

Indexes

Kathleen Durant PhD
Northeastern University
CS 3200

Outline for the day

- Index definition
- Types of tree files organization
 - B+ trees
 - ISAM
- Choosing indexed fields
- Indexes in InnoDB

Indexes

- A typical file allows us to retrieve records:
 - by specifying a file offset or a *rid*, or
 - by scanning all records sequentially
- Sometimes, we want to **retrieve records by specifying the values in one or more fields**
 - *Examples:*
 - Find all students in the “CS” department
 - Find all students with a gpa > 3
- Indexes are file structures that enable us to answer such value-based queries efficiently.

Indexes

- An index on a file speeds up selections on the search key fields for the index
 - Any subset of the fields of a relation can be the search key for an index on the relation
 - Search key is not the same as a key in the DB
- An index contains a collection of data entries, and supports efficient retrieval of all data entries with a given key value k .

Why Index?

- Database tables have many records and ..
 - linear search is very slow, complexity is $O(n)$
 - Keeping a file sorted to apply a binary search is costly
- Indexes improve search performance
 - But add extra cost to INSERT/UPDATE/ DELETE
- Many options for indexes
 - Hash Indexes (MEMORY and NDB)
 - Bitmap Indexes (not available in MySQL)
 - B-Tree Indexes and derivatives (MyISAM, InnoDB)

Index Concept

- Main idea: ***A separate data structure used to locate records***
- Most generally, index is a list of value/address pairs
 - Each pair is an index “entry”
 - Value is the index “key”
 - Address will point to a data record, or to a data page
 - The assumption is that the value/address pair will be much smaller in size than the full record
- If index is small, a copy can be maintained in memory
 - Permanent disk copy is still needed

Indexing Pitfalls

- Index itself is a data store
 - Occupies disk space
 - Must worry about maintenance, consistency, recovery, etc.
- Large indices won't fit in memory
 - May require multiple seeks to locate record entry

Essential for Multilevel Indexes

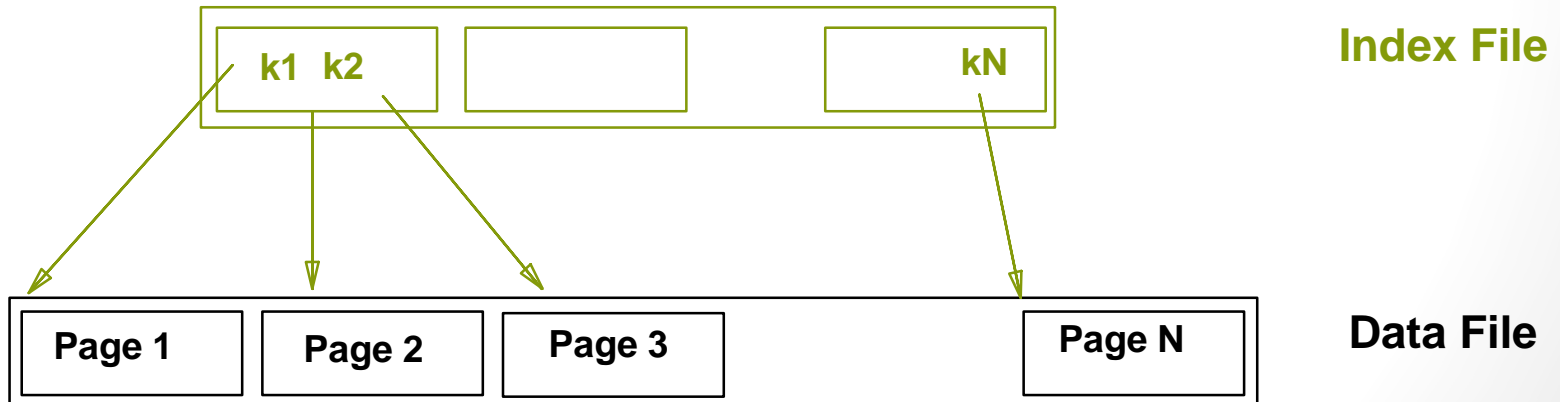
- Should support efficient random access
 - Should also support efficient sequential access, if possible
- Should have low height
 - Implies high fan out: refers to the number of children nodes for an internal node.
- Should be efficiently updatable
- Should be storage-efficient
- Top level(s) should fit in memory

Tree Structured Indexes

- Tree-structured indexing techniques support both *range searches* and *equality searches*.
- Tree structures with search keys on *value-based domains*
 - ISAM: static structure
 - B+ tree: dynamic, adjusts gracefully under inserts and deletes.

