# ER Diagram

- ## Many to Many
A can have multiple B and B can have multiple A

| A | R | B |
|---|---|---|

- ## One to Many
A can have multiple B but B can only have one A

| A | R | B |
|---|---|---|

- ## One to One
A can have only one B and B can have only one B

| A | R | B |
|---|---|---|

- ## At least one
Bold arrow specify there is at least one element.

| A — R | A have at least one |
|---|---|
| A → R | A have exactly one |

- ## Aggregation
Allows us to treat a reltionship set R as an entity set so that R can participate in other relationships

eid  name  salary
**Employees**

Monitors  until

Started_on
pid  pbudget  since  did  dname  budget
**Projects**  Sponsors  **Departments**

# Relational algebra

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

- Condition/Theta Join $Rout = Rin1 \bowtie cRin2 = \sigma c(Rin1 \times Rin2)$
- Equi Join: $Rout = Rin1 \bowtie Rin1a1 = Rin2b1, ..., Rin1an = Rin2bnRin2$ 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)
        CHECK (name != 'batman' OR name != 'joker')
)

--Drop table
DROP TABLE Students

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

--Creating index
CREATE INDEX useless_index ON table_name(column_name);
```

### 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
```

```
<!DOCTYPE DiscoverTheWorld [
<!ELEMENT DiscoverTheWorld (tour*,reservation*)>
<!ELEMENT tour (type, start-date, duration, price) >
<!ELEMENT reservation (cname, caddress, cost,
special*)>
<!ATTLIST tour TourId ID #REQUIRED >
<!ATTLIST reservation ResID ID #REQUIRED
TourID IDREF #REQUIRED>
<!ELEMENT type (#PCDATA) >
<!ELEMENT start-date (#PCDATA) >
<!ELEMENT duration (#PCDATA) >
<!ELEMENT price (#PCDATA) >
<!ELEMENT cname (#PCDATA) >
<!ELEMENT caddress (#PCDATA) >
<!ELEMENT cost (#PCDATA) >
<!ELEMENT special (#PCDATA) >
<!ATTLIST special price CDATA #REQUIRED>
]>
```

```
<!DOCTYPE Politics [ <!ELEMENT Politics (Politician*,
Province*)>
<!ELEMENT Politician ((CurrentMayor | CurrentMop)?,
address?)>
<!ELEMENT CurrentMayor (since?)>
<!ELEMENT CurrentMoP (since?)>
<!ELEMENT Province (City+, Riding+, population?)>
<!ELEMENT City (population?)>
<!ELEMENT Riding (population?)>
<!ELEMENT address (#PCDATA) >
<!ELEMENT since (#PCDATA) >
<!ELEMENT population (#PCDATA) >
<!ATTLIST Politician pname ID REQUIRED
website CDATA IMPLIED
friends IDREFS IMPLIED>
<!ATTLIST CurrentMayor cityID IDREF REQUIRED>
<!ATTLIST CurrentMoP rname IDREF REQUIRED>
<!ATTLIST Province pname ID REQUIRED>
<!ATTLIST City cityID ID REQUIRED
cname CDATA REQUIRED>
<!ATTLIST Riding rname CDATA REQUIRED>
]>
```

```xml
<bibliography>
        <books>
                <book ISBN="23456" year="1995">
                        <title> Foundations ...
</title>
                        <author> Hull </author>
                        <author> Abiteboul </author>
                        <publ> Addison Wesley </publ>
                </book>
                <book> ...</book>
        </books>
        <journals>
                <journal>
                        <title> ... </title>
                        <article> ... </article>
                        ...
                </journal>
                <journal> ... </journal>
        </journal>
</bibliography>
```
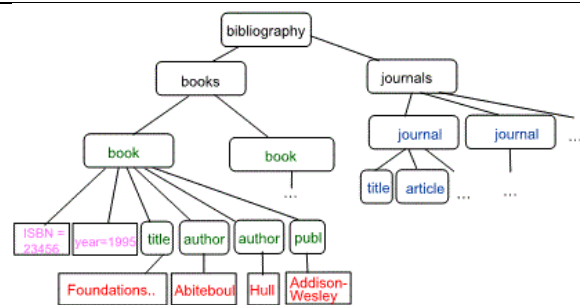
```xml
<DiscoverTheWorld>
        <tour TourId="1">
                <type> Brazil junge </type>
                <start-date> 16-April </start-date>
                <duration> 14 </duration>
                <price> 2229 </price>
        </tour>
        <tour TourId="2">
                <type> Brazil junge </type>
                <start-date> 30-April </start-date>
                <duration> 21 </duration>
                <price> 2999 </price>
        </tour>
        <tour TourId="3">
                <type> Kenia safari </type>
                <start-date> 30-April </start-date>
                <duration> 21 </duration>
                <price> 3229 </price>
        </tour>
        <reservation ResId="541" TourId="1">
                <cname> Bettina Kemme </cname>
                <caddress> Montreal </caddress>
                <cost> 2579 </cost>
                <special price="5"> vegetarian </special>
                <special price="20"> single </special>
        </reservation>
        <reservation ResId="542" TourId="2">
                <cname> Your Name </cname>
                <caddress> Your Address </caddress>
                <cost> 3105 </cost>
                <special price="5"> vegetarian </special>
        </reservation>
</DiscoverTheWorld>
```

```
<!DOCTYPE people[
    <!ELEMENT people(person*)>
    <!ELEEMNT 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-seperated 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

where: condition

| For | Let |
|---|---|
| `for $b in document("bib.xml")/bib/book`<br>return <result> $b</result> | `let  $b in document("bib.xml")/bib/book`<br>return <result> $b</result> |
| `<result><book>…</book></result>`<br>`<result><book>…</book></result>`<br>`<result><book>…</book></result>`<br>…<br><result><book>…</book></result> | `<result>`<br>`<book>…</book>`<br>`<book>…</book>`<br>`..`<br>`<book>…</book>`<br>`</result>` |

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



- **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>')
```

