



Database Design

CMSC424 Lecture Notes

Nick Roussopoulos

Department of Computer Science
University of Maryland
Spring 2017

1

© Nick Roussopoulos



Database Management System (DBMS)

- **DBMS contains information about a particular enterprise**
 - Collection of interrelated data
 - Set of programs to access and manipulate the data
 - An environment that is both *convenient* and *efficient* to use
- **Database Applications:**
 - Banking: transactions
 - Airlines: reservations, schedules
 - Universities: registration, grades
 - Sales: customers, products, purchases
 - Online retailers: order tracking, customized recommendations
 - Manufacturing: production, inventory, orders, supply chain
 - Human resources: employee records, salaries, tax deductions
- **Databases can be very large or huge (exabytes = 10^{18} Bytes)**
- **Databases touch all aspects of our lives**



© Nick Roussopoulos

2



University Database Example

- **Application program example**
 - Add new students, instructors, and courses
 - Register students for courses, and generate class rosters
 - Assign grades to students, compute grade point averages (GPA) and generate transcripts



© Nick Roussopoulos

3



Early days: Used files to store databases

- Prior to 70's, database applications were built directly on top of file systems
- Data redundancy and inconsistency
 - Multiple file copies in a variety of data structures
 - duplication of information in different files
- Difficulty in accessing data
 - Need to write a new program to carry out a new task
 - Programmers kept their own files (hard for others to access)
- Integrity problems
 - Integrity constraints (e.g., account balance > 0) become "buried" in program code rather than being stated explicitly
 - Hard to add new constraints or change existing ones



© Nick Roussopoulos

4



Drawbacks of using files to store data (cont.)

- **Atomicity of updates**
 - Failures may leave database in an inconsistent state with partial updates carried out
 - Example: Transfer of funds from one account to another should either complete or not happen at all
- **Concurrent access by multiple users**
 - Concurrent access needed for performance
 - Uncontrolled concurrent accesses can lead to inconsistencies
 - Example: Two people reading a balance (say 100) and updating it by withdrawing money (say 50 each) at the same time
- **Security problems**
 - Hard to provide user access to some, but not all, data

