

## ER Diagram

- Many to Many

A can have multiple B and B can have multiple A



- One to Many

A can have multiple B but B can only have one A



- One to One

A can have only one B and B can have only one A



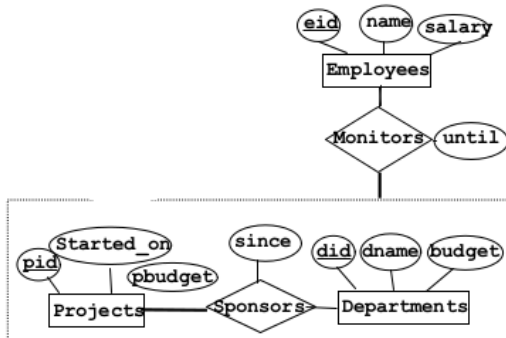
- At least one

Bold arrow specify there is at least one element.

	A have at least one
	A have exactly one

- Aggregation

Allows us to treat a relationship set R as an entity set so that R can participate in other relationships



## Relational algebra

Selection	$\sigma_{condition}(Element)$
Projection	$\pi_{attributes}(Element)$
Renaming	$\rho(R(A_1, \dots, A_n), R_{alias}(B_1, \dots, B_n))$
Cross product	$\times$
Join	$\bowtie$
Division	$\setminus$
Intersection	$\cap$
Union	$\cup$
Set different	$-$

- Condition/Theta Join  $R_{out} = R_{in1} \bowtie_{condition} R_{in2} = \sigma_{condition}(R_{in1} \times R_{in2})$
- Equi Join:  $R_{out} = R_{in1} \bowtie_{R_{in1}.a_1 = R_{in2}.b_1, \dots, R_{in1}.a_n = R_{in2}.b_n} R_{in2}$  Condition join where condition contains ONLY equalities
- Natural Join: Equijoin on all common attribute

## Sql

### Datatype

Char(n)	A character string of fixed length n
---------	--------------------------------------

VarChar(n)	Denotes a string of up to n charaters
INT or INTEGER	An integer
SHORTINT	Smaller integer
FLOAT or REAL	Float number
DOUBLE PRECISION	Double
DATE	Date format YYYY-MM-DD
TIME	Time format: hh:mm:ss

#### Table operations

```
--Create table
CREATE TABLE Students
(
    id INT NOT NULL,
    name VARCHAR(20),
    login CHAR(10),
    major VARCHAR(20) DEFAULT 'undefined',
    school_id INT,
    PRIMARY KEY(id),
    FOREIGN KEY(school_id) REFERENCES School(id)
)

--Drop table
DROP TABLE Students

--Alter table
ALTER TABLE Students ADD COLUMN firstyear:integer
```

#### Row operation

```
--INSERT
INSERT INTO Students (id, name, faculty) VALUES (8908998, 'Dupont', 'Science')

--Delete
DELETE FROM Students WHERE id = 0894984

--Update
UPDATE Students SET faculty = 'Arts' WHERE id = 9849849
```

#### Trigger

```
CREATE TRIGGER updateSkater
AFTER DELETE ON Skaters
REFERENCING OLD TABLE AS DeletedSkaters
FOR EACH STATEMENT
INSERT
INTO StatisticsTable(ModTable, ModType, Count)
SELECT 'Skaters', 'delete', COUNT(*)
FROM DeletedSkaters

Use begin/end to encapsulate more than one
action
FOR EACH ROW/STATEMENT
WHEN ...
BEGIN ATOMIC
do 1thing;
do 2nd thing;
END
```

#### XML

WHAT DAFUQ?!!@#!@#!?

#### DTD

<pre>&lt;!DOCTYPE DiscoverTheWorld [ &lt;!ELEMENT DiscoverTheWorld (tour*,reservation*)&gt; &lt;!ELEMENT tour (type, start-date, duration, price) &gt; &lt;!ELEMENT reservation (cname, caddress, cost, special*)&gt; &lt;!ATTLIST tour TourId ID #REQUIRED &gt; &lt;!ATTLIST reservation ResID ID #REQUIRED TourID IDREF #REQUIRED&gt; &lt;!ELEMENT type (#PCDATA) &gt; &lt;!ELEMENT start-date (#PCDATA) &gt; &lt;!ELEMENT duration (#PCDATA) &gt; &lt;!ELEMENT price (#PCDATA) &gt; &lt;!ELEMENT cname (#PCDATA) &gt; &lt;!ELEMENT caddress (#PCDATA) &gt; &lt;!ELEMENT cost (#PCDATA) &gt;</pre>	<pre>&lt;!DOCTYPE Politics [ &lt;!ELEMENT Politics (Politician*, Province*)&gt; &lt;!ELEMENT Politician ((CurrentMayor   CurrentMop)?, address?)&gt; &lt;!ELEMENT CurrentMayor (since?)&gt; &lt;!ELEMENT CurrentMoP (since?)&gt; &lt;!ELEMENT Province (City+, Riding+, population?)&gt; &lt;!ELEMENT City (population?)&gt; &lt;!ELEMENT Riding (population?)&gt; &lt;!ELEMENT address (#PCDATA) &gt; &lt;!ELEMENT since (#PCDATA) &gt; &lt;!ELEMENT population (#PCDATA) &gt; &lt;!ATTLIST Politician pname ID REQUIRED website CDATA IMPLIED friends IDREFS IMPLIED&gt;</pre>
--	--

