

$[R]$: the # of pages to store R

P_R : # of records per page of R

$|R|$: # of records of R ($P_R * [R]$)

$SNLJ : [R] + |R| * [S]$

R should be smaller

$PNLJ : [R] + ([R] * [S])$

$BNLJ : [R] + \lceil [R] / (B-2) \rceil * [S]$

$INLJ : |R| * (\text{lookup cost}) + [R]$

lookup = cost to get to leaf +
Alt 1: 0

Alt 2/3 clustered: 1 per matched page of S

Alt 2/3 unclustered: 1 per matched tuple of S

Naive Hash Join

Load all R into hash table

Scan S and probe R

$R < (B-2) * \text{hash_fill}$

Grace Hash Join

Good if input hashed, need output hashed

- Partition tuples from R and S by join key
- Build and probe for each partition

cost: $2 * k * ([R] + [S]) + [R] + [S]$

k = # of passes to get partitions of R to fit in mem

$R < (B-1)(B-2) \Rightarrow k=1$

Partitioning Phase: divides R into $(B-1)$ runs of size $[R]/(B-1)$

Matching Phase: requires each $R/(B-1) < (B-2)$

Does not include cost of writing matching output

Sort-Merge Join

Requires equality predicate (Equi-joins & natural joins)

cost: Sort R + Sort S + $(|R| + |S|)$

last term $([R] * [S])$ worst case (all tups match)

$B > \sqrt{\max([R], [S])} \Rightarrow$ sorted in 2 passes

$B * (B-1) > \max([R], [S])$ not sensitive to data skew

Parallel - Increase HW
- Fixed workload throughput ↑

Scale up - Increase HW
Increase HW throughput / const

Shared Nothing: runs on commodity HW
scales up w/ data
doesn't rely on HW

Inter-query - Treat each query separately → pipeline, bushy tree

Intra-query - Inter-operator, intra-operator

↳ Partition Parallelism

Data Partitioning

Range, Hash, Round-Robin

Symmetric Hash Joins

1. Each node allocates two hash tables, one for each side

2. For each tuple R

- a. build into R for join key
- b. probe into S for matches

3. For each tuple S

- a. symmetric to R

Hashes

- selection cascade and pushdown

- projection " " "

- avoid cartesian products

Physical Equivalences

- base table access - Heap scan, Index scan

- Equi-Joins - Block, Index Nested, SIM Join, GII Join

- Non-Equi → block Nested Loop

System R

Cost: # I/Os + CPU factor * # tuples

Catalog

NTuples # of tuples in table

NPages # of disk pages in a table

LowHigh min/max value in a col

Nkeys # of distinct in a col

IHeight height of an index

INPages # of disk pages in an index

Parallel Sort Merge Joins

one for each relation

Pass 0...n-1 one like parallel sorting

Pass n: merge join partitions locally on each node

Parallel Aggregation

$\text{sum}(S) = \sum_{i=1}^n S_i$

$\text{count} = \sum_i \text{count}(S_i)$

$\text{avg}(S) = \sum \sum(S_i) / \sum \text{count}(S_i)$

Asymmetric shuffle

- only partition S since R

is already partitioned properly

- If R partitioned, just partition S, then run local join at every node and union the result.

Term col = col 2 (handle for joins)

sel = 1 / Nkeys(I)

Term col1 = col 1 (handle for joins)

sel = 1 / Max(Nkeys(I1), Nkeys(I2))

Term col > value

sel = $\frac{\text{high}(I) - \text{value}}{\text{high}(I) - \text{low}(I) + 1}$

needs index on both otherwise

$\frac{1}{Nkeys(I)}$

Subtract ⋅ for OR

Broadcast Joins

- R is small, send R to every partition of S

- Do local join then union

Parallel sort

Pass 0: shuffle data across machines

Pass 1-n: done independently, single-node sorting

Pipeline Breaker

Sort, still can't start until sorted

Hash build, each hash can't start probing until hashtable built

Query Parser: checks correctness, auth generates a parse tree

Query rewriter: converts queries to canonical form (flatten views, subqueries into fewer blocks)

