Carnegie Mellon University Database Systems Relational Model & Algebra

15-445/645 SPRING 2025 **X** PROF. JIGNESH PATEL





CMU·DB 15-445/645 (Spring 2025)

WaitList

- We do **<u>not</u>** control the waitlist.
- Admins will move students off the waitlist as spots become available.
- To improve your chances of getting into the class (though not a guarantee), stay in the class and complete P0.
- This class will be offered in Fall'25 too!

COURSE OVERVIEW

This course is about the design/implementation of database management systems (DBMSs).

This is <u>**not**</u> a course about how to use a DBMS to build applications or how to administer a DBMS. \rightarrow See <u>CMU 95-703</u> (Heinz College)

COURSE LOGISTICS

Course Policies + Schedule: <u>Course Web Page</u> Discussion + Announcements: <u>Piazza</u> Homeworks + Projects: <u>Gradescope</u> Final Grades: <u>Canvas</u>

Notes: \rightarrow Do <u>not</u> post your solutions on Github. \rightarrow Do <u>not</u> email instructors / TAs for help.



TEXTBOOK

Database System Concepts 7th Edition Silberschatz, Korth, & Sudarshan

We also provide lecture notes that cover topics not found in textbook.



GRADING RUBRIC

Homeworks (15%) Projects (45%) Midterm Exam (20%) Final Exam (20%)

HOMEWORKS

Six homework assignments that cover lecture and reading material.

- \rightarrow First homework is a SQL assignment.
- \rightarrow The rest will be pencil-and-paper assignments.

Submit all assignments via Gradescope.

All homework should be done individually.

PROJECTS

All projects will use the CMU DB Group <u>BusTub</u> academic DBMS.

- \rightarrow Each project builds on the previous one.
- → We will <u>**not**</u> teach you how to write/debug C++17.
- \rightarrow See the <u>15-445/645 Bootcamp</u>.

15-445/645 (Spring 2025)

Total of <u>four</u> late days the entire semester for projects only.

We will hold an online recitation for each project after it is released.



C++ Requirement

thaisent to 6. All the projects are in C++ If you are new to C++, you must pick it up quickly... If you can take and get all the questions on the following guizzes right, you are all set: Scoping: https://www.learncpp.com/cpp-tutorial/chapter-7-summary-and-quiz/ Type Conversion: https://www.learncpp.com/cpp-tutorial/chapter-10-summary-and-quiz/ Ivalues/rvalues: https://www.learncpp.com/cpp-tutorial/chapter-12-summary-and-quiz/ Stack and heap: https://www.learncpp.com/cpp-tutorial/chapter-20-summary-and-quiz/ Move Semantics: https://www.learncpp.com/cpp-tutorial/chapter-22-summary-and-quiz/ Templates: https://www.learncpp.com/cpp-tutorial/chapter-26-summary-and-quiz/

... take it upon yourself to catch up ...

... also https://db.in.tum.de/teaching/ss23/c++praktikum/slides/lecture-10.2.pdf?lang=en

-445/645 (Spring 2025)

C++ Bootcamp: This Friday Jan. 17th from 3pm-4pm in GHC 4303

Project 0 (P0): Goals

- \rightarrow Get you started on C++, so you are not surprised later.
- \rightarrow Get you thinking about algorithms and concurrency.
- \rightarrow P0 is about building a Skip List data structure.
- \rightarrow We will discuss Skip List in more detail later in the class.
- \rightarrow <u>P0 is published</u>; due on Jan 26 @ 11:59pm.
- \rightarrow No late days allowed for PO.

If you can't score 100% on P0, you can't stay in this class, even if you are currently enrolled.

CMU-DB 15-445/645 (Spring 2025)

OFFICE HOURS

Instructors and TAs will hold office hours on weekdays (Mon-Fri) at different times.

We will also hold a TA power session on the Saturday before each project is due.

There will **<u>not</u>** be any office hours on Sundays.

