

ER Model

- overlap constraint: can an entity be in more than one subclass?
- covering constraint: Must every entity of the superclass be in one of the subclasses?
- key constraints
 - one-many → thin arrow
 - one-one → key constraints in both directions
- participation constraints
 - thick line, at least one
- combined key + participation: exactly once



Relational Model

- SQL
 - Create table X ()
 - Select ... from X where ...
 - Insert into X Values ()
 - Drop Table X
 - Delete from X where ...
 - Alter Table X ADD column ...
 - Update X set ... where ...

Integrity Constraints

- need to always be true
- NOT NULL

Primary Key Constraints

- need a unique row identifier

Referential Integrity

- all foreign key constraints are enforced
- foreign key (row) references X

ER-Relational Translation

Eg: DB2: CREATE TABLE Companies
 (name VARCHAR(30),
 Addr VARCHAR(50),
 Empl INTEGER,
 PRIMARY KEY (name))

→ Companies(name, address, empl)

- many-many relationship translated to a table

- participation constraints can't usually be reflected besides NOT NULL

- weak entity set and identifying relationship are translated into a single table

- ISA hierarchies: distribute info among relations OR sublasses have all attributes OR one big table with null values

Relational Algebra

- closed → can be composed and concatenated
- relations are sets (no two rows are identical)
- selection: σ (get rows)
- projection: Π (get cols, eliminate dups)
- renaming: ρ
- set difference: $-$ (in first but not second)
- $\sigma_{\text{condition}}(\text{relation})$, $R_1 \bowtie_c R_2$
- e.g. all makers that don't make laptops: $\Pi_{\text{maker}}(\text{computer}) - \Pi_{\text{maker}}(\text{category} = \text{laptop} \wedge \text{computer})$
- e.g. names of skaters who participated in local or regional comps:
 $\rho(\text{Temp}, \sigma_{\text{type} = \text{local}} \vee \text{type} = \text{regional})(C) \bowtie P$, $\text{Name}(\text{Temp} \bowtie S)$
- e.g. names of skaters who participated in local or regional comps:
 $P(\text{locals}, \Pi_{\text{skid}}(\text{type} = \text{local})(C) \bowtie P)$, $P(\text{regionals}, \Pi_{\text{skid}}(\text{type} = \text{regional})(C) \bowtie P)$, $\text{locals} \bowtie \text{regionals}$
- rules: $\Pi_L(\sigma_C(R)) = \sigma_C(\Pi_L(R))$ (if C considers attributes of L), $R_1 \bowtie R_2 = R_2 \bowtie R_1$,
 $\Pi_{L_2}(\Pi_{L_1}(R)) = \Pi_{L_2}(R)$ if $L_2 \subseteq L_1$, $\sigma_{C_2}(\sigma_{C_1}(R)) = \sigma_{C_1 \cup C_2}(R)$, $R_1 \bowtie (R_2 \bowtie R_3) = (R_1 \bowtie R_2) \bowtie R_3$

Quiz 1: ER

- relationships vs Entity sets:



payment as entity set allows several payments of the same company to the same party by same lobbyist (more flexibility) → can't really track a relation directly. Relationships can only have one instance amongst the exact same entities.

- Weak entity → when introducing an artificial key we should not make it a partial key – e.g. no weak entity

Quiz 2: Relational

- In a relational model: ① all rows stored must be distinct
 ② given a query, a database system can implement various ways to retrieve results ③ query languages built for relational models are data-centric
- candidate keys: all primary keys are candidate keys.
 If the combination of attributes (a1, a2) is a candidate key, then it's possible that 2 records have the same value for a1.
 → candidate keys have to be the minimal subset.
 e.g. if a1 is a candidate key alone, (a1, a2) is not one.
- foreign keys: a relation can have more than one foreign key to the same relation. Foreign keys must refer to something that exists. A foreign key by itself can also be a primary key but a foreign key cannot, by itself, be a primary key.
- can only reflect the NOT NULL participation constraint (thick arrow). A one-one key constraint implies you can merge the entity and relation.
- In the relational translation of a weak entity set, the foreign key is strictly to the relation that it is immediately dependent on.
- "at most one" constraint (\rightarrow) means that entity should be the primary key of relationship set table.
- many-many needs at least 2 tables

Quiz 3: Relational algebra

- In the union of 2 relations, the output may have same num rows as the input relations if the inputs had identical rows.
- selection and projection are commutative if the attributes in the selection are a subset of the attributes of the proj.
- be careful that projections don't eliminate necessary attributes
- non-natural join needs attribute to join on specified
- can't natural join with no common attribute.

