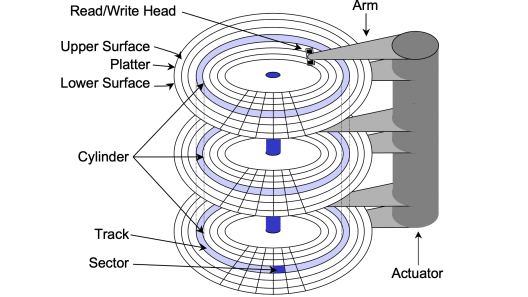
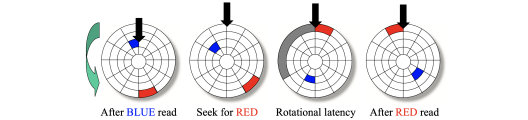


Storage

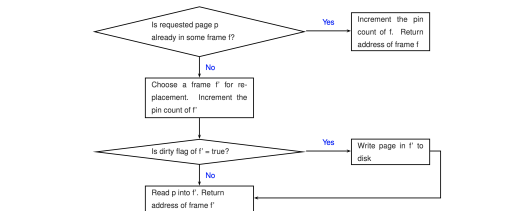
Metadata. Relation schemas, view definitions, indexes, stats.
Log files. Information maintained for data recovery.
Primary Memory. Registers, static RAM (caches), dynamic RAM (physical memory).
Secondary Memory. Magnetic disks (HDD), SSD.
Tertiary Memory. Optical disks, tapes.
Storage. DBMS stores data on non-volatile disk for persistence. Processing data in main memory (RAM).



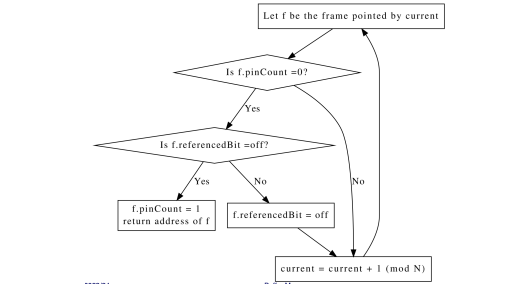
Seek time. Moving arms to position disk head on track.
Rotational Delay. Waiting for block to rotate under head. Depends on rotation speed. Average delay is time for half revolution.
Transfer time. Actually moving data to/from disk surface. number of requested sectors on same track \times time for 1 revolution divided by number of sectors per track.
Access Time. Sum of seek time, rotational delay and transfer time.
Response time. queueing delay + access time.



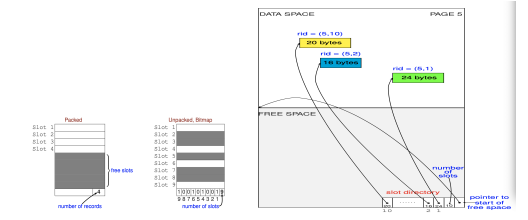
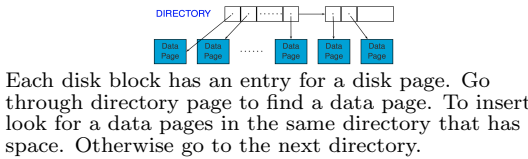
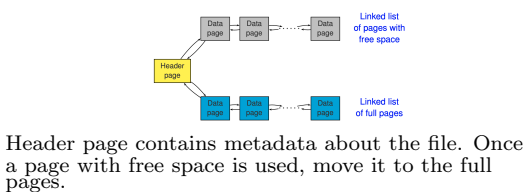
Sequential IO. Everything written on same track.
Random IO. Things may be written on different tracks. If the size of the data is grater than the track, store in the same cylinder.
SSD. No moving parts so faster random IO. Updates to page requires erasure of multiple pages before overwriting. Limited number of times a page can be erased.
Disk blocks. Data is stored/retrieved from here. Each is a sequence of ≥ 1 contiguous sectors.
File layer. Deals with organisation and retrieval of data.
Buffer Manager. Controls reading/writing of disk pages.
Disk Space Manager. Keeps track of pages used by file layer.
Buffer Pool. Main memory allocoated for DBMS.
Frames. A partition of the buffer pool into block-sized pages.
Dirty page. Page in the buffer that has been modified and not updated on disk.
Pin count. Number of clients using the page initialised to 0.
Dirty flag. Whether page is dirty initialised to false.



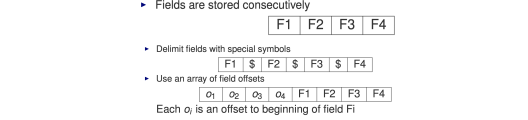
Pinning. Incrementing pin count by pinning the requested page in its frame.
Unpinning. Decreasing pin count. Dirty flag should be true if page is dirty.
Replacement. Only when pin count is 0. Before that write to disk if dirty flag is true.
Replacement Policy. Decide which unpinned page to replace.
LRU. A queue of pointers to frames with pin count 0. Exploits locality. Spatial at the page level and temporal at the record level.
Clock. Current variable points to some buffer frame. Each frame has a reference bit (on when pin count is 0). Replace a page that has reference bit off and pin count 0.



Files. Each relation is a file of records. Each record has a unique record identifier called RID.
Heap file. Unordered file.



