CSc 460 Database Design

Fall 2014

Richard Snodgrass

Why Use a Database?

CSc 460 Database Design

Fall 2014

Acknowledgements

- These slides were written by Richard T. Snodgrass (University of Arizona) and Christian S. Jensen (Aalborg University).
- Michael Soo (amazon.com) provided some of the query processing slides.
- Kristian Torp (Aalborg University) converted the slides from Island Presents to Powerpoint.
- Curtis Dyreson (Washington State University) merged the slides with those from Peter Stewart (Bond University).
- The motivational example was given by Gary Lindstrom and Wei Tao (University of Utah).

Why Use a Database?

Outline

- Why Use a Database?
- The Field

Why Use a Database

• Overview of the Course

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, M. D. Soo, and C. E. Dyreson

M-3

M-2

Prevalence of Databases

- Behind every successful website, there is a powerful database.
- Examples:
 - Amazon's website
 - American Airlinesz
 - Dell's ordering system
 - UPS / FedEx tracking
 - · Wal-Mart's inventory system

Why Use a Database?

Data Management Example

- Scenario
 - You are in a movie web startup.
 - Your customers rent DVD copies of movies.
 - Several copies of each movie.

• Needs

- · Which movies have a customer rented?
- Are any disks overdue?
- When will a disk become available?

Why Use a Database?



Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, M. D. Soo, and C. E. Dyreson



Complication: Integrity

- · Lacks data integrity, consistency
 - Clerk misspells value/field
 - Customer: Jane Doek, Rented: Eraserhead, Deu: Jan. 19, 2010
 Inputs improper value, same value differently
 - Customer: Jane Doe, Rented: The Eraserhead, Due: Feb. 29, 2010 Forgets/adds/reorders field
 - Terms: weekly special Due: Jan. 19, 2010, Rented: Eraserhead
- Requirements
 - Enforce constraints to permit only valid information to be input.

DBMS Concepts Integrity constraints Types

Why Use a Database?



м-

м-



- Two clerks edit rented.txt file at the same time.
 1) Ben starts to edit rented.txt, reads it into memory.
 2) Sarah starts to edit rented.txt.
 - 3) Ben adds a record.
 - 4) Ben saves rented.txt to disk.
 - 5) Sarah saves rented.txt to disk.
 - Ben's added record disappears!

· Requirements

Why Use a Database?

- · Must support multiple readers and writers.
- Updates to data must (appear to) occur in serial order.

DBMS Concepts Serializability Concurrency control

Complication: Crashes

- Crash during update may lead to inconsistent state.
- Crash during update may lead to inconsistent state.
 - Ben makes 250 of 500 edits to change Jane Doe to her preferred name Jan Doe.
 Before he saves it, Windows crashes!
- Before ne saves it, windows crashe

· Requirements

- Must update on all-or-none basis.
- Implemented by commit or rollback if necessary.

DBMS Concepts Transactions Commit Rollback Recovery

Why Use a Database?





M-1

M-14



Complication: Efficiency

- · Your customer list grows virally.
 - rented.txt file gets huge (gigabytes of data).
 - Slow to edit.

Why Use a Database?

Slow to query for customer information.

Requirements

- New data structures to improve query performance.
- System automatically modifies queries to improve speed.
- Ability of system to scale to handle huge datasets.

DBMS Concepts Indexes Query optimization Database tuning

Why Use a Database?



M-1

M-17

Limitations of File-based Systems

- Program must implement
 - Security
 - Concurrency control
 - Support for schema reorganization
 - · Performance enhancing data structures, e.g., indexes
- Observation
 - Many applications need these services.
- Solution
 - Build and sell a software system to provide services!

Why Use a Database?

Outline

- Why Use a Database?
- The Field
- Overview of the Course

Why Use a Database?

Major Players in the Field

- Oracle
 - 108,000 employees
 - \$38.2B (US dollars) annual revenue (FY 2014)
 - 34% market share
- Microsoft
 - Produces SQLServer
 - 89,000 employees
 - \$86.7B annual revenue (includes Office, Windows, etc.)
- IBM
 - Produces DB2
 - 105,000 employees in US, 75,000 employees in India
 - \$99.8B annual revenue

Why Use a Database?

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, M. D. Soo, and C. E. Dyreson

M-18





The Field

· Professional Publications

- ACM Transactions on Database Systems (TODS)
- quarterly, 700 pages per year
 IEEE Transactions on Knowledge and Data Engineering (TKDE)
- ▲ bimonthly, 1000 pages a year
- The VLDB Journal (VLDBJ)
- quarterly, 450 pages a year
- Information Systems (Info Sys)
- ▲ bimonthly, 600 pages a year
- Numerous other journals
- Unrefereed Technical Publications
 - ACM SIGMOD Record
 - quarterly, 250 pages a year
 - IEEE Data Engineering Bulletin

▲ quarterly, 250 pages a year

Why Use a Database?

The Field, cont.

Conferences

- ACM SIGMOD International Conference on Management of Data (SIGMOD)
- ▲ late May or early June, 500 pages a year
- ACM Principles of Database Systems (PODS)
- ▲ In conjunction with SIGMOD, 300 pages a year
- International Conference on Very Large Data Bases (VLDB)
- Mid-September, 500 pages a year
 IEEE International Conference on Data Engineering (ICDE)
- February, 600 pages a year
- Extending Data Base Technology (EDBT), International Conference on Database Theory (ICDT)
- ▲ alternate years (March and January), 400 pages a year
- 8-10 specialized conferences a year: 300 × 8 = 2500 pages a year

• 8,000 pages per year of research papers

Why Use a Database?