SQL

- $\sigma \Leftrightarrow \text{WHERE}$, $\Pi \Leftrightarrow \text{SELECT}$
- no elimination of duplicates (unless DISTINCT)
- Join, always indicate condition
- self join e.g.: `SELECT p1.sid, p2.sid`
 $\text{FROM Players p1, Players p2}$
 $\text{WHERE p1.cid=p2.cid AND p1.sid < p2.sid}$
- ^(or) ^(and) Union/Intersect/Join need set-compatible relations
 \hookrightarrow provide duplicate elimination (use UNION ALL to avoid)
- Nested Queries (IN, or NOT IN)
 \hookrightarrow e.g. names of skaters who participated in comp 101
`SELECT sname`
`FROM Skaters`
`WHERE sid IN (SELECT sid`
 $\quad \quad \quad \text{FROM Participates}$
 $\quad \quad \quad \text{WHERE cid=101})$
- \exists - true iff the nested relation is not empty
- ANY, ALL - e.g. skaters with highest rating:
`SELECT *`
`FROM Skaters`
`WHERE >= ALL (SELECT rating`
`FROM Skaters)`
- e.g. Names of skaters who participated in all competitions:
`SELECT sname`
`FROM Skaters s`
`WHERE NOT EXISTS ((SELECT c.cid`
`FROM Competition c)`
`EXCEPT`
`(SELECT p.cid`
`FROM Participates p`
`WHERE p.sid=s.sid))`
- COUNT, SUM, AVG, MAX, MIN → apply to single attribute/column
 \hookrightarrow e.g. Names of skaters w/ highest rankings
`SELECT sname`
`FROM skaters`
`WHERE rating =`
 $\quad \quad \quad (\text{SELECT MAX(rating)}$
`FROM skaters)`
- If any aggregation is used then each element in the attribute list of the select clause must be aggregated or appear in a group-by
- HAVING: e.g. for each rating, find the min age, only consider rating levels w/ at least 2 skaters:
`SELECT rating, MIN(age)`
`FROM skaters`
`GROUP BY rating`
`HAVING COUNT(*) >= 2`
- Aggregate operations cannot be nested
 \hookrightarrow View: CREATE VIEW name (a1, a2) AS SELECT ...
 \hookrightarrow like an intermediate relation
- COUNT (\leftrightarrow) counts NULLs
- NULLs → unknown logic value
 \hookrightarrow NOT(U)=U, U AND T=U, U AND F=F, U AND U=U
 $\quad \quad \quad$ U OR T=T, U OR F=U, U OR U=U

Inner Join (default)

\hookrightarrow no match in other relation \Rightarrow no output

Outer Join

\hookrightarrow no match in other relation \Rightarrow dummy record (NULLs)

\hookrightarrow LEFT \Rightarrow every tuple in the left relation appears in the output no matter what

\hookrightarrow RIGHT \Rightarrow pad dangling tuples from right relation

\hookrightarrow FULL \Rightarrow pad dangling tuples from right and left relations

If B has NOT NULL foreign key referencing A

\hookrightarrow A INNER JOIN B = A RIGHT OUTER JOIN B

WITH (like a view for the scope of the statement)

```
WITH partinfo(sid, sname, age, rank, cid) AS (
    SELECT S.sid, S.sname, S.age, P.rank, P.cid
    FROM Skaters S INNER JOIN participates P
    ON S.sid, P.pid
)
SELECT sid, name, cid
FROM partinfo
WHERE age > 7;
WITH A(...) AS
(
    SELECT ...
),B(...) AS
(
    SELECT ...
)
SELECT ...;
```

Derived Table : temp table def within FROM
`SELECT sid, sname, cid`

```
FROM( SELECT s.sid, s.sname, s.age, p.rank, p.cid
      FROM skaters s INNER JOIN participates p
      ON s.sid=p.sid
    )partinfo
WHERE age > 7;
```

Insertion/Deletion:

```
INSERT INTO skaters VALUES (68, 'Jacky', 10, 10);
INSERT INTO skaters (sid, name) VALUES (68, 'Jacky');
INSERT INTO activeSkaters (
    SELECT skaters.sid, skaters.name
    FROM skaters, participates
    WHERE skaters.sid=participates.sid);
DELETE FROM competitions WHERE cid=103;
DELETE FROM competitions;
UPDATE skaters
SET ranking = 10, age=age+1
WHERE name='debby' OR name='lily';
```

Integrity Constraints - every instance of a relation must satisfy. \rightarrow Checks: in create statement

\hookrightarrow Check (constraint)

\hookrightarrow Checks only happen with insert/update, not delete

\hookrightarrow ALTER TABLE matches

```
ADD CONSTRAINT check_dates CHECK (
    mDate >= '2023-07-20' AND mDate <= '2023-08-20');
CREATE VIEW PlayerInfo (name, number, natassoc)
AS SELECT p.name, p.number, t.natassoc
FROM players p, teams t
WHERE p.country=t.country;
```

eg.

