



MONASH  
University

MONASH  
INFORMATION  
TECHNOLOGY

# Week 3

## The Relational Database Model



# Overview

Once we have a conceptual model, it is time to move to the second stage and map this to a logical model

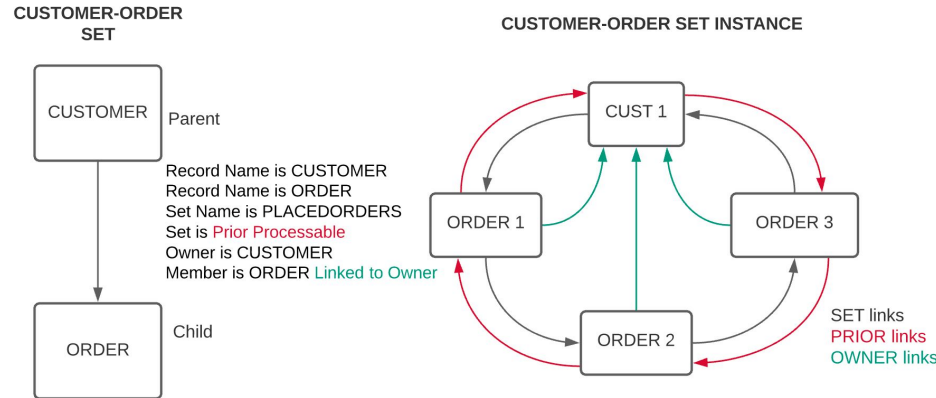
For our unit this will involve mapping to the *Relational Model* in preparation for implementation in a RDBMS. First before we consider this mapping it is necessary to have a clear understanding of the Relation Model and its use:

- Relational Model
- Relational Algebra

# Early Database Models

- Hierarchical (1970's eg. IBM Information Management System (IMS))
  - 1:M relationships, a tree of linked records, child has only one parent
- Network (1970's eg. Integrated Data Store IDS, basis for the CODASYL group)
  - child may have multiple parents
- Both Navigational - move around in data via *embedded links* (pointers)

Network:



# The Relational Model

- Introduced by CODD in 1970 - the fundamental basis for the relational DBMS
- Basic structure is the mathematical concept of a RELATION mapped to the 'concept' of a table (tabular representation of relation)
  - Relation - abstract object
  - Table - pictorial representation
  - Storage structure - "real thing" - eg. isam file of 1's and 0's
- Relational Model Terminology
  - DOMAIN - set of atomic (indivisible) values
  - Examples (name, data type, data format):
    - customer\_number domain - 5 character string of the form xxxdd
    - name domain - 20 character string
    - address domain - 30 character string containing street, town & postcode
    - credit\_limit domain - money in the range \$1,000 to \$99,999

# A Relation

- A relation consists of two parts
  - heading
  - body
- **Relation Heading**
  - Also called Relational **Schema** consists of a fixed set of attributes
    - $R(A_1, A_2, \dots, A_n)$ 
      - $R$  = relation name,  $A_i$  = attribute  $i$
  - Each attribute corresponds to one underlying domain:
    - Customer relation heading:
      - CUSTOMER (custno, custname, custadd, custcredlimit)
        - »  $\text{dom}(\text{custno}) = \text{customer\_number}$
        - »  $\text{dom}(\text{custname}) = \text{name}$
        - »  $\text{dom}(\text{custadd}) = \text{address}$
        - »  $\text{dom}(\text{custcredlimit}) = \text{credit\_limit}$

custno	custname	custadd	custcredlimit
--------	----------	---------	---------------

# Relation Body

## ▪ Relation Body

- Also called Relation Instance (state of the relation at any point in time)
  - $r(R) = \{t_1, t_2, t_3, \dots, t_m\}$
  - consists of a time-varying set of n-tuples
    - Relation R consists of tuples  $t_1, t_2, t_3 \dots t_m$
    - $m$  = number of tuples = **relation cardinality**
  - each n-tuple is an ordered list of n values
  - $t = \langle v_1, v_2, \dots, v_n \rangle$ 
    - $n$  = number of values in tuple (no of attributes) = **relation degree**
- In the tabular representation:
  - Relation heading  $\Rightarrow$  column headings
  - Relation body  $\Rightarrow$  set of data rows

custno	custname	custadd	custcredlimit
SMI13	SMITH	Wide Rd, Clayton, 3168	2000
JON44	JONES	Narrow St, Clayton, 3168	10000
BRO23	BROWN	Here Rd, Clayton, 3168	10000