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, M. D. Soo, and C. E. Dyreson

M-21

м-

M-2

M-23

The Field, cont.

• Other

- · Technical reports from university and industrial research labs
- Trade rags: Data Base Newsletter, Database Review, InfoDB, etc.: thousands
 of pages a year
- Perhaps 100 database textbooks
- · Several hundred database user guidebooks

Outline

Why Use a Database?

- Why Use a Database?
- The Field
- · Overview of the Course

The Course

Why Use a Database?

- Motivation (M)
- The Relational Model (R)
 - Definition
 - Integrity constraints
- Relational Query Languages (E)
 - Relational algebraRelational calculi: domain and tuple
 - Relational calculi: domain and tuple
- Microsoft Access (A)
 - DDL

DML
 Why Use a Database?

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, M. D. Soo, and C. E. Dyreson

M-24

| The Course, cont. | |
|--|------|
| SQL DDL and DML (S) Schema definition Querying Modifications Views Embedded SQL | |
| Transaction management Conceptual Design (C) Basic Entity-Relationship Model Constraints Specialization and generalization View integration Mapping an ER schema to tables | |
| Why Hee a Database? | M-25 |

The Course, cont.

• Web/Database Connections (W)

- HTTP
- Static Database Access
- Dynamic Database Access
- Application Design and Implementation (T)
 - Direct database interfaces
 - Indirect interfaces
 - Web interfaces

Why Use a Database?

The Course, cont.

- Logical Design (L)
 - Properties of a good design
 - Functional dependencies and keys
 - Normal forms: 3NF, BCNF
 - Decomposition algorithms
- Physical Data Organization (P)
 - Relational structures: heap, sorted, compressed
 - Indexes: primary and secondary, B-trees
- The Future (\mathbf{F})
 - Data models
 - Query languages
 - Next-generation database technologies

Why Use a Database?

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, M. D. Soo, and C. E. Dyreson

M-27

M-3

















The Relational Model

CSc 460 Database Systems

Fall 2014

Acknowledgements

- These slides were written by Richard T. Snodgrass (University of Arizona) and Christian S. Jensen (Aalborg University).
- Kristian Torp (Aalborg University) converted the slides from island presents to Powerpoint.
- Curtis E. Dyreson (Washington State University) modified the slides.

Overview

Relational Model

- The Relational Model (RM)
 Relations
 - Properties of relations
- Integrity Constraints

The Relational Model

| Course Concepts | |
|----------------------|-----|
| | - |
| | |
| | |
| Relational Model | |
| | |
| | |
| | |
| The Relational Model | R-4 |



Sets

Definition: A *set* is an *unordered* collection of *distinct* objects.

Examples

 {3, 4, a} is a set
 {a, 3, 4} is the same set as {3, 4, a}
 {4, 4} is not a set

· Set operations include

- *Intersection*, e.g., $\{3, 4, a\} \cap \{b, 4\} = \{4\}$
- Union, e.g., $\{3, 4, a\} \cup \{b, 4\} = \{3, 4, a, b\}$
- *Difference*, e.g., {3, 4, a} {b, 4} = {3, a}

e Relational Model

Domains and Attributes

Definition: A domain is a set of values of some "type".

- Positive integers = {1, 2, 3, 4, ...}
- Alphanumeric characters = {'a', 'b', ..., 'Z', '0', ..., '9'}
- A common database restriction: Domains are atomic.
 - The following programming types are atomic domains. integer, char, float, varchar
 - Structs and records are associated with a composite domain.
 struct { street: char[], city: char[], state: char[2] } address

Definition: An attribute is the name of a domain.

The Relational Model

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Relations

Definition: Given sets $A_1, A_2, ..., A_n$ a *relation r* is a subset of their Cartesian product.

 $r \subseteq A_1 \times A_2 \times \dots \times A_n$

• *r* is a set of *n*-tuples $(a_1, a_2, ..., a_n)$ where $a_i \in A_i$

Definition: $R(A_1, A_2, ..., A_n)$ is the *schema* of relation *r*.

- A_1, A_2, \dots, A_n are *domains*; their names are *attributes*.
- *R* is the name of the relation.
- *Notation:* r(R) denotes that *r* is a relation on the relation schema *R*.



ACM Turing Award winner







Tuples

Definition: An element *t* of a relation *r* is called a *tuple*.

- We refer to component values of a tuple *t* by $t[A_i] = v_i$ (the value of attribute A_i for tuple *t*).
 - Alternatively, use 'dot' notation, e.g., tA_i is the i^{th} attribute of tuple t
- $t[A_i, ..., A_k]$ refers to the subtuple of t containing the values of attributes $A_i, ..., A_k$ respectively.
- · Table metaphor
 - Tuple is a row
 - Attribute is a column

e Relational Model

Characteristics of Relations

- Tuples in a relation are *unordered*.
- Example: The following two relations have the same information content.



Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.







Characteristics of Relations (cont.)









- An *integrity constraint* is a *predicate* that is required to be true of every possible instance of the database schema.
- Purpose: To keep the database in a "good" state.
- A good state is one that is consistent with the modelled reality.How: Check to see if an update will move to a "bad" state.





- Are specified at the schema level
 - A constraint must be true for all possible instances of a schema.
 - Critical notion: *possible* instances
- Some can be specified (conceptually) with relational algebra.
- In commercial databases, specified with SQL (typically).

The Relational Model

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Enforcement of Integrity Constraints

- · Evaluate on update (insert/delete/modification)
- · Cost of enforcement
 - Only allow specification of ICs that can be tested efficiently.
 - How many to evaluate?