Query Opt

Plan space, cost estimation, search strategy

Alternative Plans

Single table

↳ selects, projects, groupBy/Agg

↳ fileScan / index

Cost Estimates for Single Plan

Index I on primary key

matches selection:

cost = $(\text{Height}(I) + 1) + 1$

for B+Tree

Clustered Index I matching

selection:

$(NPages(I) + NPages(R)) * \text{selectivity}$

Non-clustered Index I matching

selection:

$(NPages(I) + NTuples(R)) * \text{selectivity}$

Sequential Scan of file: $NPages(R)$

Recall: Must charge for dup elimination if required

Left deep plans differ in

- order of relations
- access method for each leaf operator
- the join method for each join operator

Enumerated using N passes

Pass 1: Find best 1-relation plan for each relation

Pass i : Find best way to join result of an $(i-1)$ -relation plan to the i -th relation. (i between 2 and N)

Keep cheapest plan overall and each interesting order

Interesting Orders Index scans outputs sorted on index column

Sorted by anything that we can use later

↳ ORDER BY / GROUP BY

↳ Join attributes of yet-to-be-added joins

IR's Bag of Words

Given corpus of text file, build table

Inverted file (term -> text, docID -> text)

keep posting list like A[3] sorted by DocID

near → add position to inverted file and posting list, check that positions are k off

Hellstrom Inequality

$$\frac{d_{app}}{dt} \ll \frac{d_{env}}{dt} \quad \text{Data indep. is most important when this holds.}$$

Logical data indep.

Maintain views when logical struct changes

Physical data indep.

Maintain logical struct when physical struct changes

Data model = collection of concepts for describing data

Schema = description of a particular collection of data, using a given data model

Relational model of data

- main concept = relation, rows and cols
- every relation has a schema

2PC Recovery

Coordinator recovers
wl commit and encl → nothing
wl commit only → re-run phase 2
wl "abort" → Nothing

Participant recovers:

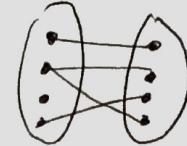
- no prepare/commit/abort? Nothing
- wl prepare & commit? send Ack to coord
- wl just prepare? send inquiry to coordinator
- wl "abort"? Nothing

Commit iff coordinator logged a commit

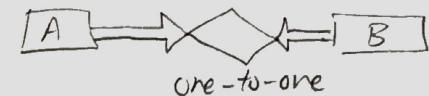
Selectivity Est

$$c >= v \Rightarrow \frac{(highkey - v)}{(highkey - lowkey + 1)} + \frac{1}{\# \text{ of distinct values of } c}$$

$$c < v \Rightarrow \frac{(v - lowkey)}{high - low + 1}$$



1-to-many | many-to-1



Inequalities & Float

$$c >= v \Rightarrow \frac{(high - v)}{high - low}$$

$$c <= v \Rightarrow \frac{(v - low)}{high - low}$$

ER

Entity - Real world obj described by attrs

Entity set - collection of the same type of entities

Each entity set has a key

Each entity has a domain

Relationships - Association among 2+ entities

Relationship set - collection of relationships involving the same entity sets

↳ An n-ary relationship set R relates n entity sets E₁, ..., E_n. Each relationship in R involves e ∈ E₁, ..., e_n ∈ E_n