# Relation Properties

- No duplicate tuples
  - by definition sets do not contain duplicate elements
    - hence tuples must be unique
- Tuples are unordered within a relation
  - by definition sets are not ordered
    - hence tuples can only be accessed by content
- No ordering of attributes within a tuple
  - by definition sets are not ordered

# Relation Properties cont'd

- Tuple values are atomic - cannot be divided
  - EMPLOYEE (eid, ename, departno, dependants)
    - not allowed: dependants (depname, depage) multivalued
  - hence no multivalued (repeating) attributes allowed, called the first normal form rule
- COMPARE with tabular representation
  - normally nothing to prevent duplicate rows
  - rows are ordered
  - columns are ordered
  - tables and relations are not the same 'thing'



# Functional Dependency

## ▪ Functional Dependency:

- A set of attributes A functionally determines an attribute B if, and only if, for each A value, there is exactly one value of B in the relation. It is denoted as  $A \rightarrow B$  (A determines B, or B depends on A)
  - $\text{orderno} \rightarrow \text{orderdate}$
  - $\text{prodno} \rightarrow \text{proddesc}$
  - $\text{orderno}, \text{prodno} \rightarrow \text{qtyordered}$

ORDERNO	ORDERDATE
10	01/MAY/19
11	02/MAY/19
12	03/MAY/19
13	04/MAY/19
14	04/MAY/19
15	05/MAY/19
16	06/MAY/19

ORDERNO	PRODNO	QTYORDERED	LINEPRICE
10	101	1	11.98
11	101	1	11.98
11	103	2	123.58
12	104	10	479.8
13	105	2	140.36
14	106	1	31.99
15	107	3	116.73

PRODNO	PRODDISC	PRODUNITPRICE
101	Salmon - Smoked, Sliced	11.98
102	Brocolinni - Gaylan, Chinese	80.75
103	Pasta - Lasagne, Fresh	61.79
104	Melon - Cantaloupe	47.98
105	Wine - Peller Estates Late	70.18
106	Peas - Pigeon, Dry	31.99
107	Pumpkin - Seed	38.91

# Relational Model Keys

- A **superkey** of a relation R is an attribute or set of attributes which exhibits only the uniqueness property
  - No two tuples of R have the same value for the superkey (Uniqueness property)
  - $t1[\text{superkey}] \neq t2[\text{superkey}]$
- A **candidate key** CK of a relation R is an attribute or set of attributes which exhibits the following properties:
  - Uniqueness property (as above), *and*
  - No proper subset of CK has the uniqueness property (Minimality or Irreducibility property) ie. a minimal superkey
- One candidate key is chosen to be the **primary key** (PK) of a relation. Remaining candidate keys are termed alternate keys (AK).

Many possible superkeys

Potentially many possible candidate keys

Only ONE primary key (may be composed of many attributes - a composite primary key)

**Free text:**

**Q1. List all the super keys for:**

Unit_code	Mark	Studid	Year
FIT2094	80	111	2016
FIT2094	20	111	2015
FIT2004	100	111	2016
FIT2004	40	222	2015
FIT2004	40	333	2015

**Treat this as only a small sample of the data**

**Free text:**

**Q2. Given the business rules:**

- a physician can only treat one patient on a particular date and time,
- a patient can only be treated by one physician on a particular date and time, and
- a patient may receive multiple treatments from the same physician on a particular date and time

**Identify the candidate key/s of this relation:**

**PATIENT\_TREATMENT (patient\_id, physician\_id, treatment\_code,  
pt\_date, pt\_time, pt\_result)**

# Selection of a Primary key

- **A primary key must be chosen considering the data that *may be added to the table in the future***
  - Names, dates of birth etc are rarely unique and as such are not a good option
  - PK should be free of 'extra' semantic meaning and security compliant, preferably a single attribute, preferably numeric (see Table 5.3 Coronel & Morris)
  - Natural vs Surrogate primary key
    - PATIENT\_TREATMENT (patient\_id, physician\_id, treatment\_code, pt\_date, pt\_time, pt\_result)
      - Superkey
      - CK
      - PK
      - Issues with PK?