<pre> &lt;!ELEMENT special (#PCDATA) &gt; &lt;!--ATTLIST special price CDATA #REQUIRED--&gt; ]&gt; </pre>	<pre> &lt;!--ATTLIST CurrentMayor cityID IDREF REQUIRED--&gt; &lt;!--ATTLIST CurrentMoP rname IDREF REQUIRED--&gt; &lt;!--ATTLIST Province pname ID REQUIRED--&gt; &lt;!--ATTLIST City cityID ID REQUIRED cname CDATA REQUIRED--&gt; &lt;!--ATTLIST Riding rname CDATA REQUIRED--&gt; ]&gt; </pre>
<pre> &lt;bibliography&gt;   &lt;books&gt;     &lt;book ISBN="23456" year="1995"&gt;       &lt;title&gt; Foundations ...     &lt;/title&gt;     &lt;author&gt; Hull &lt;/author&gt;     &lt;author&gt; Abiteboul &lt;/author&gt;     &lt;publ&gt; Addison Wesley &lt;/publ&gt;     &lt;/book&gt;     &lt;book&gt; ... &lt;/book&gt;   &lt;/books&gt;   &lt;journals&gt;     &lt;journal&gt;       &lt;title&gt; ... &lt;/title&gt;       &lt;article&gt; ... &lt;/article&gt;       ...     &lt;/journal&gt;     &lt;journal&gt; ... &lt;/journal&gt;   &lt;/journals&gt; &lt;/bibliography&gt; </pre>	<pre> &lt;DiscoverTheWorld&gt;   &lt;tour TourId="1"&gt;     &lt;type&gt; Brazil jungle &lt;/type&gt;     &lt;start-date&gt; 16-April &lt;/start-date&gt;     &lt;duration&gt; 14 &lt;/duration&gt;     &lt;price&gt; 2229 &lt;/price&gt;   &lt;/tour&gt;   &lt;tour TourId="2"&gt;     &lt;type&gt; Brazil jungle &lt;/type&gt;     &lt;start-date&gt; 30-April &lt;/start-date&gt;     &lt;duration&gt; 21 &lt;/duration&gt;     &lt;price&gt; 2999 &lt;/price&gt;   &lt;/tour&gt;   &lt;tour TourId="3"&gt;     &lt;type&gt; Kenia safari &lt;/type&gt;     &lt;start-date&gt; 30-April &lt;/start-date&gt;     &lt;duration&gt; 21 &lt;/duration&gt;     &lt;price&gt; 3229 &lt;/price&gt;   &lt;/tour&gt;   &lt;reservation ResId="541" TourId="1"&gt;     &lt;cname&gt; Bettina Kemme &lt;/cname&gt;     &lt;caddress&gt; Montreal &lt;/caddress&gt;     &lt;cost&gt; 2579 &lt;/cost&gt;     &lt;special price="5"&gt; vegetarian &lt;/special&gt;     &lt;special price="20"&gt; single &lt;/special&gt;   &lt;/reservation&gt;   &lt;reservation ResId="542" TourId="2"&gt;     &lt;cname&gt; Your Name &lt;/cname&gt;     &lt;caddress&gt; Your Address &lt;/caddress&gt;     &lt;cost&gt; 3105 &lt;/cost&gt;     &lt;special price="5"&gt; vegetarian &lt;/special&gt;   &lt;/reservation&gt; &lt;/DiscoverTheWorld&gt; </pre>

```

<!DOCTYPE people[
  <!--ELEMENT people(person*)-->
  <!--ELEMENT person(name*, (lastname|familyname)?)-->
  <!--ATTLIST person PID ID #REQUIRED
    age CDATA #IMPLIED
    children IDREFS #IMPLIED
    mother IDREF #IMPLIED
  -->
  <!--ELEMENT name(#PCDATA)-->
  <!--ELEMENT lastname(#PCDATA)-->
  <!--ELEMENT familyname (#PCDATA)-->
]>

```

Data types: PCDATA (parsed character data) or CDATA (unparsed)

#### Attributes

- ID unique identifier (similar to primary key)
- IDREF: reference to single ID
- IDREFS: space-separated list of references

#### Values

- can give a default value
- #REQUIRED must exist
- #IMPLIED optional

Specified in an XML file with <!DOCTYPE name SYSTEM "path/to/thing.dtd">

Can use regex style things too. \* is 0 or more. + is 1 or more, (a | b)? is one or the other

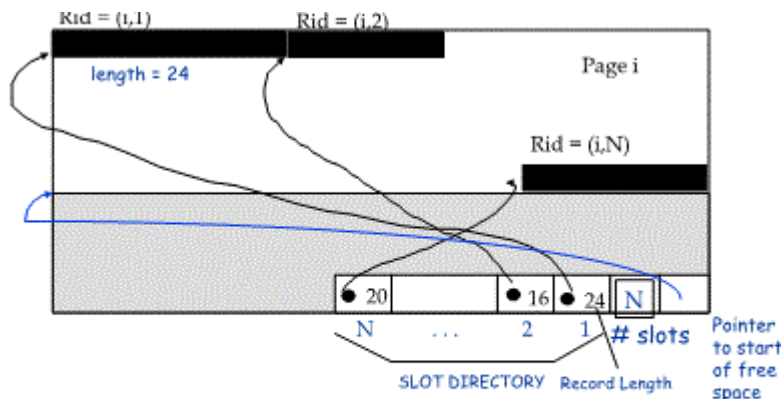
## XPATH

- /bibliography/book/author all author elements by root navigating through those elements
- /bibliography/book/@ISBN All ISBN attributes
- //title all title elements anywhere in the document
- /bibliography/\*/title titles of bibliography entries assuming that there could be books, journals, reports, etc...
- /bibliography/book[@year>1995] returns books where the year > 1995
- /bibliography/book[author='FooBar']/@Year returns the years of books written by FooBar
- /bibliography/book[count(author) <2]
- /bibliography/book/author[position()=1]/name position is the location of the node in the node set

## XQuery

For	Let
for \$b in document("bib.xml")/bib/book return <result> \$b</result>	let \$b in document("bib.xml")/bib/book return <result> \$b</result>
<result><book>...</book></result> <result><book>...</book></result> <result><book>...</book></result> ... <result><book>...</book></result>	<result> <book>...</book> <book>...</book> .. <book>...</book> </result>