RID. Tuple pair of page id and slot number.
Packed Org. Store records in contiguous slots.
Unpacked Org. Uses a bit array to maintain free slots.

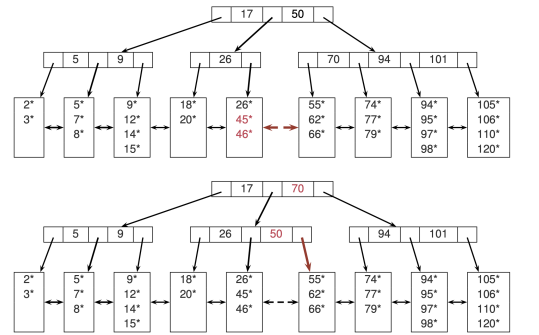


Indexing

Search Key. Sequence of k data attributes, $k \geq 1$
Composite search key. If $k > 1$
Unique index. Search key is a candidate key. Otherwise non-unique.
Data entries. Index stored as a file and records in an index file.
 B^+ tree. Leaf nodes are at the bottom most level. Root node is at level 0.
Height of tree. Leaf nodes are at level h .
Sibling nodes. Nodes at the same level if they have the same parent nodes.
Leaf Nodes. Store sorted data entries of tuple pair of search key value and RID. Doubly linked.
Internal nodes. Stored in the form p_0, k_1, p_1, \dots . Keys are sorted $k_1 < k_2 < k_3, \dots$. p_i is the disk page address. For k^* in index subtree of T_i rooted at $p_i, k^* \in [k_i, k_{i+1})$
Index entry. A k_i, p_i pair. k_i is the separator between node contents pointed to by p_{i-1}, p_i .
Order. Controls space utilisation. Each non-root node contains m entries where $m \in [d, 2d]$. Root node contains m entries, where $m \in [1, 2d]$.
Equality search. At each internal node, find the largest key of the node such that $k \geq k_i$. If k_i exists, search subtree at p_i otherwise search subtree at p_0 .
Range search. Find the node based on equality then traverse the doubly LL to find the remaining.
Format 1 data entry. k^* is an actual data record with search key value of k .
Format 2 data entry. k^* is a pair of search key value k and RID.
Format 3 data entry. k^* is a pair of search key value k and a list of RIDs.

Insertion. Find the leaf node that to insert into by performing searching.
Overflowed node. Node already has $2d$ entries and a new entry is being inserted into it.
Splitting. Get all the records in the leaf (including the one to be inserted) and sort them ascending. Then take the $d + 1$ entries and put into a new leaf node. Create a new index using the **smallest** key in the **new** leaf node. Insert it into parent node of overflowed.
Propagation of splits. Internal nodes overflow when the number of keys exceed $2d$. Get all the keys and sort them ascending. Take the middle key and push it to the parent.
Redistribution. Only at leaf nodes. Perform distribution to a non-full adjacent sibling and update the index accordingly. Try right sibling first then left otherwise use splitting.
Deletion. Find the leaf node that the record is in and remove it.
Underflowed. Non-root node. If the result of deletion causes the number of entries to be $d - 1$.
Redistribution of leaves. An underflowed nodes

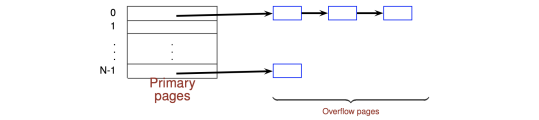
can be redistributed using an adjacent sibling entry. Try right sibling then left. Otherwise split.
Merging Condition. Both siblings have d entries. Merge with the right one.
Merging at leaves. If can merge, remove the parent key and combine with the sibling.
Propogation of nodes merging. If an internal node underflows after merging the leaves, pull down the separating key from parent node to form a merged nodes.
Balancing underflowed node. Try the right sibling first. Push the smallest one up to the parent. Then bring the middle key down.



Bulk loading. Sort all data entries by search keys and load the leaf pages with sorted entries. Construct the leaf pages with $2d$ entries. Insert pages into rightmost parent-of-leaf level page of tree.
Advantages. Efficient construction algorithm. Leaf pages are allocated sequentially.
Clustered Index. Order of data is the same or close to the order data records. Otherwise its **unclustered**. Format 1 is clustered. At most 1 clustered index for each relation.
Dense index. There is an index record for every search key value in the data. Otherwise its **sparse**. Unclustered index must be dense. Clustered can be sparse or dense.

Hash Index

Static Hashing. Hash function to identify which bucket data should be stored. 1 primary data page and overflow pages are chained to the primary data page. Perform modulo of N to decide which bucket.



Linear Hashing. Grow linearly by splitting of buckets.
Splitting a bucket. Add a new bucket (split image) and redistribute across split image and the original. Increase the value of next by 1. If the value of next becomes the same as level, reset value of next and increase level by 1.
Hash functions. Each round has 2 hash functions. If a bucket has been split, use the new hash function else use the current one.
Redistribution. To determine whether a record should move or not, and if the original table size is s , perform the record mod $2s$. If the result is 1, move it to the split image else remains the same.

$$|R| + ||R|| \times \left(\log_F(\lceil \frac{||S||}{b_d} \rceil) + \lceil \frac{||S||}{b_d ||\pi_{B_j}(S)||} \rceil \right)$$
