

# *W4111 – Introduction to Databases*

## *Lecture 9: Module II (3), NoSQL (3)*



# *Contents*

# *Module II Reminder*

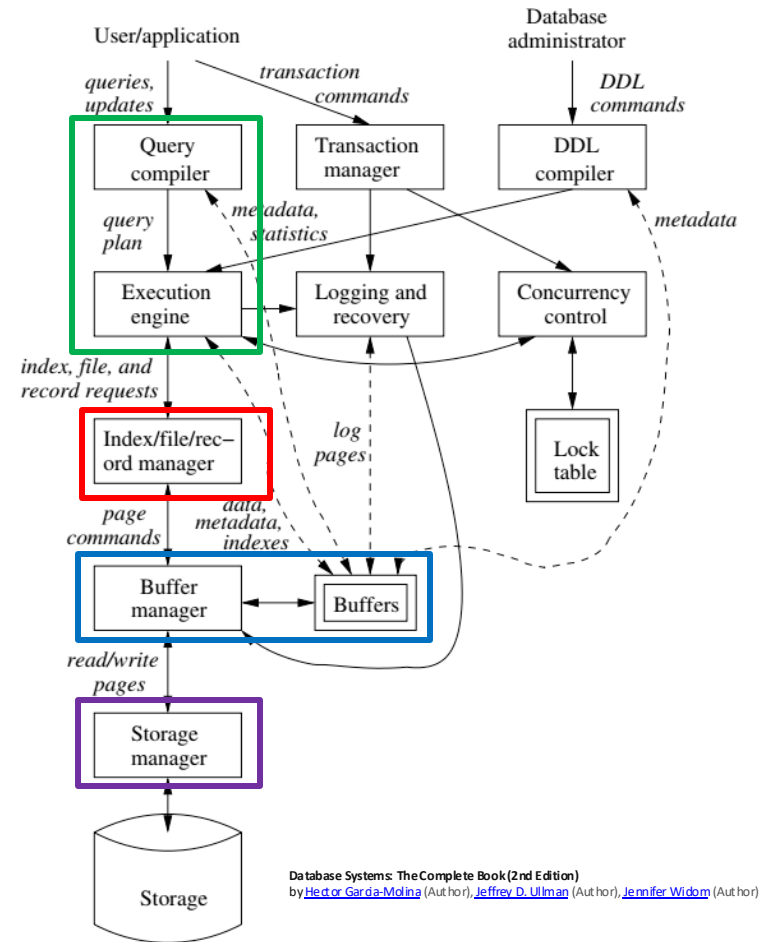
# Data Management

## Previously

- Load/save things quickly.
- Storage Mgmt. (cont)
- Access data quickly.

## Today

- Find things quickly (cont).
- Query processing, e.g. transform
  - Declarative language to
  - Procedural, functional, execution control



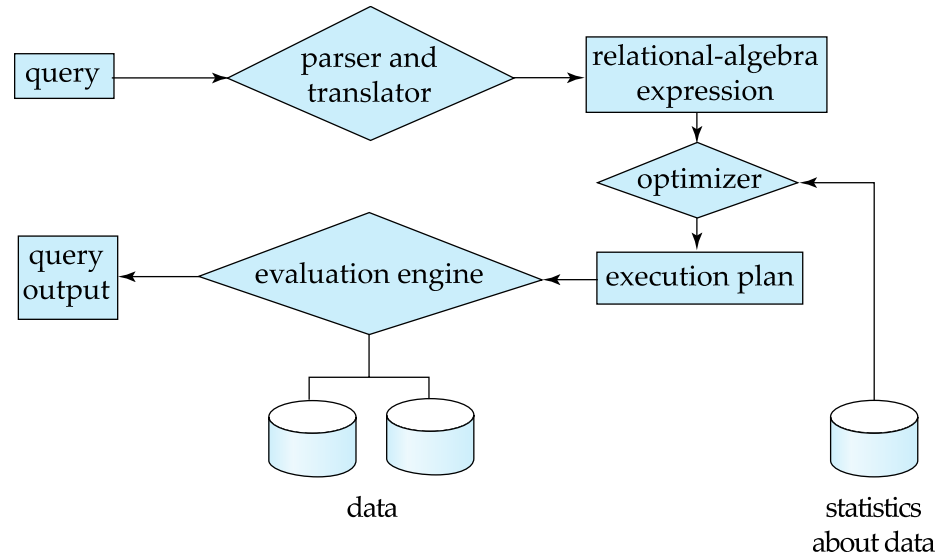
# *Query Processing*

# *Overview*

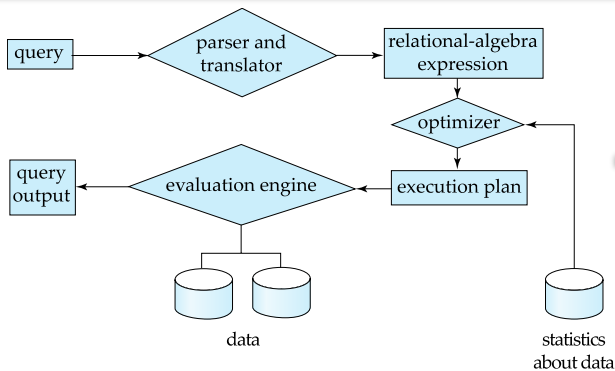


# Basic Steps in Query Processing

1. Parsing and translation
2. Optimization
3. Evaluation



# Query Processing



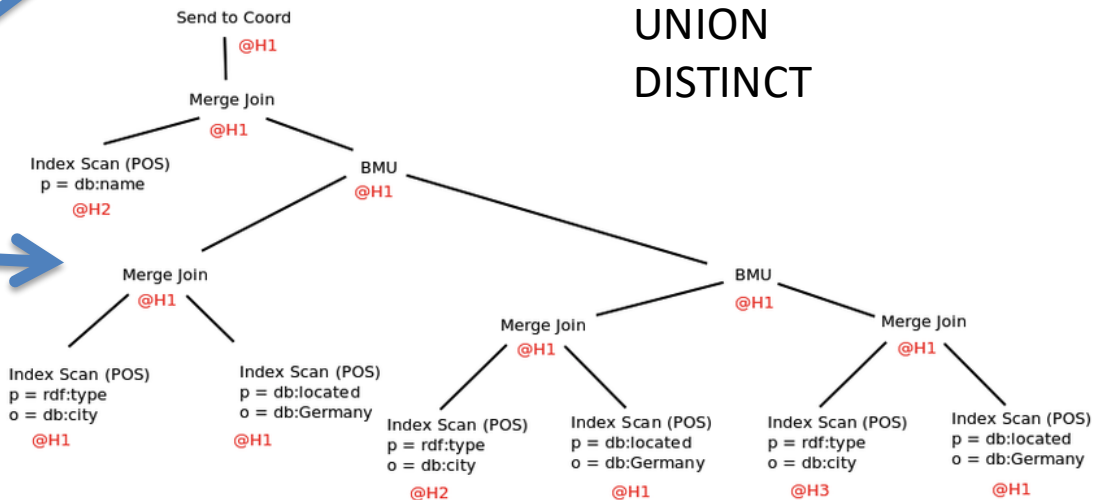
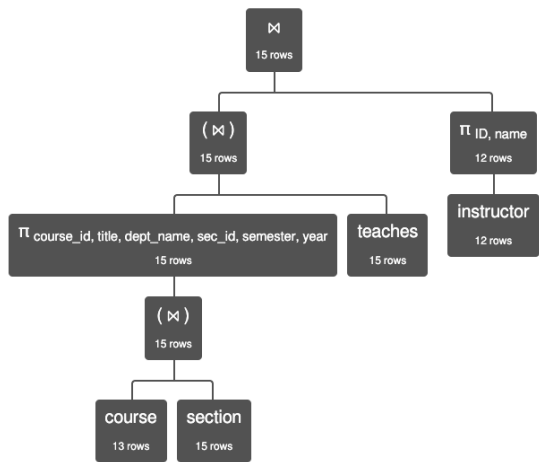
```

1 (
2   (π course_id, title, dept_name, sec_id, semester, year
3     (course ⋈ section)
4   )
5   ⋈ teaches
6 )
7
8 (π ID, name (instructor))

```

SELECT  
PROJECT

SELECT  
PROJECT  
JOIN  
UNION  
DISTINCT

$$\{SA1, SA2, \{PA1, PA2,$$




# *JOIN*

## *(Algorithm Selection Cont)*



# Join Operation

- Several different algorithms to implement joins
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- Choice based on cost estimate
- Examples use the following information
  - Number of records of *student*: 5,000    *takes*: 10,000
  - Number of blocks of *student*: 100    *takes*: 400



# Nested-Loop Join

- To compute the theta join  $r \bowtie_{\theta} s$   
    **for each** tuple  $t_r$  **in**  $r$  **do begin**  
        **for each** tuple  $t_s$  **in**  $s$  **do begin**  
            test pair  $(t_r, t_s)$  to see if they satisfy the join condition  $\theta$   
            if they do, add  $t_r \cdot t_s$  to the result.  
        **end**  
    **end**
- $r$  is called the **outer relation** and  $s$  the **inner relation** of the join.
- Requires no indices and can be used with any kind of join condition.
- Expensive since it examines every pair of tuples in the two relations.

For clarity (or confusion), I sometimes use the terms:

- Scan table for the outer relation
- Probe table for the inner relation



## Nested-Loop Join (Cont.)

- In the worst case, if there is enough memory only to hold one block of each relation, the estimated cost is
$$n_r * b_s + b_r \text{ block transfers, plus } n_r + b_r \text{ seeks}$$
- If the smaller relation fits entirely in memory, use that as the inner relation.
  - Reduces cost to  $b_r + b_s$  block transfers and 2 seeks
- Assuming worst case memory availability cost estimate is
  - with *student* as outer relation:
    - $5000 * 400 + 100 = 2,000,100$  block transfers,
    - $5000 + 100 = 5100$  seeks
  - with *takes* as the outer relation
    - $10000 * 100 + 400 = 1,000,400$  block transfers and 10,400 seeks
- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.
- Block nested-loops algorithm (next slide) is preferable.



# Indexed Nested-Loop Join

- Index lookups can replace file scans if
  - join is an equi-join or natural join and
  - an index is available on the inner relation's join attribute
    - Can construct an index just to compute a join.
- For each tuple  $t_r$  in the outer relation  $r$ , use the index to look up tuples in  $s$  that satisfy the join condition with tuple  $t_r$ .
- Worst case: buffer has space for only one page of  $r$ , and, for each tuple in  $r$ , we perform an index lookup on  $s$ .
- Cost of the join:  $b_r (t_r + t_s) + n_r * c$ 
  - Where  $c$  is the cost of traversing index and fetching all matching  $s$  tuples for one tuple of  $r$
  - $c$  can be estimated as cost of a single selection on  $s$  using the join condition.
- If indices are available on join attributes of both  $r$  and  $s$ , use the relation with fewer tuples as the outer relation.