key constraint
(at most one) →

participation constraint
(at least one) →

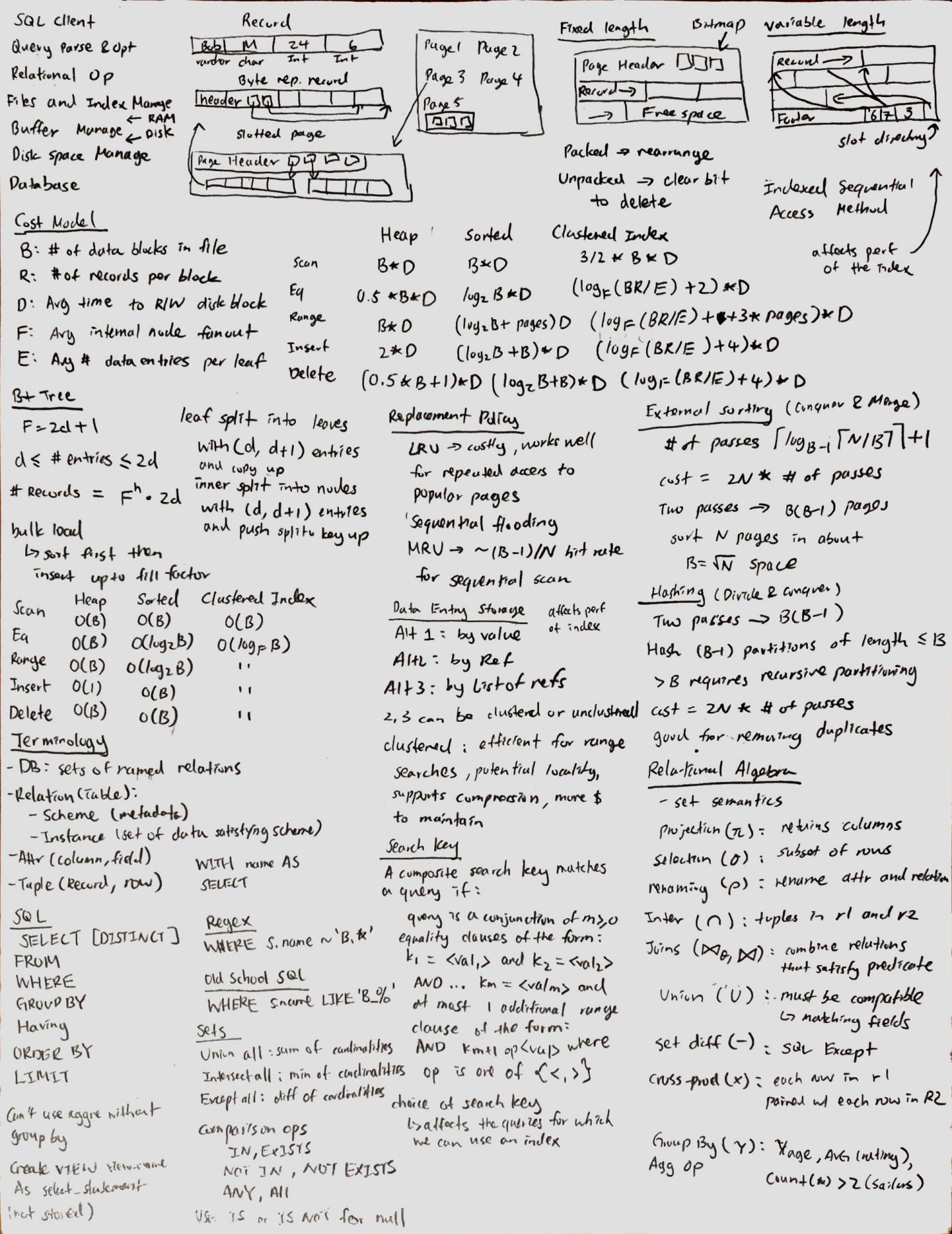
key constraint w/ total participation
(exactly one) →

Non-key partial participation
(0 or more) →

Weak Entities
- Identified uniquely only by considering primary key of another entity

↳ owner set and weak entity set must be one-to-many
(owner) (weak entities)

↳ weak entity must have total participation in the identifying relationship set



Cost of External I/O

$$B=5 \quad N=108$$

Pass 0: $\lceil 108/5 \rceil = 22$ sorted runs of 5 pages each

Pass 1: $\lceil 22/4 \rceil = 6$ sorted runs of 20 pages each

Pass 2: $\lceil 6/1 \rceil = 2$ sorted runs of 80 and 28 pages

Pass 3: sorted 108 pages

Classical IR Ranking: vector space model, similarity is the cos of the angle between two normalized vectors or dot product

$$TF \times IDF \text{ (Term Freq} \times \text{Inverse Doc Freq)} = \text{DocTermRank}$$

occurrence counts aren't good term weights, want to measure repeated or unusual word, $IDF = \log(\frac{\text{total # docs}}{\# \text{docs w/ term}})$

$$\text{Precision} = \frac{\text{TP}}{\text{Selected}} \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

FDS

Superkey: set of columns that determines all cols in table

Candidate key: a minimal set of columns that determines all cols

$K \rightarrow \{ \text{all attributes} \}$, for any $L \subset K$, $L \rightarrow \{ \text{all attributes} \}$

Update anomaly: need to update all records w/ same key

Insertion anomaly: need to insert but don't know the val?