Remedy for Violation

- · When an integrity violation occurs either
 - Cancel the operation that causes the violation.
 - Perform the operation but inform the user of the violation.
 - Trigger additional updates so the violation is corrected.
 - Execute a user-specified error-correction routine.
- · Specified for each individual constraint

Key Constraints

Definition: A *superkey* of R is a set of attributes K of R such that no two tuples in any valid relation instance r(R) will have the same value for K.

• For any distinct tuples t_1 and $t_2 \in r(R)$, $t_1 \neq t_2 \Longrightarrow t_1[K] \neq t_2[K]$

Definition: A *candidate key* of *R* is a *minimal* superkey; that is, a superkey *K* such that removal of any attribute from *K* results in a set of attributes that is not a superkey.

Definition: A *primary key* is a specified candidate key.

Example

- Consider two relation schemas (with no keys specified).
 - Customer (ID, Name)
 - Reserved (CustomerID, FilmID, ResDate)
- · Assumption: No two customers have the same ID
 - Superkeys of Customer {ID}
 - {ID, Name}
 - Candidate key of Customer (and also the primary key, since only one) $\{{\rm ID}\}$
- Assumption: Customer can reserve many films over several days.
 The superkey of Reserved
 - {CustomerID, FilmID, ResDate}
 - One candidate key (see above), primary key (see above)

Keys, Fact and Fiction

- A key is often more than a single attribute.
 - Some DB products impose a single-attribute key restriction
 - Smaller keys are sometimes more efficient.
 - This optimization detracts from semantics.

• An attribute in a key can be of any type.

- Some DB products restrict keys to integer or auto-number types.
- Easy for the software to check and maintain
- Detracts from semantics

• Choosing a key is difficult.

- Schema may change over time.
- Assume SSN is a good key, new law makes using SSN illegal.
- Converting after a "wrong choice" can be costly.
- Keys are "shared" as foreign keys in other relations

Entity Integrity

- The primary key attributes *P* = {*a*₁, ..., *a*_m} of each relation schema *R* in *S* cannot have null values in any tuple of *r*(*R*).
 - $t[P] \neq Null$ for any tuple $t \in r(R)$
 - Alternatively, $\forall k \ 1 \le k \le m \ (\forall t \in r(R) \ t[a_k] \ne Null)$
- This is because primary key values are used to identify individual tuples.

The Relational Model

Domain Constraints

- · This is similar to programming languages
 - Strong typing or type-by-name in programming languages.
 - Most DBMSs have only a restricted selection of simple, predefined types available.
 - Each domain creates a separate type.
- Consider these sample relation schemas.

Film (<u>ID</u>, Title) Customer (<u>ID</u>, Name)

- It is not meaningful to compare Film and Customer IDs
 - Create FilmID domain of type integer
 - Create CustomerID domain of type integer
 - DBMS does not allow IDs of domain FilmIDs to be compared with IDs of domain CustomerIDs.

lational Model

Referential Integrity

- Referential integrity is the "glue" that keeps the database together, ensuring the consistency of foreign keys.
 - A foreign key *points to* or *references* some relation
 - A foreign key is a key in the other relation
 - Referential integrity says that the information that is pointed to must exist.
- Example:If a CustomerID, e.g., 123456, is associated with a FilmID, e.g., 100001, in a tuple in the Reserved relation, a customer with CustomerID 123456 must exist in the Customer relation, and a film with FilmID 100001 must exist in the Film relation.

Referential Integrity, Definition

Definition: Let r_1 and r_2 be instances of schemas R_1 and R_2 , let *K* be a candidate key for R_1 , and let α be a subset of R_2 that is compatible with *K*. Then, α is a foreign key for R_2 if every value for α in a tuple in r_2 also appears as a value of *K* in a tuple in r_1 .

 $\forall r_1(R_1) \ \forall r_2(R_2) \ \exists t_1 \in r_1 \ \forall t_2 \in r_2 \ (r_2[\alpha] \neq Null \land r_2[\alpha] = r_1[K])$

• Constraints of this type are termed *referential integrity constraints* or *subset constraints*.

The Relational Model

Enforcement of Referential Integrity

- Consider the constraint that every foreign key value in r_2 is in the referenced relation r_1 .
- Insertion into *r*₂:
- Check that a corresponding tuple exists in r_1 .
- Deletion from *r_l*:
 - Check that no matching tuple existed in r_2 .
- Updates of tuples in r_2 are treated as insertions above.
- Updates of tuples in r_1 are treated as deletions above.
- Note that cascading deletes are possible.

Assertions

- *Assertions* are more general integrity constraints, expressed directly as predicates which must always be satisfied.
- Think of these as sanity checks.
- Examples
 - No DVD rents for less than \$0.00.
 - The spoken and subtitled languages for a foreign film must be different.
 - Reservations must have a date of today.
 - No DVD is being rented by more than one customer.

ne Relational Model

Relational Query Languages

CSc 460 Database Systems

Fall 2014

Acknowledgements

- These slides were written by Richard T. Snodgrass (University of Arizona) and Christian S. Jensen (Aalborg University).
- Kristian Torp (Aalborg University) converted the slides from island presents to Powerpoint.
- Curtis Dyreson (Washington State University) modified the slides.
- Rick converted to Jun Yang's dbsym font, using George Necula's TexPoint package.
- Lester McCann's article on relational division was helpful: "On Making Relational Division Comprehensible," in *Proceedings of the 33rd ASEE/IEEE Frontiers in Education Conference*, pp. F2C-6 to F2C-11, Boulder, CO, November, 2003.