PROJECT LATE POLICY

You will lose 10% of the points for a project or homework for every 24 hours it is late.

You have a total of **four** late days to be used for **projects only**.

We will grant no-penalty extensions due to extreme circumstances (e.g., medical emergencies). \rightarrow If something comes up, please contact the instructors as

soon as possible.



PLAGIARISM WARNING



The homework and projects must be your own original work. They are <u>**not**</u> group assignments. You may <u>**not**</u> copy source code from other people or the web.

Plagiarism is <u>**not**</u> tolerated. You will get lit up. \rightarrow Please ask instructors (not TAs!) if you are unsure.

See <u>CMU's Policy on Academic Integrity</u> for additional information.



TODAY'S AGENDA

Database Systems Background Relational Model Relational Algebra Alternative Data Models Q&A Session

DATABASE

Organized collection of inter-related data that models some aspect of the real-world.

Databases are the core component of most computer applications.

DATABASE EXAMPLE

Create a database that models a digital music store to keep track of artists and albums.

Information we need to keep track of in our store:

- \rightarrow Information about <u>Artists</u>
- \rightarrow The <u>Albums</u> those Artists released

FLAT FILE STRAWMAN

Store our database as comma-separated value (CSV) files that we manage ourselves in application code.

- \rightarrow Use a separate file per entity.
- \rightarrow The application must parse the files each time they want to read/update records.

Artist(name, year, country)

"Wu-Tang Clan",1992,"USA"

```
"Notorious BIG",1992,"USA"
```

```
"GZA",1990,"USA"
```

15-445/645 (Spring 2025)

Album(name, artist, year)

"Enter the Wu-Tang", "Wu-Tang Clan", 1993

```
"<u>St.Ides Mix Tape</u>", "Wu-Tang Clan", 1994
```

"Liquid Swords", "GZA", 1990

FLAT FILE STRAWMAN

Example: Get the year that GZA went solo.

Artist(name, year, country)

"Wu-Tang Clan",1992,"USA"

```
"Notorious BIG",1992,"USA"
```

"GZA",1990,"USA"



for line in file.readlines():
record = parse(line)
if record[0] == "GZA":
 print(int(record[1]))



FLAT FILES: DATA INTEGRITY

How do we ensure that the artist is the same for each album entry?

What if somebody overwrites the album year with an invalid string?

What if there are multiple artists on an album?

What happens if we delete an artist that has albums?



FLAT FILES: IMPLEMENTATION

How do you find a particular record?

What if we now want to create a new application that uses the same database? What if that application is running on a different machine?

What if two threads try to write to the same file at the same time?



FLAT FILES: DURABILITY

What if the machine crashes while our program is updating a record?

What if we want to replicate the database on multiple machines for high availability?

DATABASE MANAGEMENT SYSTEM

A <u>database management system</u> (DBMS) is software that allows applications to store and analyze information in a database.

A general-purpose DBMS supports the definition, creation, querying, update, and administration of databases in accordance with some <u>data model</u>.

DATA MODELS

A <u>data model</u> is a collection of concepts for describing the data in a database.

A <u>schema</u> is a description of a particular collection of data, using a given data model.

- \rightarrow This defines the structure of data for a data model.
- \rightarrow Otherwise, you have random bits with no meaning.

DATA MODELS

← Most DBMSs Relational Key/Value Graph Document / JSON / XML / Object Wide-Column / Column-family Array (Vector, Matrix, Tensor) Hierarchical Network Semantic **Entity-Relationship**

ECMU-DB 15-445/645 (Spring 2025)

DATA MODELS



ECMU-DB 15-445/645 (Spring 2025)

EARLY DBMSs

Early database applications were difficult to build and maintain on available DBMSs in the 1960s. \rightarrow Examples: IDS, IMS, CODASYL

 \rightarrow Computers were expensive, humans were cheap.

Tight coupling between logical and physical layers.

Programmers had to (roughly) know what queries the application would execute before they could deploy the database.