TABLE 5.3

**DESIRABLE PRIMARY KEY CHARACTERISTICS**

PK CHARACTERISTIC	RATIONALE
Unique values	The PK must uniquely identify each entity instance. A primary key must be able to guarantee unique values. It cannot contain nulls.
Nonintelligent	The PK should not have embedded semantic meaning other than to uniquely identify each entity instance. An attribute with embedded semantic meaning is probably better used as a descriptive characteristic of the entity than as an identifier. For example, a student ID of 650973 would be preferred over Smith, Martha L. as a primary key identifier.
No change over time	If an attribute has semantic meaning, it might be subject to updates, which is why names do not make good primary keys. If Vickie Smith is the primary key, what happens if she changes her name when she gets married? If a primary key is subject to change, the foreign key values must be updated, thus adding to the database work load. Furthermore, changing a primary key value means that you are basically changing the identity of an entity. In short, the PK should be permanent and unchangeable.
Preferably single-attribute	A primary key should have the minimum number of attributes possible (irreducible). Single-attribute primary keys are desirable but not required. Single-attribute primary keys simplify the implementation of foreign keys. Having multiple-attribute primary keys can cause primary keys of related entities to grow through the possible addition of many attributes, thus adding to the database workload and making (application) coding more cumbersome.
Preferably numeric	Unique values can be better managed when they are numeric, because the database can use internal routines to implement a counter-style attribute that automatically increments values with the addition of each new row. In fact, most database systems include the ability to use special constructs, such as Autonumber in Microsoft Access, sequence in Oracle, or uniqueidentifier in MS SQL Server to support self-incrementing primary key attributes.
Security-compliant	The selected primary key must not be composed of any attribute(s) that might be considered a security risk or violation. For example, using a Social Security number as a PK in an EMPLOYEE table is not a good idea.

# Null in the Relational Model Implementation

- ***NULL is a concept created and implemented by SQL***
- NULL is NOT a value - is a representation of the fact that there is *NO VALUE*
- Reasons for a NULL:
  - VALUE NOT APPLICABLE -
    - EMP relation - empno, deptno, salary, commission
      - commission only applies to staff in sales dept
  - VALUE UNKNOWN -
    - Joe's salary is NULL, Joe's salary is currently unknown
  - VALUE DOES NOT EXIST -
    - Tax File Number - is applicable to all employees but Joe may not have a number at this time
  - VALUE UNDEFINED -
    - Certain items explicitly undefined eg. divide by zero
      - Columns Number\_of\_payments, Total\_payments
      - Column Average\_payment\_made
      - If Number\_of\_payments = 0 => Average undefined

# Writing Relations

- Relations may be represented using the following notation:
  - `RELATION_NAME (attribute1, attribute2,...)`
- The primary key is underlined.
- Example:
  - **`STAFF (staff_id, staff_surname, staff_initials, staff_address, staff_phone)`**



**Q3. A well designed relational database ( a database based on the relational model) has:**

- A. No redundant data
- B. Minimal redundant data
- C. A large amount of redundant data
- D. A level of redundancy based on the vendors implementation

# Relational Database

- A relational database is a collection of normalised relations.
- Normalisation is part of the design phase of the database and will be discussed in a later lecture.

Example relational database:

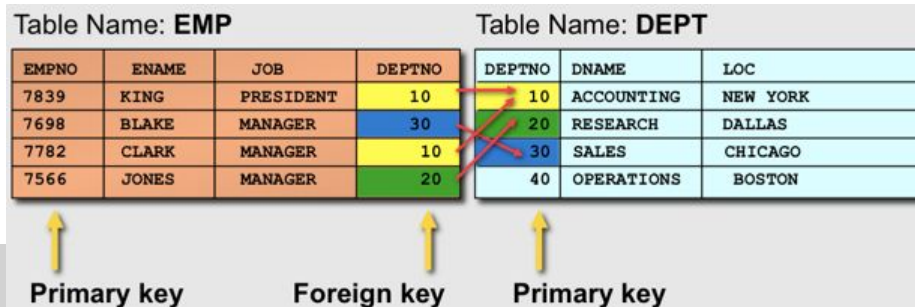
**ORDER (order\_id, orderdate)**

**ORDER\_LINE (order\_id, product\_id, ol\_quantity)**

**PRODUCT (product\_id, product\_description, product\_unitprice)**

# Foreign Key (FK) - Implementation

- FK: An attribute/s in a relation that exists in the same, or another relation as a Primary Key.
- Referential Integrity
  - A Foreign Key value must either *match* the **full primary key** in a relation **or be NULL**.
- The pairing of PK and FK creates relationships (**logical** connections) between tables when implemented in a RDBMS. Hence the abstraction away from the underlying storage model.