Overview

Relational Algebra (RA)

- Operators
 - Complete set projection, selection, Cartesian product, difference, union
 - Convenient intersection, theta join, equijoin, natural join, semijoin, antijoin, division
- Example queries
- Limitations
- Extended Relational Algebra Operators
- Tuple Relational Calculus (TRC)
- Domain Relational Calculus (DRC)

Relational Query Languages

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.





Algebra

- An algebra is a mathematical object with
 - A set of *objects*
 - A set of *operations*, each of which takes one or more objects and evaluates to an object
- Some algebras are *closed*: result of each operation on any possible input objects is a valid object.

Relational Algebra

- Five primary operators ("complete" or "query complete")
 - Selection: σ
 - Projection: π
 - Union: \cup
 - Difference: –
 - Cartesian Product: ×
- · Convenient derived operators
 - Intersection: \cap
 - Theta join n_{θ} , equijoin n_{\pm} , and natural join n
 - Semijoins: \bowtie and \bowtie
 - Antijoin: ▷
 - Relational division: ÷
- elational Query Languages





- Choose a row or rows, choice is based on a *condition* or *predicate*
- Result has same schema.
- Remember: rows are unordered so cannot specify row number, e.g., 2nd row or 4th row.

| Friends [] | V <i>ame</i> Tom Jane | Age 4 10 | Result Name Age Jane 10 |
|------------|-----------------------------|----------------|---|
| | Mel | 9 | Mel 9 |
| Selec | t rows | from . | Friends such that $Age > 8$: $\sigma_{Aee>8}(Friends)$. |

Selection, Formal Description

 $\sigma_P(r) = \{t \mid t \in r \land P(t)\}$

- *P* is a formula in *propositional calculus*, dealing with terms of the form:
 - attribute (or constant) = attribute (or constant)
 - attribute ≠ attribute
 - attribute < attribute
 - attribute ≤ attribute
 - attribute ≥ attribute
 - attribute > attribute
 torm a torm (N/);
 - term ∧ term (Note: ∧ is AND)
 term ∨ term (Note: ∨ is OR)
 - \neg (term) (*Note:* \neg *is NOT*)

Relational Query Language





































Overview

Relational Algebra (RA)

Operators

- Complete set projection, selection, Cartesian product, difference, union
- Convenient intersection, theta join, equijoin, natural join, semijoin, division
 Limitations
- Extended Relational Algebra Operators
- Tuple Relational Calculus (TRC)
- Domain Relational Calculus (DRC)

elational Query Languages



Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.





Joins

- Joins are Cartesian products coupled with selections and projections.
- Theta join
 - $r \bowtie_{\theta} s = \sigma_{\theta} (r \times s)$
 - Example: $r \bowtie_{A < E} s = \sigma_{A < E} (r \times s)$
- Equijoin
 - A theta join in which θ is an equality predicate
 - Example: $r \bowtie_{A=E} s = \sigma_{A=E} (r \times s)$
- Natural join ⋈
- Semijoins
 - Left semijoin ▷
 - Right semijoin >>>
- Antijoin ⊳















Movie Web Startup Schema (subset)

- Customer (<u>CustomerID</u>, Name, Street, City, State, Zipcode)
- Film (FilmID, Title, RentalPrice, Kind)
- Reserved (CustomerID, FilmID, ResDate)

Movie Web Queries

- List information for films with a rental price over \$4.
 - $\sigma_{RentalPrice > 4}(Film)$
- List the titles of films with a rental price over \$4.

 $\pi_{\text{Title}} (\sigma_{\text{RentalPrice} > 4}(\text{Film}))$

Movie Web Queries, cont.

• List outrageously priced films (over \$4 or under \$1).

 $\pi_{\text{Title}} \; (\sigma_{\text{RentalPrice} > 4}(Film) \cup \sigma_{\text{RentalPrice} < 1}(Film))$

 $\pi_{Title} \; (\sigma_{RentalPrice \; > \; 4 \; \lor RentalPrice \; < \; 1}(Film))$

• List the ID numbers of the films that are expensive and have been reserved.

 $\pi_{FilmId} \ (\sigma_{RentalPrice \ > \ 4}(Film)) \cap \pi_{FilmId}(Reserved)$

 $\pi_{\text{FilmId}} (\sigma_{\text{RentalPrice} > 4}(\text{Film})) \bowtie \pi_{\text{FilmId}}(\text{Reserved})$

Relational Query Languages

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Movie Web Queries, cont.

• List the IDs of the expensive films that have not been reserved.

 $\pi_{FilmId} \ (\sigma_{RentalPrice > 4}(Film)) - \pi_{FilmId}(Reserved)$

Movie Web Queries, cont.

• List the titles of all reserved films.

 $\pi_{Title} \; (\sigma_{Film.FilmID \; = \; Reserved.FilmID}(Film \times Reserved))$

 π_{Title} (Film \bowtie Reserved)

• List the customers who have reserved a film.

 $\pi_{Name} (\sigma_{Customer.CustomerID = Reserved.CustomerID} \\ (Customer \times Reserved))$

 π_{Name} (Customer \bowtie Reserved)

Relational Query Languages

Movie Web Queries, cont.

• List the customers who have reserved expensive films.

 $\pi_{Name} \ (\sigma_{RentalPrice > 4} (Customer \bowtie Reserved \bowtie Film))$

 $\pi_{Name} \ (Customer \bowtie \sigma_{RentalPrice > 4} \ (Reserved \bowtie Film))$

 $\pi_{Name} \; (Customer \bowtie Reserved \bowtie \sigma_{RentalPrice > 4} \; (Film))$

• List the streets of customers who have reserved foreign films.