EARLY DBMSs

Ted Codd was a mathematician at IBM Research in the late 1960s.

Codd saw IBM's developers rewriting database programs every time the database's schema or layout changed.

Devised the relational model in 1969.



Edgar F. Codd



DERIVABILITY, REDUNDANCY AND CONSISTENCY OF RELATIONS STORED IN LARGE DATA BANKS

E. F. Codd Research Division San Jose, California

ABSTRACT: The large, integrated data banks of the future will contain many relations of various degrees in stored form. It will not be unusual for this set of stored relations to be redundant. Two types of redundancy are defined and discussed. One type may be employed to improve accessibility of certain kinds of information which happen to be in great demand. When either type of redundantd about it and have some means of detecting any "logical" inconsistencies in the total set of stored relations. Consistency checking might be helpful in tracking down unauthorized (and possibly fraudulent) changes in the data bank should know

RJ 599(# 12343) August 19, 1969

da

da

SECMU-DB

15-445/645 (Spring 2025

LIMITED DISTRIBUTION NOTICE - This report has been submitted for publication elsewhere and has been issued as a Research Report for early dissemination of its contents. As a courtesy to the intended publisher, it should not be widely distributed until after the date of outside publication.

Copies may be requested from IBM Thomas J. Watson Research Center, Post Office Box 218, Vorktoan Heights, New York 10598 Information Retrieval

LY DBN

at

vriting

nged.

1969.

he

A Relational Model of Data for Large Shared Data Banks

E. F. CODD IBM Research Laboratory, San Jose, California

Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A prompting service which supplies such information is not a satisfactory solution. Activities of users at terminals and most application programs should remain unaffected when the internal representation of data is changed and even when some aspects of the external representation are changed. Changes in data representation will often be needed as a result of changes in query, update, and report traffic and natural growth in the types of stored information. Existing noninferential, formatted data systems provide users with tree-structured files or slightly more general network models of the data. In Section 1, inadequacies of these models are discussed. A model based on n-ary relations, a normal form for data base relations, and the concept of a universal data sublanguage are introduced. In Section 2, certain operations on relations (other than logical inference) are discussed and applied to the problems of redundancy and consistency in the user's model.

KEY WORDS AND PHRASES: data bank, data bank, data bank, data organizanian, hisrarchies of data, networks of data, relationa, derivability, relandancy, consultancy, composition, join, retrieval languaga, predicate calcula, security, data integrity CR CATECORES: 3/20, 3/23, 3/25, 4/20, 4/22, 4/29

1. Relational Model and Normal Form

1.1. INTRODUCTION

This paper is concerned with the application of elementary relation theory to systems which provide shared access to large banks of formatted data. Except for a paper by Childs [1], the principal application of relations to data systems has been to deductive question -answering systems. Levein and Maron [2] provide numerous references to work in this area.

In contrast, the problems treated here are those of data independence—the independence of application programs and terminal activities from growth in data types and changes in data representation—and extain kinds of data inconsistency which are expected to become troublesome even in nondeductive systems.

Volume 13 / Number 6 / June, 1970

P. BAXENDALE, Editor

The relational view (or model) of data described in Section 1 appears to be superior in several respects to the graph or network model [3, 4] presently in vogue for noninferential systems. It provides a means of describing data with its natural structure only—that is, without superimposing any additional structure for machine representation purposes. Accordingly, it provides a basis for a high level data language which will yield maximal independence between programs on the one hand and machine representation and organization of data on the other.

A further advantage of the relational view is that it forms a sound basis for treating derivability, redundancy, and consistency of relations—these are discussed in Section 2. The network model, on the other hand, has spawned a number of confusions, not the least of which is mistaking the derivation of connections for the derivation of relations (see remarks in Section 2 on the "connection trap").

Finally, the relational view permits a clearer evaluation of the scope and logical limitations of present formatted data systems, and also the relative merits (from a logical standpoint) of competing representations of data within a single system. Examples of this clearer perspective are eited in various parts of this paper. Implementations of systems to support the relational model are not discussed.