SQL examples from P2 & A1

• Stadium names + locations where Christine Sinclair has scored ≥ 1 goal

```
FROM( SELECT matchId
      FROM goals
      WHERE playerId = ( SELECT playerId
                           FROM players
                           WHERE name = 'Christine Sinclair')
      GROUP BY matchId
      HAVING count(*) >= 1)goalcounts, matches m, stadiums s
      WHERE goalcounts.matchId=m.matchId AND m.sName = s.sName;
```

• Info of players who played in all matches

```
SELECT p.name, p.number, p.country
FROM( SELECT playerId, COUNT(*) pcount
      FROM plays_in
      GROUP BY playerId)pc,
      (SELECT country, COUNT(*) tcount
      FROM team_plays
      GROUP BY country)tp, players p
      WHERE p.playerId=pc.playerId AND p.country=tp.country AND
      pc.pcount=tp.tcount;
```

• Country, #matches, #goals scored per team

```
SELECT t.country, COUNT(DISTINCT t.matchid) num_matches,
      COUNT(t.matchid) num_goals
FROM team_plays t LEFT OUTER JOIN goals g
      ON g.matchid = t.matchid AND t.country = g.country AND g.penalty
      = FALSE
      GROUP BY t.country;
```

• total tickets and sold tickets per match

```
SELECT m.matchid, m.sname stadium, COUNT(t.ticketid) total_tickets,
      (SELECT COUNT(ticketid)
      FROM tickets
      WHERE tickets.matchid = m.matchid AND sold = TRUE) AS
      sold_tickets
      FROM matches m LEFT JOIN tickets t ON m.matchid = t.matchid
      GROUP BY m.matchid, m.sname
      ORDER BY m.matchid;
```

• match, minute, scoring player and country for all goals scored in second half

```
SELECT g.matchid, g.minute, p.name, p.country
      FROM goals g, players p
      WHERE g.playerid = p.playerid AND g.minute >= '00:30:00' AND g.penalty
      = FALSE
      ORDER BY g.matchid;
```

• Topic Id and name of topics used by all URLs

```
SELECT topicid, name
      FROM topics t
      WHERE NOT EXISTS ((SELECT w.url
                           FROM webpages w)
                           EXCEPT (SELECT l.url
                           FROM links l
                           WHERE l.topicid=t.topicid))
      ORDER BY topicid;
```

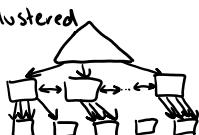
Note on Views: DB does not allow insertion
into complex views that have join conditions
on the view definition

Buffers

- DBMS stores info persistently on disks
- Pin count = # upper layer programs using page
- record ID : rid

Indexing

- helpful for equality and range queries on search key
- B+ Tree: each node = one page
 - leaf contain data entries - one data entry = pointer to one tuple
 - assume leaf and data pages are not in main mem
 - Inserting data - if not enough space, split node and redistribute (50%)
- Indirect Indexing I: $K^* = \langle K, \text{rid of data record w/ search key } k \rangle$
 - gives several entries w/ same search key side by side (more space)
- Indirect Indexing II: $\langle K, \text{list of rids of data record w/ search key } k \rangle$
 - gives a list for each search key
 - variable length data entries, can have large entries that span a page
- Direct Indexing: leaf pages contain data entries instead of pointers
 - $\langle K, \text{full record} \rangle$
 - leaves represent sorted file of data records
- Clustered Index - relation in file sorted by the search key attributes of the index.
- direct index is clustered
- indirect clustered:



- Unclustered Index - relation in heap file or sorted by an attribute other than the Search key
- Indirect unclustered:



- Primary Index: search key contains primary key
- Secondary Index: search key is not unique
- Unique index: search key is primary key or unique attribute

B+ tree cost example:

Given: R(A, B, C) = 4000 bytes
 int = 4 bytes, char[40] = 40 bytes
 size of B ~ Uniform across [0; 19999]
 200,000 tuples

data page = 4000 bytes
 each 75% full
 rid = 10 bytes, ptr = 8
 indirect index II

CALCULATE: # data pages in the heap file: - tuple size = 168,
 - # tuples/page = $4000 \cdot 0.75 / 168 \approx 18$
 - # pages = $200,000 / 18 \approx 11,111 \text{ pages}$

Leaf pages: # data entries: distinct vals of B = 20,000

- # records w/ same val / # rids in data entry: $200,000 / 20,000 = 10 \frac{\text{records}}{\text{val/6}}$
 - size of data entry: size(B) + 10 rids + size(rid) = $4 + 10 \cdot 10 = 104 \frac{\text{bytes}}{\text{entry}}$
 - # data entries per page: $4000 \cdot 0.75 / 104 \approx 29$
 - # leaf pages: $20,000 / 29 = 690 \text{ leaf pages}$

height of the tree: size of index entry in intermediate node = $4 + 8 = 12 \text{ bytes}$
 max num of index entry per intermediate page: $4000 / 12 = 333$
 single root can't hold all pointers \Rightarrow height > 1
 690 pointers to leaf pages, int. nodes 75% full \Rightarrow 3 intermediate nodes
 height = 2 (see pic \rightarrow)

DB2: CREATE INDEX ind1 ON Students (startyear);

DROP INDEX ind1;

referential integrity/ uniqueness

CREATE UNIQUE INDEX ind1 ON Students(email)

additional attributes:

CREATE UNIQUE INDEX ind1 ON Students(sid)
 INCLUDE (name)

Index on sid, entry contains sid, name, rid \downarrow doesn't need to access data pages
 SELECT name FROM Students WHERE sid = 100

clustered index:

CREATE INDEX ind1 ON Students (name) CLUSTER

multi-attribute index:

CREATE INDEX ind2 ON Students(email,sid)

Quiz 4 - Buff. Mgmt and idx:

If a frame is dirty, it means the content of the page that resides in the frame is different to the content of the page as stored on stable storage

pin counter is dictated by how many requests are using a particular page as multiple requests could be accessing the same page

page = 4000 bytes, 2 bytes per pointer/tracker, records = 140¹⁰⁰⁰, no empty slots

↳ Directam: $100 \cdot 2 + 2 = 204 \text{ bytes} \Rightarrow \text{space} = 4000 - 204 = 3796$
 $3796 / 40 = 94 \text{ records} \Rightarrow 94 \cdot 40 + 94 \cdot 2 + 2 \cdot 2 = 3950 \Rightarrow \text{space for +1}$
 $\Rightarrow 95 \cdot 40 + 95 \cdot 2 + 2 \cdot 2 = 3994 \quad \boxed{95 \text{ records max per page}}$

leaf node = 100 data entries, interm. node = 140 index entries

↳ find max # records indexed by tree of height 2: root holds 140, each interm holds 140, each leaf holds 100 $\Rightarrow 140 \cdot 140 \cdot 100 = \boxed{196,000}$

RLA, B.C (each size=4), index on A w/ attr B, page = 4000, ptr = 8
 ↳ find # index entries per node (tree 75% full)

$8 + 4 = 12 \text{ bytes per index}, 4000 \cdot 0.75 = 3000, 3000 / 12 = \boxed{250}$

Indirect Index II on A (size=10), each val of A occurs 4 times, page = 4000, ptr = 8
 ↳ find # data entries per leaf page (75% full)

$4000 \cdot 0.75 = 3000, \text{size(entry)} = 10 + 4 \cdot 10 = 50, 3000 / 50 = \boxed{60}$

Indirect Index II on A, B, A=4, B=6, 100000 records, 5 records
 ↳ find # leaf pages: size(data entry) = $10 + 5 \cdot 6 = 60, 4000 \cdot 0.75 / 60 = 50$ entries
 $100000 / 5 = 20000 \text{ data entries total}, 20000 / 60 = \boxed{400 \text{ pages}}$

A multi-attribute B+ index on attributes A,B of a relation always has a larger space requirement than an index on A with additional attribute B

You can have both an indirect, clustered index on A of relation R and an indirect non-clustered index on Attributes A,B of R.

Multi-attribute index - order matters → ordered into groups by one then ordered within groups by other. Supports queries on first attribute or both but is not useful for only second attribute (full scan)

leaves

Query Evaluation

- assume root and intermediate nodes are in main mem
- execution tree: leaf nodes are the relations, inner nodes are operators

↳ e.g.

Schema for examples:

Users(U): $\text{card}(U) = 40,000$ tuples

- 80 ~~tuples~~ data page, 500 data pages

- Index on uid has 170 leaf pages

- Index on uname has 300 leaf pages

Group members(GM): $\text{card}(GM) = 100,000$

- 100 ~~tuples~~ data page, 1,000 data pages, no index on gid (500 groups)

Selection

- reduction factor: $\text{Red}(\sigma_{\text{condition}}(R)) = |\sigma_{\text{condition}}(R)| / |R|$

- e.g. (assume 10,000 w/ exp 5) $\text{Red}(\sigma_{\text{exp=5}}(U)) = 10,000 / 40,000 = 0.25$

- unknown # tuples for condition \Rightarrow assume uniform + independent distribution

↳ e.g. $\text{Red}(\sigma_{\text{exp=5}}(U)) = \frac{1}{10} = 0.1$

- red = size of selected range / total size of domain

- e.g. $\text{Red}(\sigma_{\text{exp=5 and age <= 16}}(U)) = \text{Red}(\sigma_{\text{exp=5}}(U)) * \text{Red}(\sigma_{\text{age <= 16}}(U))$

- Size of result of Selection = (# input tuples)(reduction factor)

↳ e.g. rf = 0.1, 40,000 tuples \Rightarrow |result| = $40,000 \cdot 0.1 = 4000$ tuples

Simple Selections:

↳ no index: search on arbitrary attribute: scan relation (cost = # data pages)

↳ search on primary key attribute: Scan half the relation (cost = $\frac{\text{# data pages}}{2}$)

↳ index on selection attributes: cost = # leaf pages w/ matches + num data pages w/o match

↳ clustered B+ Tree: cost = 1 leaf page + (# matching tuples) / (tuples per page)