 π_{Street} (Customer $\bowtie \text{Reserved} \bowtie \sigma_{\text{Kind} = "F"}$ (Film))

Relational Query Languages

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Renaming Operator

- Find the film(s) with the highest rental price.
- We need a renaming operator: ρ_{Name}
- This operator doesn't really do anything, other than make it easier to do self-joins.
- Alternative formulation of the query: Find the film(s) with a rental price for which no other rental price is higher.

 $\pi_{\text{Title}}(\text{Film}) -$

 $\pi_{F2.Title} \ (Film \bowtie_{Film.RentalPrice > F2.RentalPrice} \ \rho_{F2}(Film))$





Semijoins in Distributed Query Processing

- Assume that relation *r* is very large, and resides in a database in Los Angeles.
- Relation *s* is small, and resides in a database in New York City.
- Application desires $r \bowtie s$, in New York City.
- Sending all of *r* from Los Angeles to New York City would be very expensive.
- Alternative processing strategy:
 - New York City sends $\pi_{R \cap S}(s)$ to Los Angeles.
 - Los Angeles computes *r* ⊨ *s*, and sends the result to New York City. This table is significantly smaller than *r*.
 - New York City computes $(r \bowtie s) \bowtie s = r \bowtie s$

Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.













Relational Division

- Let *r* and *s* be relations on schemes *R* and *S* respectively, where
 - $R = (A_1, A_2, ..., A_n, B_1, B_2, ..., B_n,)$
 - $S = (B_1, B_2, ..., B_n)$
- The result of *r* divided by *s* is a relation on scheme $R \div S = (A_1, A_2, ..., A_n)$
- $(r \div s) \times s \subseteq r$

al Ouery Lan

Movie Web Queries, cont.
 List the customers who have reserved *all* the foreign films.
 Identify the foreign films.
 π_{FilmID} (σ_{Kind = "F"} (Film))
 Identify those customers who have reserved all foreign films.

 $\pi_{\text{CustomerID,FilmID}}$ (Reserved) $\div \pi_{\text{FilmID}}$ ($\sigma_{\text{Kind} = \text{"F"}}$ (Film))

- Now figure out the names of those customers. $\pi_{Name} \ (Customer \bowtie$

 $\begin{array}{l}(\pi_{CustomerID,FilmID} \ (Reserved) \\ \div \ \pi_{FilmID} \ (\sigma_{Kind = ``F''} \ (Film))))\end{array}$

Hints

- 1. Make right argument one column of the stuff after "all."
- 2. Make left argument two columns, one the stuff before "all" and the other the column above.
- 3. Both arguments should be IDs of something.

Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

| Customers Who Have Reserved | All the For | eign Films | |
|---|------------------------|------------|--|
| • $R = \{$ CustomerID, FilmID $\}$ "R | eservations | " | |
| • $r = \pi_{\text{CustomerID FilmID}}$ (Reserved) <i>CustomerID Film</i> | | | |
| ; | 12345 | 1 | |
| | 12345 | 2 | |
| • $S = \{ FilmID \}$ "Foreign Films" | 12345 | 7 | |
| • $s = \pi_{\text{FilmID}} (\sigma_{\text{Kind} = \text{"F"}} (\text{Film}))))$ | 24680 | 1 | |
| Finind C Kind - F C 7777 | 24680 | 19 | |
| FilmID | 36925 | 1 | |
| 1 | 36925 | 19 | |
| 7 | 36925 | 7 | |
| Shared columns: B = { FilmII Extra columns in R: A = { Custor Reservations ÷ Foreigh Films is w |)} nerID } what? | | |
| Relational Query Languages | | E-43 | |

Relational Division

- Let *r* and *s* be relations on schemes *R* and *S* respectively, where
 - $R = (A_1, A_2, ..., A_n, B_1, B_2, ..., B_n,)$
 - $S = (B_1, B_2, ..., B_n)$
- The result of *r* divided by *s* is a relation on scheme $R \div S = (A_1, A_2, ..., A_n)$
- Can be done with just the five basic operators (!): $r \div s = \pi_{R-S} (r) - \pi_{R-S} ((\pi_{R-S} (r) \times s) - r)$

| Division Example, con | nt. | |
|--|--|-------------------------------|
| 1. $\pi_{R-S}(r)$: (the <i>A</i> columns of <i>r</i>) | CustomerID 12345 24680 36925 | |
| 2. $\pi_{R-S}(r) \times s$: (This is the maximum possible result.) | CustomerID 12345 12345 24680 24680 | <i>FilmID</i> 1 7 1 7 1 7 7 7 |
| 3. $(\pi_{R-S}(r) \times s) - r$: (the possible tuples that <i>aren't</i> in <i>r</i> , termed the | 36925 36925 <i>CustomerID</i> 24680 | 1 7 <i>FilmID</i> 7 |







Intuition for Equivalence

- $r \div s = \pi_{R-S}(r) \pi_{R-S}((\pi_{R-S}(r) \times s) r)$
- $\pi_{R-S}(r)$: the A columns $(A_1, A_2, ..., A_n)$ of r
- $\pi_{R-S}(r) \times s$: the *A* and *B* columns; the maximum possible product
- $(\pi_{R-S}(r) \times s) r$: the possible tuples that *aren't* in *r* (termed the *victims*)
- $\pi_{R-S}((\pi_{R-S}(r) \times s) r)$: the *A* columns of the victims
- π_{R-S} (r): the A columns (A₁, A₂, ..., A_n) of r: the maximum possible result of r ÷ s, note: r ÷ s ⊆ π_{R-S} (r)
- $\pi_{R-S}(r) \pi_{R-S}((\pi_{R-S}(r) \times s) r)$: the resulting tuples (just the *A* columns) that *aren't* victims

.....