1.2. DATA DEPENDENCIES IN PRESENT SYSTEMS

The provision of data description tables in recently developed information systems represents a major advance toward the goal of data independence [5, 6, 7]. Such tables facilitate changing certain characteristics of the data representation stored in a data bank. However, the variety of data representation characteristics which can be changed without logically impairing some application programs is still quite limited. Further, the model of data with which ueers interact is still cluttered with representational properties, particularly in regard to the representation of collections of data (as opposed to individual items). Three of the principal kinds of data dependencies which still need to be removed are: ordering dependence, indexing dependence, and access path dependence. In some systems these dependencies are not clearly separable from one another.

1.2.1. Ordering Dependence. Elements of data in a data bank may be stored in a variety of ways, some involving no concern for ordering, some permitting each element to participate in one ordering only, others permitting each element to participate in several orderings. Let us consider those existing systems which either require or permit data elements to be stored in at least one total ordering which is closely associated with the hardware-determinel ordering not grants might be stored in ascending order by part serial number. Such systems normally permit application programs to assume that the order of presentation of records from such a file is identical to (or is a subordering of) the

CODASYL

The Differences and Similarities Between the Data Base Set and Relational Views of Data.

→ <u>ACM SIGFIDET Workshop on Data</u> <u>Description, Access, and Control in Ann</u> <u>Arbor, Michigan, held 1–3 May 1974</u>



Codd











Stonebraker

COBOL/CODASYL camp:

- 1. The relational model is too mathematical. No mere mortal programmer will be able to understand your newfangled languages.
- 2. Even if you can get programmers to learn your new languages, you won't be able to build an efficient implementation of them.
- 3. On-line transaction processing applications want to do record-oriented operations.

Empty set

Relational camp:

- 1. Nothing as complicated as the DBTG proposal can possibly be the right way to do data management.
- 2. Any set-oriented query is too hard to program using the DBTG data manipulation language.
- 3. The CODASYL model has no formal underpinning with which to define the semantics of the complex operations in the model.

The record set, basic structure of navigational (e.g. CODASYL) databse model. A set consists of one parent record (also called "the owner"), and n child records (also called members records)

EXAMPLE DB 15-445/645 (Spring 2025) The <u>relational model</u> defines a database abstraction based on relations to avoid maintenance overhead.

Key tenets:

- \rightarrow Store database in simple data structures (relations).
- \rightarrow Physical storage left up to the DBMS implementation.
- → Access data through high-level language, DBMS figures out best execution strategy.

RELATIONAL MODEL

Structure: The definition of the database's relations and their contents independent of their physical representation.

Integrity: Ensure the database's contents satisfy constraints.

Manipulation: Programming interface for accessing and modifying a database's contents.



DATA INDEPENDENCE

Isolate the user/application from lowlevel data representation.

- → The user only worries about high-level application logic.
- → DBMS optimizes the layout according to operating environment, database contents, and workload.
- → DBMS can then re-optimize the database if/when these factors changes.



DATA INDEPENDENCE

Isolate the user/application from lowlevel data representation.

- → The user only worries about high-level application logic.
- → DBMS optimizes the layout according to operating environment, database contents, and workload.
- → DBMS can then re-optimize the database if/when these factors changes.



DATA INDEPENDENCE

Isolate the user/application from lowlevel data representation.

- → The user only worries about high-level application logic.
- → DBMS optimizes the layout according to operating environment, database contents, and workload.
- → DBMS can then re-optimize the database if/when these factors changes.



RELATIONAL MODEL

A <u>relation</u> is an unordered set that contain the relationship of attributes that represent entities.

A <u>tuple</u> is a set of attribute values (aka its <u>domain</u>) in the relation.

- \rightarrow Values are (normally) atomic/scalar.
- → The special value **NULL** is a member of every domain (if allowed).