# Merge-Join

1. Sort both relations on their join attribute (if not already sorted on the join attributes).
2. Merge the sorted relations to join them
  1. Join step is similar to the merge stage of the sort-merge algorithm.
  2. Main difference is handling of duplicate values in join attribute — every pair with same value on join attribute must be matched
  3. Detailed algorithm in book

		<i>a1</i>	<i>a2</i>
<i>pr</i> →		a	3
		b	1
		d	8
		d	13
		f	7
		m	5
		q	6
		<i>r</i>	

		<i>a1</i>	<i>a3</i>
<i>ps</i> →		a	A
		b	G
		c	L
		d	N
		m	B
		<i>s</i>	

$\#(A) = n, \#(B) = m$

$n * m$

$n * \log(n) + m * \log(m) + n + m$



## Merge-Join (Cont.)

- Can be used only for equi-joins and natural joins
- Each block needs to be read only once (assuming all tuples for any given value of the join attributes fit in memory)
- Thus the cost of merge join is:  
$$b_r + b_s \text{ block transfers} + \lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil \text{ seeks}$$

+ the cost of sorting if relations are unsorted.
- **hybrid merge-join:** If one relation is sorted, and the other has a secondary B<sup>+</sup>-tree index on the join attribute
  - Merge the sorted relation with the leaf entries of the B<sup>+</sup>-tree .
  - Sort the result on the addresses of the unsorted relation's tuples
  - Scan the unsorted relation in physical address order and merge with previous result, to replace addresses by the actual tuples
    - Sequential scan more efficient than random lookup



# Hash-Join

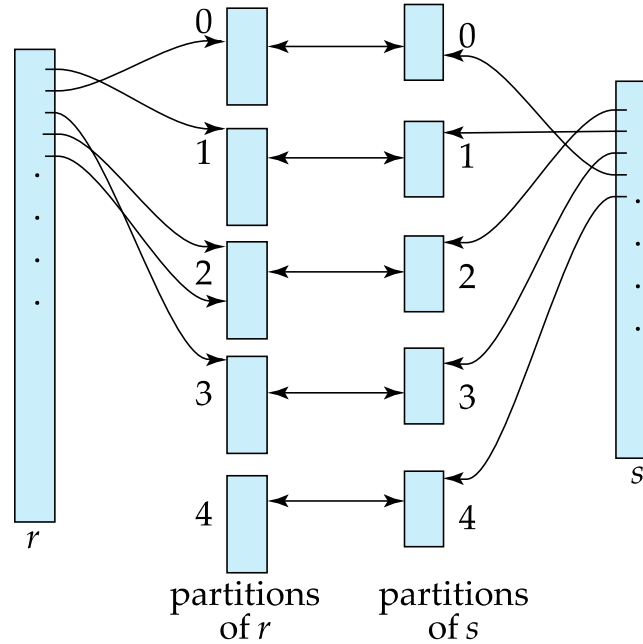
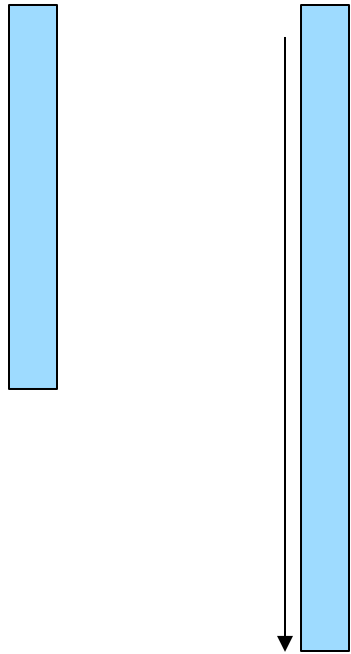
- Applicable for equi-joins and natural joins.
- A hash function  $h$  is used to partition tuples of both relations
- $h$  maps *JoinAttrs* values to  $\{0, 1, \dots, n\}$ , where *JoinAttrs* denotes the common attributes of  $r$  and  $s$  used in the natural join.
  - $r_0, r_1, \dots, r_n$  denote partitions of  $r$  tuples
    - Each tuple  $t_r \in r$  is put in partition  $r_i$  where  $i = h(t_r[\text{JoinAttrs}])$ .
  - $r_0, r_1, \dots, r_n$  denotes partitions of  $s$  tuples
    - Each tuple  $t_s \in s$  is put in partition  $s_i$ , where  $i = h(t_s[\text{JoinAttrs}])$ .
- *Note:* In book, Figure 12.10  $r_i$  is denoted as  $H_{r_i}$ ,  $s_i$  is denoted as  $H_{s_i}$  and  $n$  is denoted as  $n_h$ .





## Hash-Join (Cont.)

On  $a.X=b.Z$





# Hash-Join Algorithm

The hash-join of  $r$  and  $s$  is computed as follows.

1. Partition the relation  $s$  using hashing function  $h$ . When partitioning a relation, one block of memory is reserved as the output buffer for each partition.
2. Partition  $r$  similarly.
3. For each  $i$ :
  - (a) Load  $s_i$  into memory and build an in-memory hash index on it using the join attribute. This hash index uses a different hash function than the earlier one  $h$ .
  - (b) Read the tuples in  $r_i$  from the disk one by one. For each tuple  $t_r$  locate each matching tuple  $t_s$  in  $s_i$  using the in-memory hash index. Output the concatenation of their attributes.

Relation  $s$  is called the **build input** and  $r$  is called the **probe input**.



## Hash-Join algorithm (Cont.)

- The value  $n$  and the hash function  $h$  is chosen such that each  $s_i$  should fit in memory.
  - Typically  $n$  is chosen as  $\lceil b_s/M \rceil * f$  where  $f$  is a “**fudge factor**”, typically around 1.2
  - The probe relation partitions  $s_i$  need not fit in memory
- **Recursive partitioning** required if number of partitions  $n$  is greater than number of pages  $M$  of memory.
  - instead of partitioning  $n$  ways, use  $M - 1$  partitions for  $s$
  - Further partition the  $M - 1$  partitions using a different hash function
  - Use same partitioning method on  $r$
  - Rarely required: e.g., with block size of 4 KB, recursive partitioning not needed for relations of  $< 1$  GB with memory size of 2MB, or relations of  $< 36$  GB with memory of 12 MB

# *Other Operations (Algorithm Selection)*



# Other Operations

- **Duplicate elimination** can be implemented via hashing or sorting.
  - On sorting duplicates will come adjacent to each other, and all but one set of duplicates can be deleted.
  - *Optimization:* duplicates can be deleted during run generation as well as at intermediate merge steps in external sort-merge.
  - Hashing is similar – duplicates will come into the same bucket.
- **Projection:**
  - perform projection on each tuple
  - followed by duplicate elimination.



# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.
  - **Sorting** or **hashing** can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.
  - Optimization: **partial aggregation**
    - combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values
    - For count, min, max, sum: keep aggregate values on tuples found so far in the group.
      - When combining partial aggregate for count, add up the partial aggregates
    - For avg, keep sum and count, and divide sum by count at the end

# *Evaluation*



# Evaluation of Expressions

- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
  - **Materialization**: generate results of an expression whose inputs are relations or are already computed, **materialize** (store) it on disk. Repeat.
  - **Pipelining**: pass on tuples to parent operations even as an operation is being executed
- We study above alternatives in more detail



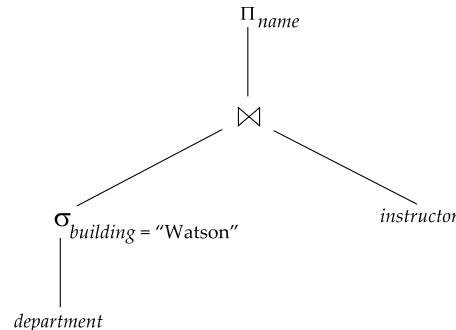


# Materialization

- **Materialized evaluation:** evaluate one operation at a time, starting at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations.
- E.g., in figure below, compute and store

$$\sigma_{building="Watson"}(department)$$

then compute the store its join with *instructor*, and finally compute the projection on *name*.





## Materialization (Cont.)

- Materialized evaluation is always applicable
- Cost of writing results to disk and reading them back can be quite high
  - Our cost formulas for operations ignore cost of writing results to disk, so
    - Overall cost = Sum of costs of individual operations + cost of writing intermediate results to disk
- **Double buffering**: use two output buffers for each operation, when one is full write it to disk while the other is getting filled
  - Allows overlap of disk writes with computation and reduces execution time



# Pipelining

- **Pipelined evaluation:** evaluate several operations simultaneously, passing the results of one operation on to the next.
- E.g., in previous expression tree, don't store result of
$$\sigma_{building="Watson"}(department)$$
  - instead, pass tuples directly to the join.. Similarly, don't store result of join, pass tuples directly to projection.
- Much cheaper than materialization: no need to store a temporary relation to disk.
- Pipelining may not always be possible – e.g., sort, hash-join.
- For pipelining to be effective, use evaluation algorithms that generate output tuples even as tuples are received for inputs to the operation.
- Pipelines can be executed in two ways: **demand driven** and **producer driven**

# *Query Optimization*



## Chapter 16: Query Optimization

**Database System Concepts, 7<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use



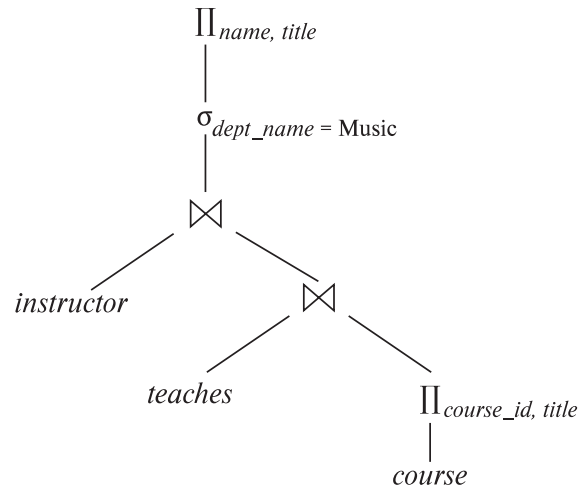
# Outline

- Introduction
- Transformation of Relational Expressions
- Catalog Information for Cost Estimation
- Statistical Information for Cost Estimation
- Cost-based optimization
- Dynamic Programming for Choosing Evaluation Plans
- Materialized views

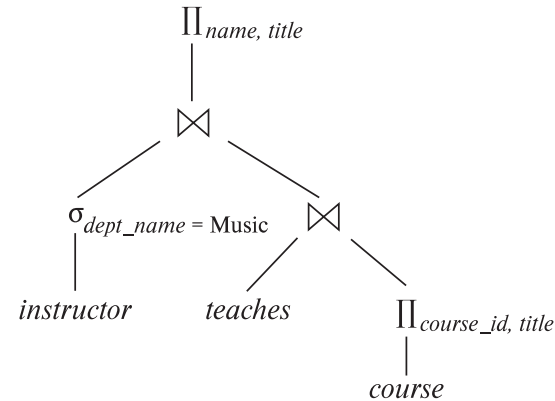


# Introduction

- Alternative ways of evaluating a given query
  - Equivalent expressions
  - Different algorithms for each operation



(a) Initial expression tree

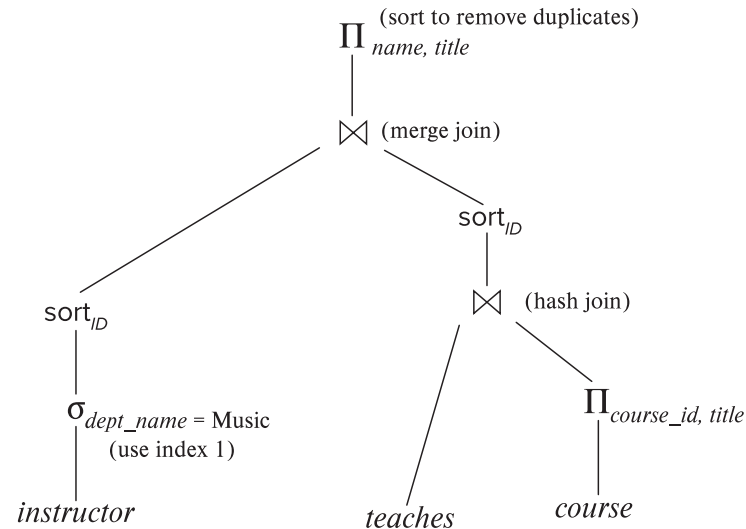


(b) Transformed expression tree



## Introduction (Cont.)

- An **evaluation plan** defines exactly what algorithm is used for each operation, and how the execution of the operations is coordinated.