- e.g. $\text{SELECT * FROM Users WHERE uid=123}$

cost = 1 leaf page + 1 data page

- e.g. $\text{SELECT * FROM Users WHERE uname LIKE 'B%'}$

cost = 1 leaf page + 2 data pages (100 tuples match)

- e.g. $\text{SELECT * FROM Users WHERE uname < 'F'}$

cost = 1 leaf page + 125 data pages (10,000 tuples match)

↳ non-clustered B+ Tree: worst case: # data pages = # matching tuples

- e.g. $\text{SELECT * FROM Users WHERE uid=123}$

cost = 1 leaf page + 1 data page (point query)

- e.g. $\text{SELECT * FROM Users WHERE uname LIKE 'B%'}$

cost = 1-2 leaf pages + 100 data pages (bc data not ordered, still better than scan)

- e.g. $\text{SELECT * FROM Users WHERE uname < 'F'}$

cost = 75 leaf pages + 10,000 pages (worse than scan)

- non-clustered only useful with very small reduction factors

- if all desired leaf pages fit in mem then can sort to speed up

Selections on multiple attributes

↳ e.g. A=100 AND B=50 (500 pages)

- no index \Rightarrow scan, cost = 500

- index on A \Rightarrow cost = query for A = 100 (get A, check value of B)

- one index per attribute \Rightarrow intersection based: find rids w/ A=100 (A index), find rids w/ B=50 (index B), intersect rids. (good if all fit into mem)

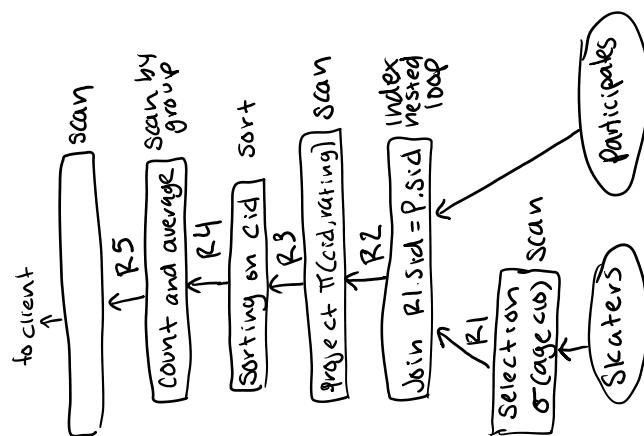
↳ use one index: use index on A if A is unique or has small rf

- one index with both attributes = good

↳ e.g. A=100 OR B=50

- no index \Rightarrow Scan, cost = 500 - 2 indexes \Rightarrow Union of rids

- index on A \Rightarrow Not useful - 1 index w/ both \Rightarrow may or may not be faster



• Sorting: $N = \# \text{ pages}, B = \# \text{ blocks}$

↳ $\# \text{ passes} = 1 + \lceil \log_2 \lceil N/B \rceil \rceil$

↳ cost = $2N(\# \text{ passes})$

Joins

- cardinality estimations

↳ $|U \bowtie GM| = |GM|$

↳ join attribute is a primary key for U, foreign key for GM \Rightarrow each GM matches with one U

↳ $|R1 \times R2| = |R1| + |R2|$

↳ joins with conditions \rightarrow reduction factor of condition * cardinality of join

- Simple Nested Loop Join: for each tuple in the outer relation, scan the entire inner relation

↳ cost = #outer pages + |outer| * #inner pages

- page Nested loop join: for each page P_o of the outer relation, get each page P_i of the inner relation

↳ cost = #outer pages + (#outer pages) * (#inner pages)

- block Nested loop join: for each block of pages b_{po} of the outer relation, get each page P_i of the inner relation

↳ b_{po} and one P_i must fit in main memory

↳ cost = #outer pages + #outer pages / # b_{po} * #inner pages

- Index nested loop join: for each tuple in the outer relation, find the matching tuples of the inner relation through an index (inner must have an index on join attr)

↳ cost = #outer pages + |outer relation| + cost finding

Block vs Index Nested Loop

- best case for block: cost = #outer pages + #inner pages

↳ index better - $|inner| > |outer| + \text{cost of matching}$

- Sort-Merge join: sort relations on join column, then scan to merge. Outer is scanned once, each inner is scanned once per matching outer

↳ if relations are already sorted:

↳ cost = #outer pages + #inner pages

↳ if relations need to be sorted:

↳ cost = $(2(\text{num passes}) + 1)(\# \text{outer pages})$

+ $(2(\text{num passes}) + 1)(\# \text{inner pages})$

↳ Sort by pipelining:

↳ cost = $(2(\text{num passes} - 1) + 1)(\# \text{outer pages})$

+ $(2(\text{num passes} - 1) + 1)(\# \text{inner pages})$

Query Evaluation cont.