<p>-Basic Queries: '/' to navigate one path at a time Example:/Bookstore/Book/Title  '/' all paths following this Example://Title When wanting to access an attribute of an element use @ Example:/Bookstore/Book/data(@ISBN) ' ' OR operator ONLY USED INSIDE CONDITIONS Example:/Bookstore/Book Magazine/Title '=' can act like == like in Self-Join Queries below Navigation accesses: Example: parent::* *::child following-silbling::*</p> <p><b><u>-Queries involving CONDITIONS</u></b> <b><u>1condition:</u></b> Example:/Bookstore/Book[@Price&lt;30] Example:/Bookstore/Book/Authors/Author[2] the 2nd author of each element</p> <p><b><u>2conditions:</u></b> To write a condition, it needs to be inside '['...]' then followed by the output that we are looking for Example:/Bookstore/Book[@Price&lt;30]/Title <b><u>Condition to find elements that contain other elements:</u></b> Example:/Bookstore/Book[Remark]/Title <b><u>Conditions &amp;&amp; Conditions + Some output:</u></b> Example: Looking for a title that has one last name =Ullman and price&lt;90 /Bookstore/Book[@Price&lt;90 and Authors/Author/Last_Name="Ullman"]/Title <b><u>Conditions Inside Conditions + Some output:</u></b> Example:Looking for a title with author ="Jeffrey Ullman" and price&lt;90 /Bookstore/Book[@Price&lt;90 and Authors/Author/[Last_Name="Ullman" and First_Name="Jeffrey"]]/Title <b><u>Conditions &amp;&amp; !Conditions + Some output:</u></b> Example:Looking for a title with author ="Ullman" and NOT author="Widom" /Bookstore/Book[/Authors/Author/Last_Name="Ullman" and count(/Authors/Last_Name="Widom"=0)]/Title <b><u>The condition 'contains':</u></b> /Bookstore/Book[contains(Remark, "great")]/Title</p>	<p><b><u>Self-Join Query</u></b> Querying two instances of the database at one and joining them together.Trying to find the magazines wheres theres a book with the same title. Example: doc("BookstoreQ.xml")/Bookstore/Magazine[Title=doc("BookstoreQ.xml")/Bookstore/Book/Title]</p> <p><b><u>Navigation Accesses</u></b> The name() function returns the name of a tag or element To find all elements whose parent is not "Bookstore" or "Book" <b><u>/Bookstore/*[name(parent::*)!="Bookstore" and name(parent::*)!="Book"]</u></b></p>
--	---



- ➔ **Record id (rid)** = internal identifier of a record:   
 <page id, slot #>.
- ➔ Can move records on page without changing rid;

XML in DB@^%\$^\$#

```
INSERT INTO MyXML(id, INFO) VALUES (1000,
'<customerinfo cid="1000">
<name>Kathy Jones</name>
<addr country =Canada">
  <street>123 fake</street>
  <city>Ottawa</city>
  <prov-state>Ontario</prov-state>
  <pcode-zip>H0H 0H0</pcode-zip>
</addr>
</customerinfo>')
```

Buffer

<p>DBMS stores information persistently on ( "hard" ) disks.</p> <ul style="list-style-type: none"> <li>❑ Unit of transfer main-memory/disk: disk blocks or pages.</li> <li>❑ Timing: <ul style="list-style-type: none"> <li>☆ 2- 20 msec for random data block (bad seek time)</li> <li>☆ If blocks are sequentially on disk, only +1ms per block</li> <li>☆ Compare main memory access: in nanoseconds</li> </ul> </li> <li>❑ Basic operations (READ/WRITE from/to disk)</li> <li>❑ Why disks? <ul style="list-style-type: none"> <li>☆ Cheaper than Main Memory</li> <li>☆ Higher Capacity</li> <li>☆ Main Memory is volatile</li> </ul> </li> </ul>	<p>When loading a page from disk:</p> <ul style="list-style-type: none"> <li>☆ Replacement frame must have "pin counter" of 0</li> <li>❑ When requesting a page that is in the buffer <ul style="list-style-type: none"> <li>☆ Increment pin counter</li> </ul> </li> <li>❑ After operation has finished <ul style="list-style-type: none"> <li>☆ Decrement pin counter</li> <li>☆ Set dirty bit if page has been modified:</li> </ul> </li> <li>❑ Frame is chosen for replacement by a replacement policy: <ul style="list-style-type: none"> <li>☆ Only unpinned page can be chosen (pin count = 0)</li> <li>☆ Least-recently-used (LRU), Clock, MRU etc.</li> </ul> </li> </ul>	<p>If requested is not in pool:</p> <ul style="list-style-type: none"> <li>☆ If there is an empty frame, use it</li> <li>☆ Else choose an empty frame for replacement. If the frame is dirty (page was modified), write it to disk</li> <li>☆ Read requested page into chosen frame</li> </ul> <p>Buffer management in DBMS requires ability to:</p> <ul style="list-style-type: none"> <li>☆ <b>pin a page</b> in buffer pool, <b>force a page to disk</b> (important for implementing CC &amp; recovery),</li> <li>☆ adjust <b>replacement policy</b>, and <b>pre-fetch pages</b> based on access patterns in typical DB operations.</li> </ul>
---	--	---

Indexing

COST MODEL	HEAP FILES	SORTED FILES
<p>Measure performances by simplifying the parameters (10 focused):</p> <ul style="list-style-type: none"> <li>☆ only consider disk reads (ignore writes)</li> <li>☆ only consider number of I/Os and not the individual time for each</li> </ul>	<ul style="list-style-type: none"> <li>☆ Linked, unordered list of all pages of the file</li> <li>☆ Is it good for: <ul style="list-style-type: none"> <li>● scan retrieving all records (SELECT *)?</li> <li>▲ yes, you have to retrieve all</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>☆ Records are ordered according to one or more attributes of the relation</li> <li>☆ Is it good for: <ul style="list-style-type: none"> <li>● scan retrieving all records (SELECT *)?</li> </ul> </li> </ul>

read (ignores page pre-fetch) ☆ Average-case analysis: based on several simplistic assumptions. ● delete/update  ▲ depends on where	pages anyway ● equality search on primary key ▲ not great: have to read on avg half the pages for 1 record ● range search or equality search on non-primary key ▲ not great, all pages need to be read ● insert ▲ yes, can insert anywhere ● delete/update ▲ depends on where	▲ yes, you have to retrieve all pages anyway ● equality search on sort attribute ▲ good: find first qualifying page with binary search (log2) ● range search on sort attribute ▲ good: find first qualifying page with binary search (log2): adjacent pages might have additional matching records
---	---	--