- Find out how to view query execution plans on your favorite database





## Introduction (Cont.)

- Cost difference between evaluation plans for a query can be enormous
  - E.g., seconds vs. days in some cases
- Steps in **cost-based query optimization**
  1. Generate logically equivalent expressions using **equivalence rules**
  2. Annotate resultant expressions to get alternative query plans
  3. Choose the cheapest plan based on **estimated cost**
- Estimation of plan cost based on:
  - Statistical information about relations. Examples:
    - number of tuples, number of distinct values for an attribute
  - Statistics estimation for intermediate results
    - to compute cost of complex expressions
  - Cost formulae for algorithms, computed using statistics



# Viewing Query Evaluation Plans

- Most database support **explain** <query>
  - Displays plan chosen by query optimizer, along with cost estimates
  - Some syntax variations between databases
    - Oracle: **explain plan for** <query> followed by **select \* from** table (*dbms\_xplan.display*)
    - SQL Server: **set showplan\_text on**
- Some databases (e.g. PostgreSQL) support **explain analyse** <query>
  - Shows actual runtime statistics found by running the query, in addition to showing the plan
- Some databases (e.g. PostgreSQL) show cost as *f..l*
  - *f* is the cost of delivering first tuple and *l* is cost of delivering all results



# Generating Equivalent Expressions

**Database System Concepts, 7<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Transformation of Relational Expressions

- Two relational algebra expressions are said to be **equivalent** if the two expressions generate the same set of tuples on every *legal* database instance
  - Note: order of tuples is irrelevant
  - we don't care if they generate different results on databases that violate integrity constraints
- In SQL, inputs and outputs are multisets of tuples
  - Two expressions in the multiset version of the relational algebra are said to be equivalent if the two expressions generate the same multiset of tuples on every legal database instance.
- An **equivalence rule** says that expressions of two forms are equivalent
  - Can replace expression of first form by second, or vice versa



# Equivalence Rules

1. Conjunctive selection operations can be deconstructed into a sequence of individual selections.

$$\sigma_{\theta_1 \wedge \theta_2}(E) \equiv \sigma_{\theta_1}(\sigma_{\theta_2}(E))$$

2. Selection operations are commutative.

$$\sigma_{\theta_1}(\sigma_{\theta_2}(E)) \equiv \sigma_{\theta_2}(\sigma_{\theta_1}(E))$$

3. Only the last in a sequence of projection operations is needed, the others can be omitted.

$$\Pi_{L_1}(\Pi_{L_2}(\dots(\Pi_{L_n}(E))\dots)) \equiv \Pi_{L_1}(E)$$

where  $L_1 \subseteq L_2 \dots \subseteq L_n$

4. Selections can be combined with Cartesian products and theta joins.

a.  $\sigma_{\theta}(E_1 \times E_2) \equiv E_1 \bowtie_{\theta} E_2$

b.  $\sigma_{\theta_1}(E_1 \bowtie_{\theta_2} E_2) \equiv E_1 \bowtie_{\theta_1 \wedge \theta_2} E_2$



## Equivalence Rules (Cont.)

5. Theta-join operations (and natural joins) are commutative.

$$E_1 \bowtie E_2 \equiv E_2 \bowtie E_1$$

6. (a) Natural join operations are associative:

$$(E_1 \bowtie E_2) \bowtie E_3 \equiv E_1 \bowtie (E_2 \bowtie E_3)$$

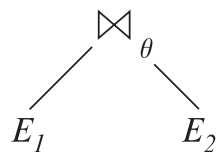
- (b) Theta joins are associative in the following manner:

$$(E_1 \bowtie_{\theta_1} E_2) \bowtie_{\theta_2 \wedge \theta_3} E_3 \equiv E_1 \bowtie_{\theta_1 \wedge \theta_3} (E_2 \bowtie_{\theta_2} E_3)$$

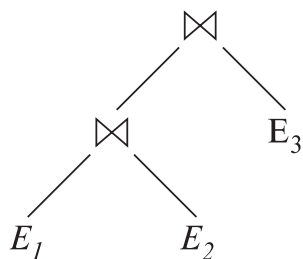
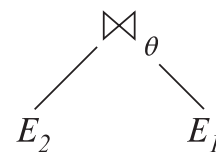
where  $\theta_2$  involves attributes from only  $E_2$  and  $E_3$ .



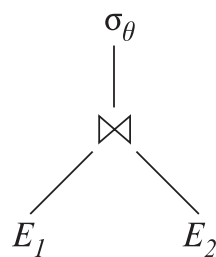
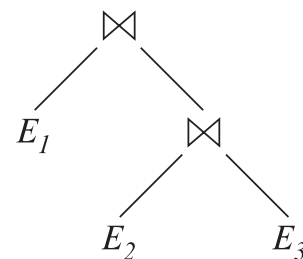
# Pictorial Depiction of Equivalence Rules



Rule 5  
 $\longleftrightarrow$

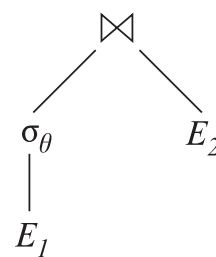


Rule 6.a  
 $\longleftrightarrow$



Rule 7.a  
 $\longleftrightarrow$

If  $\theta$  only has  
 attributes from  $E_1$





## Equivalence Rules (Cont.)

7. The selection operation distributes over the theta join operation under the following two conditions:
- (a) When all the attributes in  $\theta_0$  involve only the attributes of one of the expressions ( $E_1$ ) being joined.

$$\sigma_{\theta_0}(E_1 \bowtie_{\theta} E_2) \equiv (\sigma_{\theta_0}(E_1)) \bowtie_{\theta} E_2$$

- (b) When  $\theta_1$  involves only the attributes of  $E_1$  and  $\theta_2$  involves only the attributes of  $E_2$ .

$$\sigma_{\theta_1 \wedge \theta_2}(E_1 \bowtie_{\theta} E_2) \equiv (\sigma_{\theta_1}(E_1)) \bowtie_{\theta} (\sigma_{\theta_2}(E_2))$$





## Equivalence Rules (Cont.)

8. The projection operation distributes over the theta join operation as follows:

(a) if  $\theta$  involves only attributes from  $L_1 \cup L_2$ :

$$\Pi_{L_1 \cup L_2}(E_1 \bowtie_{\theta} E_2) \equiv \Pi_{L_1}(E_1) \bowtie_{\theta} \Pi_{L_2}(E_2)$$

(b) In general, consider a join  $E_1 \bowtie_{\theta} E_2$ .

- Let  $L_1$  and  $L_2$  be sets of attributes from  $E_1$  and  $E_2$ , respectively.
- Let  $L_3$  be attributes of  $E_1$  that are involved in join condition  $\theta$ , but are not in  $L_1 \cup L_2$ , and
- let  $L_4$  be attributes of  $E_2$  that are involved in join condition  $\theta$ , but are not in  $L_1 \cup L_2$ .

$$\Pi_{L_1 \cup L_2}(E_1 \bowtie_{\theta} E_2) \equiv \Pi_{L_1 \cup L_2}(\Pi_{L_1 \cup L_3}(E_1) \bowtie_{\theta} \Pi_{L_2 \cup L_4}(E_2))$$

Similar equivalences hold for outerjoin operations:  $\bowtie$ ,  $\ltimes$ , and  $\rhd$



## Equivalence Rules (Cont.)

9. The set operations union and intersection are commutative

$$E_1 \cup E_2 \equiv E_2 \cup E_1$$

$$E_1 \cap E_2 \equiv E_2 \cap E_1$$

(set difference is not commutative).

10. Set union and intersection are associative.

$$(E_1 \cup E_2) \cup E_3 \equiv E_1 \cup (E_2 \cup E_3)$$

$$(E_1 \cap E_2) \cap E_3 \equiv E_1 \cap (E_2 \cap E_3)$$

11. The selection operation distributes over  $\cup$ ,  $\cap$  and  $-$ .

a.  $\sigma_{\theta}(E_1 \cup E_2) \equiv \sigma_{\theta}(E_1) \cup \sigma_{\theta}(E_2)$

b.  $\sigma_{\theta}(E_1 \cap E_2) \equiv \sigma_{\theta}(E_1) \cap \sigma_{\theta}(E_2)$

c.  $\sigma_{\theta}(E_1 - E_2) \equiv \sigma_{\theta}(E_1) - \sigma_{\theta}(E_2)$

d.  $\sigma_{\theta}(E_1 \cap E_2) \equiv \sigma_{\theta}(E_1) \cap E_2$

e.  $\sigma_{\theta}(E_1 - E_2) \equiv \sigma_{\theta}(E_1) - E_2$

preceding equivalence does not hold for  $\cup$

12. The projection operation distributes over union

$$\Pi_L(E_1 \cup E_2) \equiv (\Pi_L(E_1)) \cup (\Pi_L(E_2))$$



## Equivalence Rules (Cont.)

13. Selection distributes over aggregation as below  
$$\sigma_{\theta}(G\gamma_A(E)) \equiv G\gamma_A(\sigma_{\theta}(E))$$
  
provided  $\theta$  only involves attributes in  $G$
14. a. Full outerjoin is commutative:  
$$E_1 \bowtie E_2 \equiv E_2 \bowtie E_1$$
  
b. Left and right outerjoin are not commutative, but:  
$$E_1 \bowtie E_2 \equiv E_2 \bowtie E_1$$
15. Selection distributes over left and right outerjoins as below, provided  $\theta_1$  only involves attributes of  $E_1$   
a.  $\sigma_{\theta_1}(E_1 \bowtie_{\theta} E_2) \equiv (\sigma_{\theta_1}(E_1)) \bowtie_{\theta} E_2$   
b.  $\sigma_{\theta_1}(E_1 \ltimes_{\theta} E_2) \equiv E_2 \bowtie_{\theta} (\sigma_{\theta_1}(E_1))$
16. Outerjoins can be replaced by inner joins under some conditions  
a.  $\sigma_{\theta_1}(E_1 \bowtie_{\theta} E_2) \equiv \sigma_{\theta_1}(E_1 \ltimes_{\theta} E_2)$   
b.  $\sigma_{\theta_1}(E_1 \ltimes_{\theta} E_1) \equiv \sigma_{\theta_1}(E_1 \bowtie_{\theta} E_2)$   
provided  $\theta_1$  is null rejecting on  $E_2$



## Equivalence Rules (Cont.)

Note that several equivalences that hold for joins do not hold for outerjoins

- $\sigma_{\text{year}=2017}(\text{instructor} \bowtie \text{teaches}) \neq \sigma_{\text{year}=2017}(\text{instructor} \bowtie \text{teaches})$
- Outerjoins are not associative
$$(r \bowtie s) \bowtie t \neq r \bowtie (s \bowtie t)$$
  - e.g. with  $r(A,B) = \{(1,1)\}$ ,  $s(B,C) = \{(1,1)\}$ ,  $t(A,C) = \{\}$