Range Searches

- “Find all students with $gpa > 3.0$ ”
 - If data is in a sorted file, do binary search to find first such student, then scan to find others.
 - Cost of binary search can be quite high.
- Simple idea: Create an ‘index’ file.



➤ *Can do a binary search on (smaller) index file!*

ISAM

- = Indexed Sequential Access Method
 - IBM terminology
 - “Indexed Sequential” more general term (non-IBM)
 - ISAM as described in textbook is very close to B+ tree
 - simpler versions exist
- Main idea: ***maintain sequential ordered file but give it an index***
 - Sequentiality for efficient “batch” processing
 - Index for random record access

ISAM Technique

- Build a dense index of the pages (1st level index)
 - Sparse from a record viewpoint
- Then build an index of the 1st level index (2nd level index)
- Continue recursively until top level index fits on 1 page
- Some implementations may stop after a fixed # of levels

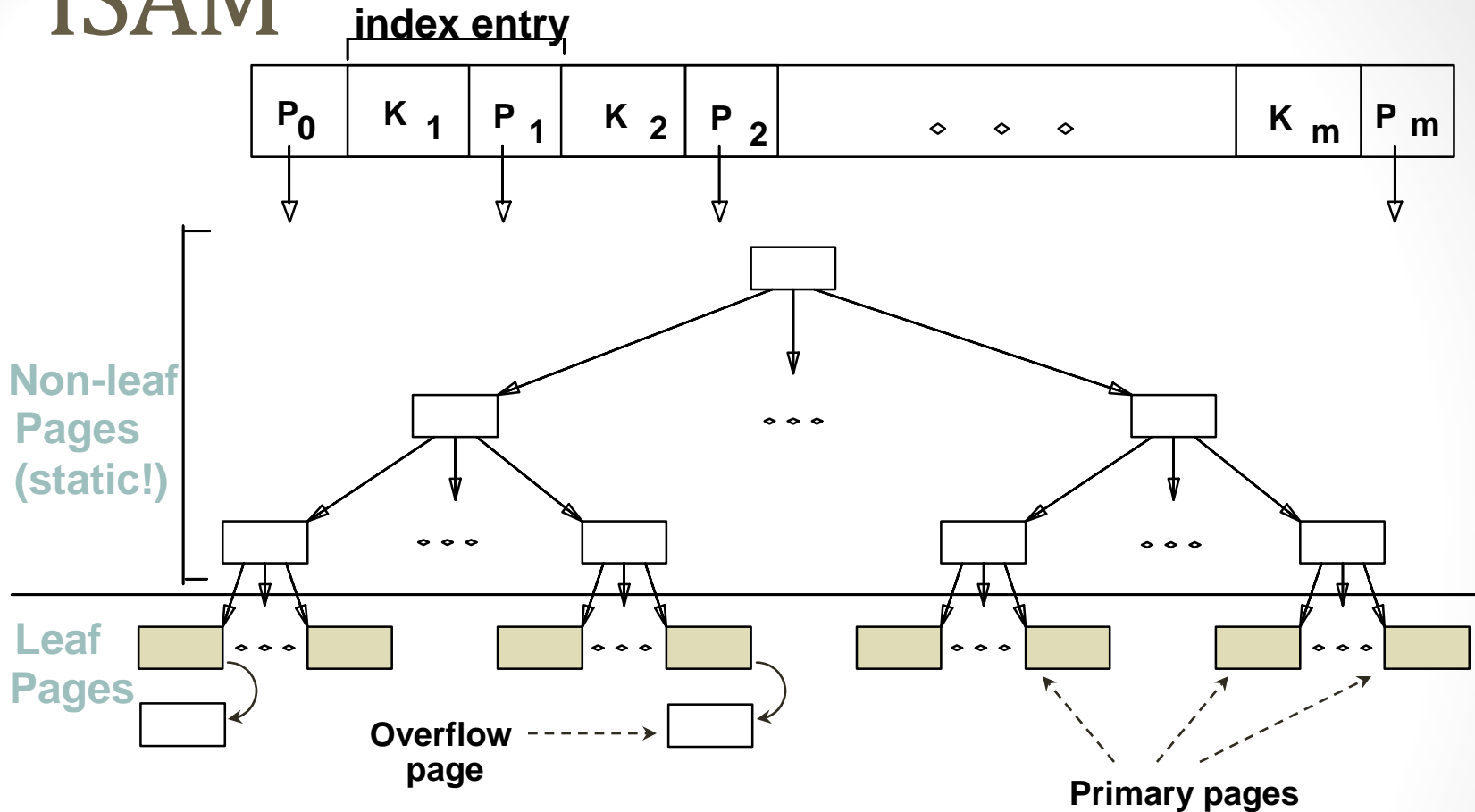
Updating an ISAM File

- Data set must be kept sequential
 - So that it can be processed without the index
 - May have to rewrite entire file to add records
 - Could use overflow pages
 - chained together or in fixed locations (overflow area)
- Index is usually NOT updated as records are added or deleted
- Once in a while the whole thing is “reorganized”
 - Data pages recopied to eliminate overflows
 - Index recreated

ISAM Pros, Cons

- Pro
 - Relatively simple
 - Great for true sequential access
- Cons