- Projections: done on the fly with another operation
- Union: sort \rightarrow scan \rightarrow merge
- Aggregate operations: usually done last
 - \hookrightarrow without grouping \Rightarrow scan
 - \hookrightarrow with grouping \rightarrow sort on group-by attribute then scan relation and compute aggregate for each group
- Pipelining: as soon as a tuple is determined, it is forwarded to next operation
 - \hookrightarrow parallel execution of operators
 - \hookrightarrow no materialization of intermediate relations if it can be avoided
- Algebraic Optimization (use a tree)
 - \hookrightarrow perform operations that eliminate a lot of tuples
 - push down projections and selections (careful with projections)
 - \hookrightarrow consider # tuples that flow between operators
 - \hookrightarrow e.g. $\Pi_{\text{uname}, \text{exp}}(\sigma_{\text{uid}=\text{G1} \wedge \text{age} > 20}(\text{U} \bowtie \text{GM}))$
 $\Rightarrow \Pi_{\text{uname}, \text{exp}}(\sigma_{\text{age} > 20}(\text{U}) \bowtie \Pi_{\text{uid}=\text{G1}}(\text{GM}))$
- Cost-based Optimization
 - \hookrightarrow alternative plans created by algebraic optimization
 - \hookrightarrow considers how operators are executed
 - \hookrightarrow process:
 - Pass 1: Find best 1-relation plan for each relation
 - Pass 2: Find best way to join results of each 1-relation to another relation (consider inner and outer)
 - Pass N: Find best way to join results of (N-1)-relation plan to the Nth relation
 - prune high cost alternatives
 - bottom-up calculation each possible alternative; estimate the cost

optimization example:

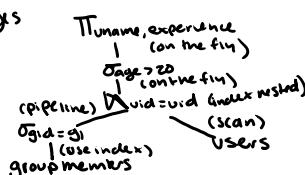
```
SELECT u.uname, u.experience
FROM Users u, Groupmembers g
WHERE u.uid = g.uid AND g.gid='G1'
AND u.age>20
```

Pass 1:

- users: scan is cheapest since index is non-clustered and age>20 has high RF
 - \hookrightarrow cost = 500, output = 20,000 tuples
- groupmembers: index on gid matches gid=g1 \Rightarrow $\frac{100000 \text{ records}}{500 \text{ groups}} = 200 \text{ tuples}$
 - \hookrightarrow worst case cost = (1-2 leaf pages) + 200 data pages \times 201-202
 - \hookrightarrow Scan \rightarrow cost = 10000
 - \hookrightarrow output 200 tuples = 2 pages (leave in mem)

Pass 2:

- users as outer:
 - merge sort join: sort 200 tuples of GM (in main mem), sort 500 pages of U and combine merge join w/ pass 1 of sort: 2 passes $\Rightarrow 3 \cdot 500 = 1500$
 - Block nested loop join: inner relation = 200 tuples = 2 pages in main mem
 read in U once: 500 pages
 - Index nested loop join on uid not possible
- Group member as outer:
 - Index nested loop: cost = $200 \cdot (1 \text{ leaf} + 1 \text{ data}) = 400 \text{ pages}$ \leftarrow best
 - Block nested loop: read in U once: 500 pages
 - Merge sort same as above
- Selection + projection on the fly



Assignment 2

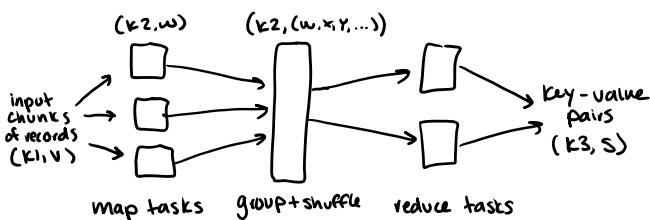
- # Data Pages required for a relation:
 - \hookrightarrow # tuples/page = (data page size) * (fill) / (size of tuple)
 - \hookrightarrow # pages = (total # tuples) / (# tuples per page)
- # leaf pages for an index (indirect II)
 - \hookrightarrow # data entries/page = (data page size) * (fill) / (data entry size)
 - \hookrightarrow # leaf pages = (total # data entries) / (# data entries per page)
- I/O for a conditional selection w/ index ...
 - \hookrightarrow non-clustered indirect type II index on desired attribute: $\underbrace{(\# \text{leaf pages})(\% \text{ matching tuples})}_{\text{worst case}} + (\text{total } \# \text{ tuples})(\% \text{ matching})$
 - \hookrightarrow non-clustered indirect (II) index on 2 attributes
 - (#leaf pages matching attr 1) + (#leaf pages matching Attr2) + $\frac{(\# \text{tuples})}{(\% \text{ matching})}$
- Group by / Selection query
 - \hookrightarrow no index: Scan tuples, only desired attributes (tuple size \cdot # tuples/page)
 - (if it all fits in main mem), sort, group by, select: $\text{I/O} = \# \text{ tuples}$
 - \hookrightarrow Index on group by attribute: Scan all leaf pages + select ($\text{I/O} = \# \text{ leaf pages}$)
- Index nested loop join (clustered)
 - \hookrightarrow $\text{I/O} = \# \text{outer pages} + |\text{outer relation}| \times (\text{cost of finding matching})$
- Block nested loop join
 - \hookrightarrow use smaller relation as outer.
 - $\text{I/O} = \# \text{outer pages} + (\# \text{outer pages} / \text{block size}) \times (\# \text{inner pages})$
- Sort merge join