Artist(name, year, country)

name	year	country
Wu-Tang Clan	1992	USA
Notorious BIG	1992	USA
GZA	1990	USA

= Table with *n* columns

n-ary Relation
RELATIONAL MODEL: PRIMARY KEYS

A relation's **primary key** uniquely identifies a single tuple.

Some DBMSs automatically create an internal primary key if a table does not define one.

DBMS can auto-generation unique
primary keys via an <u>identity column</u>:
 → IDENTITY (SQL Standard)
 → SEQUENCE (PostgreSQL / Oracle)
 → AUTO_INCREMENT (MySQL)

Artist(name, year, country)

name	year	country
Wu-Tang Clan	1992	USA
Notorious BIG	1992	USA
GZA	1990	USA

ECMU-DB 15-445/645 (Spring 2025)

RELATIONAL MODEL: PRIMARY KEYS

A relation's **primary key** uniquely identifies a single tuple.

Some DBMSs automatically create an internal primary key if a table does not define one.

DBMS can auto-generation unique
primary keys via an <u>identity column</u>:
 → IDENTITY (SQL Standard)
 → SEQUENCE (PostgreSQL / Oracle)
 → AUTO_INCREMENT (MySQL)

Artist(id, name, year, country) id country name year 101 Wu-Tang Clan 1992 USA 102 Notorious BIG 1992 USA 103 GZA 1990 IUSA

ECMU-DB 15-445/645 (Spring 2025)

A <u>foreign key</u> specifies that an attribute from one relation maps to a tuple in another relation.



Artist(id, name, year, country)

id	name	year	country
101	Wu-Tang Clan	1992	USA
102	Notorious BIG	1992	USA
103	GZA	1990	USA

Album(id, name, artists, year)

id	name	artists	year
11	<u>Enter the Wu-Tang</u>	101	1993
22	<u>St.Ides Mix Tape</u>	???	1994
33	Liquid Swords	103	1995

Artist(id, name, year, country)

id	name	year	country
101	Wu-Tang Clan	1992	USA
102	Notorious BIG	1992	USA
103	GZA	1990	USA

Album(id, name, artists, year)

id	name	artists	year
11	Enter the Wu-Tang	101	1993
22	<u>St.Ides Mix Tape</u>	???	1994
33	Liquid Swords	103	1995

ECMU·DB 15-445/645 (Spring 2025)

id	name	year	country
101	Wu-Tang Clan	1992	USA
102	Notorious BIG	1992	USA
103	GZA	1990	USA

Artist(id, name, year, country)

Album(id, name, year)

id	name	year
11	Enter the Wu-Tang	1993
22	<u>St.Ides Mix Tape</u>	1994
33	<u>Liquid Swords</u>	1995

ArtistAlbum(artist_id, album_id)

artist_id	album_id
101	11
101	22
103	22
102	22

RELATIONAL MODEL: CONSTRAINTS

User-defined conditions that must hold for <u>any</u> instance of the database.

- \rightarrow Can validate data within a single tuple or across entire relation(s).
- → DBMS prevents modifications that violate any constraint.

Unique key and referential (fkey) constraints are the most common.

SQL:92 supports global asserts but these are rarely used (too slow).

Artist(id, name, year, country)

id	name	year	country
101	Wu-Tang Clan	1992	USA
102	Notorious BIG	1992	USA
103	GZA	1990	USA

```
CREATE TABLE Artist (
  name VARCHAR NOT NULL,
  year INT,
  country CHAR(60),
  CHECK (year > 1900)
);
```

CREATE ASSERTION myAssert
CHECK (<SQL>);

DATA MANIPULATION LANGUAGES (DML)

The API that a DBMS exposes to applications to store and retrieve information from a database.

Procedural:

→ The query specifies the (high-level) strategy to find the desired result based on sets / bags.

Non-Procedural (Declarative):

 \rightarrow The query specifies only what data is wanted and not how to find it.

← Relational Calculus

Relational Algebra

RELATIONAL ALGEBRA

Fundamental operations to retrieve and manipulate tuples in a relation.

