

# CS 186 FINAL

## Cheat Sheet

### SQL:

- WHERE auto-filters NULL entries in col.
- SELECT <cols> FROM <table> AS <name> WHERE <predicate> GROUP BY <cols> HAVING <predicate> ORDER BY <cols> LIMIT <num>

- Only select grouped/aggregated cols after GROUP BY.
- WHERE <string col> LIKE "... compares to regex string.  
["%"] = 0+ characters and "\_" = 1 character]
- FROM A, B is cross product join condition
- FROM A INNER JOIN B ON <predicate> [cross-join with join cond. in WHERE = inner join]
- A LEFT OUTER JOIN B ON <...> preserves A's cols. [similar RIGHT OUTER JOIN, FULL OUTER JOIN]
- A NATURAL JOIN B = inner join with implied condition on columns with same name

```

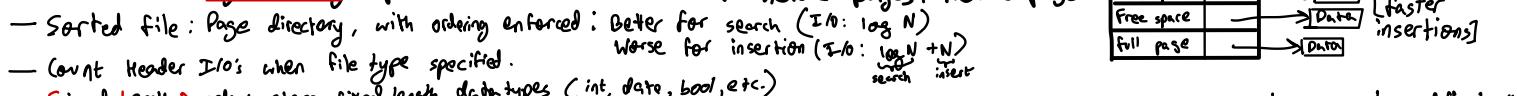
Subqueries
- SELECT *
  FROM classes
  WHERE EXISTS(
    SELECT *
      FROM enrollment
      WHERE classes.num
        = enrollment.num
  );
  ) WITH <name> AS (
    <query>
  ) ...

```

### Disk & Files:

- Disk : R/W arm assembly on platters spinning 15k rpm. SSD is stationary, faster alternative.
- Heap file : Linked-list implementation: Header Page → Full pg 1 → Full pg 2 → ... full pages [Remember to update header on insertions]

Page directory implementation: Linked list of header pages. Header page =



Sorted file: Page directory, with ordering enforced: Better for search (I/O: log N) Worse for insertion (I/O: log N + N)

Count header I/O's when file type specified.

Fixed Length Records: store fixed length data types (int, date, bool, etc.)

Variable Length Records: store fixed first then variable length (varchar) with pointers to end of varchar. VLR in some schema can have diff length

Packed FLR pages: records stacked, store pointer to end of last record. — max capacity of VLR page:

Unpacked FLR pages: slots stacked store bitmap of what slots are free.

Varchar (n) can have 0-X bytes size

VLR pages have footer (slot directory, slot count, free space ptr)

— Ptrs are int-size(offsets) [ptr, record len] — Each VARCHAR has a 4-byte ptr: slot count + 4 + (4 + 4) # records

— slot dir size :

page size - 8 [with int = 4 bytes]  
record size + 8

slot count + 4 + (4 + 4) # records  
free space ptr per len

### B+ Trees

- Each node has [d, 2d] entries assuming no deletes (so [d+1, 2d+1] child ptrs) [- Leaf nodes contain nextLeaf ptr.]
- If insertion overflow: split node into d, d+1 records. copy if leaf, else push.
- Only delete within leaf nodes.
- Alt 1 By value: leaf nodes contain records. (key, record) [No support for multiple index] [file must be sorted]
- Alt 2 By ref: leaf nodes point to records. (key, [Page num, Record num]). [leaf nodes are data pages]
- Alt 3 By list of ref: list of ptrs with each key. (key, list of [""]). MAX LEAVES =  $(2d+1)^h$
- Clustered: Data pages sorted same as in index. (Alt 1 has to be clustered, Alt 2/3 don't care) MAX RECORDS =  $2d(2d+1)^h$  [Alt 1 data pages for
- Bulk loading: fill leaf nodes upto 2f.d. fill inner nodes completely.

### Buffer Management

- Frame ID: RAM Address of frame - Page ID: Memory Addr of Page in frame - Dirty Bit: modified? - Pin #: number of threads using page at same time. [set by requestor using page at same time.]
- Least Recently Used: costly, sequential scanning means 0 hit rate of page i.e. file/index [decremented by requestor of page]
- Most Recently Used: good for sequential scan, bad for most use.
- Clock Policy: approximates LRU
  - set clock to first unpinned frame on start. Set ref bit to 1 when page is initially read in.
  - If hit, set ref bit to 1 without moving clock hand.
  - To evict: traverse cache until first ref=0 is found, setting 1's to 0 when passed.
  - skip pinned pages when ref bit 0 found, evict, replace, set to 1, move clock to next frame.
- Sequential flooding: when an sequential scan has 0 hit rate.