| | | |
|---|---|---|
| DBMS stores information persistently on ( "hard" ) disks.<br>❑ Unit of transfer main-memory/disk: disk blocks or pages.<br>❑ Timing:<br>☆ 2- 20 msec for random data block (bad seek time)<br>☆ If blocks are sequentially on disk, only +1ms per block<br>☆ Compare main memory access: in nanoseconds<br>❑ Basic operations (READ/WRITE from/to disk)<br>❑ Why disks?<br>☆ Cheaper than Main Memory<br>☆ Higher Capacity<br>☆ Main Memory is volatile | When loading a page from disk:<br>☆ Replacement frame must have "pin counter" of 0<br>❑ When requesting a page that is in the buffer<br>☆ Increment pin counter<br>❑ After operation has finished<br>☆ Decrement pin counter<br>☆ Set dirty bit if page has been modified:<br>❑ Frame is chosen for replacement by a replacement policy:<br>☆ Only unpinned page can be chosen (pin count = 0)<br>☆ Least-recently-used (LRU), Clock, MRU etc. | If requested is not in pool:<br>☆ If there is an empty frame, use it<br>☆ Else choose an empty frame for replacement. If the frame is dirty (page was modified), write it to disk<br>☆ Read requested page into chosen frame Buffer management in DBMS requires ability to:<br>☆ **pin a page** in buffer pool, **force a page to disk** (important for implementing CC & recovery),<br>☆ adjust **replacement policy**, and **pre-fetch pages** based on access patterns in typical DB operations. |