→ Based on set algebra (unordered lists with no duplicates).

Each operator takes one or more relations as its inputs and outputs a new relation.

→ We can "chain" operators together to create more complex operations.

15-445/645 (Spring 2025

σ Select

π Projection

- **U** Union
- ∩ Intersection
- Difference
- × Product

🛛 Join

45

- Choose a subset of the tuples from a relation that satisfies a selection predicate.
- → Predicate acts as a filter to retain only tuples that fulfill its qualifying requirement.
- → Can combine multiple predicates using conjunctions / disjunctions.

Syntax: $\sigma_{\text{predicate}}(R)$

R	(a_id	,b_id	•
	a_id	b_id	
	a1	101	
	a2	102	
	a2	103	
	a3	104	

- Choose a subset of the tuples from a relation that satisfies a selection predicate.
- → Predicate acts as a filter to retain only tuples that fulfill its qualifying requirement.
- → Can combine multiple predicates using conjunctions / disjunctions.

Syntax: opredicate(R)

R	(a_id	,b_id
	a_id	b_id
	a1	101
	a2	102
	a2	103
	a3	104





15-445/645 (Spring 2025)

- Choose a subset of the tuples from a relation that satisfies a selection predicate.
- → Predicate acts as a filter to retain only tuples that fulfill its qualifying requirement.
- → Can combine multiple predicates using conjunctions / disjunctions.

Syntax: $\sigma_{\text{predicate}}(R)$

15-445/645 (Spring 2025)

R	(a_id	,b_id	
	a_id	b_id	
	a1	101	
	a2	102	
	a2	103	
	a3	104	

ס _{a_i}	d='	a2'	(R))
a_i		b_	id	

a2 a2 102

103

a_id=	'a2'۸	b_	id>10	2(R	
	id	h	id		

/n

ā	a_id	b_id
ā	a2	103

- Choose a subset of the tuples from a relation that satisfies a selection predicate.
- → Predicate acts as a filter to retain only tuples that fulfill its qualifying requirement.
- → Can combine multiple predicates using conjunctions / disjunctions.

Syntax: $\sigma_{\text{predicate}}(R)$

-445/645 (Spring 2025

R	(a_id	,b_id)
	a_id	b_id	
	a1	101	
	a2	102	
	a2	103	
	a3	104	

$\sigma_{a_{id}=a_{2}}(\kappa)$				
a_id	b_id			
a2	102			
a2	103			

$\sigma_{a_{ic}}$	= '	a2'۸	b_	id>102	(R)
	a	id	h	id	

а	_id	b_id
а	2	103

SELECT	*	FROM	R			
WHERE	a_	id='a	a2'	AND	b_	_id>102;

- Choose a subset of the tuples from a relation that satisfies a selection predicate.
- → Predicate acts as a filter to retain only tuples that fulfill its qualifying requirement.
- → Can combine multiple predicates using conjunctions / disjunctions.

Syntax: $\sigma_{\text{predicate}}(R)$

R	(a_id	,b_id)
	a_id	b_id	
	a1	101	
	a2	102	
	a2	103	
	a3	104	

>

{a2'} (K)	$\sigma{a_{a_{a_{a_{a_{a_{a_{a_{a_{a_{a_{a_{a_$	d='a2'∧	b_id>10	2 (R)
b_id		a_id	b_id	
102		a2	103	
103				
CT * F	FROM R			
RE a_i	id='a2'	AND	o_id>1	02
	a2'(R) b_id 102 103 CT * F RE a_i	a2'(R) σ _{a_i} b_id 102 103 CT * FROM R RE a_id='a2'	a2'(R) σ _{a_id='a2'Λ} b_id 102 103 CT * FROM R RE a_id='a2' AND I	a2'(R) b_id 102 103 CT ★ FROM R RE a_id='a2' ∧ b_id>10 a_id b_id a2 103 CT ★ FROM R

RELATIONAL ALGEBRA: PROJECTION

Generate a relation with tuples that contains only the specified attributes.