### Relational Algebra:

- $\pi$  [Select] -- [set minus] -  $\times$  [cross product] -  $\sigma$  [where for equality predicates with  $\wedge, \vee$ ] -  $\bowtie$  [Inner Join] -  $\cup$  [Union] -  $\cap$  [Intersection]
- $\bowtie$  [Natural Join] -  $\rho_{ab}$  [rename col a to b]

### Hashing:

- Divide phase:
  - Input values, stream in using 1 input buffer, hash using  $h_{p1}$ , stream output through  $B-1$  output buffers, to  $B-1$  partitions.
  - $r_2 = \text{sum of pages in partitions w/ } > B \text{ pages each. Use } h_{p2} \text{ to rehash. So on until all buffers } < B.$
- Converge phase:
  - For each partition, use all  $B$  input buffers to load in, hash using  $h_r$  to re-order within partition, flush all  $B$  to memory.
- I/O cost:  $\sum_{i=1}^{\# \text{ of passes}} (\text{Input pages for hash} + \text{Output pages for hash}) + 2 \text{ (Final Output pages)}$

(in-memory) sorted

- First sort each page, then pass 0: use all  $B$  as input buffers and merge  $B$  pages at a time, creating runs of size  $B$ .
- Pass 1: merge  $B-1$  runs at a time. use 1 as output buffer,  $B-1$  as input. [me diff is bc merge is out-of-memory]
- I/Os:  $2N(1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil)$

## Joins:

- $[R]$  is pages in R,  $|R|$  is records in R
- Simple nested loop join:  $[R] + |R|[S]$
- Block nested loop join:  $[R] + \lceil \frac{|R|}{B-2} \rceil [S]$ . (use B-2 input buff for R, 1 for S, 1 output)
- Index nested LJ:  $[R] + |R|(\text{cost to look up records in } S)$
- Grace Hash Join: Partition R and S until  $\min(|\text{records}|, |\text{Part}(R)|) \leq B-2 + i$ . Then,  $i$ , rehash smaller partition and naive hash join. (I/O: ~~hashing + hashed(R) + S~~) doesn't work with key skew.
- Sort Merge SMJ:  $2[R](1 + \lceil \log_{B-2} |R| \rceil) + 2[S](1 + \lceil \log_{B-2} \lceil \frac{|S|}{B} \rceil \rceil) + [S] + [R]$  or  $[R] + |R|[S]$  in worst case
- advance R or S pointers until match
- mark the match in S, check further until non-match
- go back to Mark in S, repeat advancement until match
- Optimized SMJ: sort both R and S until  $\max(\text{runs}(R), \text{runs}(S)) \leq B-1$ . If  $\text{runs}(R) + \text{runs}(S) \leq B-1$ , then merge all in one pass. I/O:  $SMJ - 2E] - 2[R]$   
Else if  $\max(\text{runs}(R), \text{runs}(S)) \leq B-2$ , sort larger one, merge with runs of smaller. I/O:  $SMJ - 2\max([R], [S])$   
Else if  $\min(\text{runs}(R), \text{runs}(S)) \leq B-2$ , sort smaller... I/O:  $SMJ - 2\min([R], [S])$ . Else SMJ.