Relational Completeness

- All the operators can be expressed in terms of the five basic operators: σ, π, −, ×, ∪ (along with ρ).
- This set is called a *complete set* of relational algebraic operators.
- Any query language that is at least as powerful as these operators is termed (query) relationally complete.

Relational Query Languages

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Limitations of the Algebra

- Can't do arithmetic.
- Find the rental price assuming a 10% increase.Can't do aggregates.
 - How many films has each customer reserved?
- Can't handle "missing" data.
- Make a list of the films, along with who reserved it, if applicable.
- Can't perform transitive closure.
 - For a partof(Part, ConstituentPart) relation, find all parts in the car door.
- Can't sort, or print in various formats.
 Print a reserved summary, sorted by customer name.
- Can't modify the database.
 - Increase all \$3.25 rentals to \$3.50.
- tional Query Languages

Overview

Relational Algebra (RA)

- Operators
 - Complete set projection, selection, Cartesian product, difference, union
 Convenient intersection, theta join, equijoin, natural join, semijoin,
- division

 Example queries
- Example que
 Limitations
- Extended Relational Algebra Operators
- Tuple Relational Calculus (TRC)
- Domain Relational Calculus (DRC)

telational Query Language

Extending the Projection Operator

- The *generalized* projection operator allows expressions in addition to column names in the subscript.
- Find the rental price, assuming a 10% increase.

 $\pi_{\text{Title, RentalPrice}*1.1}(\text{Film})$

Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.



- Appends new attributes.
- Each aggregate function agg_i decides which attribute to aggregate over, e.g., agg₂ could be the sum of the third attribute, or the average of the fifth attribute.
 Note: Text does not retain original columns.
- elational Overy Languages





Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Aggregates Functions, cont.

• How many reservations are there at each rental price?

 $\pi_{RentalPrice, \ count}(_{RentalPrice}\mathsf{F}_{count} \ (Reserved \bowtie Film))$

Outer Joins

- In a regular equijoin or natural join, tuples in *r* or *s* that do not having matching tuples in the other relation do not appear in the result.
- The outer joins retain these tuples, and place nulls in the missing attributes.
- Left outer join:

 $r \bowtie s = r \bowtie s \cup ((r - (r \bowtie s)) \times (null, ..., null))$ • Right outer join:

- $r \bowtie s = r \bowtie s \cup ((null, ..., null) \times (s (s \bowtie r)))$
- Full outer join:

 $r \Leftrightarrow s = r \bowtie s \cup ((r - (r \bowtie s)) \times \{(null, ..., null)\})$ $\cup (\{(null, ..., null)\} \times (s - (s \bowtie r)))$

Note: text uses a slightly different notation.

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Movie Web Query

• Make a list of the films, along with who reserved it, *if applicable*.

 $\pi_{\text{Title, Name}}(\text{Film} \Join (\text{Reserved} \Join \text{Customer}))$

Overview

• Relational Algebra (RA)

- Operators Limitations
- Extended Relational Algebra Operators
- Tuple Relational Calculus (TRC)
- Domain Relational Calculus (DRC)

Relational Calculus

- The relational algebra is *procedural*, specifying *how*, that is, a sequence of operations to derive the desired results.
- Relational calculus is based on first-order predicate calculus.
- Relational calculus is more *declarative*, specifying *what* is desired.
- The expressive power of the two languages is identical.
 This implies that relational calculus is relationally complete.
- Many commercial relational query languages are based on the relational calculus.
- The implementations are based on the relational algebra.

 Retired Very Languages

Tuple Relational Calculus

- A tuple variable ranges over tuples of a particular relation.
- Example: List the information about expensive films. {*t* / *Film*(*t*) ∧ *t*.*RentalPrice* > 4}
- *Film(t)* specifies the range relation *Film* for the tuple variable *t*.
- Each tuple satisfying *t*.*RentalPrice* > 4 is retrieved.
- The entire tuple is retrieved.
- List the titles of expensive films.

{ $t.Title | Film(t) \land t.RentalPrice > 4$ }

Relational Query Languag

Syntax of Tuple Relational Calculus

```
\{t_1.A_1, t_2.A_2, \dots, t_j.A_j \mid P(t_1, t_2, \dots, t_j)\}
```

- t_1, t_2, \dots, t_j are tuple variables.
- Each *t_k* is an tuple of the relation over which it ranges.
- *P* is a predicate, which is made up of atoms, of the following types.
 - $r(t_k)$ identifies r as the range of t_k
 - Atoms from predicate calculus: \land , \lor , \neg , \forall , \exists

Relational Query Languag

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Movie Web Queries

 List the outrageously priced films (over \$4 or under \$1). {t / Film(t) ∧ (t.RentalPrice > 4 ∨

t.RentalPrice < 1)}

 $\{t.FilmID, t.Title, t.RentalPrice, t.Kind \mid Film(t) \land \\ (t.RentalPrice > 4 \lor t.RentalPrice < 1) \}$

• List the ID numbers of the films that are expensive and have been reserved.

 $\begin{aligned} & \{t.FilmID \mid Film(t) \land (t.RentalPrice > 4 \\ & \land (\exists r \ (Reserved(r) \land t.FilmID = r.FilmID)) \} \end{aligned}$

Movie Web Queries, cont.

• List the IDs of the expensive films that have *not* been reserved.

