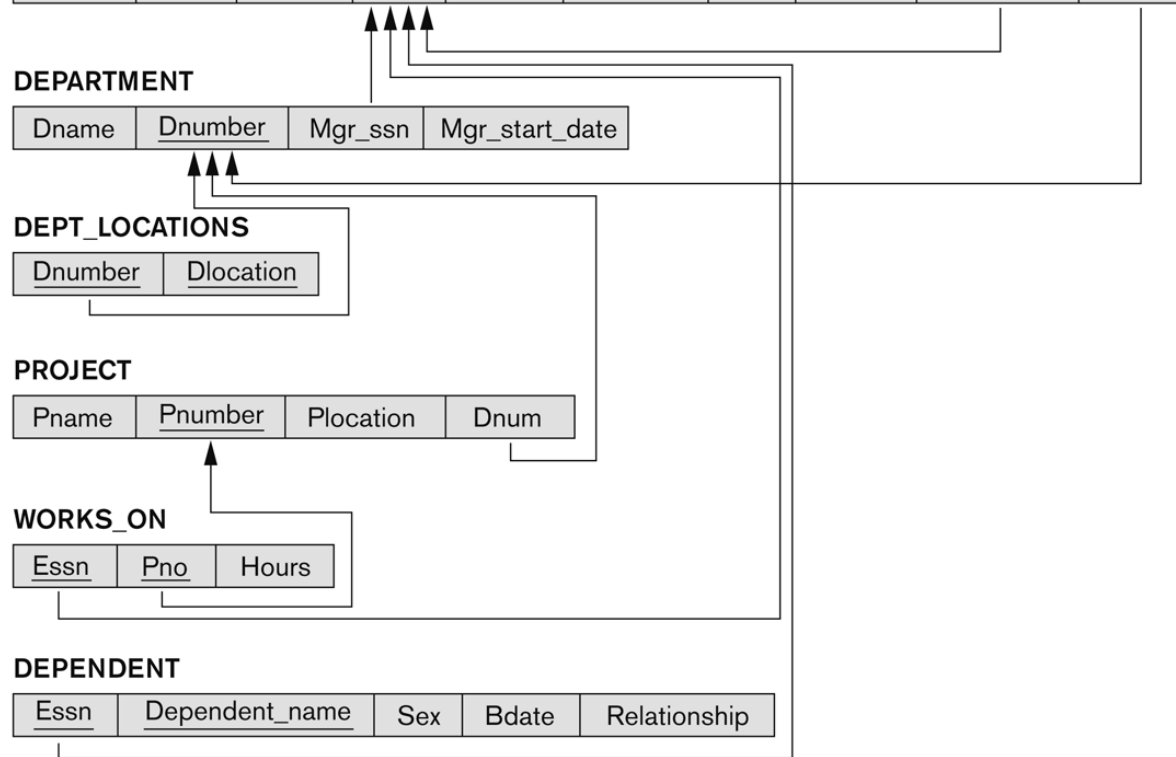
Pname	<u>Pnumber</u>	Plocation	Dnum
-------	----------------	-----------	------

WORKS_ON

<u>Essn</u>	<u>Pno</u>	Hours
-------------	------------	-------

DEPENDENT

<u>Essn</u>	<u>Dependent_name</u>	Sex	Bdate	Relationship
-------------	-----------------------	-----	-------	--------------



Other Types of Constraints

- ✓ Semantic Integrity Constraints:
 - based on application semantics and cannot be expressed by the model per se
 - Example: “the max. no. of hours per employee for all projects he or she works on is 56 hrs per week”
- ✓ A **constraint specification** language may have to be used to express these
- ✓ SQL-99 allows triggers and **ASSERTIONS** to express for some of these

Populated database state

- ✓ Each *relation* will have many tuples in its current relation state
- ✓ The *relational database state* is a union of all the individual relation states
- ✓ Whenever the database is changed, a new state arises
- ✓ Basic operations for changing the database:
 - INSERT a new tuple in a relation
 - DELETE an existing tuple from a relation
 - MODIFY an attribute of an existing tuple
- ✓ Next slide shows an example state for the COMPANY database

Populated database state for COMPANY

Figure 5.6

One possible database state for the COMPANY relational database schema.