#### Q4. Business rules:

Runners may form a team, the runner who registers the team is recorded as the team leader. Each team can have up to 5 members (runners).

Identify the FK(s):

**TEAM(team\_id, team\_name, team\_leader)**

**RUNNER(runner\_id, runner\_name, team\_id)**

- A. team\_leader in TEAM
- B. team\_id in TEAM
- C. runner\_id in RUNNER
- D. team\_id in RUNNER

# Data Integrity - Implementation

- Entity integrity
  - Primary key value must not be NULL.
    - No duplicate tuple property then ensures that each primary key must be unique
    - Implemented in the RDBMS via a unique index on the PK
- Referential integrity
  - The values of FK must either match a value of a full PK in the related relation or be NULL.
- Column/Domain integrity
  - All values in a given column must come from the same domain (the same data type and range).

**Q5. The following set of relations:**

**HOSPITAL (hosp\_id, hosp\_name, hosp\_phone)**

**DOCTOR (hosp\_id, dr\_id, dr\_name, dr\_mobile)**

**PATIENT (pat\_id, pat\_name, pat\_dob, dr\_id)**

- A. have no integrity issues
- B. violate entity integrity
- C. violate referential integrity
- D. violate column/domain integrity

*Multiple responses allowed*

# Relational DMLs

- Relational Calculus
- Relational Algebra
- Transform Oriented Languages (e.g. SQL)
- Graphical Languages
- Exhibit the “closure” property - queries on relations produce relations

# Relational Calculus

- Based on mathematical logic.
- Non-procedural.
- Primarily of theoretical importance.
- May be used as a yardstick for measuring the power of other relational languages (“relational completeness”).
- Operators may be applied to any number of relations.



# RELATIONAL ALGEBRA

Manipulation of relational data

# Relational Algebra

- Relationally complete.
- Procedural.
- Operators only apply to at most two relations at a time.
- 8 basic operations:
  - single relation: selection, projection
  - cartesian product, join
  - union
  - intersection
  - difference
  - division
- *Pure form has no concept of NULL*

# Relational Operation PROJECT

PRDETAIL (project\_code, project\_manager, project\_bid\_price)

$\pi$

PROJECT_CODE	PROJECT_MANAGER	PROJECT_BID_PRICE
21-5Z	Holly B. Parker	\$16,833,460.00
25-2D	Jane D. Grant	\$12,500,000.00
25-5A	George F. Dorts	\$32,512,420.00
25-9T	Holly B. Parker	\$21,563,234.00
27-4Q	George F. Dorts	\$10,314,545.00
29-2D	Holly B. Parker	\$25,559,999.00
31-7P	William K. Moor	\$56,850,000.00

# Relational Operation SELECT

PRDETAIL (project\_code, project\_manager, project\_bid\_price)