 $\{t.FilmID \mid Film(t) \land (t.RentalPrice > 4 \\ \land \neg \exists r (Reserved(r) \land t.FilmID = r.FilmID)) \}$

Movie Web Queries, cont.

• List the titles of all reserved films.

· List the customers who have reserved a film.

{c.Name | $\exists r (Customer(c) \land Reserved(r) \land c.CustomerID = r.CustomerID)$ }

Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Movie Web Queries, cont.

• List the customers who have reserved expensive films.

 $\{c.Name \mid \exists r \exists f (Customer(c) \land Reserved(r) \\ \land Film(f) \land f.RentalPrice > 4 \\ \land c.CustomerID = r.CustomerID \\ \land r.FilmID = f.FilmID) \}$

• List the streets of customers who have reserved foreign films.

```
 \{ c.Street \mid \exists r \exists f (Customer(c) \land Reserved(r) \\ \land Film(f) \land f.Kind = "F" \\ \land c.CustomerID = r.CustomerID \\ \land r.FilmID = f.FilmID) \}
```


- List the customers who have reserved all the foreign films.
- Rephrased: List the customers for which, for all films, if that film was a foreign film, then that film must have been reserved by that customer.

 $\{c.Name \mid Customer(c) \\ \land \forall f (Film(f) \land f.Kind = "F" \Rightarrow \\ (\exists r (Reserved(r) \land \\ r.FilmID = f.FilmID \land \\ r.CustomerID = c.CustomerID))) \}$

 $a \Rightarrow b$ is shorthand for " $\neg a \lor b$ " or "**if** *a* **then** *b*"

Equivalences• $\forall x (P(x)) = \neg \exists x (\neg P(x))$ • $\exists x (P(x)) = \neg \forall x (\neg P(x))$ • $\exists x (P(x)) = \neg \exists x (\neg P(x) \lor \neg Q(x))$ • $\forall x (P(x) \lor Q(x)) = \neg \exists x (\neg P(x) \land \neg Q(x))$ • $\exists x (P(x) \lor Q(x)) = \neg \forall x (\neg P(x) \land \neg Q(x))$ • $\exists x (P(x) \land Q(x)) = \neg \forall x (\neg P(x) \lor \neg Q(x))$ • $\exists x (P(x) \land Q(x)) = \neg \forall x (\neg P(x) \lor \neg Q(x))$ • $\forall x (P(x)) \Rightarrow \exists x (P(x))$ • $\neg \exists x (P(x)) \Rightarrow \neg \forall x (P(x))$

Movie Web Queries, cont.

- Find the film(s) with the highest rental price.
- Rephrased: Find the film(s) for which all films have a rental price that is not higher.

 $\{ f.FilmID \mid Film(f) \land \forall f_2 (Film(f_2) \Rightarrow \}$

 f_2 .RentalPrice $\leq f$.RentalPrice)

- Alternative formulation: Find the film(s) with a rental price for which no other rental price is higher.
- Rephrased: Find the film(s) for which there does not exist a film with a higher rental price.

 $\{ f.FilmID \mid Film(f) \land \neg \exists f_2 (Film(f_2)) \}$

 $\land f_2$. Rental Price > f. Rental Price)

Overview

- Relational Algebra (RA)
 - Operators Limitations
 - Extended Relational Algebra Operators
- Tuple Relational Calculus (TRC)
- Domain Relational Calculus (DRC)

Domain Relational Calculus

• Each query is an expression of the form

 $\{x_1, x_2, \dots, x_n \mid P(x_1, x_2, \dots, x_n)\}$

- *P* is a formula similar to the formula in the *predicate calculus*.
- List the information on the expensive films.

 $\{F, T, R, K | Film(F, T, R, K) \land R > 4\}$

• List the titles of the expensive films.

 $\{ T \mid \exists F \exists R \exists K (Film(F, T, R, K) \land R > 4) \}$

Relational Query Language

Movie Web Queries

al Query La

• List the outrageously priced films (over \$4 or under \$1).

 $\{T \mid \exists F \exists R \exists K (Film(F, T, R, K) \land (R > 4 \lor R < 1))\}$

• List the ID numbers of the films that are expensive and have been reserved.

 $\{F \mid \exists T \exists R \exists K (Film(F, T, R, K) \land R > 4 \\ \land \exists I \exists D (Reserved(I, F, D))) \}$

Movie Web Queries, cont.

• List the IDs of the expensive films that have not been reserved.

 $\{F \mid \exists T, R, K \ (Film(F, T, R, K) \land R > 4 \\ \land \neg(\exists I \exists D \ (Reserved(I, F, D)))) \}$

Movie Web Queries, cont.

• List the titles of all reserved films.

 $\{T \mid \exists F, R, K \ (Film(F, T, R, K) \\ \land \exists I, D \ (Reserved(I, F, D))) \}$

• List the customers who have reserved a film.

 $\{N \mid \exists I, S, C, ST, Z \ (Customer(I, N, S, C, ST, Z) \\ \land \exists F, D \ (Reserved(I, F, D))) \}$

Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Movie Web Queries, cont.

• List the customers who have reserved expensive films.

 $\begin{aligned} & \{N \mid \exists I, S, C, ST, Z \ (Customer(I, N, S, C, ST, Z) \\ & \land \exists F, D \ (Reserved(I, F, D) \\ & \land (\exists T, R, K \ (Film(F, T, R, K) \land R > 4)) \end{aligned}$

• List the streets of customers who have reserved foreign films.