| COST MODEL | HEAP FILES | SORTED FILES |
|---|---|---|
| Measure performances by simplifying the parameters (IO focused):<br>☆ only consider disk reads (ignore writes)<br>☆ only consider number of I/Os and not the individual time for each read (ignores page pre-fetch)<br>☆ Average-case analysis; based on several simplistic assumptions. ● delete/update<br><br>▲ depends on where | ☆ Linked, unordered list of all pages of the file<br>☆ Is it good for:<br>● scan retrieving all records (SELECT *)?<br>▲ yes, you have to retrieve all pages anyway<br>● equality search on primary key<br>▲ not great: have to read on avg half the pages for 1 record<br>● range search or equality search on non-primary key<br>▲ not great, all pages need to be read<br>● insert<br>▲ yes, can insert anywhere<br>● delete/update<br>▲ depends on where | ☆ Records are ordered according to one or more attributes of the relation<br>☆ Is it good for:<br>● scan retrieving all records (SELECT *)?<br>▲ yes, you have to retrieve all pages anyway<br>● equality search on sort attribute<br>▲ good: find first qualifying page with binary search (log2)<br>● range search on sort attribute<br>▲ good: find first qualifying page with binary search (log2):<br>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 | $$\frac{number\ of\ tuples\ *\ tuple\ size}{fill\ rate\ *\ page\ size}$$ | $$\frac{172*200000}{40000*0.80}=10750$$ |
| Index entry size in root and intermediate pages | $size\ of\ key + size\ of\ pointer$ | $6+8=14\ bytes$ |
| Average number of rids per data entry | $$\frac{number\ of\ tuples}{different\ values\ (if\ uniform)}$$ | $$\frac{200,000}{20,000}=10$$ |
| Average length per data entry | $size\ of\ key\ +\ (number\ of\ rids\ *\ size\ of\ rid)$ | $6\ +\ 10*10\ =\ 106$ |
| Average number of data entries per leaf page | $$\frac{fillrate\ *\ page\ size}{length\ of\ data\ entry}$$ | $$\frac{0.75*4000}{106}=28\ entries\ per\ page$$ |
| Estimate number of leaf page | $$\frac{number\ of\ different\ values}{number\ of\ entrier\ per\ page}$$ | $$\frac{20,000}{28}=715$$ |
| Number of entries in intermediate pages | $$\frac{fillrate\ *\ page\ size}{lenght\ of\ index\ enty}$$ | $$min=\frac{0.5*4000}{14}=143,\max=\frac{1*4000}{14}=285$$ |
| Height of tree | $(nb\ of\ entry\ in\ intermediate\ page)^{h-1}>nb\ of\ leaf\ page$ | 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

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

| | | |
|---|---|---|
| **Force Flush strategy**<br>• All changes are flush to disk BEFORE commit | • Completed transaction need not action<br>• Active transaction might have partial changes on disk(Need undone) | • Append to log file log record before flushing<br>• At commit/abort append to log file commit/abort log record<br>• When recovering from crash: Scan log backward for each record if commited ignore otherwise install Before-Image of the record |
| **No force flush strategy**<br>• Changes might be flushed at any time(BEFORE/AFTER commit) | • Done transaction might have missing changes (must be redone)<br>• Active/Aborted transaction might have been flushed before crash(Must be undone) | • 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)<br>• Flush before-image to disk before flushing the P<br>• Flush after-image to disk before commit of T<br>• At commit/abort append commit/abort record to log file and flush |

| | |
|---|---|
| • **Unrepeatable read:** If T1 read twice the same data item but T2 change its value between the first and the second<br>• **Dirty read:** If T2 read from T1 before T1 commit.<br>• **Lost update:** If T2 modify a data item modified by T1 without taking in account the value modified by T1. | |

| Schedule | Schedule examples: |
|---|---|
| • <u>Serial schedule:</u> All transaction one after the other<br>• <u>Non-serial schedule:</u> Transaction overlap<br>- <u>Serializable:</u> Dependency graph has no cycle(T1 always does action before T2)<br>- <u>Recoverable schedule:</u> If transaction $T_i$ reads a value written by transaction $T_j$ then $T_i$ commit only after $T_j$ committed<br>- <u>Avoiding cascading aborts</u>: A transaction reads only values written by committed transactions.<br>- <u>Strict:</u> A transaction only read or overwrite value written by committed transaction | • Strict and serializable<br><br>$$r1(x), w2(x), c2, w1(x), c1$$<br>• Avoids cascading aborts, non-strict, serializable<br>• Recoverable, not avoiding cascade aborts, serializable<br><br>$$r1(x), w2(y), w2(x), r1(y), c2, c1$$<br>• Not recoverable, serializable<br><br>$$r1(x), w2(y), r1(x), c1, c2$$<br>• Not recoverable, Not-serializable |