EMPLOYEE

Fname	Minit	Lname	Ssn	Bdate	Address	Sex	Salary	Super_ssn	Dno
John	B	Smith	123456789	1965-01-09	731 Fondren, Houston, TX	M	30000	333445555	5
Franklin	T	Wong	333445555	1955-12-08	638 Voss, Houston, TX	M	40000	888665555	5
Alicia	J	Zelaya	999887777	1968-01-19	3321 Castle, Spring, TX	F	25000	987654321	4
Jennifer	S	Wallace	987654321	1941-06-20	291 Berry, Bellaire, TX	F	43000	888665555	4
Ramesh	K	Narayan	666884444	1962-09-15	975 Fire Oak, Humble, TX	M	38000	333445555	5
Joyce	A	English	453453453	1972-07-31	5631 Rice, Houston, TX	F	25000	333445555	5
Ahmad	V	Jabbar	987987987	1969-03-29	980 Dallas, Houston, TX	M	25000	987654321	4
James	E	Borg	888665555	1937-11-10	450 Stone, Houston, TX	M	55000	NULL	1

DEPARTMENT

Dname	Dnumber	Mgr_ssn	Mgr_start_date
Research	5	333445555	1988-05-22
Administration	4	987654321	1995-01-01
Headquarters	1	888665555	1981-06-19

DEPT_LOCATIONS

Dnumber	Dlocation
1	Houston
4	Stafford
5	Bellaire
5	Sugarland
5	Houston

WORKS_ON

Essn	Pno	Hours
123456789	1	32.5
123456789	2	7.5
666884444	3	40.0
453453453	1	20.0
453453453	2	20.0
333445555	2	10.0
333445555	3	10.0
333445555	10	10.0
333445555	20	10.0
999887777	30	30.0
999887777	10	10.0
987987987	10	35.0
987987987	30	5.0
987654321	30	20.0
987654321	20	15.0
888665555	20	NULL

PROJECT

Pname	Pnumber	Plocation	Dnum
ProductX	1	Bellaire	5
ProductY	2	Sugarland	5
ProductZ	3	Houston	5
Computerization	10	Stafford	4
Reorganization	20	Houston	1
Newbenefits	30	Stafford	4

DEPENDENT

Essn	Dependent_name	Sex	Bdate	Relationship
333445555	Alice	F	1986-04-05	Daughter
333445555	Theodore	M	1983-10-25	Son
333445555	Joy	F	1958-05-03	Spouse
987654321	Abner	M	1942-02-28	Spouse
123456789	Michael	M	1988-01-04	Son
123456789	Alice	F	1988-12-30	Daughter
123456789	Elizabeth	F	1967-05-05	Spouse

Update Operations on Relations

- ✓ INSERT a tuple.
- ✓ DELETE a tuple.
- ✓ MODIFY a tuple.
- ✓ Integrity constraints should not be violated by the update operations.
- ✓ Updates may **propagate** to cause other updates automatically. This may be necessary to maintain integrity constraints.

Update Operations on Relations

- ✓ In case of integrity violation, several actions can be taken:
 - Cancel the operation that causes the violation (RESTRICT or REJECT option)
 - Perform the operation but inform the user of the violation
 - Trigger additional updates so the violation is corrected (CASCADE option, SET NULL option)
 - Execute a user-specified error-correction routine

Possible violations for each operation

- ✓ INSERT may violate any of the constraints:
 - Domain constraint:
 - if one of the attribute values provided for the new tuple is not of the specified attribute domain
 - Key constraint:
 - if the value of a key attribute in the new tuple already exists in another tuple in the relation
 - Referential integrity:
 - if a foreign key value in the new tuple references a primary key value that does not exist in the referenced relation
 - Entity integrity:
 - if the primary key value is null in the new tuple

Possible violations for each operation

- ✓ DELETE may violate only referential integrity:
 - If the primary key value of the tuple being deleted is referenced from other tuples in the database
 - Can be remedied by several actions: RESTRICT, CASCADE, SET NULL
 - RESTRICT option: reject the deletion
 - CASCADE option: propagate the new primary key value into the foreign keys of the referencing tuples
 - SET NULL option: set the foreign keys of the referencing tuples to NULL
 - One of the above options must be specified during database design for each foreign key constraint

Possible violations for each operation

- ✓ UPDATE may violate domain constraint and NOT NULL constraint on an attribute being modified
- ✓ Any of the other constraints may also be violated, depending on the attribute being updated:
 - Updating the primary key (PK):
 - Similar to a DELETE followed by an INSERT
 - Need to specify similar options to DELETE
 - Updating a foreign key (FK):
 - May violate referential integrity
 - Updating an ordinary attribute (neither PK nor FK):
 - Can only violate domain constraints

Summary

- ✓ Relational Model Concepts
 - Definitions
 - Characteristics of relations
- ✓ Discussed Relational Model Constraints and Relational Database Schemas
 - Domain constraints'
 - Key constraints
 - Entity integrity
 - Referential integrity
- ✓ Described the Relational Update Operations and Dealing with Constraint Violations

The image features a spiral-bound notebook with a cream-colored page. A horizontal line is drawn across the upper portion of the page. The title 'Relational Algebra' is centered on the page in a bold, brown font. The spiral binding is visible on the left side of the notebook.