$\sigma$

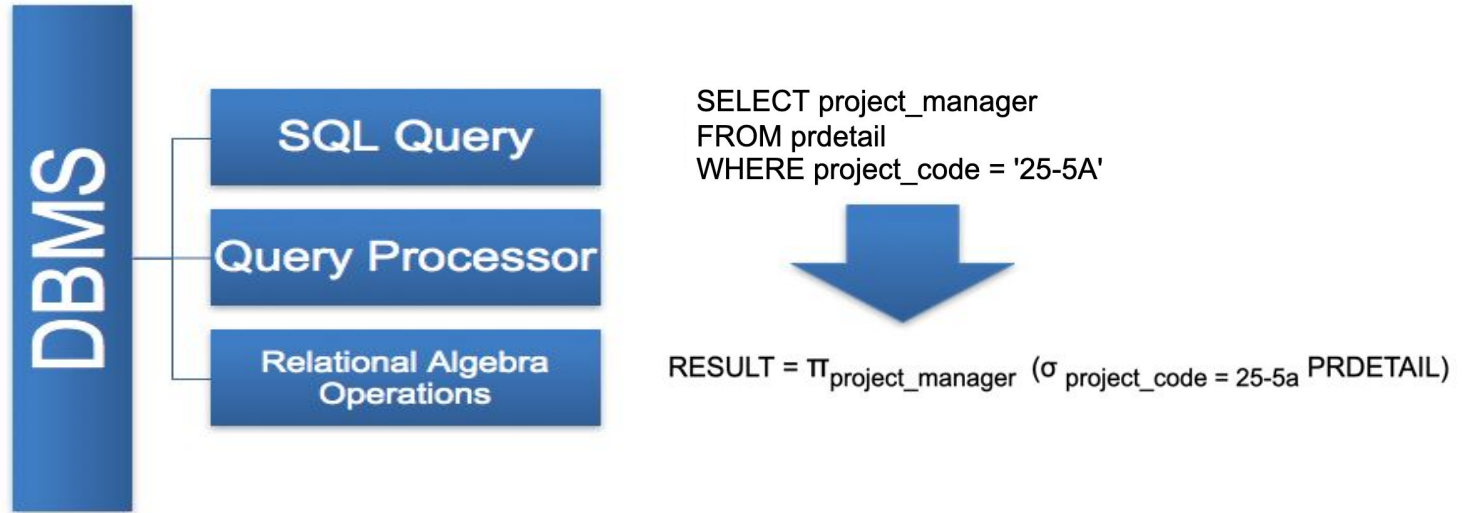
PROJECT_CODE	PROJECT_MANAGER	PROJECT_BID_PRICE
21-5Z	Holly B. Parker	\$16,833,460.00
25-2D	Jane D. Grant	\$12,500,000.00
25-5A	George F. Dorts	\$32,512,420.00
25-9T	Holly B. Parker	\$21,563,234.00
27-4Q	George F. Dorts	\$10,314,545.00
29-2D	Holly B. Parker	\$25,559,999.00
31-7P	William K. Moor	\$56,850,000.00

## Relational Operation Multiple Actions

PROJECT_CODE	PROJECT_MANAGER	PROJECT_BID_PRICE
21-5Z	Holly B. Parker	\$16,833,460.00
25-2D	Jane D. Grant	\$12,500,000.00
25-5A	George F. Dorts	\$32,512,420.00
25-9T	Holly B. Parker	\$21,563,234.00
27-4Q	George F. Dorts	\$10,314,545.00
29-2D	Holly B. Parker	\$25,559,999.00
31-7P	William K. Moor	\$56,850,000.00

RESULT =  $\Pi_{\text{project\_manager}}$  ( $\sigma_{\text{project\_code} = 25-5A}$  PRDETAIL)

# SQL vs Relational Algebra in the Database



# Relational Algebra *select* and *project*

The following relations represent a karate dojo member training attendance:

**SENSEI** (sensei\_id, sensei\_name)

**TRAINING\_SCHEDULE** (training\_day, training\_time, group\_id, sensei\_id)

**ATTENDANCE** (training\_day, training\_time, member\_id, attendance\_date)

**MEMBER** (member\_id, member\_name, member\_dob, member\_belt, group\_id)

**GROUP** (group\_id, group\_name, group\_age\_range)

- Primary keys are underlined
- A karate member falls into one of the age level groups: Tiny Tiger (for 4-7 year old), Young Dragon (for 8-14 years old), or Adult (for 14+ years old) and owns a certain color of belt (e.g. white, green, brown or black)
- Sensei (Karate teachers) are scheduled to train an age level group of karate members in a particular day and time (e.g. Sensei Luke Nakamura trains Tiny Tiger members every Tuesday 5pm)
- A karate member may attend more than one training schedule of their age level group in a given week.