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

| Lock **request:** | Lock release: |
|---|---|
| • If lock is S, no X lock is active and the request queue is empty:<br>- Add the lock to the granted lock queue and set the lock type to S<br>• If lock is X and no lock active(request queue is also empty):<br>- Add the lock to the granted lock queue and set the lock type to X<br>• Otherwise<br>- Add the lock to the request lock queue | • Remove the lock from the granted lock queue<br>• If this was the only lock granted on this object:<br>- Grant one X lock(If the first of the request is a X lock)<br>- Grant n S lock(If the first n element are S lock) |
| **Deadlocks:** | **Solve Deadlock:** |
| • Make the wait-for graph($T_i$ need ressource lock by $T_j$)<br>• If cycle then we have a deadlock (Noooooooooooooo…) | • Add a timeout for each transaction and abort if transaction timeout. Problem on what timeout value to choose<br>• Request all the lock at the beginning of the transcation |

| Predicate locking: | Predicate looking example: |
|---|---|

**Predicate locking:**
- Grant lock on all records that satisfies logical predicates(e.g. depid>5, age > 2*salary)
- More bullshit

**Predicate looking example:**
- Assume 2 tranascrtions:
- *UPDATE Skaters set rating = 7 WHERE sid = 123*
- *SELECT max(age) FROM Skaters WHERE rating = 5*
- Assume: T1 execute first then it has a X-lock on Skaters with sid=123
- Assume: T2 has to scan the entire table to get skater with rating=5
  - For each tuple
    - set S-lock on tuple
    - Check condition
    - If condition TRUE keep lock and return value
    - If condition FALSE release lock
  - It need to read the tuple where sid=123 and rating= 5 but block has T1 has a lock on it.
  - T2 is block by T1 although there is no conflict

---

**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:** | **Parallel Query Evaluation:** |

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

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

**Horizontal data Partitioning:**

- Data
  - Large table $R(K, A, B, C)$
  - Key value store $KV(K, V)$
- Goal
  - Partition into chunks $c_1, \ldots c_n$ 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) mod P+1
- Range partitioned on attribute X
  - Partition range of X into: $-\infty = v_0 < v_1 < \cdots < v_p = \infty$
  - Record r goes to chunk i, if $v_{i-1} < t.X < v_i$

**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

**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 $(k, v1), (k, v2), \ldots (k, vn)$ with same key k to one key/value-list pair $(k, (v1, v2, \ldots vn))$
  - For Word count: all ('and', 1), ('and', 1), ('and', 1) … are transformed into one ('and', (1,1,1,….))
- 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 $(k, (v1, v2, \ldots vn)$
    Output (k, n)