 - Not very dynamic
 - Inefficient if lots of overflow pages
 - Can only be one ISAM index per file

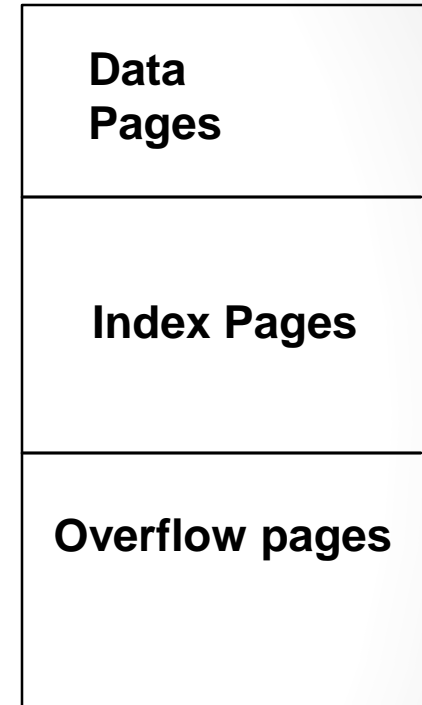
ISAM



- Leaf pages contain sorted data records
- Non-leaf part directs searches to the data records; **static once built**
- Inserts/deletes: use **overflow pages**, bad for frequent inserts.

Comments on ISAM

- *File creation*: Leaf (data) pages allocated sequentially, sorted by search key; then index pages allocated, then space for overflow pages.
 - *Index entries*: <search key value, page id>; they `direct` search for *data entries*, which are in leaf pages.
 - Search: Start at root; use key comparisons to go to leaf. Cost $\log_F N$; $F = \# \text{ entries/index pg}$, $N = \# \text{ leaf pgs}$
 - Insert: Find leaf data entry belongs to, and put it there.
 - Delete: Find and remove from leaf; if empty overflow page, de-allocate.
- *Static tree structure: inserts/deletes affect only leaf pages.*



Definition of B+ tree

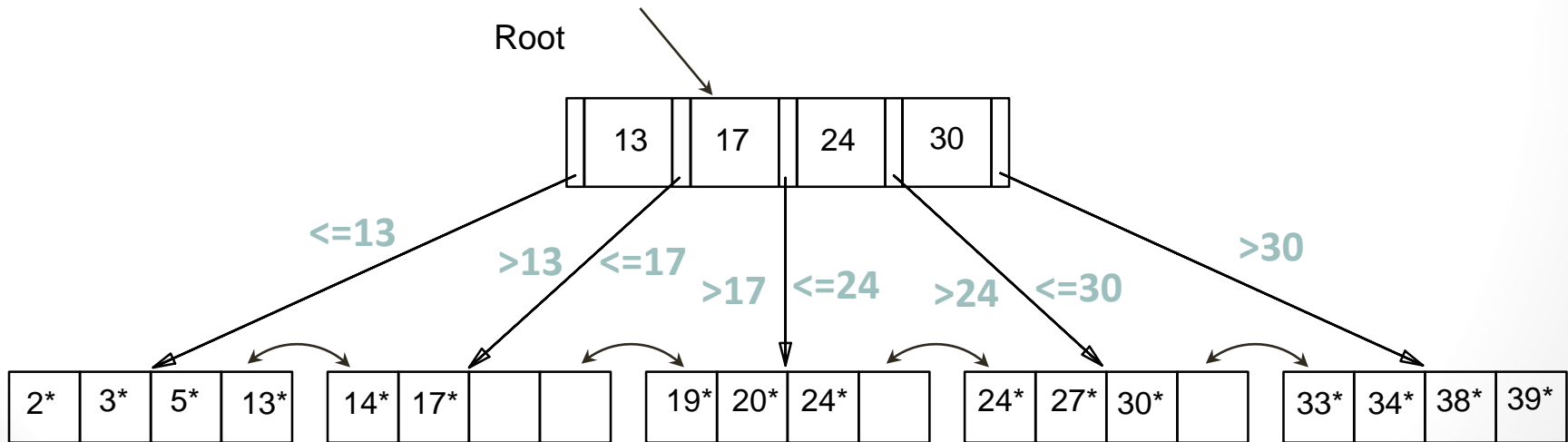
- A B+ tree of order n is a height-balanced tree , where each node may have up to n children, and in which:
 - All leaves (leaf nodes) are on the same level of the tree
 - No node can contain more than n children
 - Size limitation
 - All nodes except the root have at least $n/2$ children
 - The root is either a leaf node, or it has at least $n/2$ children
- Ensures that a fixed maximum number of reads would be required to access any data requested, based on the height of the tree