Large Scale Data Processing

- throughput: transactions/queries per time unit
- response time: time for execution of a single transaction
- OLTP (online transaction processing): workload of short, update intensive queries
- OLAP (online analytical processing): workload of complex queries
- Inter-query parallelism: different queries run in parallel on different processors
- Inter-operator parallelism: different operators w/in same execution tree run on different processors. Pipelining leads to parallelism
- Intra-operator parallelism: single operator runs on many processors (scan, join)
- Data partitioning: partition into chunks of records stored at nodes
 - ↳ Hash function or range partition
- Execution path: push operator to nodes w/ chunks
 - execute locally at each node in parallel
 - send results to coordinating node
 - merge results
 - minimize I/O and communication costs
- Vertical partitioning/column stores \Rightarrow less I/O

Map-reduce

- 1) Distribute data (key-value pairs)
- 2) Map tasks: extract info and output new key-value pair
- 3) Shuffle + sort
- 4) Reduce tasks: aggregate, summarize, filter
- 5) Write results



- input/output are key/value pairs
- map + reduce are written by us
- group + shuffle are done by the system
- example: word count, DS(K, doc)

wordcount map:

For each input key/value pair (dkey, dtext)
 For each word w of dtext
 output key-value pair (w, 1)

System groups eg (K, (v₁, v₂, v₃, ..., v_n))

WordCount Reduce:

For each input key/value-list pair (K, (v₁, ..., v_n))
 output (K, n)

Phases:

- ↳ map tasks: record reader \rightarrow map function \rightarrow write to local file
- ↳ group + shuffle: group keys and aggregate value-lists \rightarrow copy from map to reduce
- ↳ reduce tasks: reduce \rightarrow write to file system

Implementation:

- ↳ Master node controls execution, partitions file into m parts, assigns workers to m map tasks
- ↳ Map workers execute map tasks and write to local disk
- ↳ Master assigns workers to r reduce tasks
- ↳ Reduce workers implement group and shuffle and execute reduce tasks
- master handles failures (assigns new workers)

• Selection w/ MapReduce

$\hookrightarrow R(A, B, C)$, selection w/ condition c on R

- map: for each tuple t of R for which c holds, output (A, (B, C))
- reduce: identity (output (A, (B, C)))

• Join w/ Map-reduce

\hookrightarrow natural join $R(A, B, C)$ with $Q(C, D, E)$

- map: For each tuple (a, b, c) of R, output (c, (R, (a, b)))
 For each tuple (c, d, e) of Q, output (c, (Q, (d, e)))

- group + shuffle aggregates key-value pairs

- reduce: For each tuple (c, value-list)

$$R_t = Q_t = \text{empty}$$

for each v = (val, tuple) in value-list

if v.ref = R; insert tuple into R_t

else insert tuple into Q_t

for v1 in R_t, for v2 in Q_t, output (c, (v₁, v₂))

• Projection with map-reduce

\hookrightarrow projection on B, C of R

- map: for each tuple t = (a, (b, c)) of R, let t' = (b, c)
 output (t', 0)

- group + shuffle

- reduce: for each tuple (t', (0, 0, 0, ...)), output (t', 0)

• Group by w/ Map-reduce

\hookrightarrow grouping: SELECT b, max(c) GROUP BY b

- map: for each tuple (a, (b, c)) of R, output (b, c)

- group + shuffle produces key-value list

- reduce: for each (b, (c₁, c₂, ...)) perform aggregation

Declarative Languages

- pig latin - run complex queries that require several map-reduce jobs

• E.g.: Users = load 'users' as (name, age);

Fltrd = filter Users by age >= 18 and age <= 15;

Pages = load 'pages' as (uname, url);

Jnd = join Fltrd by name, Pages by uname;

Grpd = group Jnd by url;

Smmld = foreach Grpd generate (\$0), COUNT(\$1) as clicks;

Srted = order Smmld by clicks desc;

Top5 = limit Srted 5; store Top5 into 'top5sites';

• Load: read info from a file into a temp relation

• Store: write relation into file

• Selection: Res = filter R1 by ...