| | |
|---|---|
| **Relational Operators with Map/ reduce**<br><br>• Assume R(A, B, C) relation (no duplicates)<br>• **Selection with condition c on R**<br>- for each tuple t of R for which condition c holds, output $(t,t)$<br>- Reduce: identity, that is output $(t,t)$<br>• **Projection on A, B, of R**<br>- Map: transform each tuple $t = (a,b,c)$ of R into tuple $t' = (a,b)$ of R, and output $(t',t')$<br>- There might now be duplicates, that is several $(t',t')$ tuples, the group function will aggregate them to $(t',(t',...,t'))$<br>- Reduce: for each tuple $(t',(t',...,t'))$ output $(t',t')$ | • **Grouping: SELECT a, sum(b) GROUP by (a)**<br>- Map: for each tuple (a, b,c) of R output (a,b)<br>- Group and shuffle will create for each value $a$ a key/value-list $(a,(b1,b2,...)$<br>- Reduce: for each $(a,(b1,b2,...))$ perform aggregation $(e.g., b1 + b2,...)$ |
| • **Natural Join R(A,B,C) with Q(C,D,E) via hash-join(SELECT FROM R1,R2)**<br>- Map:<br>    For each tuple $(a,b,c)$ of R, output $(c,(R,(a,b)))$<br>    For each tuple $(c,d,e)$ of Q, output $(c,(Q,(d,e)))$<br>- Group and shuffle will aggregate all key/value pairs with same c-value<br>- Reduce:<br>    For each tuple (c, value-list), example:<br>    $(value-list = (R,(a1,b1)),(R,(a2,b2)),...(Q,(d1,e1)),...))$<br>      Rt = Qt = empty;<br>      for each v=(rel,tuple) in value-list<br>        if v.rel = R: insert tuple into Rt else insert tuple into Qt<br>      for v1 in Rt, for v2 in Qt, output(c,v1,v2)<br>      Basically produces all combinations (c, ai,bi,dj,ej) | Pig latin:<br><br>• Users = **LOAD** 'users' **AS** (name,age);<br>• Filtered = **FILTER** Users **BY** age >= 18 **AND** age <= 25;<br>• Pages = **LOAD** 'pages' **AS** (uname, url);<br>• Joined = **JOIN** Fltrd **BY** name, Pages **BY** uname;<br>• Grouped = **GROUP** Jnd **BY** url;<br>• Smmd = **FOREACH** Grpd **GENERATE** ($0), COUNT($1) **AS** clicks;<br>• Srtd = **ORDER** Smmd **BY** clicks desc;<br>• Top5 = **LIMIT** Srtd 5;<br>• **STORE** Top5 **INTO** 'top5sites'<br><br> |
| **Pig examples:**<br><br>```<br>raw = LOAD ….<br>-- filter per percent<br>fltrd = FILTER raw by percent >= 60;<br><br>gen = foreach fltrd generate CONCAT(firstname, CONCAT(' ', lastname));<br>results = DISTINCT gen;<br>STORE results INTO 's3n://comp421-h4/q1_results';<br>``` | ```<br>raw = LOAD ………<br><br>--some data entries use the middle name as well, so this way we will catch all of them<br>fltrd = FILTER raw by votes >= 100;<br>SPLIT fltrd INTO winners IF elected == 1, losers IF elected == 0;<br><br>elections = JOIN winners BY (date, type, parl, prov, riding), losers BY (date, type, parl, prov, riding);<br><br>vote_differences = foreach elections generate winners::lastname as winner, losers::lastname as loser, (winners::votes-losers::votes) as vote_difference:int;<br><br>results = FILTER vote_differences by vote_difference < 10;<br>--print the result tuple to the screen<br><br>STORE results INTO 's3n://comp421-h4/q2_results';<br>``` |

```
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

| | |
|---|---|
| **Examples for flowing problems:**<br>• Participates<br>  - 100 000 tuples<br>  - 1000 pages<br>  - 100 tuples per page<br>• Skaters<br>  - 40 000 tuples<br>  - 500 pages<br>  - 80 tuples per page<br>  - Index on sid has 170 leaf page<br>  - Index on names has 300 leaf page | • **Reduction factor of a condition defined as**<br>  - $Red(\sigma_{condition}(R)) = \|\sigma_{condition}(R)\|/\|R\|$<br>  - $Red\left(\sigma_{rating=5}(Skaters)\right) = \frac{\|\sigma_{rating=5}(Skaters)\|}{\|Skaters\|} = \frac{15}{100} = 0.15$<br>• **If not known, DBMS makes simple assumptions**<br>  - $Red\left(\sigma_{rating=5}(Skaters)\right) = \frac{1}{\|diff\ rating\|} = 0.1$<br>  - $Red\left(\sigma_{age\leq10}(Skaters)\right) = \frac{10-min(age)}{max(age)-min(age)} = \frac{10-4}{30-4} = \frac{6}{26} = 0.025$<br>• **Result size:** *number of input tuples * reduction factor*<br>• How to know the number of different values, max, min,<br>  - through indices, heuristics, separate statistics (histograms) |