Let suppose we have a relation R (A, B, C, D, F) such that:

- A and B are int (6 byte)
- C-F are char [40] (10 byte per char).
- Tuple = 172 bytes. 200,000 tuples
- Each data page has 4000 bytes and is around 80% full
- B values are uniformly distributed
- Rid = 10 bytes
- Size of pointer in intermediate page = 8 bytes
- Index pages are 4K and between 50%-100% full

Goal	Formula	With this example
Number of pages	$\frac{\text{number of tuples} * \text{tuple size}}{\text{fill rate} * \text{page size}}$	$\frac{172 * 200000}{40000 * 0.80} = 10750$
Index entry size in root and intermediate pages	$\text{size of key} + \text{size of pointer}$	$6 + 8 = 14 \text{ bytes}$
Average number of rids per data entry	$\frac{\text{number of tuples}}{\text{different values (if uniform)}}$	$\frac{200,000}{20,000} = 10$
Average length per data entry	$\text{size of key} + (\text{number of rids} * \text{size of rid})$	$6 + 10 * 10 = 106$
Average number of data entries per leaf page	$\frac{\text{fillrate} * \text{page size}}{\text{length of data entry}}$	$\frac{0.75 * 4000}{106} = 28 \text{ entries per page}$
Estimate number of leaf page	$\frac{\text{number of different values}}{\text{number of entrier per page}}$	$\frac{20,000}{28} = 715$
Number of entries in intermediate pages	$\frac{\text{fillrate} * \text{page size}}{\text{lenght of index enty}}$	$\min = \frac{0.5 * 4000}{14} = 143, \max = \frac{1 * 4000}{14} = 285$
Height of tree	$(\text{nb of entry in intermediate page})^{h-1} > \text{nb of leaf page}$	3

Non-clustered index B-tree with <k, list of rid>

Height of tree = Number of leaf pages / (min | max)? number of entries in intermediate pages

Give the pids of all projects within department D2 that started in 2014.

$$\pi_{pid} \left( \sigma_{dep_{id}=D2 \wedge start_{date}=2014} (Project) \right)$$

Give the pids of all projects that have at least one excellent evaluation

$$\pi_{Project.pid} \left( \sigma_{Evaluation.grade='excellent'} (Project \bowtie Evaluation) \right)$$

<b><u>Force Flush strategy</u></b> <ul style="list-style-type: none"> <li>All changes are flush to disk BEFORE commit</li> </ul>	<ul style="list-style-type: none"> <li>Completed transaction need not action</li> <li>Active transaction might have partial changes on disk(Need undone)</li> </ul>	<ul style="list-style-type: none"> <li>Append to log file log record before flushing</li> <li>At commit/abort append to log file commit/abort log record</li> <li>When recovering from crash: Scan log backward for each record if committed ignore otherwise install Before-Image of the record</li> </ul>
<b><u>No force flush strategy</u></b> <ul style="list-style-type: none"> <li>Changes might be flushed at any time(BEFORE/AFTER commit)</li> </ul>	<ul style="list-style-type: none"> <li>Done transaction might have missing changes (must be redone)</li> <li>Active/Aborted transaction might have been flushed before crash(Must be undone)</li> </ul>	<ul style="list-style-type: none"> <li>For each write(x) of a transaction T with x being on page P: Log record with before AND after image of x(Before so you can undo changes, After so you can redo changes)</li> <li>Flush before-image to disk before flushing the P</li> <li>Flush after-image to disk before commit of T</li> <li>At commit/abort append commit/abort record to log file and flush</li> </ul>

<ul style="list-style-type: none"> <li><b><u>Unrepeatable read:</u></b> If T1 read twice the same data item but T2 change its value between the first and the second</li> <li><b><u>Dirty read:</u></b> If T2 read from T1 before T1 commit.</li> <li><b><u>Lost update:</u></b> If T2 modify a data item modified by T1 without taking in account the value modified by T1.</li> </ul>	
---	--

Schedule <ul style="list-style-type: none"> <li><u>Serial schedule:</u> All transaction one after the other</li> <li><u>Non-serial schedule:</u> Transaction overlap <ul style="list-style-type: none"> <li><u>Serializable:</u> Dependency graph has no cycle(T1 always does action before T2)</li> <li><u>Recoverable schedule:</u> If transaction <math>T_i</math> reads a value written by transaction <math>T_j</math> then <math>T_i</math> commit only after <math>T_j</math> committed</li> <li><u>Avoiding cascading aborts:</u> A transaction reads only values written by committed transactions.</li> <li><u>Strict:</u> A transaction only read or overwrite value written by committed transaction</li> </ul> </li> </ul>	Schedule examples: <ul style="list-style-type: none"> <li>Strict and serializable <math>r1(x), w2(x), c2, w1(x), c1</math></li> <li>Avoids cascading aborts, non-strict, serializable</li> <li>Recoverable, not avoiding cascade aborts, serializable <math>r1(x), w2(y), w2(x), r1(y), c2, c1</math></li> <li>Not recoverable, serializable <math>r1(x), w2(y), r1(x), c1, c2</math></li> <li>Not recoverable, Not-serializable</li> </ul>
---	---