Database systems offer solutions to all the above problems



© Nick Roussopoulos

5



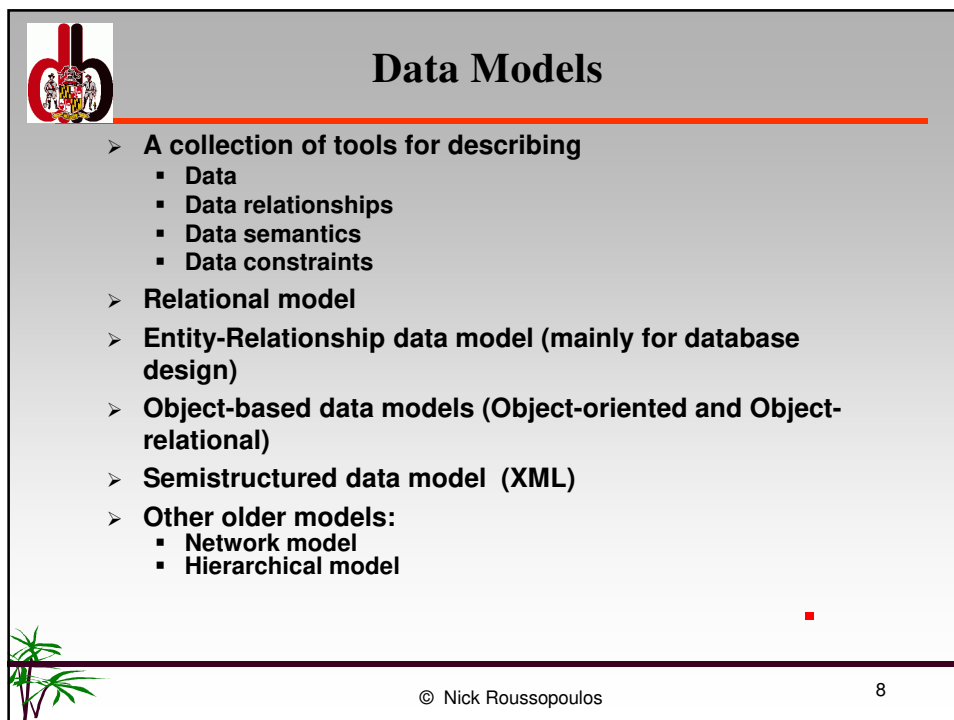
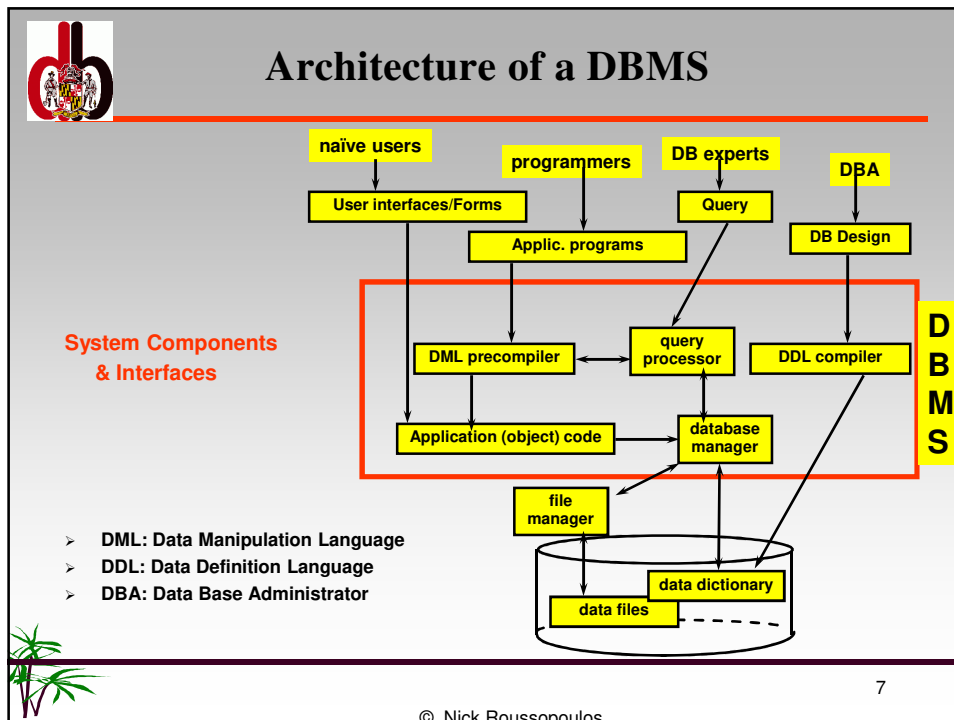
Reasons for Using a DBMS

- **provides data independence**
 - avoids hard wiring
- **reduces data redundancy and inconsistency**
- **allows concurrent multi-user access**
- **provides crash recovery**
- **provides security of data**
- **protects the data from being corrupted**
 - maintains integrity



© Nick Roussopoulos

6





What are models for?

- A model represents a perception of reality



- The modeling process is to fix a perception of reality
- During the modeling process we
 - select aspects
 - abstract to form a simple comprehension of the phenomena in the real world
- A Database is a Model of Reality
- A Database System is software to support the definition & use of the database



9

© Nick Roussopoulos



What are models for?

- useful when we want to
 - Examine
 - Evaluate
 - Manage a subset of the real world
- the cost of using the models is considerably lower than the cost of doing experiments in the real world

Examples:

- airplane simulator
- nuclear plant simulator
- economic model
- data model
- a map



10

© Nick Roussopoulos



The Relational Model

wrote	author	ISBN	publish-city	year
	Brown	56789	New York	1988
	Jeff	Tim	Washington	2001
	Jim	John	Clarcksville	2012