Write the relational algebra for the following query (**your answer must show an understanding of query efficiency**):

**(1) Show the name and dob of all black belt members.**

# JOIN

- Join operator used to combine data from two or more relations, based on a common attribute or attributes.
- Different types:
  - theta-join
  - equi-join
  - natural join
  - outer join



# THETA JOIN (Generalised join)

$$(\text{Relation\_1}) \bowtie_F (\text{Relation\_2})$$

- $F$  is a predicate (i.e. truth-valued function) which is of the form  $\text{Relation\_1.a}_i \theta \text{Relation\_2.b}_i$ 
  - $\text{CUSTOMER.cust\_no} \theta \text{ORDER.cust\_no}$
- $\theta$  is one of the standard arithmetic comparison operators,  
 $<, \leq, =, \geq, >$
- Most commonly,  $\theta$  is equals ( $=$ ), but can be *any* of the operators
  - $\text{EMPLOYEE.emp\_sal} > \text{SALARYSCALE.step\_5}$

# NATURAL JOIN

STUDENT		MARK		
ID	Name	ID	Subj	Marks
1	Alice	1	1004	95
2	Bob	2	1045	55
		1	1045	90

**Step 1: STUDENT X MARK**

**Step 2: delete rows where IDs do not match (select =)**

**Result at Step 2 is an Equijoin**

STUDENT.ID	Name	MARK.ID	Subj	Marks
1	Alice	1	1004	95
<del>1</del>	<del>Alice</del>	<del>2</del>	<del>1045</del>	<del>55</del>
1	Alice	1	1045	90
<del>2</del>	<del>Bob</del>	<del>1</del>	<del>1004</del>	<del>95</del>
2	Bob	2	1045	55
<del>2</del>	<del>Bob</del>	<del>1</del>	<del>1045</del>	<del>90</del>

# NATURAL JOIN

STUDENT		MARK		
ID	Name	ID	Subj	Marks
1	Alice	1	1004	95
2	Bob	2	1045	55
		1	1045	90

**Step 1: STUDENT X MARK**

**Step 2: delete rows where IDs do not match (select =)**

**Step 3: delete duplicate columns (project away)**

**Result at Step 3 is a Natural Join**

STUDENT.ID	Name	MARK.ID	Subj	Marks
1	Alice	1	1004	95
1	Alice	1	1045	90
2	Bob	2	1045	55

# NATURAL JOIN

STUDENT			MARK		
ID	Name		ID	Subj	Marks
1	Alice	⋈	1	1004	95
2	Bob		2	1045	55
			1	1045	90

**Step 1: STUDENT X MARK**

**Step 2: delete rows where IDs do not match (select =)**

**Step 3: delete duplicate columns (project away)**

ID	Name	Subj	Marks
1	Alice	1004	95
1	Alice	1045	90
2	Bob	1045	55

**A natural join of STUDENT and MARK**

# LEFT OUTER JOIN (SQL Implementation)

STUDENT			MARK		
ID	Name		ID	Subj	Marks
1	Alice	⋈	1	1004	95
2	Bob		2	1045	55
3	Chris		1	1045	90
			4	1004	100

***Get (incomplete) information of only Chris***

ID	Name	Subj	Marks
1	Alice	1004	95
1	Alice	1045	90
2	Bob	1045	55
3	Chris	Null	Null

**A left outer join of STUDENT and MARK**

**All data from the LEFT is included, nulls fill the right side**

# RIGHT OUTER JOIN (SQL Implementation)

STUDENT			MARK		
ID	Name		ID	Subj	Marks
1	Alice	⋈	1	1004	95
2	Bob		2	1045	55
3	Chris		1	1045	90
			4	1004	100

***Get (incomplete) information of the student with ID 4***

ID	Name	Subj	Marks
1	Alice	1004	95
1	Alice	1045	90
2	Bob	1045	55
4	Null	1004	100

**A right outer join of STUDENT and MARK.  
All data from the RIGHT is included, nulls  
fill the left side**

# FULL OUTER JOIN (SQL Implementation)

## STUDENT

ID	Name
1	Alice
2	Bob
3	Chris



## MARK

ID	Subj	Marks
1	1004	95
2	1045	55
1	1045	90
4	1004	100

***Get (incomplete) information of both Chris and student with ID 4***

ID	Name	Subj	Marks
1	Alice	1004	95
1	Alice	1045	90
2	Bob	1045	55
3	Chris	Null	Null
4	Null	1004	100