Unrecoverable	Recoverable schedule with cascading abort	Recoverable schedule with commit	Avoids cascading	Non strict	Strict	Strict and serializable
T1      T2 R(A) W(A)  R(A) commit commit	T1      T2 R(A) W(A)  R(A) abort abort	T1      T2 R(A) W(A)  R(A) commit commit	T1      T2 R(A) W(A) abort  R(A) commit	T1      T2 W(A)  W(A) abort  commit	T1      T2 W(A) abort  W(A) commit	T1      T2 R(x)  W(x) commit  W(x) commit

<b><u>Lock request:</u></b> <ul style="list-style-type: none"> <li>If lock is S, no X lock is active and the request queue is empty: <ul style="list-style-type: none"> <li>Add the lock to the granted lock queue and set the lock type to S</li> </ul> </li> <li>If lock is X and no lock active(request queue is also empty): <ul style="list-style-type: none"> <li>Add the lock to the granted lock queue and set the lock type to X</li> </ul> </li> <li>Otherwise <ul style="list-style-type: none"> <li>Add the lock to the request lock queue</li> </ul> </li> </ul>	<b><u>Lock release:</u></b> <ul style="list-style-type: none"> <li>Remove the lock from the granted lock queue</li> <li>If this was the only lock granted on this object: <ul style="list-style-type: none"> <li>Grant one X lock(If the first of the request is a X lock)</li> <li>Grant n S lock(If the first n element are S lock)</li> </ul> </li> </ul>
<b><u>Deadlocks:</u></b> <ul style="list-style-type: none"> <li>Make the wait-for graph(<math>T_i</math> need resource lock by <math>T_j</math>)</li> <li>If cycle then we have a deadlock (Noooooooooooooooooo...)</li> </ul>	<b><u>Solve Deadlock:</u></b> <ul style="list-style-type: none"> <li>Add a timeout for each transaction and abort if transaction timeout. Problem on what timeout value to choose</li> <li>Request all the lock at the beginning of the transaction</li> </ul>

<b><u>Predicate locking:</u></b> <ul style="list-style-type: none"> <li>Grant lock on all records that satisfies logical predicates(e.g. <math>depid &gt; 5, age &gt; 2 * salary</math>)</li> <li>More bullshit</li> </ul>	<b><u>Predicate locking example:</u></b> <ul style="list-style-type: none"> <li>Assume 2 transactions: <ul style="list-style-type: none"> <li><math>UPDATE Skaters \text{ set rating} = 7 \text{ WHERE sid} = 123</math></li> <li><math>SELECT max(age) \text{ FROM Skaters WHERE rating} = 5</math></li> </ul> </li> <li>Assume: T1 execute first then it has a X-lock on Skaters with <math>sid=123</math></li> <li>Assume: T2 has to scan the entire table to get skater with <math>rating=5</math> <ul style="list-style-type: none"> <li>For each tuple <ul style="list-style-type: none"> <li>set S-lock on tuple</li> <li>Check condition</li> <li>If condition TRUE keep lock and return value</li> <li>If condition FALSE release lock</li> </ul> </li> <li>It need to read the tuple where <math>sid=123</math> and <math>rating=5</math> but block has T1 has a lock on it.</li> <li>T2 is block by T1 although there is no conflict</li> </ul> </li> </ul>
--	--

**Problems of strict 2PL locking:**

- Very restrictive, low concurrency, problem with long query
- More and more exception

**In order to allow for more concurrency, SQL2 defines various levels of isolation**

- Assumed to be implemented by different forms of locking
- Avoid different levels of anomalies
- Used for non-critical transactions or read-only transactions
- Lower levels of isolation do NOT provide serializability

**Problem**

- Definitions are no more appropriate if systems do not use locking but other forms of concurrency control
- For instance, Oracle’s “serializable” level does not provide serializable schedule as defined in the literature

Isolation level:

- In principle isolation levels are independent of concurrency control mechanics
- In reality they were defined with locking in mind

Isolation level/Anomaly	Dirty read	Unrepeatable read	Phantom
Read uncommitted	Maybe	Maybe	Maybe
Read committed	No	Maybe	Maybe
Repeated reads	No	No	Maybe
Serializable	No	No	No

- Read uncommitted:**
  - Read op. do not set locks; can read not-committed updates
- Read committed:**
  - Read op. set short S locks; have to wait for X locks to be released
  - release lock immediately after execution of op
- Repeated reads:**
  - Read operations set standard lock S locks; standard 2PL
- Serializable:**
  - Read op. must set S locks that cover all objects that are read
  - predicate locks or coarse locks (e.g., lock on entire relation)





# Big data

Some bullshit info:

- Hardware
  - CPU does not increase
  - Instead muticode
- Usage
  - Astronomy: high-resolution, high-frequency sky surveys
  - Medicine: digital records, MRI, ultrasound
  - Biology: sequencing data
  - User behavior data: click streams, search logs

## Horizontal data Partitioning:

- Data
  - Large table  $R(K, A, B, C)$
  - Key value store  $KV(K, V)$
- Goal
  - Partition into chunks  $c_1, \dots, c_n$  of records stored at N nodes
- Range partition
  - Equal size of each chunk
- Hash partitioned on attribute X
  - Record  $r$  goes to chunk  $i$ , according to hash function
  - xample: hash function  $H(r.X) \bmod P+1$
- Range partitioned on attribute X
  - Partition range of X into:  $-\infty = v_0 < v_1 < \dots < v_p = \infty$
  - Record  $r$  goes to chunk  $i$ , if  $v_{i-1} < t.X < v_i$