## Query Optimization:

- Streaming Operators: small cost to get next tuple (eg- SELECT)
- Blocking Operators: can't produce any output before consuming full input (eg- SORT)
- Selectivity estimation: what % of pages make it through operator
  - $X = a : \frac{1}{\text{unique}(x)}$
  - $X = Y : \frac{1}{\max(1, \text{unique}(x), 1, \text{unique}(y))}$
  - $X > a : \frac{\max(X) - a}{\max(X) - \min(a) + 1}$
  - cond1 AND cond2 : Selectivity 1 \* selectivity 2
- Join Selectivity:  $\frac{1}{\max(\text{unique}(A.id), \text{unique}(B.id))}$ . Output tuples =  $|A| \cdot |B| \cdot \text{selectivity}$ .
- Heuristics for optimization: Push down selects( $\sigma$ ) and projects( $\pi$ ); only consider left-deep plans; cross joins <sup>only if</sup> no other option.
- System R (Gellinger Optimizer):
  - Pass 1: scan each table using full/index scan. advance optimal + interesting access plans.
  - Pass 2-n: do same thing for joins and not scans. [only SMJ can produce interesting orders. SNLS, INLJ preserve left relation order].
  - Join I/O cost depends on:
    - whether intermediate relations (from prev. operator) are materialized or streamed. System R doesn't materialize.
    - interesting orders from prev. operators. <sup>for NLT, only materialize right relation</sup>

## Transactions & Concurrency

- Inconsistent Reads (Write-Read Conflict): reads part of what was updated (eg- read R pre-update)
- Lost Update (Write-Write Conflict): overrides someone else's update (eg-  $P \leftarrow P+5$ ,  $P \leftarrow 2P$  instead of  $P \leftarrow 2P+5$ ) but S post-update
- Dirty Reads (Write-Read Conflict): reads an update that isn't committed.
- Unrepeatable Reads (Read-Write): transaction has multiple reads, but another transaction modifies record b/w reads.
- Atomicity: transaction either commits or aborts. Either way, either all changes happen or none.
- Consistency: if DB starts out consistent, it ends up consistent after Xaction.
- Isolation: each Xact should execute as if it ran by itself. No other Xact should affect it.
- Durability: if Xact commits, its changes persist despite failures.
- Concurrent execution increases throughput and reduces latency.
- Serial schedule = NO interleaved transactions.
- Schedule serializable iff equiv. to serial schedule.
- Conflicts occur when: 2 operations from diff Xactions operate on same resource, and  $\geq 1$  is write.
- Conflict-equivalent: schedules that order conflicts the same.  $\rightarrow$  (d/w about seq. reads)
- Conflict-equiv  $\Rightarrow$  equiv  $\circ$  conflict-serializable  $\Rightarrow$  serializable
- Conflict serializable iff dependency graph acyclic.
- View-equivalent: some initial reads, some dependent reads, some winning writes.
- View serializable finds all conflict-serializable schedules plus some serializable schedules w/ blind writes
- Blind writes: sequential writes with no interleaved reads.
- Two-Phase Locking (2PL): ensures conflict-serializable schedules.
  - Xact must acquire Shared lock before read, Exclusive lock before write
  - Xact cannot acquire new locks after releasing any lock.
  - 2PL is prone to cascading aborts ( $T_1$  write + release,  $T_2$  read,  $T_1$  abort,  $T_2$  forced to abort).
  - Strict 2PL ensures all locks released together. (Prevents cascading aborts).
  - Lock Management
- Deadlock avoidance: assign priority based on age.
- Wait-Die: If  $(T_i > T_j)$   $T_i$  waits:  $T_i$  aborts. [In this example,  $T_k$  will only wait for  $T_i$  on locks before  $T_i$  has  $T_j$ ]
- Wound-Wait: If  $(T_i > T_j)$   $T_j$  aborts:  $T_i$  waits. [ $T_i$  wants what  $T_j$  has]

Resource	Granted	Mode	Owner
A	$T_1, T_2$	S	$T_3, T_4, T_5$
B	$T_6$	X	$T_2, T_3, T_5$

- IS: intent to place S on finer granularity

- IX: " " X " "

- SIX: S and IX together. S at this level, so no other Xact can get X on anything in that subtree.

- If on S's own, can't acquire new locks while waiting.

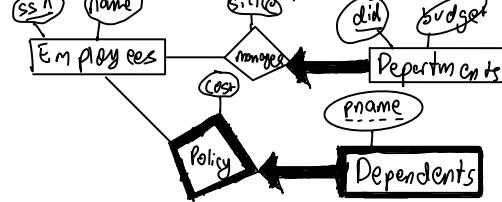
	IS	IX	S	SIX	X
IS	✓	✓	✓	✓	✗
IX	✓	✓	✗	✗	✗
S	✗	✗	✓	✗	✗
SIX	✗	✗	✓	✗	✗
X	✗	✗	✗	✗	✗

- To get S/IS on node, hold IS/IX on parent.

