• Given domains  $D_1$ ,  $D_2$ , ...,  $D_n$  a **relation** *r* is a subset of  $D_1 \times D_2 \times ... \times D_n$  (*cartesian product*)

Thus, a relation is a set of tuples (a1, a2, ..., an) where each ai  $\in$  Di

- Schema of a relation consists of
  - attribute definitions
    - name
    - type/domain
  - integrity constraints

Tuple $\longrightarrow$ RowRelation $\longrightarrow$ TableMathGeneral

## 

#### account(Account\_schema) Relation from a schema

account_number	branch_name	balance
A-101	Downtown	500
A-102	Perryridge	400
A-201	Brighton	900
A-215	Mianus	700
A-217	Brighton	750
A-222	Redwood	700
A-305	Round Hill	350

#### Relation instance

# **Relation Instance**

- The current values (relation instance) of a relation are specified by a table
- An element t of r is a tuple, represented by a row in a table
- Order of tuples is irrelevant (tuples may be stored in an arbitrary order)





- A database consists of multiple relations
- Information about an enterprise is broken up into parts, with each relation storing one part of the information
- E.g.

account : information about accounts
depositor : which customer owns which account
customer : information about customers

### The customer Relation

customer_name	customer_street	customer_city
Adams	Spring	Pittsfield
Brooks	Senator	Brooklyn
Curry	North	Rye
Glenn	Sand Hill	Woodside
Green	Walnut	Stamford
Hayes	Main	Harrison
Johnson	Alma	Palo Alto
Jones	Main	Harrison
Lindsay	Park	Pittsfield
Smith	North	Rye
Turner	Putnam	Stamford
Williams	Nassau	Princeton

Customer\_schema = (customer\_name, customer\_street, customer\_city)

# The depositor Relation

customer_name	account_number
Hayes	A-102
Johnson	A-101
Johnson	A-201
Jones	A-217
Lindsay	A-222
Smith	A-215
Turner	A-305

Depositor\_schema = (customer\_name, account\_number)

# Why Split Information Across Relations?

• Storing all information as a single relation such as

```
Bank_schema = (account_number, balance, customer_name, ..)
```

- Results in
  - repetition of information
    - e.g., if two customers own an account (What gets repeated?)
  - the need for null values
    - e.g., to represent a customer without an account
- Normalization theory (Chapter 7) deals with how to design relational schemas

- Reflect constraints in the real-world enterprise
- Let  $K \subseteq R$
- K is a **superkey** of R if values for K are sufficient to identify a unique tuple of each possible relation r(R)
  - by "possible r " we mean a relation r that could exist in the enterprise we are modeling.
  - Example: {customer\_name, customer\_street} and {customer\_name} are both superkeys of Customer, if no two customers can possibly have the same name
    - In real life, an attribute such as customer\_id would be used instead of customer\_name to uniquely
      identify customers, but we omit it to keep our examples small, and instead assume customer names are
      unique.

Keys (Cont.)

- K is a **candidate key** if K is *minimal*: no subset of K is a superkey Example: {customer\_name} is a candidate key for Customer
- Primary key: a candidate key chosen as the principal means of identifying tuples within a relation
  - Should choose an attribute whose value never, or very rarely, changes.
  - E.g. email address is unique, but may change
  - Others?
  - Generate your own

#### R: relational schema

# K: superkey of R $\Rightarrow$ restriction on relations r(R)t<sub>1</sub>, t<sub>2</sub> $\in$ r and t<sub>1</sub> $\neq$ t<sub>2</sub> $\Rightarrow$ t<sub>1</sub>[K] $\neq$ t<sub>2</sub>[K]



- A relation schema may have an attribute that corresponds to the primary key of another relation. The attribute is called a **foreign key**.
  - E.g. customer\_name and account\_number attributes of depositor are foreign keys to customer and account respectively.
  - Only values occurring in the primary key attribute of the referenced relation may occur in the foreign key attribute of the referencing relation.



# Schema Diagram





# Query Languages

- Language in which user requests information from the database.
- Categories of languages
  - Procedural
  - Non-procedural, or **declarative**
- "Pure" languages:
  - Relational algebra
  - Tuple relational calculus
  - Domain relational calculus
- Pure languages form underlying basis of query languages that people use.