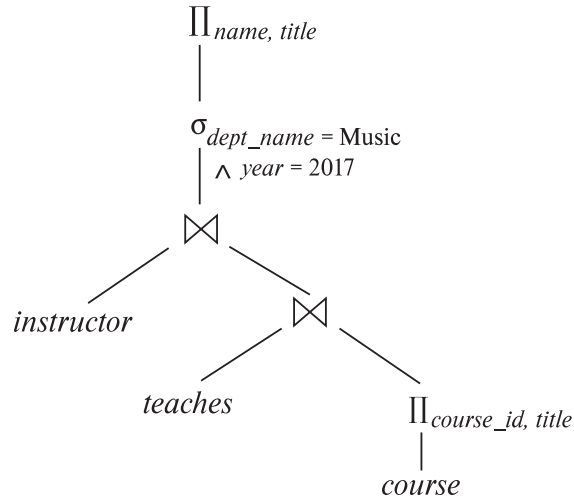


# Transformation Example: Pushing Selections

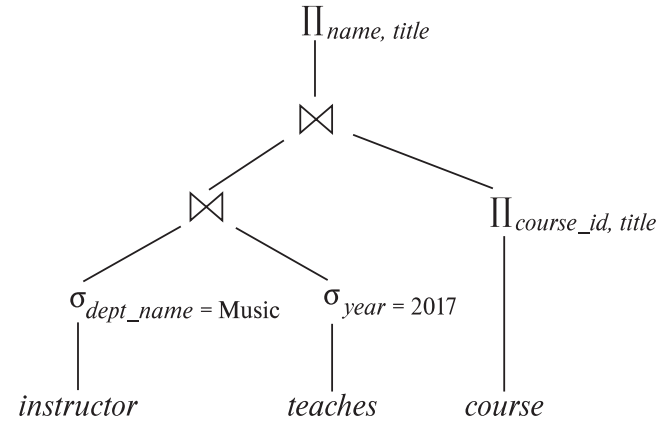
- Query: Find the names of all instructors in the Music department, along with the titles of the courses that they teach
  - $\Pi_{name, title}(\sigma_{dept\_name = 'Music'}(instructor \bowtie (teaches \bowtie \Pi_{course\_id, title}(course))))$
- Transformation using rule 7a.
  - $\Pi_{name, title}((\sigma_{dept\_name = 'Music'}(instructor)) \bowtie (teaches \bowtie \Pi_{course\_id, title}(course)))$
- Performing the selection as early as possible reduces the size of the relation to be joined.



# Multiple Transformations (Cont.)



(a) Initial expression tree



(b) Tree after multiple transformations



# Join Ordering Example

- For all relations  $r_1$ ,  $r_2$ , and  $r_3$ ,

$$(r_1 \bowtie r_2) \bowtie r_3 = r_1 \bowtie (r_2 \bowtie r_3)$$

(Join Associativity)  $\bowtie$

- If  $r_2 \bowtie r_3$  is quite large and  $r_1 \bowtie r_2$  is small, we choose

$$(r_1 \bowtie r_2) \bowtie r_3$$

so that we compute and store a smaller temporary relation.



## Join Ordering Example (Cont.)

- Consider the expression

$$\Pi_{name, title}(\sigma_{dept\_name = \text{"Music"}}(instructor) \bowtie teaches) \\ \bowtie \Pi_{course\_id, title}(course)))$$

- Could compute  $teaches \bowtie \Pi_{course\_id, title}(course)$  first, and join result with

$$\sigma_{dept\_name = \text{"Music"}}(instructor)$$

but the result of the first join is likely to be a large relation.

- Only a small fraction of the university's instructors are likely to be from the Music department
  - it is better to compute

$$\sigma_{dept\_name = \text{"Music"}}(instructor) \bowtie teaches$$

first.





# Enumeration of Equivalent Expressions

- Query optimizers use equivalence rules to **systematically** generate expressions equivalent to the given expression
- Can generate all equivalent expressions as follows:
  - Repeat
    - apply all applicable equivalence rules on every subexpression of every equivalent expression found so far
    - add newly generated expressions to the set of equivalent expressions

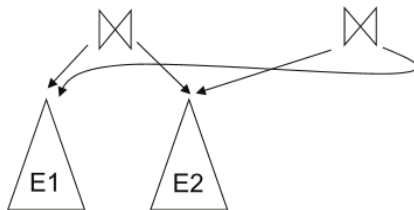
Until no new equivalent expressions are generated above

- The above approach is very expensive in space and time
  - Two approaches
    - Optimized plan generation based on transformation rules
    - Special case approach for queries with only selections, projections and joins



# Implementing Transformation Based Optimization

- Space requirements reduced by sharing common sub-expressions:
  - when E1 is generated from E2 by an equivalence rule, usually only the top level of the two are different, subtrees below are the same and can be shared using pointers
    - E.g., when applying join commutativity



- Same sub-expression may get generated multiple times
  - Detect duplicate sub-expressions and share one copy
- Time requirements are reduced by not generating all expressions
  - Dynamic programming
    - We will study only the special case of dynamic programming for join order optimization



# Cost Estimation

- Cost of each operator computer as described in Chapter 15
  - Need statistics of input relations
    - E.g., number of tuples, sizes of tuples
- Inputs can be results of sub-expressions
  - Need to estimate statistics of expression results
  - To do so, we require additional statistics
    - E.g., number of distinct values for an attribute
- More on cost estimation later



# Choice of Evaluation Plans

- Must consider the interaction of evaluation techniques when choosing evaluation plans
  - choosing the cheapest algorithm for each operation independently may not yield best overall algorithm. E.g.
    - merge-join may be costlier than hash-join, but may provide a sorted output which reduces the cost for an outer level aggregation.
    - nested-loop join may provide opportunity for pipelining
- Practical query optimizers incorporate elements of the following two broad approaches:
  1. Search all the plans and choose the best plan in a cost-based fashion.
  2. Uses heuristics to choose a plan.



# Cost-Based Optimization

- Consider finding the best join-order for  $r_1 \bowtie r_2 \bowtie \dots \bowtie r_n$ .
- There are  $(2(n-1))/(n-1)!$  different join orders for above expression. With  $n = 7$ , the number is 665280, with  $n = 10$ , the number is greater than 176 billion!
- No need to generate all the join orders. Using dynamic programming, the least-cost join order for any subset of  $\{r_1, r_2, \dots, r_n\}$  is computed only once and stored for future use.



# Statistics for Cost Estimation

**Database System Concepts, 7<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Statistical Information for Cost Estimation

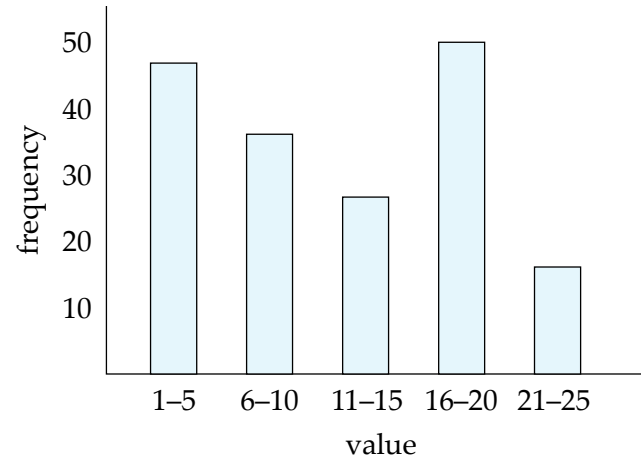
- $n_r$ : number of tuples in a relation  $r$ .
- $b_r$ : number of blocks containing tuples of  $r$ .
- $l_r$ : size of a tuple of  $r$ .
- $f_r$ : blocking factor of  $r$  — i.e., the number of tuples of  $r$  that fit into one block.
- $V(A, r)$ : number of distinct values that appear in  $r$  for attribute  $A$ ; same as the size of  $\Pi_A(r)$ .
- If tuples of  $r$  are stored together physically in a file, then:

$$b_r = \frac{n_r}{f_r}$$



# Histograms

- Histogram on attribute *age* of relation *person*



- **Equi-width** histograms
- **Equi-depth** histograms break up range such that each range has (approximately) the same number of tuples
  - E.g. (4, 8, 14, 19)
- Many databases also store  $n$  **most-frequent values** and their counts
  - Histogram is built on remaining values only





## Histograms (cont.)

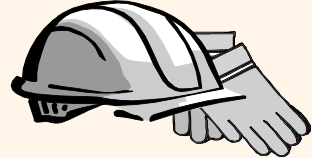
- Histograms and other statistics usually computed based on a **random sample**
- Statistics may be out of date
  - Some database require a **analyze** command to be executed to update statistics
  - Others automatically recompute statistics
    - e.g., when number of tuples in a relation changes by some percentage



# Selection Size Estimation

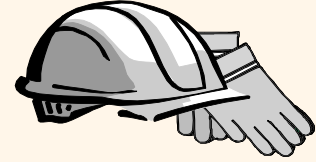
- $\sigma_{A=v}(r)$ 
  - $n_r / V(A, r)$  : number of records that will satisfy the selection
  - Equality condition on a key attribute: *size estimate* = 1
- $\sigma_{A \leq v}(r)$  (case of  $\sigma_{A \geq v}(r)$  is symmetric)
  - Let  $c$  denote the estimated number of tuples satisfying the condition.
  - If  $\min(A, r)$  and  $\max(A, r)$  are available in catalog
    - $c = 0$  if  $v < \min(A, r)$
    - $c = n_r \cdot \frac{v - \min(A, r)}{\max(A, r) - \min(A, r)}$
  - If histograms available, can refine above estimate
  - In absence of statistical information  $c$  is assumed to be  $n_r/2$ .

*Normalization (I) →*  
*Module IV*  
*Big Data, Decision Support, ... ..*



# The Evils of Redundancy

- ❖ *Redundancy* is at the root of several problems associated with relational schemas:
  - **redundant storage, insert/delete/update anomalies**
- ❖ Integrity constraints, in particular *functional dependencies*, can be used to identify schemas with such problems and to suggest refinements.
- ❖ Main refinement technique: *decomposition* (replacing ABCD with, say, AB and BCD, or ACD and ABD).
- ❖ Decomposition should be used judiciously:
  - Is there reason to decompose a relation?
  - What problems (if any) does the decomposition cause?



## Example (Contd.)

❖ Problems due to  $R \rightarrow W$  :

- Update anomaly: Can we change W in just the 1st tuple of SNLRWH?
- Insertion anomaly: What if we want to insert an employee and don't know the hourly wage for his rating?
- Deletion anomaly: If we delete all employees with rating 5, we lose the information about the wage for rating 5!

Will 2 smaller tables be better?