• Join: Res = join R1 by a1, R2 by a2

• Order by: Res = order R1 by a1 desc

• Group by: Grpd = group Rel by A

• For each:

Smmld = foreach Grpd generate (\$0), COUNT(\$1) as c

Smmld = foreach Grpd generate group, Count(Rel)

• Projection: Rel = foreach R1 generate A, B;

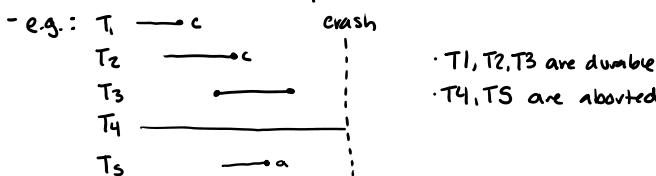
• Flattening: Res = foreach R generate \$0, flatten(\$1)

↳ make \$1 appear as part of list, not nested

• Dump to screen

Transactions

- Database transaction: series of reads and writes on objects
- OLTP: update intensive (15-50% of SQL statements are updates)
- ACID:
 - ↳ Atomicity - a transaction is atomic if it is executed all or none
 - ↳ Consistency - transaction preserves the consistency of the data
 - ↳ Isolation - in case that transactions are executed concurrently, the effect must be the same as if one at a time
 - ↳ Durability - changes made by a transaction must be permanent
- Atomicity: after a commit, all changes are installed in the DB.
 - after an abort → local recovery, eliminate partial results / undo all modifications so far
- Durability: there must be a guarantee that the changes will last - i.e. survive failures / crashes. Persistence.



- System crash recovery
 - ↳ restart server and perform global recovery:
 - transactions that committed before crash: changes in DB
 - transactions that aborted: no changes in DB
 - transactions that were active: no changes in DB
- Each update changes the records corresponding page
 - ↳ dirty pages are written to disk at buffer manager's discretion
- To undo/redo in case of crash - write log info to a log disk
- Write ahead logging (WAL):
 - ↳ for each write(x), append log record with before and after image of x to log tail. Undo changes of active transactions w/ before image
 - ↳ flush log to disk before flushing the page. (so you can undo in case of abort)
 - ↳ flush log to disk before commit (so you can redo when something committed)
 - ↳ log is an append; flush means all of log is flushed in one shot
- Recovery
 - ↳ Undo: for each transaction that did not commit, find log records and install before images.
 - ↳ Redo: for each transaction that committed, find all log records and install after images.

Isolation and Concurrency Control

- Isolation: transactions execute concurrently but each runs in isolation - and do not affect each other.
 - ↳ enforced by concurrency control
- Concurrency control provides serializable executions: affects of concurrent execution is the same as running one after the other.
- Schedule: sequence of actions (w, r, c, a) from a set of transactions
- Complete schedule: contains commit/abort for each transaction
- Serial schedule: transactions executed fully one after the other
- Dirty reads: T2 reads from T1 before T1 commits (inconsistent data)
 - ↳ e.g. T1: R1(A) W1(A) R1(B) W1(B) C
T2: R2(A) R3(B) C3
 - like T2 executes after T1 (reads what T1 wrote)
 - like T2 executes before T1 (reads B before T1 wrote B)
- Unrepeatable reads
 - ↳ T1 reads an item twice but T2 changes the value inbetween
 - ↳ e.g. T1 R1(A) R1(A) C1
T2 W2(A) C2
 - T1 conceptually executes before and after T2
- Lost update
 - ↳ e.g. T1 R1(A) W1(A) C1
T2 R2(A) W2(A) C2
 - as if T2 never happened
- Conflicting operations: 2 operations (from different transactions) if they access the same object and both operations are writes, or one is write and one is read
- Conflict Serializable Schedules
 - ↳ two schedules are conflict equivalent if they ① involve the same actions of the same committed transactions and ② every pair of conflicting actions of committed transactions is ordered the same way
 - ↳ a schedule is conflict serializable if it is conflict equivalent to some serial schedule
- e.g.: T1 r1(x) w1(x) w1(y) C1
T2 r2(x) r2(y) w2(x) C2
is conflict serializable with
T1 r1(x) w1(x) w1(y) C1
T2 r2(x) r2(y) w2(x) C2
- Serializability + conflict graphs:
 - Let S be a schedule (T, O, C)
 - each transaction T_i in T is represented by a node
 - there is an edge from T_i to T_j if an operation of T_i precedes and conflicts with an operation of T_j
 - e.g. T1: R1(A) W1(A) → R1(B) W1(B) C1
T2: R2(A) W2(A) → R2(B) W2(B) C2
 $T_1 \rightarrow T_2$
 - T1: R1(A) W1(A) → R1(B) W1(B) C1
T2: R2(A) W2(A) R2(B) W2(B) C2
 $T_1 \leftrightarrow T_2$
- Dependency Graphs: schedule is conflict serializable iff its dependency graph is acyclic
 - ↳ $T_1 \rightarrow T_2 \rightarrow T_4 \rightarrow T_3$
 $T_1 \rightarrow T_4 \rightarrow T_2 \rightarrow T_3$
- Concurrency control: during execution, take measures such that a non-serializable execution can never happen
- Locking: each transaction must obtain an S (shared) lock on object before reading and an X (exclusive) lock on object before writing. If an X lock is granted to an object, no other lock can be granted on the object at the same time. If an S lock is granted, no X lock can be granted on the same object at the same time. If a conflicting lock is active, the transaction must wait until the lock is released at the end of the other transaction.
 - ↳ ZPL
- Phase 1: Growing + locking, Phase 2: Shrinking + unlocking at end