Deletion anomaly: If we delete all record w/ a value, we lose info

F^+ = closure of F : the set of all FDs implied by F (includes trivial dependency)

Reflexivity: If $X \geq Y$, then $X \rightarrow Y$ Augmentation: If $X \rightarrow Y$, $XZ \rightarrow YZ$

Transitivity: $X \rightarrow Y, Y \rightarrow Z$, then $X \rightarrow Z$ Union: $X \rightarrow Y, X \rightarrow Z$, then $X \rightarrow YZ$ Decomp: $X \rightarrow YZ$, then $X \rightarrow Y$ and $X \rightarrow Z$

Find closure:

If you can guess missing val, then not BCNF

decomp of R into X and Y is lossless w.r.t. F iff the closure of F contains $X \cap Y \rightarrow X$

$$X+ = X$$

Repeat until no change:

for $V \rightarrow V \subseteq F$,

if $V \subseteq X+$, then add V to $X+$

Loss less if joining decomposed relations gets back original

Decompose R into X and Y is dep preserving if $(F_X \cup F_Y)^+ = F^+$

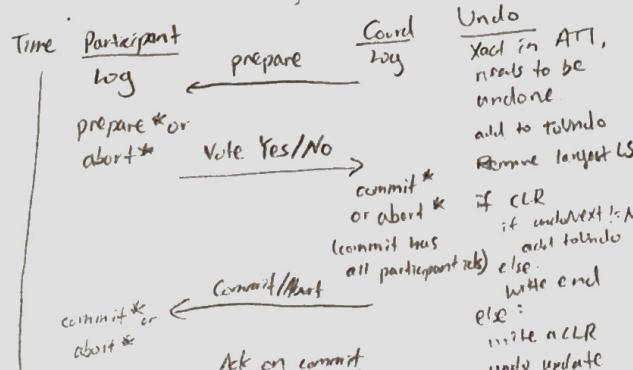
Sort-Merge Refinement

pass 0 - k-2 external sort R and S until almost sorted

pass k-1 merge R and S

$$2 * (k-1) * ([R] + [S]) + [R] + [S]$$

2 pass requires $B \geq TR + JS$



* wait for log flush before sending next msg

Insert/Delete

- $\hookrightarrow 1$ to read record
- $\hookrightarrow 1$ to write record
- $\hookrightarrow 1$ to update leaf node

Eq Search

- $\hookrightarrow 1$ to read the leaf
- \hookrightarrow add 1 to read actual records

Range Search

- $\hookrightarrow \frac{3}{2}$ to read records > 3 pages

- $\hookrightarrow \frac{3}{2}$ to read leaf
(2/3 Full Factor)

If BCNF \rightarrow 3NF
but not other way around

SQL Junks

```
SELECT H.home_id, T.sale_price
FROM Homes H
LEFT OUTER JOIN Transactions T
ON H.home_id = T.home_id
WHERE H.city = "Berkeley"
```

```
SELECT T1.dcid
FROM Trip AS T1
WHERE 5.0 = All()
```

```
SELECT T2.driver_rating
FROM Trip AS T2
WHERE T1.dcid = T2.dcid AND DriverRating IS NOT NULL
```

May not be a dep preserving decomp into BCNF

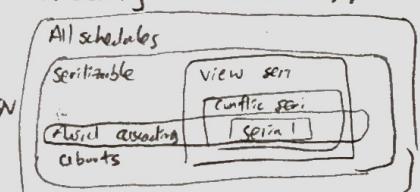
	IS	IX	S	SIX	X
IS	✓	✓	✓	✓	
IX	✓	✓			
S	✓			✓	
SIX	✓				
X					

Atomicity: All or None

Consistency: DB starts and ends consistent

Isolation: Exec of each transaction is isolated

Durability: if Xact commits, persist



$T_i > T_j$? wait : die

$T_i > T_j$? wound : wait

No Steal Steal

No Undo Redo	UNDO REDO
No Undo No Redo	UNDO NO Redo

Redo unless any of this is true

- Affected page not in DPT

- Affected is in DPT, but recLSN > LSN

- PageLSN (in DB) > LSN