Wages

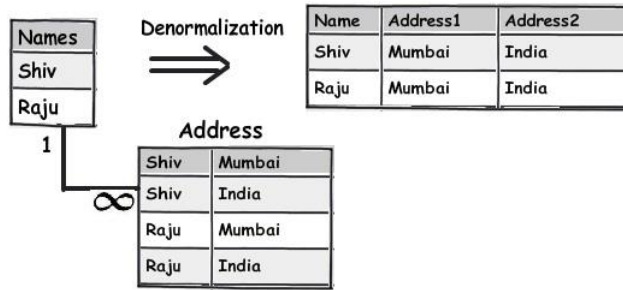
R	W
8	10
5	7

Hourly\_Emps2

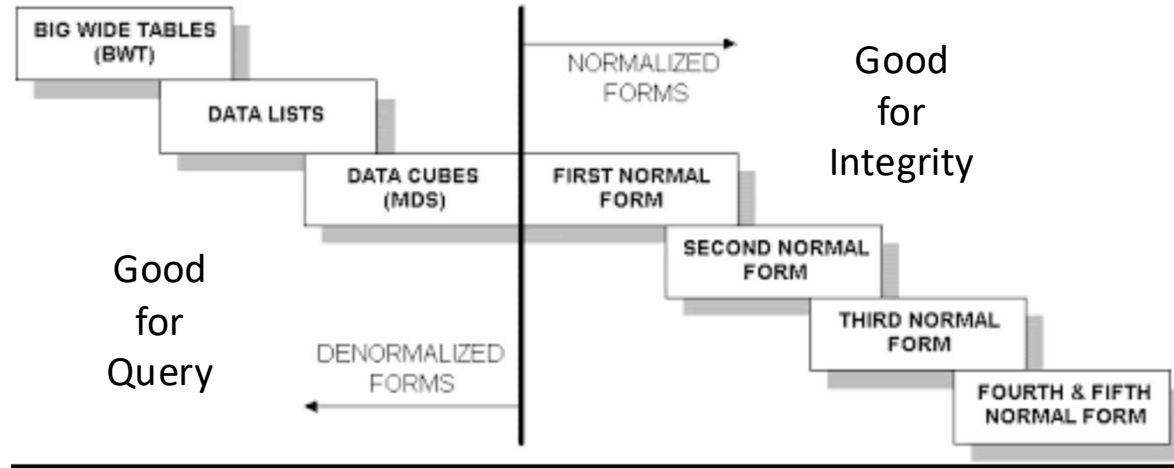
S	N	L	R	H
123-22-3666	Attishoo	48	8	40
231-31-5368	Smiley	22	8	30
131-24-3650	Smethurst	35	5	30
434-26-3751	Guldu	35	5	32
612-67-4134	Madayan	35	8	40

S	N	L	R	W	H
123-22-3666	Attishoo	48	8	10	40
231-31-5368	Smiley	22	8	10	30
131-24-3650	Smethurst	35	5	7	30
434-26-3751	Guldu	35	5	7	32
612-67-4134	Madayan	35	8	10	40

# Wide Flat Tables



- Improve query performance by precomputing and saving:
  - JOINS
  - Aggregation
  - Derived/computed columns
- One of the primary strength of the relational model is maintaining “integrity” when applications create, update and delete data. This relies on:
  - The core capabilities of the relational model, e.g. constraints.
  - A well-design database (We will cover a formal definition – “normalization” in more detail later.
- Data models that are well designed for integrity are very bad for read only analysis queries. We will build and analyze wide flat tables as part of the analysis tasks in HW3, HW4 as projects.





# Overview of Normalization



# Features of Good Relational Designs

- Hypothetically, assume our schema originally had one relation that provided information about faculty and departments.
- Suppose we combine *instructor* and *department* into *in\_dep*, which represents the natural join on the relations *instructor* and *department*

ID	name	salary	dept_name	building	budget
22222	Einstein	95000	Physics	Watson	70000
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000
83821	Brandt	92000	Comp. Sci.	Taylor	100000
15151	Mozart	40000	Music	Packard	80000
33456	Gold	87000	Physics	Watson	70000
76543	Singh	80000	Finance	Painter	120000

Instructor(id, name, salary, dept\_name)  
Department(dept\_name, building, budget)

- There is repetition of information
- Need to use null values (if we add a new department with no instructors)





# Decomposition

- The only way to avoid the repetition-of-information problem in the *in\_dep* schema is to decompose it into two schemas – *instructor* and *department* schemas.
- Not all decompositions are good. Suppose we decompose

*employee*(*ID*, *name*, *street*, *city*, *salary*)

into

*employee1* (*ID*, *name*)

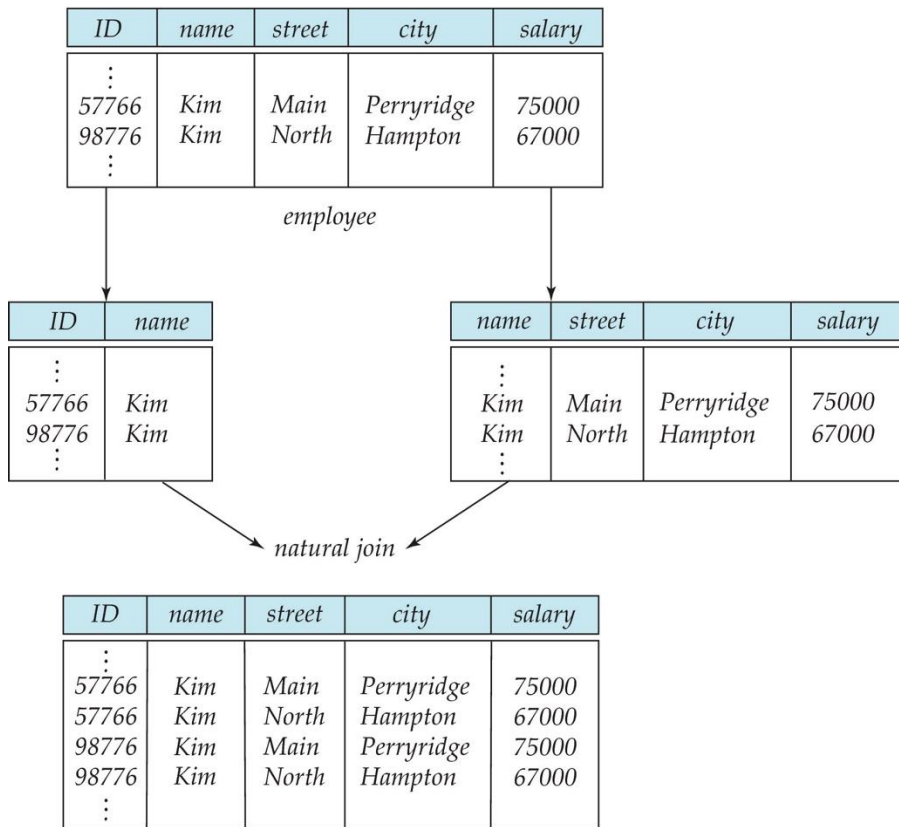
*employee2* (*name*, *street*, *city*, *salary*)

The problem arises when we have two employees with the same name

- The next slide shows how we lose information -- we cannot reconstruct the original *employee* relation -- and so, this is a **lossy decomposition**.



# Lossy Decomposition





# Normalization Theory

- Decide whether a particular relation  $R$  is in “good” form.
- In the case that a relation  $R$  is not in “good” form, decompose it into set of relations  $\{R_1, R_2, \dots, R_n\}$  such that
  - Each relation is in good form
  - The decomposition is a lossless decomposition
- Our theory is based on:
  - Functional dependencies
  - Multivalued dependencies



# Functional Dependencies

- There are usually a variety of constraints (rules) on the data in the real world.
- For example, some of the constraints that are expected to hold in a university database are:
  - Students and instructors are uniquely identified by their ID.
  - Each student and instructor has only one name.
  - Each instructor and student is (primarily) associated with only one department.
  - Each department has only one value for its budget, and only one associated building.



## Functional Dependencies (Cont.)

- An instance of a relation that satisfies all such real-world constraints is called a **legal instance** of the relation;
- A legal instance of a database is one where all the relation instances are legal instances
- Constraints on the set of legal relations.
- Require that the value for a certain set of attributes determines uniquely the value for another set of attributes.
- A functional dependency is a generalization of the notion of a *key*.



# Functional Dependencies Definition

- Let  $R$  be a relation schema

$$\alpha \subseteq R \text{ and } \beta \subseteq R$$

- The **functional dependency**

$$\alpha \rightarrow \beta$$

**holds on**  $R$  if and only if for any legal relations  $r(R)$ , whenever any two tuples  $t_1$  and  $t_2$  of  $r$  agree on the attributes  $\alpha$ , they also agree on the attributes  $\beta$ . That is,

$$t_1[\alpha] = t_2[\alpha] \Rightarrow t_1[\beta] = t_2[\beta]$$

- Example: Consider  $r(A,B)$  with the following instance of  $r$ .

1	4
1	5
3	7

- On this instance,  $B \rightarrow A$  hold;  $A \rightarrow B$  does **NOT** hold,



# Closure of a Set of Functional Dependencies

- Given a set  $F$  set of functional dependencies, there are certain other functional dependencies that are logically implied by  $F$ .
  - If  $A \rightarrow B$  and  $B \rightarrow C$ , then we can infer that  $A \rightarrow C$
  - etc.
- The set of **all** functional dependencies logically implied by  $F$  is the **closure** of  $F$ .
- We denote the *closure* of  $F$  by  $F^+$ .



# Keys and Functional Dependencies

- $K$  is a superkey for relation schema  $R$  if and only if  $K \twoheadrightarrow R$
- $K$  is a candidate key for  $R$  if and only if
  - $K \twoheadrightarrow R$ , and
  - for no  $a \subset K$ ,  $a \twoheadrightarrow R$
- Functional dependencies allow us to express constraints that cannot be expressed using superkeys. Consider the schema:

*in\_dep* (ID, name, salary, dept\_name, building, budget).

We expect these functional dependencies to hold:

*dept\_name*  $\twoheadrightarrow$  *building*

*ID*  $\twoheadrightarrow$  *building*

but would not expect the following to hold:

*dept\_name*  $\twoheadrightarrow$  *salary*

DFF Note:

- In the current data:  $ID \rightarrow dept\_name \rightarrow building$
- But
- In the schema, *dept\_name* may be NULL..





# Use of Functional Dependencies

- We use functional dependencies to:
  - To test relations to see if they are legal under a given set of functional dependencies.
    - ▶ If a relation  $r$  is legal under a set  $F$  of functional dependencies, we say that  $r$  **satisfies**  $F$ .
  - To specify constraints on the set of legal relations
    - ▶ We say that  $F$  **holds on**  $R$  if all legal relations on  $R$  satisfy the set of functional dependencies  $F$ .
- Note: A specific instance of a relation schema may satisfy a functional dependency even if the functional dependency does not hold on all legal instances.
  - For example, a specific instance of *instructor* may, by chance, satisfy  $name \rightarrow ID$ .



# Normal Forms



# Boyce-Codd Normal Form

- A relation schema  $R$  is in BCNF with respect to a set  $F$  of functional dependencies if for all functional dependencies in  $F^+$  of the form

$$\alpha \rightarrow \beta$$

where  $\alpha \subseteq R$  and  $\beta \subseteq R$ , at least one of the following holds:

- $\alpha \rightarrow \beta$  is trivial (i.e.,  $\beta \subseteq \alpha$ )
- $\alpha$  is a superkey for  $R$

## DFF Note:

- The theoretical treatment is not conveying the practical intent.
- If  $\alpha$  is a superkey, I can set a primary key/unique constraint on  $\alpha$ .
- Consider tables with address info in the rows.



# Decomposing a Schema into BCNF

- Let  $R$  be a schema  $R$  that is not in BCNF. Let  $\alpha \rightarrow \beta$  be the FD that causes a violation of BCNF.
- We decompose  $R$  into:
  - $(\alpha \cup \beta)$
  - $(R - (\beta - \alpha))$
- In our example of *in\_dep*,
  - $\alpha = \text{dept\_name}$
  - $\beta = \text{building, budget}$and *in\_dep* is replaced by
  - $(\alpha \cup \beta) = (\text{dept\_name, building, budget})$
  - $(R - (\beta - \alpha)) = (\text{ID, name, dept\_name, salary})$

DFF Note – again this is baffling

- $R = (\text{id, name\_last, name\_first, street, city, state, zipcode})$
- $\alpha \rightarrow \beta$  means  $\text{zipcode} \rightarrow (\text{city, state})$
- $(\alpha \cup \beta) = (\text{zipcode, city, state})$ , which we call Address
- $R - (\beta - \alpha) = R - \beta + \alpha = (\text{id, name\_last, name\_first, zipcode})$ , which we call Person
- Setting primary key ID in Address and zipcode on Address preserves the dependency and is lossless.