## Execution steps:

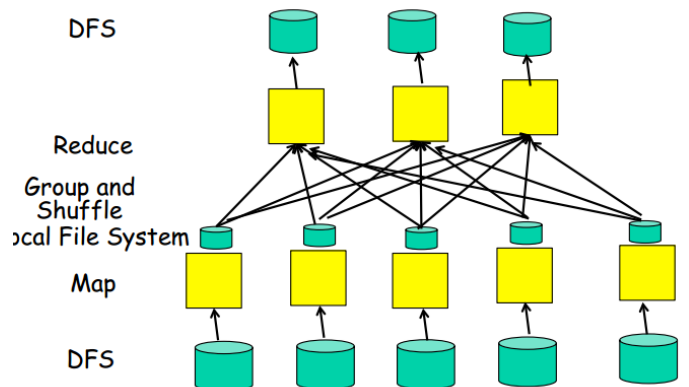
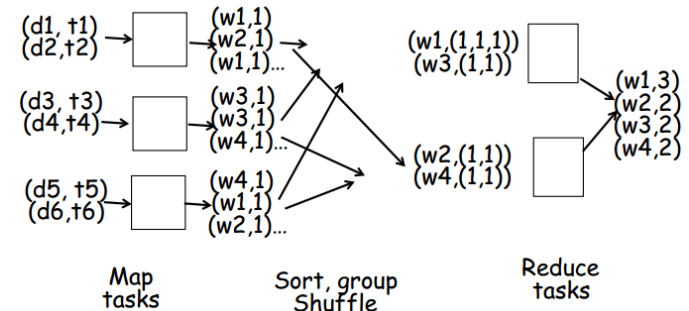
- User indicates  $m$  (number of map tasks),  $r$  (number of reduce tasks),  
key/value set = document Set DS
- System creates  $m$  map tasks and splits input set of 1key/value pairs into  $m$  partitions and gives each map task one partition as input
- Each Map task executes user written map function
  - WordCountMap:
    - For each input key/value pair  $(dkey, dtext)$
    - For each word  $w$  of  $dtext$
    - Output key-value pair  $(w, 1)$
- Next step only completes once all map tasks have completed
- System sorts map outputs by key and transforms all key/value pairs  $(k, v_1), (k, v_2), \dots, (k, v_n)$  with same key  $k$  to one key/value-list pair  $(k, (v_1, v_2, \dots, v_n))$ 
  - For Word count: all  $(\text{'and'}, 1), (\text{'and'}, 1), (\text{'and'}, 1) \dots$  are transformed into one  $(\text{'and'}, (1, 1, 1, \dots))$
- System partitions output by key into  $r$  partitions and assigns these partitions as inputs to the  $r$  reduce tasks
- Each reduce task executes user written reduce function
  - WordCountReduce:
    - For each input key/value-list pair  $(k, (v_1, v_2, \dots, v_n))$
    - Output  $(k, n)$

Parallel Query Evaluation:


- Inter-query parallelism
  - Different queries run in parallel on different processors; each query is executed sequentially
- Inter-operator parallelism
  - Different operator within the same execution tree run on different processors
- Intra-operator parallelism
  - A single operator(JOIN, GROUP, ...) runs on many processor

Vertical Partitioning:

- Column stores
- Data: relation  $R(K, A, B, C)$
- Partition into  $RA(K, A), RB(K, B), RC(K, C)$
- Query:
  - SELECT A FROM R
- Query only needs to access partition RA
- Much less IO