- To get X/IX on node, hold IX/SIX on parent.

- Locks are released in bottom-up order.

- **Observe skipping:** when a lock is released, allow all compatible locks, checking in order of wait queue.



- F+ is the set of all FDs implied by F. It is the set of all table relationships in closure of F.

-  $X \rightarrow A \in F^+$ , either  $A \subseteq X$ , or X is a superkey of R.

- Lossless decompositions are reversible. (e.g. BCNF).

-  $R \rightarrow X, Y$  is lossless wrt F iff  $F^+$  contains  $X \cap Y \rightarrow X$  or  $X \cap Y \rightarrow Y$  [ $X \cap Y$  is the set of cols in both X and Y]

-  $R \rightarrow X, Y$  is dependency preserving wrt F iff  $(F \cup UF_Y)^+ = F^+$ .

## DB Design

### Entity-Relationship Models

- Key constraints: arrows represent one-to-many (from 1 to many).

- Participation constraints: thick lines represent at-least one.

- Weak entities have partial key instead of primary key. They also have many-to-one with an 'owner' entity. Owner's primary + weak's partial can uniquely identify each weak entity.

- Functional Dependencies:  $X \rightarrow Y$  implies X determines Y.

- X is superkey iff  $X \rightarrow R$  (where R is full relation).

- Candidate key is group of cols that form superkey (can't remove any cols and still form).

- Boyce-Codd Normal Form: R is in BCNF if  $\forall X \rightarrow A \in F^+$ , either  $A \subseteq X$ , or X is a superkey of R.

- Input: R, F (FDs for R)

- Algo:  $\forall r \in R$  where r not BCNF:

a) Pick f:  $X \rightarrow A$ ;  $X, A \in r$ .

b) Let  $R_1 = X^+, R_2 = X \cup (r - X^+)$

c) Remove r from R.

d) Insert  $R_1, R_2$  into R. (as indep tables)

e) Recompute F as FDs overall r/r.

$X \in R$   
 $X \in R_2$

## Recovery

Copying CLRs is optional, for convenience)

- Force Policy: force all modified pages to disk before Xact commits.
- No-force Policy: only write back to disk when page needs to be evicted from buffer pool.
- Steal Policy: allow writing to disk before Xact finishes.
- No-steal Policy: pages can't be evicted from buffer pool until Xact commits, fastest, undo/reco possible.
- Write-Ahead Logging: log records written to disk before data pages written. (for atomicity)
- Log records can be "start, commit, Abort, or Update (<XID, pageID, offset, length, old-data, new-data>)
- Log Sequence Numbers:
  - prevLSN: stores LSN of last action from same Xact in log records
  - flushedLSN: stores LSN of last log record flushed to disk.
  - Page LSN: stores LSN of last log record where this page was modified.  
[Page i can only be flushed if pageLSN\_i ≤ flushed LSN. log flushed before page]
- Abort Process: Write ABORT record → undo actions bottom up and write Compensation Log Records for each.
- Transaction Table: X ID, status, lastLSN (most recent action from this Xact).
- Dirty Page Table: page ID, RecLSN (first operation to dirty page).
- UPDATE records are written before dirty page. COMMIT records written after all pages written to disk.
- Undo Logging: after crash, scan log from bottom. UPDATE log records represent pre-modification value)
  - for COMMIT/ABORT T: mark T as completed
  - for UPDATE T, x, v: if (T not completed) write x=v to disk: ignore [undo incomplete Xacts]
  - for START T: ignore.
- Redo Logging: no-force, no-steal [UPDATE log records rep. post modification values]
  - (Redo updates of committed Xacts) → COMMIT records written before dirty pages flushed
  - AIRES recovery Algo: Analysis → Redo → Undo
- Analysis Phase: to rebuild DPT and Xact table. Scan logs from beginning and:
  - if not END: add to Xact table, set last LSN.
  - if COMMIT/ABORT: update status in Xact table.
  - if UPDATE: if page not in DPT, add and set RecLSN.
  - if END: remove from Xact table.
  - after analysis phase, for any committing Xacts, add END log record, and remove from Xact table.
- Checkpointing writes Xact table and DPT to log to optimize recovery.
  - fuzzy Checkpointing has <begin-checkpoint> and <end-checkpoint> logs, and tables can be of any time b/w these. Thus, start recovery at begin-checkpoint.
  - Then, redo all Update/CLR operations where: Page in DPT, RecLSN ≤ LSN, pageLSN < LSN. (regardless of commit)
  - CLR records have undonextLSN = prevLSN of operation being undone.