- Similar to regular algebra (3\*x + 2\*y)
- Relations instead of numbers
- Why study?
  - foundation of low-level DBMS operations
  - in order to understand query execution
  - Build sophisticated SQL queries
- More procedural than SQL (which is declarative)



- A set: unordered collection of distinct objects
- Elements of a set
- Subset, proper subset, superset
- Union
- Intersection
- set difference
- Cartesian product (set of ordered pairs)

{0, 20, 12, 60}
{Canada, U.S.A.}

# **Relational Operators**

- Six basic operators
  - select:  $\sigma$
  - project: ∏
  - **union:** U
  - set difference: -
  - Cartesian product: x
  - rename: ρ

# Select Operation – Example

All tuples in relation loan where branch is "Perryridge"

$\sigma$ branch_name="Perryridge"(loan)					
Predicate					
loan_number	branch_name	amount			
L-15	Perryridge	1500			
L-16	Perryridge	1300			

loan_number	branch_name	amount
L-11	Round Hill	900
L-14	Downtown	1500
L-15	Perryridge	1500
L-16	Perryridge	1300
L-17	Downtown	1000
L-23	Redwood	2000
L-93	Mianus	500

All tuples with amount lent is more than \$1200  $\sigma_{amount} > 1200(loan)$ 

# Select Operation

Comparators: =,  $\neq$ , <,  $\leq$ , >,  $\geq$ 

Connectives: and  $(\land)$ , or  $(\lor)$ , not  $(\neg)$ 

 $\sigma_{branch_name="Perryridge" \land amount > 1350(loan)}$ 

loan_number	branch_name	amount	
L-15	Perryridge	1500	

loan_number	branch_name	amount
L-11	Round Hill	900
L-14	Downtown	1500
L-15	Perryridge	1500
L-16	Perryridge	1300
L-17	Downtown	1000
L-23	Redwood	2000
L-93	Mianus	500

# Project Operation – Example

#### List only part of a relation

#### $\pi_{loan_number}$ , amount(loan)

•	
loan_number	amount
L-11	900
L-14	1500

loan_number	branch_name	amount	
L-11	Round Hill	900	
L-14	Downtown	1500	
L-15	Perryridge	1500	
L-16	Perryridge	1300	
L-17	Downtown	1000	
L-23	Redwood	2000	
L-93	Mianus	500	

# Result of a relational operation is a relation

 $\pi_{loan_number}(\sigma_{branch_name="Perryridge"(loan))}$ 



loan_number	branch_name	amount
L-11	Round Hill	900
L-14	Downtown	1500
L-15	Perryridge	1500
L-16	Perryridge	1300
L-17	Downtown	1000
L-23	Redwood	2000
L-93	Mianus	500

# Union operation

• Combine the result of two operations

#### Query: customers with an account or a loan (or both)

	customer_name	account_number		customer_name	account_number		
Denositor	Hayes	A-102		Adams	L-16		
	Johnson	A-101		Curry	L-93	Borrow	or
relation	Johnson	A-201		Hayes	L-15	relation	
		A 017		Jackson	L-14	rolation	1
	JOHES	A-217		Jones	L-17		
	Lindsay	A-222		Smith	L-11		customer_name
	Smith	A-215		Smith	L-23		Adams
	Turner	A-305		Williams	L-17		Curry
	$\longleftrightarrow$					Hayes	
							Jackson
							Jones
	Smith						
$\Pi$ customer_name(depositor) U $\Pi$ customer_name(borrower) $\longrightarrow$							Williams
						Lindsay	
					Johnson		

Turner

- For (r U s) to work
  - r and s should have **same arity** (same number of attributes)
  - corresponding attributes should have same domain

# Set-difference operations

• Find tuples in one that are not present in another

 $\pi$ customer\_name(depositor) -  $\pi$ customer\_name(borrower)



• Same rules for compatibility apply as in Union

# Cartesian-Product operation

- Combine information,  $r_1 \times r_2$
- Remember: a relation is a subset of a Cartesian product
- Naming scheme to differentiate attributes: relation.attribute
  - only for non-distinct attributes
- What tuples appear in  $r_1 \times r_2$ ?
- tuples in  $r_1 \times r_2$ : all possible combinations of tuples in  $r_1$  and  $r_2$
- if  $r_1$  has  $n_1$ , and  $r_2$  has  $n_2$ , then  $r_1 \times r_2$  has  $n_{1*}n_2$  tuples

# Cartesian-Product

- For relations  $r_1(R_1)$ ,  $r_2(R_2)$ :
  - r<sub>1</sub> x r<sub>1</sub> concatenation of R<sub>1</sub> and R<sub>2</sub>
- For tuple  $t \in r_1 \times r_1$ , then:
  - $t[R_1] = t_1[R_1]$  and  $t[R_2] = t_2[R_2]$