## Map reduce

<p><b>Relational Operators with Map/ reduce</b></p> <ul style="list-style-type: none"> <li>Assume <math>R(A, B, C)</math> relation (no duplicates)</li> <li><b>Selection with condition <math>c</math> on <math>R</math></b> <ul style="list-style-type: none"> <li>for each tuple <math>t</math> of <math>R</math> for which condition <math>c</math> holds, output <math>(t, t)</math></li> <li>Reduce: identity, that is output <math>(t, t)</math></li> </ul> </li> <li><b>Projection on <math>A, B</math>, of <math>R</math></b> <ul style="list-style-type: none"> <li><u>Map</u>: transform each tuple <math>t = (a, b, c)</math> of <math>R</math> into tuple <math>t' = (a, b)</math> of <math>R</math>, and output <math>(t', t')</math></li> <li>There might now be duplicates, that is several <math>(t', t')</math> tuples, the group function will aggregate them to <math>(t', (t', \dots, t'))</math></li> <li><u>Reduce</u>: for each tuple <math>(t', (t', \dots, t'))</math> output <math>(t', t')</math></li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li><b>Grouping: <math>\text{SELECT } a, \text{sum}(b) \text{ GROUP by } (a)</math></b> <ul style="list-style-type: none"> <li><u>Map</u>: for each tuple <math>(a, b, c)</math> of <math>R</math> output <math>(a, b)</math></li> <li>Group and shuffle will create for each value <math>a</math> a key/value-list <math>(a, (b_1, b_2, \dots))</math></li> <li><u>Reduce</u>: for each <math>(a, (b_1, b_2, \dots))</math> perform aggregation (e.g., <math>b_1 + b_2, \dots</math>)</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li><b>Natural Join <math>R(A, B, C)</math> with <math>Q(C, D, E)</math> via hash-join(<math>\text{SELECT FROM } R_1, R_2</math>)</b> <ul style="list-style-type: none"> <li><u>Map</u>:           <ul style="list-style-type: none"> <li>For each tuple <math>(a, b, c)</math> of <math>R</math>, output <math>(c, (R, (a, b)))</math></li> <li>For each tuple <math>(c, d, e)</math> of <math>Q</math>, output <math>(c, (Q, (d, e)))</math></li> </ul> </li> <li>Group and shuffle will aggregate all key/value pairs with same <math>c</math>-value</li> <li><u>Reduce</u>:           <ul style="list-style-type: none"> <li>For each tuple <math>(c, \text{value-list})</math>, example:               <ul style="list-style-type: none"> <li><math>(\text{value-list} = (R, (a_1, b_1)), (R, (a_2, b_2)), \dots (Q, (d_1, e_1)), \dots)</math></li> <li><math>R_t = Q_t = \text{empty}</math>;</li> <li>for each <math>v = (\text{rel}, \text{tuple})</math> in value-list</li> <li>if <math>v.\text{rel} = R</math>: insert tuple into <math>R_t</math> else insert tuple into <math>Q_t</math></li> <li>for <math>v_1</math> in <math>R_t</math>, for <math>v_2</math> in <math>Q_t</math>, output <math>(c, v_1, v_2)</math></li> <li>Basically produces all combinations <math>(c, a_i, b_j, d_k, e_l)</math></li> </ul> </li> </ul> </li> </ul> </li></ul>	<p>Pig latin:</p> <ul style="list-style-type: none"> <li>Users = <b>LOAD</b> 'users' <b>AS</b> (name,age);</li> <li>Filtered = <b>FILTER</b> Users <b>BY</b> age &gt;= 18 <b>AND</b> age &lt;= 25;</li> <li>Pages = <b>LOAD</b> 'pages' <b>AS</b> (uname, url);</li> <li>Joined = <b>JOIN</b> Fltrd <b>BY</b> name, Pages <b>BY</b> uname;</li> <li>Grouped = <b>GROUP</b> Jnd <b>BY</b> url;</li> <li>Smmmd = <b>FOREACH</b> Grpd <b>GENERATE</b> (\$0), COUNT(\$1) <b>AS</b> clicks;</li> <li>Srtd = <b>ORDER</b> Smmmd <b>BY</b> clicks desc;</li> <li>Top5 = <b>LIMIT</b> Srtd 5;</li> <li><b>STORE</b> Top5 <b>INTO</b> 'top5sites'</li> </ul> 
<p><b>Pig examples:</b></p> <pre>raw = LOAD .... -- filter per percent fltrd = FILTER raw by percent &gt;= 60;  gen = foreach fltrd generate CONCAT(firstname, CONCAT(' ', lastname)); results = DISTINCT gen; STORE results INTO 's3n://comp421-h4/q1_results';</pre>	<pre>raw = LOAD .....  --some data entries use the middle name as well, so this way we will catch all of them fltrd = FILTER raw by votes &gt;= 100; SPLIT fltrd INTO winners IF elected == 1, losers IF elected == 0;  elections = JOIN winners BY (date, type, parl, prov, riding), losers BY (date, type, parl, prov, riding);  vote_differences = foreach elections generate winners::lastname as winner, losers::lastname as loser, (winners::votes-losers::votes) as vote_difference:int;  results = FILTER vote_differences by vote_difference &lt; 10; --print the result tuple to the screen  STORE results INTO 's3n://comp421-h4/q2_results';</pre>

```
raw = LOAD ...
--some data entries use the middle name as well, so this way
we will catch all of them
fltrd = FILTER raw by type == 'Gen' and elected == 1;

parl_group = GROUP fltrd BY parl;

parl_count = FOREACH parl_group GENERATE ($0) as parl,
COUNT($1) as count;
parl_count_before = FOREACH parl_count GENERATE ($0+1) as
parl, $1 as count;
parl_join = JOIN parl_count BY parl, parl_count_before BY
parl;

parl_diff = FOREACH parl_join GENERATE parl_count::parl as
parl, parl_count::count, parl_count::count -
parl_count_before::count;

results = ORDER parl_diff BY parl;

dump results;
```

```
raw = LOAD ...
parl_group = GROUP raw BY parl;

parl_count = FOREACH parl_group GENERATE ($0) as parl,
COUNT($1) as parl_count;

party_group = GROUP raw BY (parl, party);

party_count = FOREACH party_group GENERATE FLATTEN($0) as
(parl, party), COUNT($1) as party_count;

joined = JOIN party_count BY parl, parl_count BY parl;

results = FOREACH joined GENERATE parl_count::parl as parl,
party_count::party as party, party_count::party_count as
party_count, parl_count::parl_count as parl_count;

store results into '/user/hadoop/q4output.csv' using
PigStorage('\t', '-schema');
```