- \rightarrow Rearrange attributes' ordering.
- \rightarrow Remove unwanted attributes.
- → Manipulate values to create derived attributes.

Syntax: $\Pi_{A1,A2,...,An}(R)$

R	(a_id	l,b_id	
	a_id	b_id	
	a1	101	
	a2	102	
	a2	103	
	a3	104	

I _{b_id-}	100,a_id (0 ;	a_id='a2	2' (R))
	b_id-100	a_id	
	2	a2	
	3	a2	

SELECT	b_	_id-100),	a_i	İd	
FROM	R	WHERE	a_	_id	=	'a2';



RELATIONAL ALGEBRA: UNION

Generate a relation that contains all tuples that appear in either only one or both input relations.

Syntax: (R U S)

R	(a_ic	l,b_ic	1)
	a_id	b_id	
	a1	101	
	a2	102	
	a3	103	

5	(a_id	,b_id	
	a_id	b_id	
	a3	103	
	a4	104	
	a5	105	

(R (JS)
a_id	b_id
a1	101
a2	102
a3	103
a4	104
a5	105

(SELECT	*	FROM	R)
U	NI	ON	
(SELECT	*	FROM	S);

RELATIONAL ALGEBRA: INTERSECTION

Generate a relation that contains only the tuples that appear in both of the input relations.

Syntax: ($R \cap S$)

R	(a_io	l,b_ic	ľ
	a_id	b_id	
	a1	101	
	a2	102	
	a3	103	

S(a_id,b_id) a_id b_id

103

104

105

a3

a4

a5

(R ſ	ר S)
a_id	b_id
a3	103

(SELECT * FROM R) INTERSECT (SELECT * FROM S);



RELATIONAL ALGEBRA: DIFFERENCE

Generate a relation that contains only the tuples that appear in the first and not the second of the input relations.

Syntax: (R - S)

R(a_id,b_id)a_idb_ida1101a2102a3103

S	(a_id	,b_id)
	a_id	b_id	
	a3	103	
	a4	104	
	a5	105	



(SELECT * FROM R) EXCEPT (SELECT * FROM S);



RELATIONAL ALGEBRA: PRODUCT

Generate a relation that contains all possible combinations of tuples from the input relations.

Syntax: (R × S)

15-445/645 (Spring 2025)

SELECT * FROM R **CROSS JOIN** S;

SELECT * FROM R, S;

R(a_id,b_id) a_id b_id a1 101 a2 102 a3 103

S(a_id,b_id)

a_id	b_id
a3	103
a4	104
a5	105

(R × S)

R.a_id	R.b_id	S.a_id	S.b_id
a1	101	a3	103
a1	101	a4	104
a1	101	a5	105
a2	102	a3	103
a2	102	a4	104
a2	102	a5	105
a3	103	a3	103
a3	103	a4	104
a3	103	a5	105

Generate a relation that contains all tuples that are a combination of two tuples (one from each input relation) with a common value(s) for one or more attributes.

Syntax: (R ⋈ S)



Generate a relation that contains all tuples that are a combination of two tuples (one from each input relation) with a common value(s) for one or more attributes.

Syntax: (R ⋈ S)



 (R ⋈ S)

 a_id b_id val

 a3
 103



Generate a relation that contains all tuples that are a combination of two tuples (one from each input relation) with a common value(s) for one or more attributes.



 $(R \bowtie S)$

			-	-					- /
		R.a_id	R.b_id	S.a_id	S.b_id	S.val	a_id	b_id	val
Syntax: (R 🕨	4 S)	a3	103	a3	103	XXX	a3	103	XXX



Generate a relation that contains all tuples that are a combination of two tuples (one from each input relation) with a common value(s) for one or more attributes.



⊠ S)

(R

Syntax: (R 🖂 S) R.a_id R.b_id S 3 d S 5 d S.val a3 103 3 XXX a3 103 XXX a3 103 XXX



Generate a relation that contains all tuples that are a combination of two tuples (one from each input relation) with a common value(s) for one or more attributes.