B+ tree properties

- The **depth** of a tree is the maximum number of levels between the root node and a leaf node in the tree.
- The **degree** or **order**, of a tree is the maximum number of children allowed per parent.
 - Large degrees create broader, shallower (shorter) trees.
 - Because access time in a tree structure depends upon depth than on breadth, it is advantageous to have “bushy,” shallow trees.
- The number of key values contained in a nonleaf node is 1 less than the number of pointers (the order).
- Leaf nodes are linked in order of key values.
 - Provides a mechanism for retrieving a range of records

Example B+ Tree

- Search begins at root, and key comparisons direct it to a leaf (as in ISAM).
- Search for 5*, 15*, all data entries $\geq 24^*$...

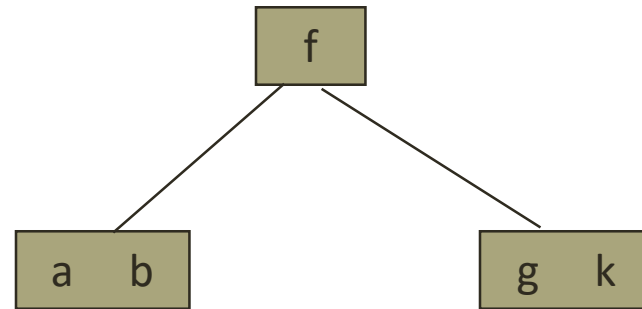
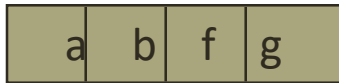


B+ Trees in Practice

- Typical **order**: 200. Typical **fill-factor**: 67%.
 - Average fan-out for internal nodes = 133
- Typical capacities:
 - Height 4: $133^4 = 312,900,700$ records
 - Height 3: $133^3 = 2,352,637$ records
- Can often hold top levels in buffer pool:
 - Level 1 = 1 page = 8 Kbytes
 - Level 2 = 133 pages = 1 Mbyte
 - Level 3 = 17,689 pages = 133 MBytes

Insertion in B Tree

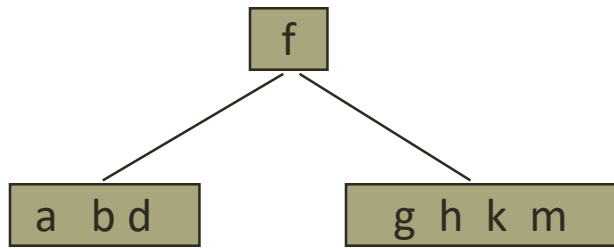
- 1. a, g, f, b:
- 2. k:



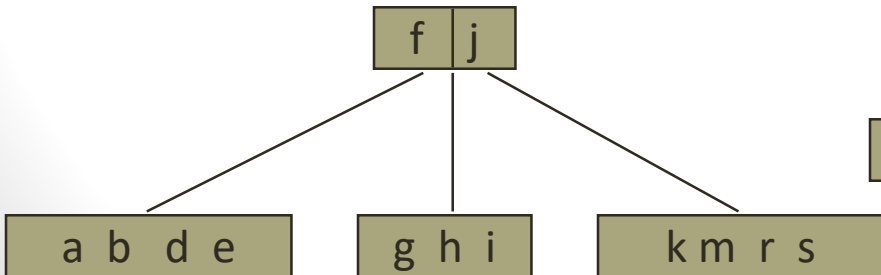
B tree: Data record stored in the tree

Insertion (cont.)