## Query evaluation

<p><b>Examples for flowing problems:</b></p> <ul style="list-style-type: none"> <li>Participates <ul style="list-style-type: none"> <li>100 000 tuples</li> <li>1000 pages</li> <li>100 tuples per page</li> </ul> </li> <li>Skaters <ul style="list-style-type: none"> <li>40 000 tuples</li> <li>500 pages</li> <li>80 tuples per page</li> <li>Index on sid has 170 leaf page</li> <li>Index on names has 300 leaf page</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li><b>Reduction factor of a condition defined as</b> <ul style="list-style-type: none"> <li><math>Red(\sigma_{condition}(R)) =  \sigma_{condition}(R) / R </math></li> <li><math>Red(\sigma_{rating=5}(Skaters)) = \frac{ \sigma_{rating=5}(Skaters) }{ Skaters } = \frac{15}{100} = 0.15</math></li> </ul> </li> <li><b>If not known, DBMS makes simple assumptions</b> <ul style="list-style-type: none"> <li><math>Red(\sigma_{rating=5}(Skaters)) = \frac{1}{ diff\ rating } = 0.1</math></li> <li><math>Red(\sigma_{age \leq 10}(Skaters)) = \frac{10 - \min(age)}{\max(age) - \min(age)} = \frac{10 - 4}{30 - 4} = \frac{6}{26} = 0.025</math></li> </ul> </li> <li><b>Result size: number of input tuples * reduction factor</b></li> <li>How to know the number of different values, max, min, <ul style="list-style-type: none"> <li>through indices, heuristics, separate statistics (histograms)</li> </ul> </li> </ul>
<p><b>Simple selection:</b></p> <ul style="list-style-type: none"> <li>No index: <ul style="list-style-type: none"> <li>Search on arbitrary attributes: scan the entire relation e.g. <math>cost = page(Skaters)</math></li> <li>Search on primary key attributes: scan on average half of S e.g. <math>cost = \frac{page(Skaters)}{2} = 250</math></li> </ul> </li> <li>Index on selection attribute <ul style="list-style-type: none"> <li>Use index to find qualifying data entries, then retrieve corresponding data records.</li> </ul> </li> </ul>	<p><b>Clustered B+tree</b></p> <ul style="list-style-type: none"> <li>Costs: <ul style="list-style-type: none"> <li><math>\#LeafPages + \#dataPageToRetrieve</math></li> </ul> </li> <li>Example 1 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name sid = 5</li> <li>1 tuple match</li> <li><math>cost = 1\ leaf\ page + 1\ data\ page = 2</math></li> </ul> </li> <li>Example 2 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name LIKE 'Z%'</li> <li>System estimate the number of matching tuples(Around 100 match on 2 data page as its clustered)</li> <li><math>cost = 1\ leaf\ page + 2\ data\ page = 3</math></li> </ul> </li> <li>Example 3 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name &lt; 'F%'</li> <li>Assume around 10 000 tuples match(On 125 datapage)</li> <li><math>cost = 1\ leaf\ page + 125\ data\ pages = 126</math></li> </ul> </li> </ul>
<p><b>Unclustered B+tree</b></p> <ul style="list-style-type: none"> <li>Costs: <ul style="list-style-type: none"> <li><math>\#dataPageToRetrieve = \#tuples</math></li> <li><math>\#LeafPages + \#dataPageToRetrieve</math></li> </ul> </li> <li>Example 1 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name sid = 5</li> <li>Same</li> </ul> </li> <li>Example 2 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name LIKE 'Z%'</li> <li>Assume around 100 match on 80 data page but we need to retrieve some data page twice</li> <li><math>cost = 1\ leaf\ page + 101\ data\ page = 101</math></li> </ul> </li> <li>Example 3 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name &lt; 'F%'</li> <li>Assume around 10 000 tuples match</li> <li><math>cost = 75\ leaf\ page + 10\ 000\ data\ pages = 10\ 001</math></li> </ul> </li> </ul>	<p><b>Unclustered B+tree with sorting:</b></p> <ul style="list-style-type: none"> <li>Sort matching data entries (rid=pid,slot-id) in leaf-pages by page-id</li> <li>Only fast if the 75 leaf pages with matching entries fit in main memory</li> <li>Retrieve each page only once and get all matching tuples</li> <li><math>\#data\ pages = \#data\ pages\ that\ have\ at\ least\ one\ matching\ tuple;</math></li> <li>worst case is total # of data pages</li> <li>Example 1 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name sid = 5</li> <li>Same</li> </ul> </li> <li>Example 2 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name LIKE 'Z%'</li> <li>Assume around 100 match on 80 data page</li> <li><math>cost = 1\ leaf\ page + 80\ data\ page = 81</math></li> </ul> </li> <li>Example 3 <ul style="list-style-type: none"> <li>SELECT * FROM Skaters WHERE name &lt; 'F%'</li> <li>Assume around 10 000 tuples match(assume thataround 490 data pages)</li> <li><math>cost = 75\ leaf\ page + 490\ data\ pages = 565</math></li> </ul> </li> <li>Note: sorting expensive if leaf-pages do not fit in main-memory</li> </ul>
<p><b>Sort:</b></p> <ul style="list-style-type: none"> <li>Sometimes a pass 2 is needed <ul style="list-style-type: none"> <li>Pass 0 created more runs than there are main memory buffers</li> <li>Therefore Pass 1 produces more than one run</li> <li>Pass 2 takes the runs of Pass 1 and merges them</li> </ul> </li> <li>Cost <ul style="list-style-type: none"> <li>SELECT sname, age FROM Skaters ORDER BY age</li> <li>If everything fits into main memory (Only pass 0 needed): <ul style="list-style-type: none"> <li>Read number of data pages</li> <li>sort and pipeline result into next operator (project)</li> </ul> </li> <li>Pass 0 + pass 1 needed <ul style="list-style-type: none"> <li>Pass 0: read # pages, write # pages (have to write temp. results!)</li> <li>Pass 1: read # pages, sort and pipeline result into next operator</li> <li>3 * #pages</li> </ul> </li> <li>Pass 0 + pass1 + pass2 needed <ul style="list-style-type: none"> <li>5 * #pages</li> </ul> </li> </ul> </li> </ul>	<p><b>Join cost estimation:</b></p> <ul style="list-style-type: none"> <li><math> Skaters Participates  =  Participates </math> <ul style="list-style-type: none"> <li>Join attribute is primary key for Skaters</li> <li>Each Participates tuple matches exactly with one Skaters tuple</li> </ul> </li> <li><math> Skaters \times Participates  =  Participates  *  Skaters </math> <ul style="list-style-type: none"> <li>Cross product is always the product of individual relation sizes</li> </ul> </li> <li>For other joins more difficult to estimate (Continues in next episode...)</li> </ul>

**Nested loop joins:**

- Simple nested loop join: For each tuple in the outer relation P, we scan the entire inner relation S
  - $cost = PartPages + CARD(P) * SkaterPages = 1000 + 100000 * 500$
- Page-oriented Nested Loops join: For each page of P, get each page of S, and write out matching pairs of tuples <p, s>, where p is in P-page and s is in S-page.
  - $cost = PartPage + PartPage * SkaterPage = 1000 + 1000 * 500$

**Join cost on relation R1(outer) and R2(inner):****Block oriented nested loop join:**

- Smaller relation fits in main memory+2extra buffer page:
$$cost = page(R1) + page(R2)$$
- No relation fits in main memory(B join frame):
$$cost = page(R1) + \frac{page(R2) * page(R1)}{B - 2}$$

**Index nested loop join**

- Index on the join column of one of the relation(R2):
$$cost = page(R1) + card(R1) * cost\_finding\_index(R2)$$
- If the join attribute is primary key in inner relation

**Sort merge join**

- Sort P and S on the join column, then scan them to do a "merge" (on join col.), and output result tuples
  - Advance scan of P until current P-tuple >= current S tuple, then advance scan of S until current S-tuple >= current P tuple; do this until current P tuple = current S tuple.
  - At this point, all P tuples with same value in Pi (current P group) and all S tuples with same value in Sj (current S group) match; output <p, s> for all pairs of such tuple 1s.
- P is scanned once; each S group is scanned once per matching P tuple. (Multiple scans of an S group are likely to find needed pages in buffer.)