### Redo Phase

All updates with  $LSN > pageLSN$  are Redone

## Parallel Query Processing

- Shared memory: every CPU shares memory + disk
- Shared disk: shared disk, individual RAMs
- Shared nothing: (used in rest of note) CPUs only communicate via messages.
- **Intra Query Parallelism**
  - Intra-Operator: divide data up and run one operator across many machines.
  - Inter-Operator: different op. diff. machines. e.g. SMJ(R,S) →  $\text{sort}(R)$  →  $\text{Project}(q)$
  - Pipeline parallelism: record passed up to parent as soon as child is done w/ it.  $\text{Filter}(x=3)$
  - Bushy tree parallelism: different branches of operator run parallel.  $\text{Scan}(S)$
- **Partitioning**
  - Sharding: each data page on 1 machine
  - Range partitioning: each machine gets a range to store.
  - Hash partitioning: self-explanatory.
  - Round-robin partitioning: only benefit is uniform split → may parallelization.  $I/O = (1 + \lceil \frac{1}{I} + \log_{\frac{1}{I}} \frac{N}{MB} \rceil) 2N$  in best case. (depends on range dist.)
  - Parallel Sorting: Range Partition → local sort on each machine. Range partition also works but high easier.
  - Parallel hashing: Hash Partition → local hash on each machine. Hash partition also works but high easier.
  - Parallel SMJ: Range Part. Rounds using some ranges → Local SMJ.  $I/O = [R] + [S] + \text{Sort}(S) + \text{Sort}(R) + [R] + [S]$
  - Parallel GHJ: Hash " " " " hash fn → GHJ.  $\text{Receive data on machines}$   $\text{stream joined}$
  - Broadcast Join: send smaller relation to every machine and do local joins.
  - All above joins need full input relations to produce output.
  - Only pipeline-friendly join is symmetric hash join.
    - start with 2 empty hash tables.
    - when record from R arrives, probe S, generate outputs, add record to table for R. VV for S.
  - Hierarchical Aggregation: parallelized SUM, COUNT, AVG, etc.

## Distributed Transactions

- Every machine has its own lock table.
- Superpose all waits for graphs to check for deadlock.
- **2 Phase Commit**: need consensus to commit.
  - Prepare phase:
    - coord → participants. participants generate PREPARE/ABORT record. • vote yes if PREPARE else no.
    - Commit/Aabort phase:
      - coord → part. [based on yes/no votes] • " " COMMIT / " " " " • send ACK (Acknowledgment) back to coord.
      - coord generates END record after receiving all ACKs. (all records must be flushed before any network messages)
  - Recovery:
    - Participant sees no PREPARE log → aborts
    - " " PREPARE log → ask coord. to repeat COMMIT/ABORT decision.
    - Coord " " no COMMIT log → aborts
    - " " COMMIT log → resend commit message to parts.
    - Part. " " " " → resend ACK
    - Coord " " END " " → no recovery to do.
  - **Presumed Abort**: Optimization for 2PC where no log record means abort. Hence, abort records never have to be flushed.
  - NO votes occur when part. fails, avoids deadlock, etc.
  - Even with another part. of higher priority, YES votes must be honored.
    - If vote Yes and coord. asks to commit, must commit.

## NoSQL

- **Online Transaction Processing**: Transactions ↑, users ↑, many updates.
- **Online Analytical Processing**: Read only workloads, data ↑, more joins/aggregation.
- Migration from OLTP to OLAP via Exact Transform Load.
- **Consistency**: multiple clients viewing should get same view of data. (Don't stop / wait)
  - Avalability**: every request must receive response that isn't error.
  - Partition tolerance**: continue operation despite dropped/delayed messages or disconnection.
- **Brewer's CAP theorem**: distributed systems can ensure 2/3 of these at most.
- 2PC sacrifices availability.
- Alternative: Eventual consistency (once communication is good, consistency ensured).