Note:

- This is an example only.
- Zip code does not imply city or state in real world.



# BCNF and Dependency Preservation

- It is not always possible to achieve both BCNF and dependency preservation

- Consider a schema:

*dept\_advisor(s\_ID, i\_ID, department\_name)*

- With function dependencies:

$i\_ID \rightarrow dept\_name$

$s\_ID, dept\_name \rightarrow i\_ID$

- *dept\_advisor* is not in BCNF

- *i\_ID* is not a superkey.

- Any decomposition of *dept\_advisor* will not include all the attributes in

$s\_ID, dept\_name \rightarrow i\_ID$

- Thus, the composition is NOT be dependency preserving



# Third Normal Form

- A relation schema  $R$  is in **third normal form (3NF)** if for all:

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- $\alpha \rightarrow \beta$  is trivial (i.e.,  $\beta \in \alpha$ )
- $\alpha$  is a superkey for  $R$
- Each attribute  $A$  in  $\beta - \alpha$  is contained in a candidate key for  $R$ .  
(**NOTE:** each attribute may be in a different candidate key)
- If a relation is in BCNF it is in 3NF (since in BCNF one of the first two conditions above must hold).
- Third condition is a minimal relaxation of BCNF to ensure dependency preservation (will see why later).



## 3NF Example

- Consider a schema:

*dept\_advisor(s\_ID, i\_ID, dept\_name)*

- With function dependencies:

$i\_ID \rightarrow dept\_name$

$s\_ID, dept\_name \rightarrow i\_ID$

- Two candidate keys =  $\{s\_ID, dept\_name\}, \{s\_ID, i\_ID\}$
- We have seen before that *dept\_advisor* is **not** in BCNF
- *R*, however, is in 3NF
  - $s\_ID, dept\_name$  is a superkey
  - $i\_ID \rightarrow dept\_name$  and  $i\_ID$  is NOT a superkey, but:
    - ▶  $\{dept\_name\} - \{i\_ID\} = \{dept\_name\}$  and
    - ▶  $dept\_name$  is contained in a candidate key



# Comparison of BCNF and 3NF

- Advantages to 3NF over BCNF. It is always possible to obtain a 3NF design without sacrificing losslessness or dependency preservation.
- Disadvantages to 3NF.
  - We may have to use null values to represent some of the possible meaningful relationships among data items.
  - There is the problem of repetition of information.





# Goals of Normalization

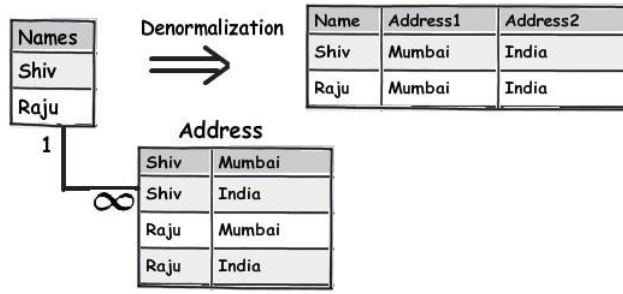
- Let  $R$  be a relation scheme with a set  $F$  of functional dependencies.
- Decide whether a relation scheme  $R$  is in “good” form.
- In the case that a relation scheme  $R$  is not in “good” form, need to decompose it into a set of relation scheme  $\{R_1, R_2, \dots, R_n\}$  such that:
  - Each relation scheme is in good form
  - The decomposition is a lossless decomposition
  - Preferably, the decomposition should be dependency preserving.



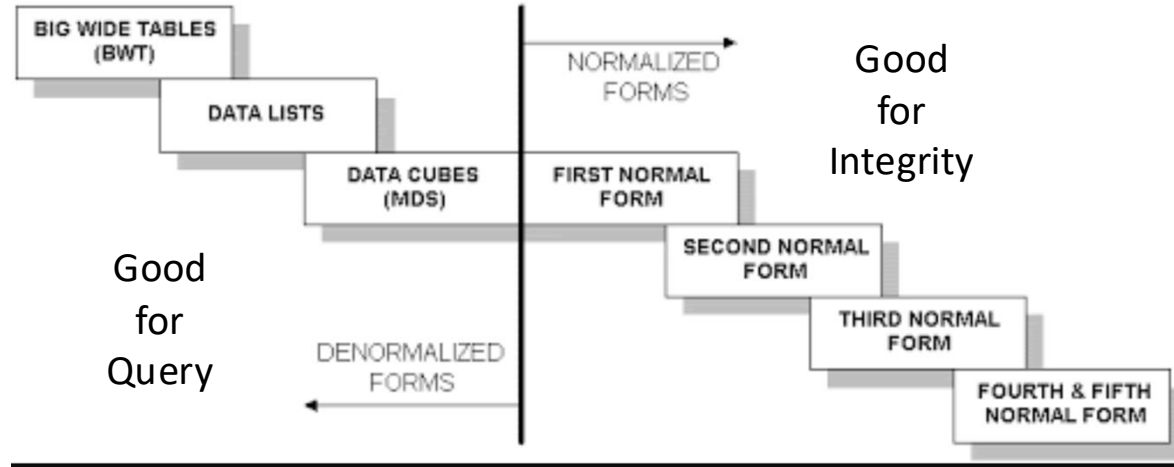
# Denormalization for Performance

- ❑ May want to use non-normalized schema for performance
- ❑ For example, displaying *prereqs* along with *course\_id*, and *title* requires join of *course* with *prereq*
- ❑ Alternative 1: Use denormalized relation containing attributes of *course* as well as *prereq* with all above attributes
  - ❑ faster lookup
  - ❑ extra space and extra execution time for updates
  - ❑ extra coding work for programmer and possibility of error in extra code
- ❑ Alternative 2: use a materialized view defined a *course* ⋈ *prereq*
  - ❑ Benefits and drawbacks same as above, except no extra coding work for programmer and avoids possible errors

# Wide Flat Tables



- Improve query performance by precomputing and saving:
  - JOINS
  - Aggregation
  - Derived/computed columns
- One of the primary strength of the relational model is maintaining “integrity” when applications create, update and delete data. This relies on:
  - The core capabilities of the relational model, e.g. constraints.
  - A well-design database (We will cover a formal definition – “normalization” in more detail later.
- Data models that are well designed for integrity are very bad for read only analysis queries. We will build and analyze wide flat tables as part of the analysis tasks in HW3, HW4 as projects.





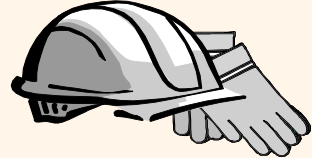
# First Normal Form

- Domain is **atomic** if its elements are considered to be indivisible units
  - Examples of non-atomic domains:
    - ▶ Set of names, composite attributes
    - ▶ Identification numbers like CS 101 that can be broken up into parts
- A relational schema R is in **first normal form** if the domains of all attributes of R are atomic
- Non-atomic values complicate storage and encourage redundant (repeated) storage of data
  - Example: Set of accounts stored with each customer, and set of owners stored with each account
  - We assume all relations are in first normal form (and revisit this in Chapter 22: Object Based Databases)



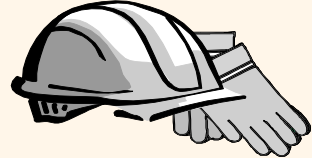
## First Normal Form (Cont.)

- Atomicity is actually a property of how the elements of the domain are used.
  - Example: Strings would normally be considered indivisible
  - Suppose that students are given roll numbers which are strings of the form *CS0012* or *EE1127*
  - If the first two characters are extracted to find the department, the domain of roll numbers is not atomic.
  - Doing so is a bad idea: leads to encoding of information in application program rather than in the database.



## Reasoning About FDs

- ❖ Given some FDs, we can usually infer additional FDs:
  - $ssn \rightarrow did, did \rightarrow lot$  implies  $ssn \rightarrow lot$
- ❖ An FD  $f$  is implied by a set of FDs  $F$  if  $f$  holds whenever all FDs in  $F$  hold.
  - $F^+ =$  *closure of  $F$*  is the set of all FDs that are implied by  $F$ .
- ❖ Armstrong's Axioms ( $X, Y, Z$  are sets of attributes):
  - Reflexivity: If  $X \subseteq Y$ , then  $Y \rightarrow X$
  - Augmentation: If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$  for any  $Z$
  - Transitivity: If  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$
- ❖ These are *sound* and *complete* inference rules for FDs!



## Lossless Decomposition

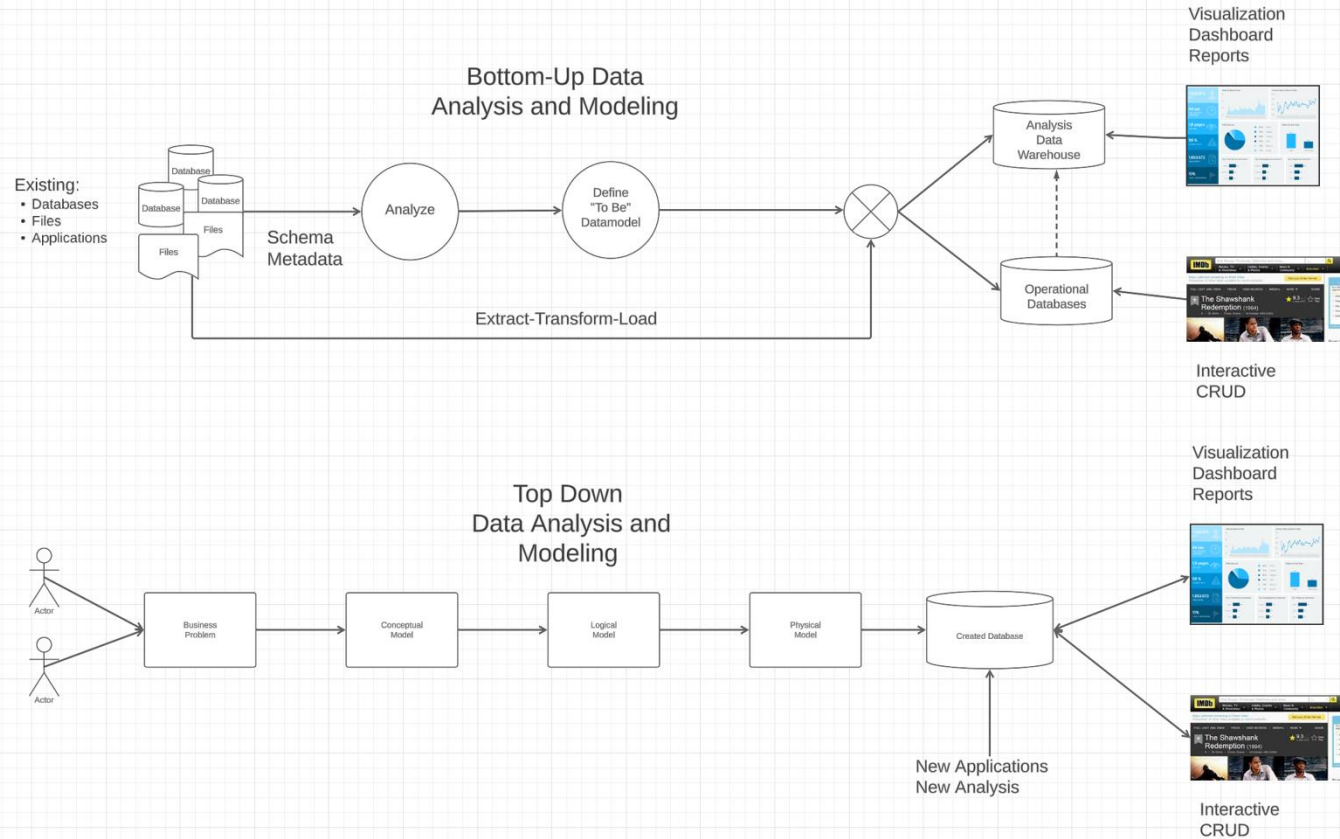
- ❖ Let  $R$  be a relation schema and let  $R_1$  and  $R_2$  form a decomposition of  $R$ .
  - That is  $R = R_1 \cup R_2$
  - *Note: This notation is confusing. This is a statement about the schema, not about the data in relations.*
- ❖ We say that the decomposition is a **lossless decomposition** if there is no loss of information by replacing  $R$  with the two relation schemas  $R_1 \cup R_2$
- ❖ Formally,
$$\Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$$
- ❖ And, conversely a decomposition is lossy if
$$r \subset \Pi_{R_1}(r) \bowtie \Pi_{R_2}(r) = r$$