Relational Algebra

Relational Query Languages

- ✓ Query languages: Allow manipulation and **retrieval of data** from a database.
- ✓ Relational model supports simple, powerful QLs:
 - Strong formal foundation based on logic.
 - Allows for much optimization.
- ✓ Query Languages **!=** programming languages!
 - QLs not expected to be “Turing complete”.
 - QLs not intended to be used for complex calculations.
 - QLs support easy, efficient access to large data sets.

Formal Relational Query Languages

- ✓ Two mathematical Query Languages form the basis for “real” languages (e.g. SQL), and for implementation:
- Relational Algebra: More **operational(procedural)**, very useful for representing execution plans.
 - Relational Calculus: Lets users describe what they want, rather than how to compute it. (**Non-operational, declarative.**)

Preliminaries

- ✓ A query is applied to relation instances, and the result of a query is also a relation instance.
 - *Schemas of input* relations for a query are **fixed**
 - The **schema for the result** of a given query is also **fixed**
Determined by definition of query language constructs.
- ✓ Positional vs. named-field notation:
 - Positional notation easier for formal definitions, named-field notation more readable.
 - Both used in SQL

Example Instances

✓ “Sailors” and “Reserves” relations for our examples. “bid”= boats.

“sid”: sailors

✓ We’ll use positional or named field notation

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/96
58	103	11/12/96

R1

S1

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

Relational Algebra

✓ Basic operations:

- Selection (σ) Selects a subset of rows from relation.
- Projection (π) Deletes unwanted columns from relation.
- Cross-product (\times) Allows us to combine two relations.
- Set-difference ($-$) Tuples in reln. 1, but not in reln. 2.
- Union (\cup) Tuples in reln. 1 and in reln. 2.

✓ Additional operations:

- Intersection, join, division, renaming: Not essential, but (very!) useful.

✓ Since each operation returns a relation, **operations can be composed!** (Algebra is “closed”.)

Selection

- ✓ Selects rows that satisfy *selection condition*.
- ✓ *Schema* of result identical to schema of (only) input relation.
- ✓ *Result* relation can be the *input* for another relational algebra operation (*Operator composition*.)

sid	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

$$\sigma_{rating > 8}^{(S2)}$$

sname	rating
yuppy	9
rusty	10

$$\pi_{sname, rating}(\sigma_{rating > 8}^{(S2)})$$

Projection

- ✓ Deletes attributes that are not in *projection list*.
- ✓ *Schema* of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- ✓ Projection operator has to eliminate *duplicates*! (Why?? what are the consequences?)
 - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

$\pi_{sname, rating}(S2)$

age
35.0
55.5

$\pi_{age}(S2)$

Union, Intersection, Set-Difference

- ✓ All of these operations take two input relations, which must be union-compatible:
 - Same number of fields.
 - `Corresponding' fields have the same type.
- ✓ What is the *schema* of result?

sid	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
44	guppy	5	35.0
28	yuppy	9	35.0

$S1 \cup S2$

sid	sname	rating	age
22	dustin	7	45.0

$S1 - S2$

sid	sname	rating	age
31	lubber	8	55.5
58	rusty	10	35.0

$S1 \cap S2$

Cross-Product

- ✓ Each row of S1 is paired with each row of R1.
- ✓ *Result schema* has one field per field of S1 and R1, with field names 'inherited' if possible.
 - *Conflict*: Both S1 and R1 have a field called *sid*.

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	22	101	10/10/96
31	lubber	8	55.5	58	103	11/12/96
58	rusty	10	35.0	22	101	10/10/96
58	rusty	10	35.0	58	103	11/12/96

▪ Renaming operator: $\rho (C(1 \rightarrow sid1, 5 \rightarrow sid2), S1 \times R1)$

Joins

✓ Condition Join: $R \bowtie_c S = \sigma_c (R \times S)$

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	58	103	11/12/96

✓ *Result schema* same as that of cross-product.

✓ Fewer tuples than cross-product. Filters tuples not satisfying the join condition.