- **Data Entities (objects) are represented by relations (alias tables)**
 - each data entity is a single tuple (alias row) in a relation
 - each property of the entity is an attribute (alias column) in the tuple
- ***Relationships among data entities are also represented by relations***
 - each relationship is a tuple in a relation
 - each property of the relationship is an attribute (column)
- **Attributes take values from a specific set called the domain of the column**
 - the domain defines the type
- ➔ **relation schema**: the name of a relation and the names of its attributes (along with their underlying domains)
- ➔ **database schema**: the set of all relation schemas

δ



11

© Nick Roussopoulos



The Relational Model (cont)

- **intension and extension of a relation**

- intension =====>
- extension =====>

author	ssn	name	city	phone
	123	Sue	New York	301-444-4444
	124	Tim	Washington	202-666-7666
	125	John	Clarcksville	301-888-9999

attributes
(or columns)

tuples
(or rows)

- **generalities**

- a relation is a set of tuples
 - order is not important and not guaranteed
 - tuples are unique (set)- however, in real life, RDBMSs allow pseudo sets
- order of attributes is no longer important as we use the attribute names
 - author.name specifies the 2nd column
 - if t in author, then t.name specifies the value of the second column
- ➔ the number of attributes in a relation determines the arity or degree of it
 - author 4th degree, book 3rd degree (ternary), library 2nd degree (binary)
- ➔ the number of tuples in a relation is called the cardinality of it.
- ➔ cardinality is a property of a relation- NOT the relation scheme

δ



12

© Nick Roussopoulos



Keys of Relations

- **Keys are defined to distinguish among rows of the table**
 - note that keys are properties of the relation scheme- NOT the relation. By looking at the relation we cannot tell what is a key but may be able to tell what isn't.
- 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 $r(R)$
 - Example: $\{ID\}$ and $\{ID, name\}$ are both superkeys of *instructor*.
- Superkey K is a **candidate key** if K is minimal
Example: $\{ID\}$ is a candidate key for *Instructor*.
- One of the candidate keys is selected to be the **primary key**.
 - which one?
- **Foreign key** constraint: Values in one relation must appear in the primary key values of another relation
 - **Referencing** relation
 - **Referenced** relation

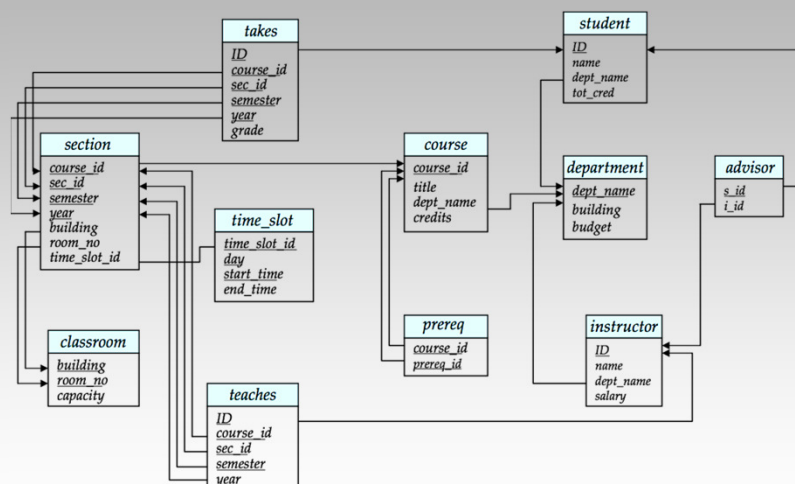


13

© Nick Roussopoulos



Schema Diagram for University Database



14

© Nick Roussopoulos



Relational Query Languages

- Procedural vs. non-procedural, or declarative
- “Pure” languages:
 - Relational algebra
 - Tuple relational calculus
 - Domain relational calculus
- Embedded in host languages or environments

δ



15

© Nick Roussopoulos



Relational Algebra

- set of operators that map one or more relations into another relation
- closed algebraic system
 - the best feature- operations on operations
 - form relational algebraic expressions
- two types of operation: database specific and set theoretic operators
 - database specific
 - [horizontal] selection (σ) (or restriction)
 - projection (π)
 - join
 - outer join
 - semijoin
 - division
 - set operators
 - union
 - difference
 - intersection
 - cartesian (cross) product