# *HW3 and HW4, “The Project” Discussion*

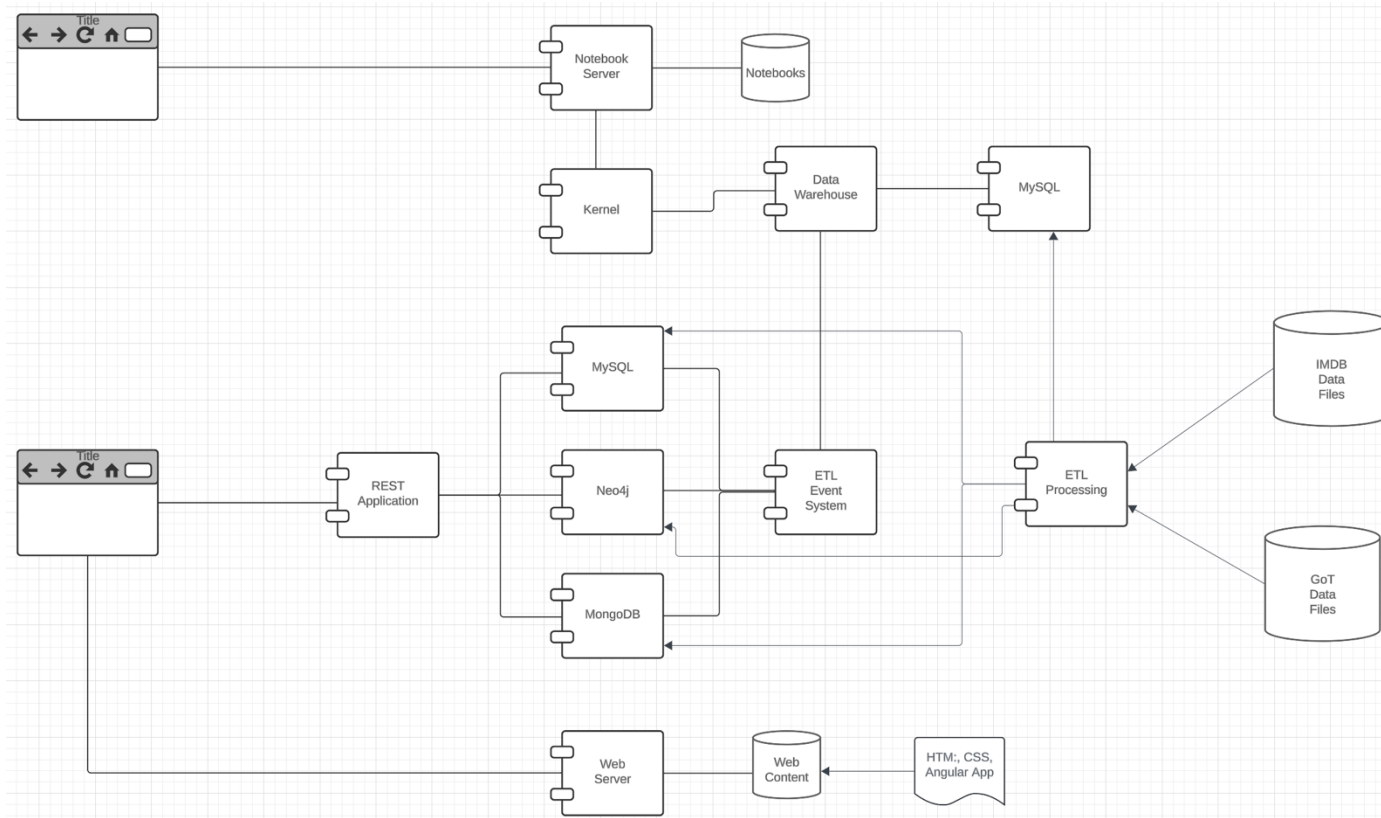


# *Overview*

# Vision

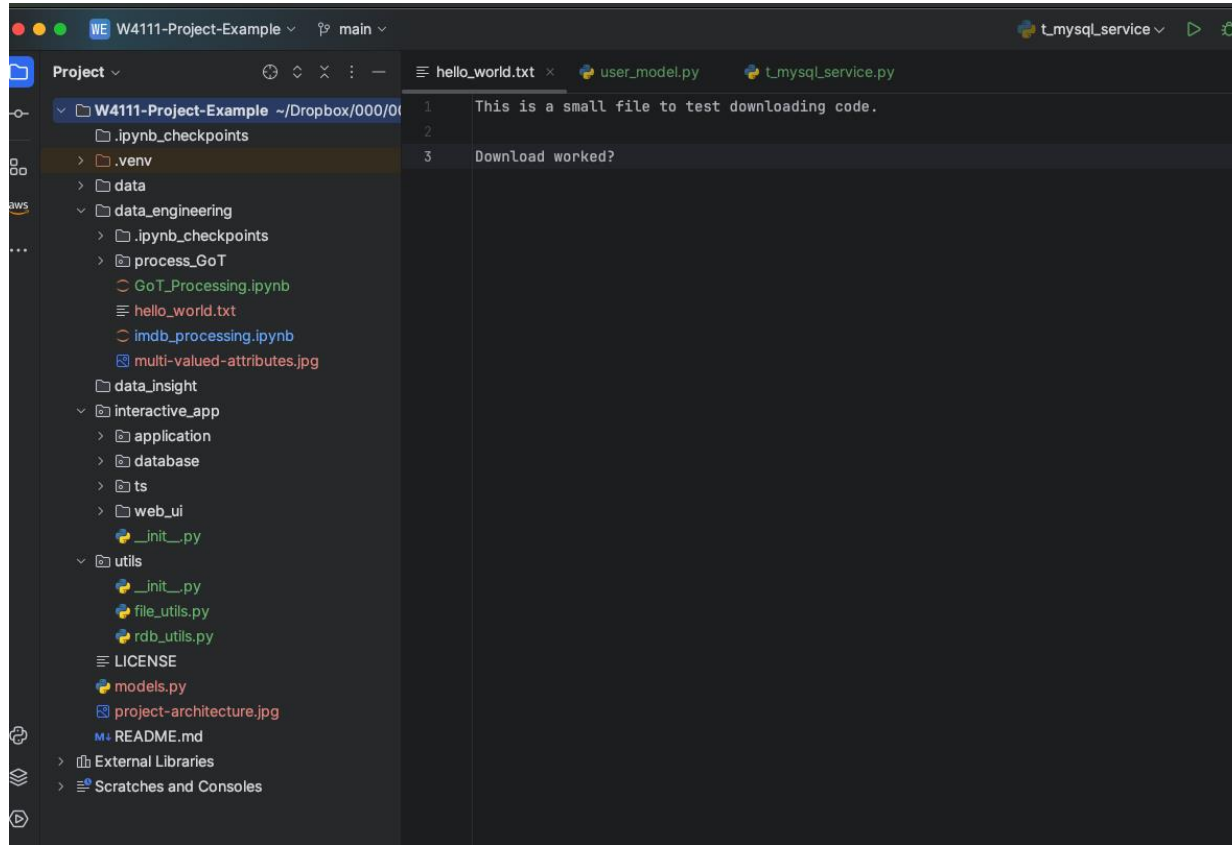


# Vision



- Both programming and non-programming implement a data engineering Jupyter notebook.
- The programming track builds a simple REST application/microservice for transformed data.
- Non-programming implements more complex transformation, and queries for visualization.
- I will provide starter projects with examples.

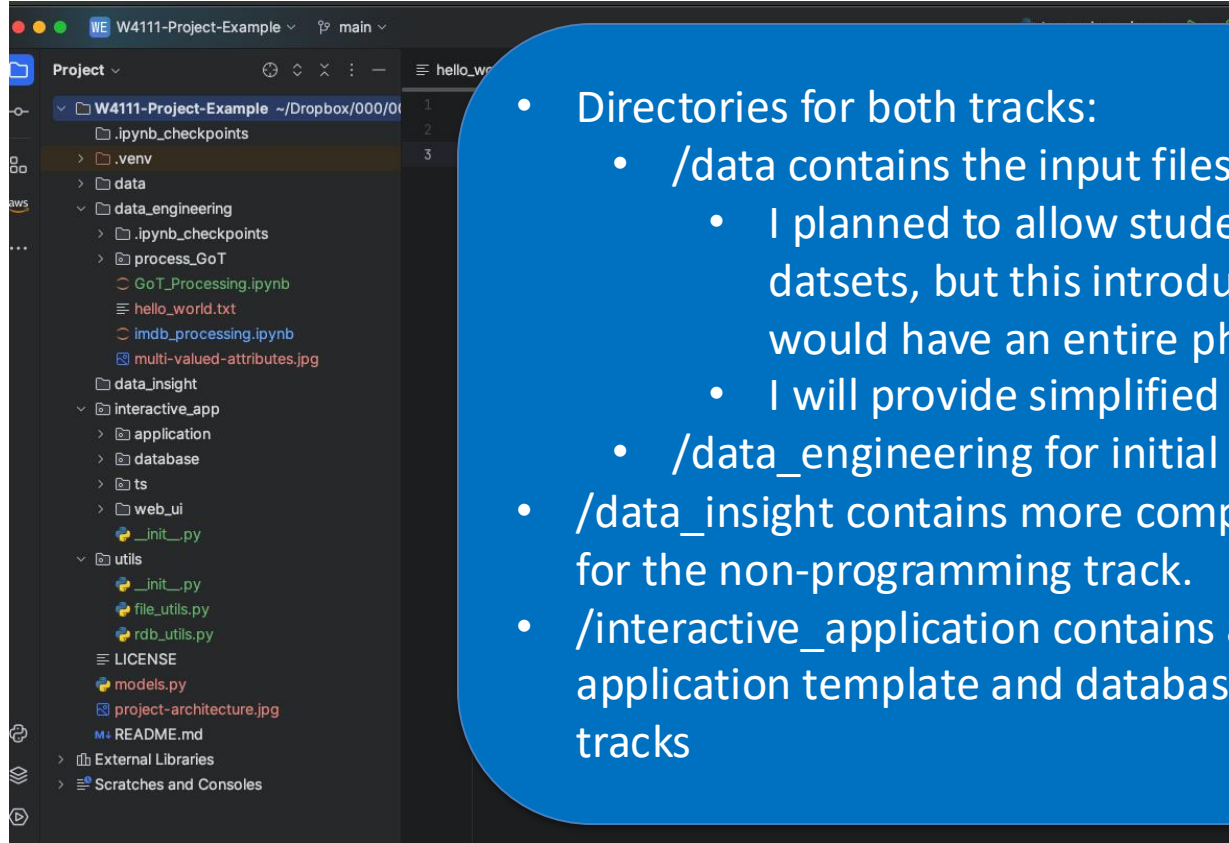
# Hw3, HW4 and Project



The screenshot shows a code editor interface with a dark theme. The top bar displays the project name 'W4111-Project-Example' and the file 'main'. The left sidebar shows the project structure, including folders like '.ipynb\_checkpoints', '.venv', 'data', 'data\_engineering', 'data\_insight', 'interactive\_app', 'utils', and files like 'hello\_world.txt', 'user\_model.py', 't\_mysql\_service.py', 'LICENSE', 'models.py', 'project-architecture.jpg', and 'README.md'. The main editor area shows the content of 'hello\_world.txt' with three lines of text:

```
1 This is a small file to test downloading code.  
2  
3 Download worked?
```