| **Simple selection:**<br>• No index:<br>  - Search on arbitrary attributes: scan the entire relation<br>    e.g. $cost = page(Skaters)$<br>  - Search on primary key attributes: scan on average half of S<br>    e.g. $cost = \frac{page(Skaters)}{2} = 250$<br>• Index on selection attribute<br>  - Use index to find qualifying data entries, then retrieve corresponding data records. | **Clustered B+tree**<br>• Costs:<br>  - $\#LeafPages + \#dataPageToRetreive$<br>• Example 1<br>  - SELECT * FROM Skaters WHERE name sid = 5<br>  - 1 tuple match<br>  - $cost = 1\ leaf\ page + 1\ data\ page = 2$<br>• Example 2<br>  - SELECT * FROM Skaters WHERE name LIKE 'Z%'<br>  - System estimate the number of matching tuples(Around 100 match on 2 data page as its clustered)<br>  - $cost = 1\ leaf\ page + 2\ data\ page = 3$<br>• Example 3<br>  - SELECT * FROM Skaters WHERE name < 'F%'<br>  - Assume around 10 000 tuples match(On 125 datapage)<br>  - $cost = 1\ leaf\ page + 125\ data\ pages = 126$ |
| **Unclustered B+tree**<br>• Costs:<br>  - $\#dataPageToRetreive = \#tuples$<br>  - $\#LeafPages + \#dataPageToRetreive$<br>• Example 1<br>  - SELECT * FROM Skaters WHERE name sid = 5<br>  - Same<br>• Example 2<br>  - SELECT * FROM Skaters WHERE name LIKE 'Z%'<br>  - Assume around 100 match on 80 data page but we need to retrieve some data page twice<br>  - $cost = 1\ leaf\ page + 101\ data\ page = 101$<br>• Example 3<br>  - SELECT * FROM Skaters WHERE name < 'F%'<br>  - Assume around 10 000 tuples match<br>  - $cost = 75\ leaf\ page + 10\ 000\ data\ pages = 10\ 001$ | **Unclustered B+tree with sorting:**<br>• Sort matching data entries (rid=pid,slot-id) in leaf-pages by page-id<br>• Only fast if the the 75 leaf pages with matching entries fit in main memory<br>• Retrieve each page only once and get all matching tuples<br>• #data pages = #data pages that have at least one matching tuple;<br>• worst case is total # of data pages<br>• Example 1<br>  - SELECT * FROM Skaters WHERE name sid = 5<br>  - Same<br>• Example 2<br>  - SELECT * FROM Skaters WHERE name LIKE 'Z%'<br>  - Assume around 100 match on 80 data page<br>  - $cost = 1\ leaf\ page + 80\ data\ page = 81$<br>• Example 3<br>  - SELECT * FROM Skaters WHERE name < 'F%'<br>  - Assume around 10 000 tuples match(assume thataround 490 data pages)<br>  - $cost = 75\ leaf\ page + 490\ data\ pages = 565$<br>• Note: sorting expensive if leaf-pages do not fit in main-memory |
| **Sort:**<br>• Sometimes a pass 2 is needed<br>  - Pass 0 created more runs than there are main memory buffers<br>  - Therefore Pass 1 produces more than one run<br>  - Pass 2 takes the runs of Pass 1 and merges them<br>• Cost<br>  - SELECT sname, age FROM Skaters ORDER BY age<br>  - If everything fits into main memory (Only pass 0 needed):<br>      Read number of data pages<br>      sort and pipeline result into next operator (project)<br>  - Pass 0 + pass 1 needed<br>      Pass 0: read # pages, write # pages (have to write temp. results!)<br>      Pass 1: read # pages, sort and pipeline result into next operator<br>      3 * #pages<br>  - Pass 0 + pass1 + pass2 needed<br>      5 * #pages | **Join cost estimation:**<br>• $\|Skaters\ Participates\| = \|Participates\|$<br>  - Join attribute is primary key for Skaters<br>  - Each Participates tuple matches exactly with one Skaters tuple<br>• $\|Skaters \times Participates\| = \|Participates\| * \|Skaters\|$<br>  - Cross product is always the product of individual relation sizes<br>• For other joins more difficult to estimate (Continues in next episode…) |

**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.)



"This problem had been my life's work. I planned to devote my remaining years to it. It's just been solved in four seconds."



THE #1 PROGRAMMER EXCUSE FOR LEGITIMATELY SLACKING OFF: "MY CODE'S COMPILING."

HEY! GET BACK TO WORK!

Pig executing!

OH. CARRY ON.

I THINK WE SHOULD BUILD AN SQL DATABASE.

UH-OH

DOES HE UNDERSTAND WHAT HE SAID OR IS IT SOMETHING HE SAW IN A TRADE MAGAZINE AD?

WHAT COLOR DO YOU WANT THAT DATABASE?

I THINK MAUVE HAS THE MOST RAM.