δ



16

© Nick Roussopoulos



Employee-Department Database

Emp	ename	sal	dept
	Gary	30K	toy
	Shirley	35K	candy
	Christos	37K	shoe
	Robin	22K	toy
	Uma	30K	shoe
	Tim	12K	

Dept	dept	floor	mgr
	candy	1	Irene
	toy	2	Jim
	men	2	John
	shoe	1	George



17

© Nick Roussopoulos



Relational Operators

- [horizontal] selection (σ): picks a subset of the rows

σ (Emp)
dept = 'toy'

Result	ename	sal	dept
	Gary	30K	toy
	Robin	22K	toy

- [vertical] projection (π): picks a subset of columns

π (Dept)
floor

Result	floor
	1
	2



18

© Nick Roussopoulos



Relational Algebra (more operators)

- Cartesian product: $r \times s$

A	B
α	1
β	2

r

C	D	E
α	10	a
β	10	a
β	20	b
γ	10	b

s

A	B	C	D	E
α	1	α	10	a
α	1	β	10	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



© Nick Roussopoulos

19



Joins

Emp	ename	sal	dept
Gary	30K	toy	
Shirley	35K	candy	
Christos	37K	shoe	
Robin	22K	toy	
Uma	30K	shoe	
Tim	12K		

Dept	dept	floor	mgr
	candy	1	Irene
	toy	2	Jim
	men	2	John
	shoe	1	George

- Equijoin (based on equality)
 $\text{Emp} \bowtie \text{Dept}$
 $\text{Emp.dept} = \text{Dept.dept}$

Result	ename	sal	dept	dept	floor	mgr
	Gary	30K	toy	toy	2	Jim
	Shirley	35K	candy	candy	1	Irene
	Christos	37K	shoe	shoe	1	George
	Robin	22K	toy	toy	2	Jim
	Uma	30K	shoe	shoe	1	George

- Natural join (also based on equality)
 $\text{Emp} \ltimes \text{Dept}$

Result	ename	sal	dept	floor	mgr
	Gary	30K	toy	2	Jim
	Shirley	35K	candy	1	Irene
	Christos	37K	shoe	1	George
	Robin	22K	toy	2	Jim
	Uma	30K	shoe	1	George

- θ join (θ is an arbitrary boolean expression)
 $\text{Emp} \ltimes_{\theta} \text{Dept}$

When θ is true $r \ltimes_{\text{true}} s = r \times s$



© Nick Roussopoulos

20



Joins (cont)

Emp	ename	sal	dept
Gary	30K	toy	
Shirley	35K	candy	
Christos	37K	shoe	
Robin	22K	toy	
Uma	30K	shoe	
Tim	12K		

Dept	dept	floor	mgr
	candy	1	Irene
	toy	2	Jim
	men	2	John
	shoe	1	George

- Semijoin $Emp \bowtie Dept = \Pi_{Emp.*} (Emp \bowtie Dept)$

Result	ename	sal	dept
	Gary	30K	toy
	Shirley	35K	candy
	Christos	37K	shoe
	Robin	22K	toy
	Uma	30K	shoe

- Antijoin $Emp \not\supset Dept = \Pi_{Emp.*} Emp - (Emp \bowtie Dept)$

Emp	ename	sal	dept
	Tim	12K	



© Nick Roussopoulos

21



Outerjoins

Emp	ename	sal	dept
Gary	30K	toy	
Shirley	35K	candy	
Christos	37K	shoe	
Robin	22K	toy	
Uma	30K	shoe	
Tim	12K		

Dept	dept	floor	mgr
	candy	1	Irene
	toy	2	Jim
	men	2	John
	shoe	1	George

- Left outerjoin $Emp \ltimes Dept$

Result	ename	sal	dept	dept	floor	mgr
	Gary	30K	toy	toy	2	Jim
	Shirley	35K	candy	candy	1	Irene
	Christos	37K	shoe	shoe	1	George
	Robin	22K	toy	toy	2	Jim
	Uma	30K	shoe	shoe	1	George
	Tim	12K		NULL	NULL	NULL