**A full outer join of STUDENT and MARK**  
**All data from the RIGHT and LEFT is included, nulls fill where necessary**

# UNION, INTERSECT, DIFFERENCE

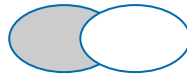
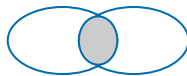
## STUDENT

student_id	student_name
1	Anne
2	Bruce
3	Claire



## STAFF

staff_id	staff_name
1	Anne
2	Jane
33	Claire



## UNION (STUDENT $\cup$ STAFF)

id	name
1	Anne
2	Bruce
3	Claire
2	Jane
33	Claire

## INTERSECT (STUDENT $\cap$ STAFF)

id	name
1	Anne

## DIFFERENCE (STUDENT - STAFF)

id	name
2	Bruce
3	Claire

**Union compatible relations required**



# Relational Algebra

The following relations represent a karate dojo member training attendance:

**SENSEI** (sensei\_id, sensei\_name)

**TRAINING\_SCHEDULE** (training\_day, training\_time, group\_id, sensei\_id)

**ATTENDANCE** (training\_day, training\_time, member\_id, attendance\_date)

**MEMBER** (member\_id, member\_name, member\_dob, member\_belt, group\_id)

**GROUP** (group\_id, group\_name, group\_age\_range)

- Primary keys are underlined
- A karate member falls into one of the age level groups: Tiny Tiger (for 4-7 year old), Young Dragon (for 8-14 years old), or Adult (for 14+ years old) and owns a certain color of belt (e.g. white, green, brown or black)
- Sensei (Karate teachers) are scheduled to train an age level group of karate members in a particular day and time (e.g. Sensei Luke Nakamura trains Tiny Tiger members every Tuesday 5pm)
- A karate member may attend more than one training schedule of their age level group in a given week.

Write the relational algebra for the following query (**your answer must show an understanding of query efficiency**):

**(2) Show the name, belt colour and attendance dates of the member with an id of 12345**

# ANSWER

$$R2 = \pi \text{ member\_name, member\_belt, attendance\_date } ( \text{MEMBER} \bowtie \text{ATTENDANCE} )$$

- this is the CANONICAL QUERY - not technically incorrect, but very inefficient, say member 12345 has only attended once in say 1000 tuples in ATTENDANCE. Your solution must demonstrate an understanding of efficiency:

$$A2a = \pi \text{ member\_id, member\_name, member\_belt } ( \sigma \text{ member\_id} = 12345 \text{ MEMBER} )$$

$$A2b = \pi \text{ member\_id, attendance\_date } ( \sigma \text{ member\_id} = 12345 \text{ ATTENDANCE} )$$

$$R2 = \pi \text{ member\_name, member\_belt, attendance\_date } ( A2a \bowtie A2b )$$

# Relational Algebra continues

The following relations represent a karate dojo member training attendance:

**SENSEI** (sensei\_id, sensei\_name)

**TRAINING\_SCHEDULE** (training\_day, training\_time, group\_id, sensei\_id)

**ATTENDANCE** (training\_day, training\_time, member\_id, attendance\_date)

**MEMBER** (member\_id, member\_name, member\_dob, member\_belt, group\_id)

**GROUP** (group\_id, group\_name, group\_age\_range)

- Primary keys are underlined
- A karate member falls into one of the age level groups: Tiny Tiger (for 4-7 year old), Young Dragon (for 8-14 years old), or Adult (for 14+ years old) and owns a certain color of belt (e.g. white, green, brown or black)
- Sensei (Karate teachers) are scheduled to train an age level group of karate members in a particular day and time (e.g. Sensei Luke Nakamura trains Tiny Tiger members every Tuesday 5pm)
- A karate member may attend more than one training schedule of their age level group in a given week.

Write the relational algebra for the following query (**your answer must show an understanding of query efficiency**):

- (3) Show the id, name and age level group name of members who were absent (did not attend any training) between 01-03-2021 and 31-03-2021 (inclusive).**

## ANSWER

A3a =  $\pi$  member\_id, member\_name, group\_id (MEMBER)

A3b =  $\pi$  member\_id ( $\sigma$  attendance\_date  $\geq$  01-03-2021 and  
attendance\_date  $\leq$  31-03-2021 (ATTENDANCE))

A3c = A3a  $\bowtie$  A3b

A3d = A3a - A3c

R =  $\pi$  member\_id, member\_name, group\_name (A3d  $\bowtie$  ( $\pi$  group\_id,  
group\_name(GROUP)))