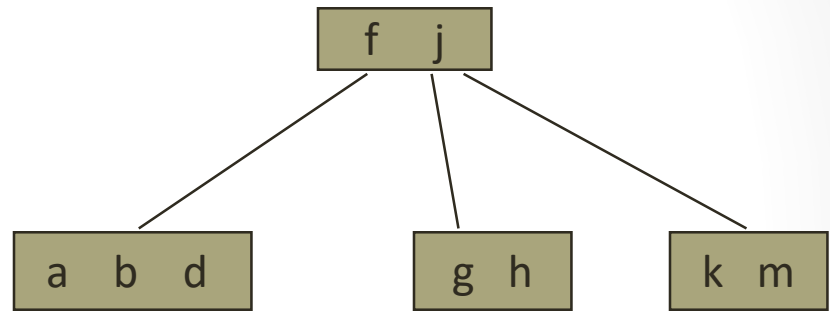
- 3.
- d, h, m:



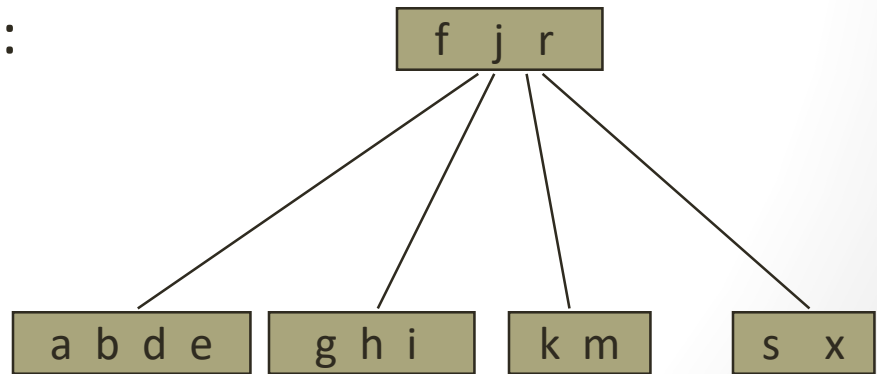
- 5.
- e, s, i, r:
-



- 4.
- j:



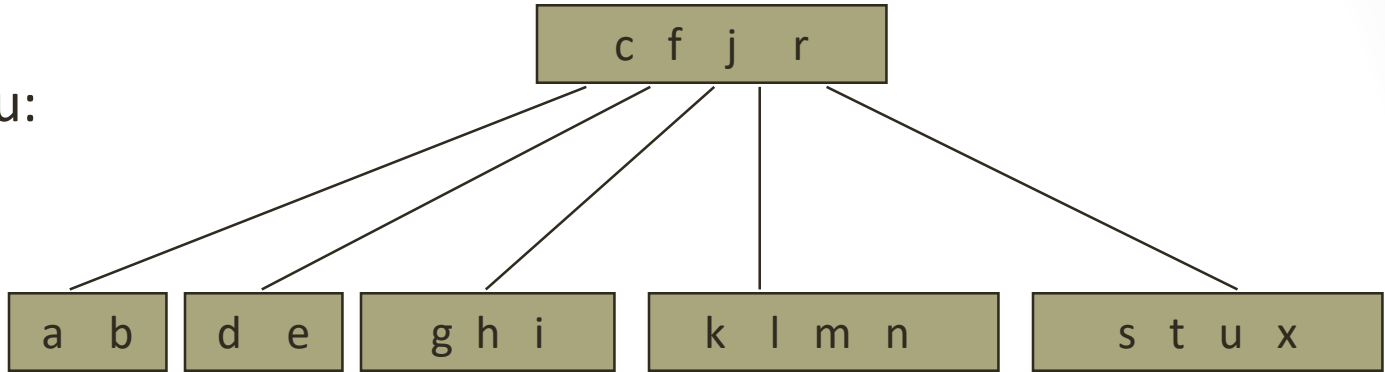
- 6.
- x:



Insertion (cont.)