- Right outerjoin $Emp \rtimes Dept$

Result	ename	sal	dept	dept	floor	mgr
	Gary	30K	toy	toy	2	Jim
	Shirley	35K	candy	candy	1	Irene
	Christos	37K	shoe	shoe	1	George
	Robin	22K	toy	toy	2	Jim
	Uma	30K	shoe	shoe	1	George
	NULL	NULL	NULL	men	2	John

- Full outerjoin $Emp \ltimes\rtimes Dept$

Result	ename	sal	dept	dept	floor	mgr
	Gary	30K	toy	toy	2	Jim
	Shirley	35K	candy	candy	1	Irene
	Christos	37K	shoe	shoe	1	George
	Robin	22K	toy	toy	2	Jim
	Uma	30K	shoe	shoe	1	George
	Tim	12K		NULL	NULL	NULL
	NULL	NULL	NULL	men	2	John



© Nick Roussopoulos

22



Do we need all these joins?

- More examples on Relational Algebra can be found in the book and
- https://en.wikipedia.org/wiki/Relational_algebra



© Nick Roussopoulos

23



Two More Relational Operators

- Rename: $\rho_{\text{new_name}}$ (E) - where E is a relational expression
- Generalized projection $\pi_{F1, F2, \dots}$ (E)
- where F1, F2, ... Are functions on columns



© Nick Roussopoulos

24



Comparing Tuples in the same relation

➤ Example: Find the smallest salary employee

- we need to compare every tuple with every other tuple in Emp
- bring them together by taking the Cartesian product Emp X Emp to compare the two sal columns
- how can we distinguish the two sal columns?

Rename one Emp with $\rho_{Emp2}(Emp)$

- Join the two $Emp \bowtie Emp2$
- and compare $Emp.sal > Emp2.sal$
- this will give all the employees that have at least another employee with smaller salary
- subtract from Emp the above to find Tim

Emp	ename	sal	dept
	Gary	30K	toy
	Shirley	35K	candy
	Christos	37K	shoe
	Robin	22K	toy
	Uma	30K	shoe
	Tim	12K	

Emp2	ename	sal	dept
	Gary	30K	toy
	Shirley	35K	candy
	Christos	37K	shoe
	Robin	22K	toy
	Uma	30K	shoe
	Tim	12K	

$\pi_{Emp.ename}(\rho_{Emp2}(Emp)) - \pi_{Emp.ename, Emp2.ename, Emp2.sal > Emp.sal}(Emp \bowtie \rho_{Emp2}(Emp))$



25

© Nick Roussopoulos



Relational Algebra - Set Operators

- union: $r \cup s$ - both relations must be union-compatible
i.e. same degree and the same domains

A	B
α	1
α	2
β	1

r

A	B
α	2
β	3

s

A	B
α	1
α	2
β	1
β	3

- difference: $r - s$ - both relations must be union compatible

A	B
α	1
α	2
β	1

r

A	B
α	2
β	3

s

A	B
α	1
β	1

- intersection: $r \cap s$ - both relations must be union compatible

A	B
α	1
α	2
β	1

r

A	B
α	2
β	3

s

A	B
α	2



26

© Nick Roussopoulos



Relational Calculi

- based on predicate calculus
- non-procedural (descriptive) as opposed to procedural (prescriptive)

Tuple Calculus

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

P: is a predicate associated with some relation R

t: tuple variable ranging over the relation R

t[A]: value of attribute A in t

- students in CMSC 424
 - $\{t \mid t \in \text{enroll} \wedge (t[\text{course\#}] = \text{CMSC424})\}$
- students in CMSC-424 conforming with the CMSC-351 prerequisite
 - $\{t \mid t \in \text{enroll} \wedge s \in \text{enroll} \wedge (t[\text{course\#}] = \text{CMSC424}) \wedge (s[\text{course\#}] = \text{CMSC351}) \wedge (t[\text{ss\#}] = s[\text{ss\#}])\}$