Syntax: (R ⋈ S)

15-445/645 (Spring 2025)



(R ⋈ S) a_id b_id val a3 103 XXX

Generate a relation that contains all tuples that are a combination of two tuples (one from each input relation) with a common value(s) for one or more attributes.

Syntax: (R ⋈ S)

R	(a_io	l,b_id	d)	S(a_id	,b_i	d,va	1)
	a_id	b_id			a_id	b_id	val	
	a1	101			a3	103	XXX	
	a2	102			a4	104	YYY	
	a3	103			a5	105	ZZZ	

(R ⋈ S)		
a_id	b_id	val
a3	103	XXX

SELECT * FROM R **NATURAL JOIN** S;

SELECT * FROM R **JOIN** S **USING** (a_id, b_id);

SELECT * FROM R JOIN S
ON R.a_id = S.a_id AND R.b_id = S.b_id;

ECMU-DB 15-445/645 (Spring 20

RELATIONAL ALGEBRA: EXTRA OPERATORS

Rename (ρ) Assignment (R←S) Duplicate Elimination (δ) Aggregation (γ) Sorting (τ) Division (R÷S)

OBSERVATION

Relational algebra defines an ordering of the highlevel steps of how to compute a query. \rightarrow Example: $\sigma_{b_{id=102}}(R \bowtie S)$ vs. $(R \bowtie (\sigma_{b_{id=102}}(S))$

A better approach is to state the high-level answer that you want the DBMS to compute.
→ Example: Retrieve the joined tuples from **R** and **S** where **b_id** equals 102.



RELATIONAL MODEL: QUERIES

The relational model is independent of any query language implementation.

SQL is the *de facto* standard (many dialects).

```
for line in file.readlines():
    record = parse(line)
    if record[0] == "GZA":
        print(int(record[1]))
```

SELECT year FROM artists
WHERE name = 'GZA';

ECMU-DB 15-445/645 (Spring 2025)



15-445/645 (Spring 2025)

A collection of record documents containing a hierarchy of named field/value pairs.

- \rightarrow A field's value can be either a scalar type, an array of values, or another document.
- \rightarrow Modern implementations use JSON. Older systems use XML or custom object representations.

Avoid "relational-object impedance mismatch" by tightly coupling objects and database.

























Application Code class Artist { int id; String name; int year; Album albums[]; } class Album { int id; String name;

int year;

ECMU-DB 15-445/645 (Spring 2025)
DOCUMENT DATA MODEL



Application Code

class Artist {
 int id;
 String name;
 int year;
 Album albums[];
}
class Album {
 int id;
 String name;
 int year;
}

```
"name": "GZA",
"year": 1990,
"albums": [
   "name": "Liquid Swords",
   "year": 1995
  },
   "name": "Beneath the Surface",
   "year": 1999
```

ECMU-DB 15-445/645 (Spring 2025)

DOCUMENT DATA MODEL



Application Code class Artist { int id; String name; int year; Album albums[]; class Album { int id; String name; int year;



ECMU-DB 15-445/645 (Spring 2025)

One-dimensional arrays used for nearest-neighbor search (exact or approximate).

- \rightarrow Used for semantic search on embeddings generated by ML-trained transformer models (think ChatGPT).
- → Native integration with modern ML tools and APIs (e.g., LangChain, OpenAI).

At their core, these systems use specialized indexes to perform NN searches quickly.



Album(id, name, year)id nameyear11Enter the Wu-Tang199322St.Ides Mix Tape199433Liquid Swords1995













ECMU-DB 15-445/645 (Spring 2025)



SCMU-DB 15-445/645 (Spring 2025)









86

CONCLUSION

Databases are ubiquitous.

Relational algebra defines the primitives for processing queries on a relational database.

We will see relational algebra again when we talk about query optimization + execution.

NEXT CLASS

Modern SQL \rightarrow Make sure you understand basic SQL before the lecture.