- **NOSQL Data Models**
  - **Key Value Store**: key is a string/int identifier, value is byte-array or record.  
(eg - AWS DynamoDB) Only supports get(key) and put(key, value).
  - **Wide-Column**: key = (rowID, columnID). value = record.  
(eg - Cassandra) supports get(key, [columns]), put(key, value).
  - **Document**: KVS with values that are semi-structured (eg - JSON / XML)
    - [Relational DB stores tuples in tables.]
    - [Document stores documents in collections.]
- **JSON** supports Arrays, Atomics (string/boolean/null/number (64-bit int)), Objects (dict of k-v pairs)
  - flexible - Schema non-enforcing - text-based not binary-based representation
  - Relational DBs use binary data formats that are strict incompatible w/ each other.
  - JSON is optimized to be shared by different apps and languages, not for storage and retrieval from disk
- **MongoDB**: get good.

## Map Reduce [for selection, group by, summation, etc.]

- **Distributed File Systems** partition data into **chunks** (usually 64 MB)
- func **map**(key1, value1) → bag((key2, value2)) . [bag is a collection (array)]
- func **reduce**((key2, bag(value2)) → bag((key3, value3)) . [works on output of map].
- **leader** splits M partitions (by key1) and assigns M **workers**
  - ↳ workers write map output to R regions (by key2). (written to disk for fault tolerance)
  - ↳ leader then assigns R workers, who write reduce output to disk.
- **Straggler** is a machine taking unusually long on task.
- MapReduce preemptively starts new (backup) executions of last few remaining tasks, and writes whichever finishes first, killing rest.

## Spark

- Uses semi-structured data objects called **Resilient Distributed Datasets** (immutable)
- Intermediate RDDs are not flushed to disk unless `.persist()` is called.
- **Lineage** (sequence of operations to construct RDDs) is logged and flushed.
- On failure, all non-persisted RDDs are recomputed from lineage.
- **actions** are **eager**: they are immediately executed (count, reduce, save, etc.)
- **transformations** are **lazy**: logged and built into optimization tree. Executed on collect.

A Partial List of Spark Transformations and Actions

Transformations:	
<code>map(f : T -&gt; U):</code>	<code>RDD&lt;T&gt; -&gt; RDD&lt;U&gt;</code>
<code>flatMap(f: T -&gt; Seq(U)):</code>	<code>RDD&lt;T&gt; -&gt; RDD&lt;U&gt;</code>
<code>filter(f:T-&gt;Bool):</code>	<code>RDD&lt;T&gt; -&gt; RDD&lt;T&gt;</code>
<code>groupByKey():</code>	<code>RDD&lt;(K,V)&gt; -&gt; RDD&lt;(K,Seq[V])&gt;</code>
<code>reduceByKey(F:(V,V)-&gt; V):</code>	<code>RDD&lt;(K,V)&gt; -&gt; RDD&lt;(K,V)&gt;</code>
<code>union():</code>	<code>(RDD&lt;T&gt;,RDD&lt;T&gt;) -&gt; RDD&lt;T&gt;</code>
<code>join():</code>	<code>(RDD&lt;(K,V)&gt;,RDD&lt;(K,W)&gt;) -&gt; RDD&lt;(K,(V,W))&gt;</code>
<code>cogroup():</code>	<code>(RDD&lt;(K,V)&gt;,RDD&lt;(K,W)&gt;) -&gt; RDD&lt;(K,(Seq[V],Seq[W]))&gt;</code>
<code>crossProduct():</code>	<code>(RDD&lt;T&gt;,RDD&lt;U&gt;) -&gt; RDD&lt;(T,U)&gt;</code>

Actions:	
<code>count():</code>	<code>RDD&lt;T&gt; -&gt; Long</code>
<code>collect():</code>	<code>RDD&lt;T&gt; -&gt; Seq&lt;T&gt;</code>
<code>reduce(f:(T,T)-&gt;T):</code>	<code>RDD&lt;T&gt; -&gt; T</code>
<code>save(path:String):</code>	Outputs RDD to a storage system e.g., HDFS

- **B** = # Blocks, **R** = # Records per Block
- Equality Search = Binary Search  $\lceil \log_2(BR) \rceil \text{ I/O}$
- Insertion:  $\log(BR) + 1$  + 1 + 2, I/O  
find page to mod. Read page to mod. write modified page, update index.