✓ Sometimes called a *theta-join*.

$$S1 \bowtie_{S1.sid < R1.sid} R1$$

Joins

- ✓ Equi-Join: A special case of condition join where the condition c contains only **equalities**.

sid	sname	rating	age	bid	day
22	dustin	7	45.0	101	10/10/96
58	rusty	10	35.0	103	11/12/96

- ✓ Result schema similar to cross-product, but only one copy of fields for which equality is specified.
- ✓ Natural Join: Equijoin on *all* common fields.

$$\pi_{sid, \dots, age, bid, \dots} (S1 \bowtie_{sid} R1)$$

Division

- ✓ Not supported as a primitive operator, but useful for expressing queries like:

Find sailors who have reserved all boats.

- ✓ Precondition: in A/B , the attributes in B must be included in the schema for A . Also, the result has attributes $A-B$.
 - SALES(supId, prodId);
 - PRODUCTS(prodId);
 - Relations SALES and PRODUCTS must be built using projections.
 - SALES/PRODUCTS: the ids of the suppliers supplying ALL products.

Examples of Division A/B

sno	pno
s1	p1
s1	p2
s1	p3
s1	p4
s2	p1
s2	p2
s3	p2
s4	p2
s4	p4

A

pno
p2

B1

sno
s1
s2
s3
s4

A/B1

pno
p2
p4

B2

sno
s1
s4

A/B2

pno
p1
p2
p4

B3

sno
s1

A/B3

Expressing A/B Using Basic Operators

- ✓ Division is not essential op; just a useful shorthand.
 - (Also true of joins, but joins are so common that systems implement joins specially. Division is NOT implemented in SQL).
- ✓ *Idea:* For *SALES/PRODUCTS*, compute all products such that there exists at least one supplier not supplying it.
 - *x* value is *disqualified* if by attaching *y* value from *B*, we obtain an *xy* tuple that is not in *C*.

$$C = \pi_{sid}((\pi_{sid}(Sales) \times Products) - Sales)$$

The answer is $\pi_{sid}(Sales) - C$

Find names of sailors who've reserved boat #103

✓ Solution 1: $\pi_{sname}((\sigma_{bid=103} Reserves) \bowtie Sailors)$

✓ Solution 2: $\rho(Temp1, \sigma_{bid=103} Reserves)$

$\rho(Temp2, Temp1 \bowtie Sailors)$

$\pi_{sname}(Temp2)$

✓ Solution 3: $\pi_{sname}(\sigma_{bid=103}(Reserves \bowtie Sailors))$

Find names of sailors who've reserved a red boat

- ✓ Information about boat color only available in Boats; so need an extra join:

$$\pi_{sname}((\sigma_{color='red'} Boats) \bowtie Reserves \bowtie Sailors)$$

- ✓ A more efficient solution:

$$\pi_{sname}(\pi_{sid}((\pi_{bid} \sigma_{color='red'} Boats) \bowtie Res) \bowtie Sailors)$$

A query optimizer can find this, given the first solution!

Find sailors who've reserved a red or a green boat

- ✓ Can identify all red or green boats, then find sailors who've reserved one of these boats:

$$\rho \text{ (Tempboats, } (\sigma_{color='red' \vee color='green'} \text{ Boats}))$$

$$\pi_{sname}(\text{Tempboats} \bowtie \text{Reserves} \bowtie \text{Sailors})$$

- ❖ Can also define Tempboats using union! (How?)

- ❖ What happens if \vee is replaced by \wedge in this query

Find sailors who've reserved a red and a green boat

- ✓ Previous approach won't work! Must identify sailors who've reserved red boats, sailors who've reserved green boats, then find the intersection (note that *sid* is a key for *Sailors*):

$$\rho (Tempred, \pi_{sid}((\sigma_{color='red'} Boats) \bowtie Reserves))$$

$$\rho (Tempgreen, \pi_{sid}((\sigma_{color='green'} Boats) \bowtie Reserves))$$

$$\pi_{sname}((Tempred \cap Tempgreen) \bowtie Sailors)$$

Find the names of sailors who've reserved all boats

- ✓ Uses division; schemas of the input relations to / must be carefully chosen:

$$\rho \text{ (} Temp\text{sids, } (\pi_{sid,bid} Reserves) / (\pi_{bid} Boats) \text{)}$$

$$\pi_{sname} (Temp\text{sids} \bowtie Sailors)$$

- ✓ To find sailors who've reserved all 'Interlake' boats:

$$\dots / \pi_{bid} (\sigma_{bname='Interlake'} Boats)$$

Summary

- ✓ The relational model has rigorously defined query languages that are simple and powerful.
- ✓ Relational algebra is more operational; useful as internal representation for query evaluation plans.
- ✓ Several ways of expressing a given query; a query optimizer should choose the most efficient version.