# Hw3, HW4 and Project



- Directories for both tracks:
  - /data contains the input files in CSV and JSON formats
    - I planned to allow students to choose their own datasets, but this introduces too much complexity. We would have an entire phase of “Is this a goo dataset?”
    - I will provide simplified IMDB and GoT data.
  - /data\_engineering for initial ETL
- /data\_insight contains more complex queries and visualization for the non-programming track.
- /interactive\_application contains a simple web UI, REST application template and database schema for the programming tracks

# *REST*

# Data Modeling Concepts and REST

Almost any data model has the same core concepts:

- Types and instances:
  - Entity Type: A definition of a type of thing with properties and relationships.
  - Entity Instance: A specific instantiation of the Entity Type
  - Entity Set Instance: An Entity Type that:
    - Has properties and relationships like any entity, but ...
    - Has at least one *special relationship* – ***contains***.
- Operations, minimally CRUD, that manipulate entity types and instances:
  - Create
  - Retrieve
  - Update
  - Delete
  - Reference/Identify/... ..
  - Host/database/table/pk



## What is REST architecture?

REST stands for REpresentational State Transfer. REST is web standards based architecture and uses HTTP Protocol. It revolves around resource where every component is a resource and a resource is accessed by a common interface using HTTP standard methods. REST was first introduced by Roy Fielding in 2000.

In REST architecture, a REST Server simply provides access to resources and REST client accesses and modifies the resources. Here each resource is identified by URIs/ global IDs. REST uses various representation to represent a resource like text, JSON, XML. JSON is the most popular one.

## HTTP methods

Following four HTTP methods are commonly used in REST based architecture.

- **GET** – Provides a read only access to a resource.
- **POST** – Used to create a new resource.
- **DELETE** – Used to remove a resource.
- **PUT** – Used to update a existing resource or create a new resource.

## Introduction to RESTful web services

A web service is a collection of open protocols and standards used for exchanging data between applications or systems. Software applications written in various programming languages and running on various platforms can use web services to exchange data over computer networks like the Internet in a manner similar to inter-process communication on a single computer. This interoperability (e.g., between Java and Python, or Windows and Linux applications) is due to the use of open standards.

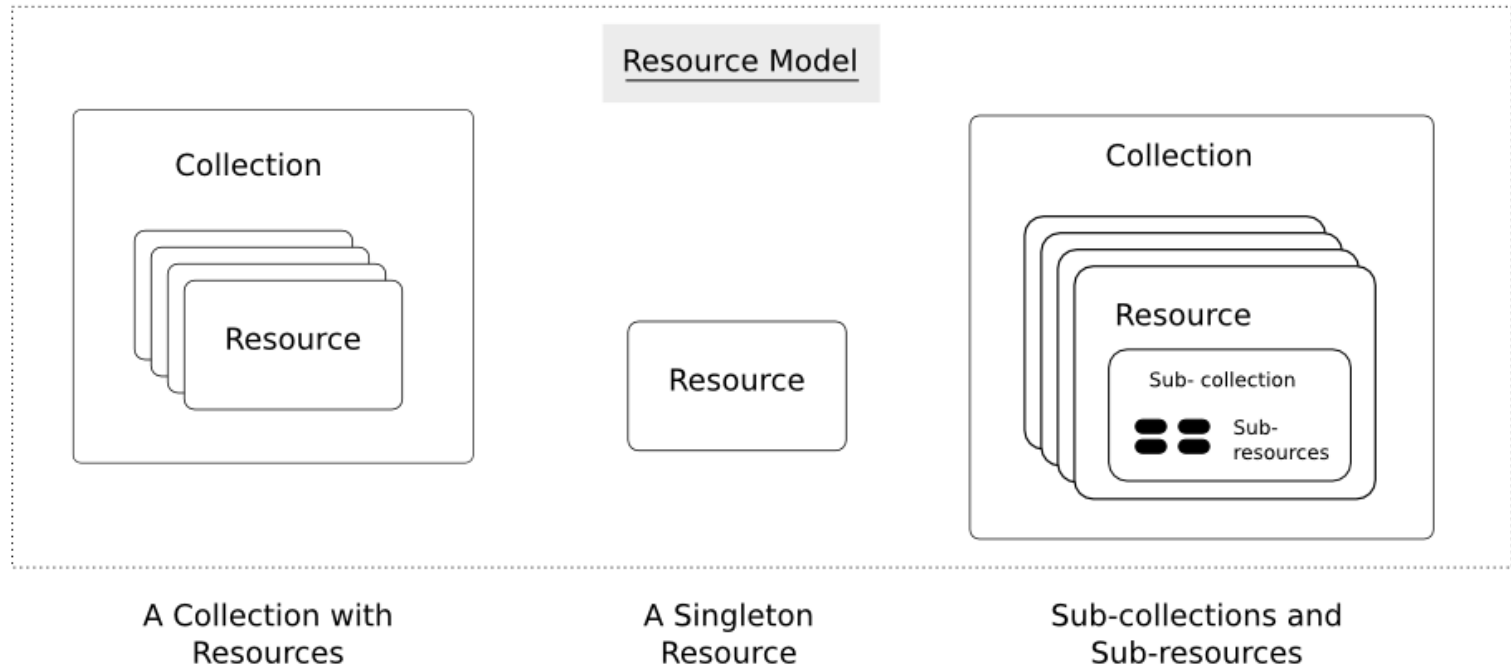
Web services based on REST Architecture are known as RESTful web services. These webservices uses HTTP methods to implement the concept of REST architecture. A RESTful web service usually defines a URI, Uniform Resource Identifier a service, provides resource representation such as JSON and set of HTTP Methods.

## Creating RESTful Webservice

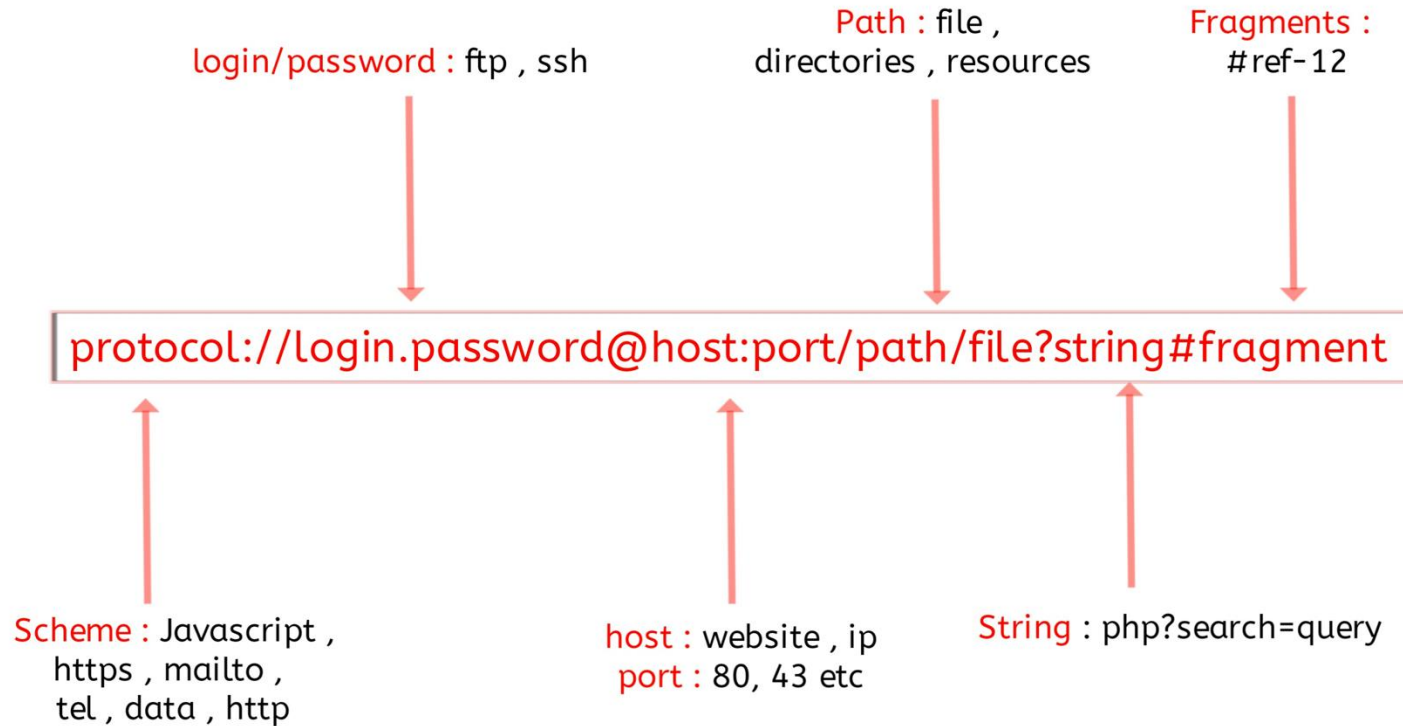
In next chapters, we'll create a webservice say user management with following functionalities –

Sr.No.	URI	HTTP Method	POST body	Result
1	/UserService/users	GET	empty	Show list of all the users.
2	/UserService/addUser	POST	JSON String	Add details of new user.
3	/UserService/getUser/:id	GET	empty	Show details of a user.

# REST and Resources



# URLs



jdbc:mysql://columbia-examples.ckkqqktwkcji.us-east-1.rds.amazonaws.com:3306

GET [http://localhost:5001/f23\\_imdb\\_clean/name\\_basics/nm0000158](http://localhost:5001/f23_imdb_clean/name_basics/nm0000158)

GET [http://localhost:5001/f23\\_imdb\\_clean/name\\_basics?deathYear=2023&birthyear=1960](http://localhost:5001/f23_imdb_clean/name_basics?deathYear=2023&birthyear=1960)

```
select * from f23_imdb_clean.name_basics where  
    deathYear=2023 AND birthyear=1960
```

PUT [http://localhost:5001/f23\\_imdb\\_clean/name\\_basics ?deathYear=2023&birthyear=1960](http://localhost:5001/f23_imdb_clean/name_basics ?deathYear=2023&birthyear=1960)

```
Body {'primaryName': 'Does not matter cause is dead.'}
```

```
update f23_imdb_clean.name_basics  
    set  
    where deathYear=2023 AND birthyear=1960
```