## K-D trees

- An example of **Spatial Indexing**
- Binary Search tree split on k different cols.
- **construct(Data)**
  - if  $|Data| = R$ : return new `Block(Data)`
  - else divide Data by thresholding on some col
  - return `(col, threshold)`

`construct(Data1)`      `construct(Data2)`

	Heap File	Sorted File	Clustered Index	K-d Trees
Scan all records	$B^*D$	$B^*D$	$3/2 * B * D$	$B * D$
Equality Search	$0.5*B*D$	$(\log_2 B)*D$	$(\log_p(BR/E)+2)*D$	$(\log_p(BR) + 2) * D$
Range Search	$B^*D$	$((\log_2 B) + \text{pages}) * D$	$(\log_p(BR/E) - 3 * \text{pages}) * D$	$(\log_p(BR) + 2 * \text{pages}) * D$
Insert	$2^*D$	$((\log_2 B) + B)*D$	$(\log_p(BR/E)+4)*D$	$(\log_p(BR) + 4) * D$
Delete	$(0.5*B+1)*D$	$((\log_2 B) + B)*D$	$(\log_p(BR/E)+4)*D$	$(\log_p(BR) + 4) * D$

- B:** Number of data blocks
- R:** Number of records per block
- D:** Average time to read/write disk block
- F:** Average internal node fanout
- E:** Average # data entries per leaf

## Eventual Consistency

### BASE Semantics

- Eventually consistent systems opt for the 3 BASE guarantees instead of ACID guarantees:
  - Basic Availability:** reads and writes are available as much as possible; however, they are not guaranteed to be consistent (i.e. a read may not get the latest data and a write might not persist after all updates propagate). This is up to the application to fix.
  - Soft State:** a database can change even without inputs (e.g. as updates propagate), so the application only has a probability of knowing its state.
  - Eventually Consistent:** given enough time (after all updates propagate) all reads will be consistent.

## Paxos Process

- **Election:** leader is elected.
  - candidate machines generate ballots with ballot-id. ballot-ids must be unique, increasing  $\geq$  MachineID + counter
  - proposer sends PREPARE(ballot-id) to majority participants.
  - if participant votes no, it ignores PREPARE and sends nothing.
  - if " " " yes, it sends PROMISE(ballot-id, old-ballot-id, old-v). this is a promise to ignore future ballots with smaller id.
- **Bill:** Proposer produces bill.
  - if proposer receives majority vote, it becomes leader.
  - if proposer receives value from some people it sent last PROPOSE to. leader sends PROPOSE(ballot-id, v) to same people. if it received old-v earlier, it must re-propose  $v = \text{old-}v$  with highest corr. old-ballot-id.
  - if participant has responded to PROPOSE from higher ballot-id, they ignore this. else, can ACCEPT(ballot-id, v)
- **Law:** Bill signed into law.
  - if majority, leader adds (ballot-id, v) to log and flushes.
  - Sends COMMIT(ballot-id, v) to all participants.

## Miscellaneous

- CLRs are counted when computing last LSN.
- All operations except COMMIT, END, ABORT are redone (in ARIES).
- Don't undo already CLR-ed operations in undo phase.

Nearest Neighbor Search: (for some point  $B$ )  
 - traverse tree down to find leaf node for  $B$ .  
 - current best estimate is best est.,  $q$ , in that block.  
 - calculate distance ( $B$ , threshold) for all boundaries of the partition. traverse all neighbors s.t.  $\text{dist} < \text{dist}(B, q)$

## Consensus (Strong Consistency)

- consensus protocols satisfy
- Termination: All non-faulty processes eventually decide on value.
- Agreement: All processes decide same val.
- Validity: Chosen value is proposed by some participant of protocol.

### Paxos

- messages eventually arrive (no bound on delay)
- messages not corrupted (failed machines just stop)
- to tolerate  $F$  failures,  $2F+1$  replicas must exist  
 $\downarrow$   
 (i.e. majority of machines must be non-faulty)

only if sent ACCEPT  
 for some ballot before