 $\begin{aligned} \{S \mid \exists I, N, C, ST, Z \ (Customer(I, N, S, C, ST, Z) \\ & \land \exists F, D \ (Reserved(I, F, D) \\ & \land \exists T, R \ (Film(F, T, R, "F")))) \end{aligned}$

Movie Web Queries, cont.

al Ouery La

• List the customers who have reserved all the foreign films.

 $\begin{array}{l} \{N \mid \exists I, S, C, ST, Z \ (Customer(I, N, S, C, ST, Z) \\ \land \forall F \ (\exists T, R \ (Film(F, T, R, ``F'') \Rightarrow \\ \exists D \ (Reserved(I, F, D)))) \} \end{array}$

Movie Web Queries, cont.

• Find the film(s) with the highest rental price.

• Alternative formulation: Find the film(s) with a rental price for which no other rental price is higher.

 $\{T \mid \exists F, R, K (Film(F, T, R, K))\}$ $\land \neg(\exists F_2, T_2, R_2, K_2(Film(F_2, T_2, R_2, K_2) \land R_2 > R)))\}$

Relational Query Languages

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Assertions

- Assertions are more general integrity constraints, expressed directly as predicates which must always be satisfied.
- These can be expressed in the algebra or calculi, of the form: "There does not exist an offending tuple."

No video rents for less than \$0.

• Algebra

- $\sigma_{RentalPrice < 0}(Film) = \{\}$
- Tuple Relational Calculus
 - $\forall f(Film(f) \Longrightarrow f.RentalPrice \ge 0)$
 - $\neg \exists f(Film(f) \land f.RentalPrice < 0)$
- Domain Relational Calculus
 - $\forall F (\exists T, \exists R \exists K(Film(F,T,R,K)) \Longrightarrow R \ge 0)$
 - $\neg \exists F, T, R, K(Film(F, T, R, K) \land R < 0)$

Another Assertion Example

- The spoken and subtitled languages for a film must be different.
- Algebra $\sigma_{spokenLanguage} = subtitleLanguage}(Film) = \{ \}$
- Tuple Relational Calculus
 - $\neg \exists f(Film(f) \land f.spokenLanguage = f.subtitleLanguage)$
- Domain Relational Calculus

```
\neg (\exists F, T, R, K, SP, SL (Film(F, T, R, K, SP, SL) \land SP = SL))
```

Relational Query Language

Summary

- · Relational algebra
 - Objects are relations (sets of n-tuples).
- Tuple relational calculus
 - Variables range over relations (sets of tuples).
 - Each variable is associated with an individual tuple.
- · Domain relational calculus
 - Variables range over domains (sets of values).
 - Each variable is associated with an individual value.
- Misnamed
 - Tuple relational calculus and Value relational calculus
 - Relational relational calculus and Domain relational calculus

Relational Query Languages

Expressive Power

- Theorem: The following three languages define *exactly* the same class of functions.
 - Relational algebra expressions
 - Safe relational tuple calculus formulas
 - Safe relational domain calculus formulas

• Proof:

- relational algebra \leq (safe) tuple relational calculus
 - Show via induction that every relational algebra expression has a counterpart
- in the tuple relational calculus. • tuple relational calculus \leq domain relational calculus
- domain relational calculus \leq relational algebra
- Corollary: All three languages are relationally complete.

elational Query Language

Completeness

Theorem: The tuple relational calculus is complete.

Proof: First convert algebraic expression to use just the five basic operators. Then, if algebraic expression is just one operator, substitute the calculus form.

- Selection $\sigma_P(r) \Rightarrow \{ t / r(t) \land P(t) \}$
- Projection
- Union

 $\pi_{X}(r) \Longrightarrow \{ t[X] / r(t) \}$ $r \cup s \Longrightarrow \{ t / r(t) \lor s(t) \}$

Difference

 $r-s \Longrightarrow \{ t / r(t) \land \neg s(t) \}$

- Cartesian product $r \times s \Rightarrow \{ t \circ q \mid r(t) \land s(q) \}$
- For algebraic expressions with more than one operator, induct on length of the algebraic expression.

Relational Query Language

Copyright © 2014 R. T. Snodgrass, C. S. Jensen, K. Torp, and C. E. Dyreson.

Completeness Example, Cont.

• Replace selection:

Replace the Cartesian product:

 $\{t.Title \mid \exists f (Film(f) \land \exists r (Reserved(r) \land f.FilmID = r.FilmID \land t = f \circ r)) \}$

• Simplify {*f.Title* | *Film(f)* ∧ ∃*r* (*Reserved(r*)

 $\land f.FilmID = r.FilmID$) }

Safety

- It is possible to write tuple calculus expressions that generate infinite relations.
- { $t/ \neg r(t)$ } results in an infinite relation if the domain of any attribute of relation *r* is infinite.
- We wish to ensure that an expression in relational calculus yields only a *finite* number of tuples.
- The *domain of a tuple relational calculus expression* is the set of all values that either appear as constant values in the expression or that exist in any tuple of the relations referenced in the expression.
- An expression is *safe* if all values in its result are from the domain of the expression.

Relational Query Languages

Equivalence of Expressive Power

- Theorem: The relational algebra is as expressive as the (safe) tuple relational calculus.
- (Informal) Proof: by induction on the number of operators in the calculus predicate
 - $\exists X(...) \Rightarrow \pi(...)$
 - $X \lor Y \implies \pi(X \times Y)$
 - $P(r) \implies \sigma_P(...)$
 - $\neg P(r) \implies U-r$
 - $X \wedge Y \implies \neg(\neg X \vee \neg Y)$
 - $\forall X(P(r)) \Rightarrow \neg \exists X(\neg P(r))$

telational Query Languages