- **Quantifiers, Bound and Free Variables**

- $\exists s$ (or $\forall s$) quantify the variable s following them-
 s is bound by the quantifiers
- queries have free variables that take independent values: t above is free



27

© Nick Roussopoulos



Atoms & Formulas

- **Atoms**
 - $R(t)$ where t is a tuple variable
 - $t[i] \theta s[j]$ where t, s are tuple variables and $\theta \in \{<, >, =, \neq, \leq, \geq\}$
- **Formulas are:**
 - an atom is a formula
 - if P and Q are formulas, then so are (P) , $\neg P$, $P \wedge Q$, $P \vee Q$
 - if $P(t)$ is a formula and t is free, then $\exists t P(t)$ and $\forall t P(t)$ are also formulas
- **Equivalences**
 - $\neg(P \wedge Q) \equiv \neg P \vee \neg Q$
 - $\neg(P \vee Q) \equiv \neg P \wedge \neg Q$
 - $\forall t \in r(P(t)) \equiv \neg \exists t \in r(\neg P(t))$
 - $P \Rightarrow Q \equiv \neg(P) \vee Q$



28

© Nick Roussopoulos



Tuple Calculus (cont)

- Problems with negation
 - $\{t \mid t \notin \text{Enroll}\}$ - Unsafe

Restriction: put negation as a second operand of a conjunction. The first operand of the conjunction defines the scope and the second filters this scope.

$\{t \mid t \in \text{Enroll} \wedge t[\text{course\#}] \neq \text{CMSC-420}\}$ - Safe

$\xleftrightarrow{\text{positive scope}} \quad \xleftrightarrow{\text{negative clause}}$



29

© Nick Roussopoulos



Domain Calculus

- $\text{dom}(X)$ - set of values of a domain; more than one attribute may receive values from $\text{dom}(X)$
- Atoms
 - $\langle x_1, x_2, \dots, x_N \rangle \in R$ - set of domain variables receiving values from the attributes of R
 - $x \theta y$ where x, y are domain variables and $\theta \in \{<, >, =, \neq, \leq, \geq\}$
 - $x \theta c$ where c is a constant
- Formulas are:
 - an atom is a formula
 - if P and Q are formulas, then so are (P) , $\neg P$, $P \wedge Q$, $P \vee Q$
 - if $P(x)$ is a formula and x is a domain variable, then $\exists x P(x)$ and $\forall x P(x)$ are also formulas
- Queries:
 - $\{\langle x_1, x_2, \dots, x_N \rangle \mid \Psi(x_1, x_2, \dots, x_N)\}$
 - Examples:
 - $\{\langle \text{ss\#}, \text{course\#}, \text{semester} \rangle \mid \text{Enroll}(\text{ss\#}, \text{course\#}, \text{semester})\}$
 - $\{\langle x, y, z \rangle \mid \text{Enroll}(x, y, z) \wedge y = \text{CMSC-424}\}$



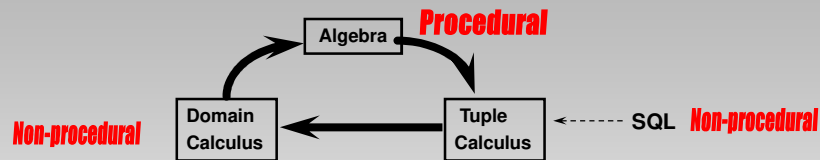
- Domain calculus expressions are safe

30

© Nick Roussopoulos



Reductions of Relational Algebra & Calculi



- the reduction proves equivalence in expressing power
- these languages are used as benchmarks for other languages
- a relational language L is relationally complete iff statements in one of the above can be mapped to statements in L and vice versa



31

© Nick Roussopoulos



Relational Algebra (one last operator)

- division R / S : given $R(A,B)$ and $S(B)$ t is in R/S if for all s in S there exists an r in R such that $r.B=s.B$ and $t.A=r.A$

Construction:

$$R/S = \pi_A(R) - \pi_A((\pi_A(R) \times S) - R)$$



<== all possible pairs including the desired



<== not desired tuples



32

© Nick Roussopoulos