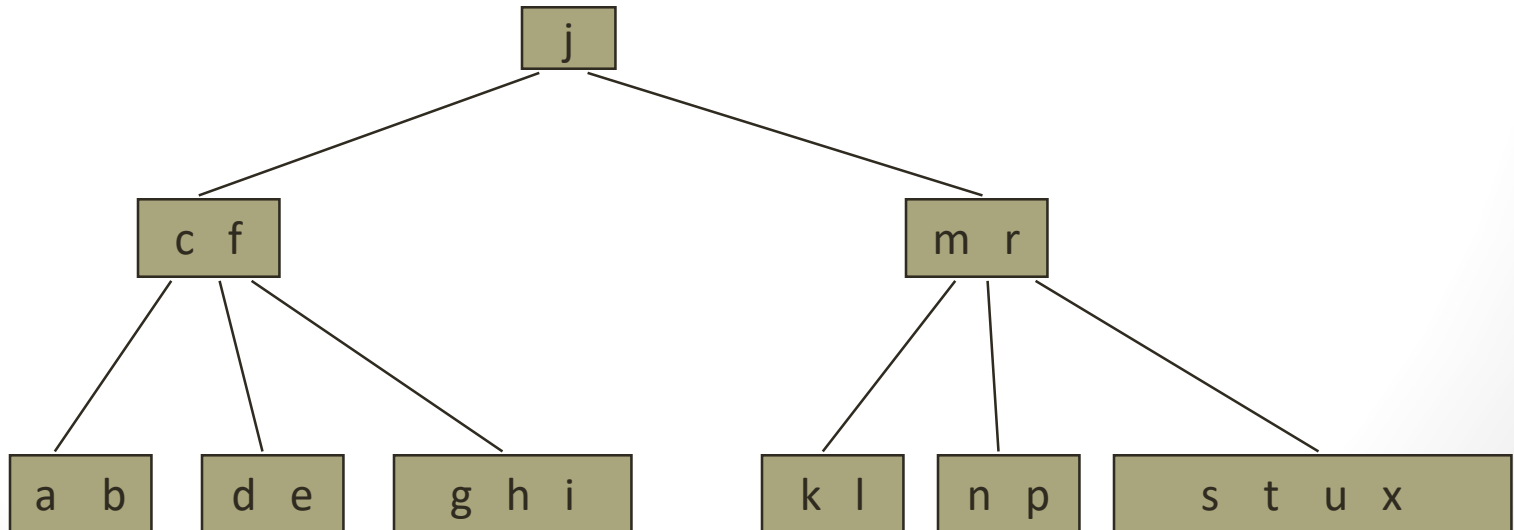
7.

c, l, n, t, u:



8.

p:



B+ Insertion algorithm

- Descend the tree where the leaf fits
- If the node has an empty spot insert the key/reference pair into the node
- If the node is already full, split into 2 nodes, distributing the keys evenly between the 2 nodes
- If the node is a leaf, take a copy of the minimum value in the second of these 2 nodes and repeat this insertion algorithm to insert it into the parent node
- If the parent is a non-leaf, exclude the middle value (median) during the split and repeat this insertion algorithm to insert this excluded value into the parent node.

B+ Deletion Algorithm

- Descend to the leaf where the key exists
- Remove the required record from the node
- If the node still has enough records then stop
- If the node does not have enough records then check its neighbors and distribute the keys. Need to fix the upper index node to represent the new split value
- If the nodes does not have enough records and cannot redistribute records from a neighbor, then merge with a neighbor node and remove the old split value from the parent node

Summary: B+ trees

- Typically, *67%* occupancy on average.
- Usually preferable to ISAM, modulo *locking* considerations; adjusts to growth gracefully.
- Key compression increases fan-out, reduces height.
- Most widely used index in database management systems because of its versatility. One of the most optimized components of a DBMS.

Process: Choice of indexes

- One approach:
 - Consider the most important queries in turn.
 - Consider the best plan using the current indexes, and see if a better plan is possible with an additional index. If so, create it.
- Must understand how a DBMS evaluates queries and creates query evaluation plans.
- Before creating an index, must also consider the impact on updates in the workload.
- Trade-off: Indexes can make select queries go faster, updates slower. Require disk space, too.

Index selection guideline

- Attributes in WHERE clause are candidates for index keys.
 - Exact match condition suggests hash index.
 - Range query suggests tree index.
- Clustering is especially useful for range queries; can also help on equality queries if there are many duplicates.
- Multi-attribute search keys should be considered when a WHERE clause contains several conditions.
 - Order of attributes is important for range queries.
 - Such indexes can sometimes enable index-only strategies for important queries: when only indexed attributes are needed.
 - For index-only strategies, clustering is not important.
- Try to choose indexes that benefit many queries.
 - Since only one index can be clustered per relation, choose it based on important queries that would benefit the most from clustering.

Indexes in InnoDB

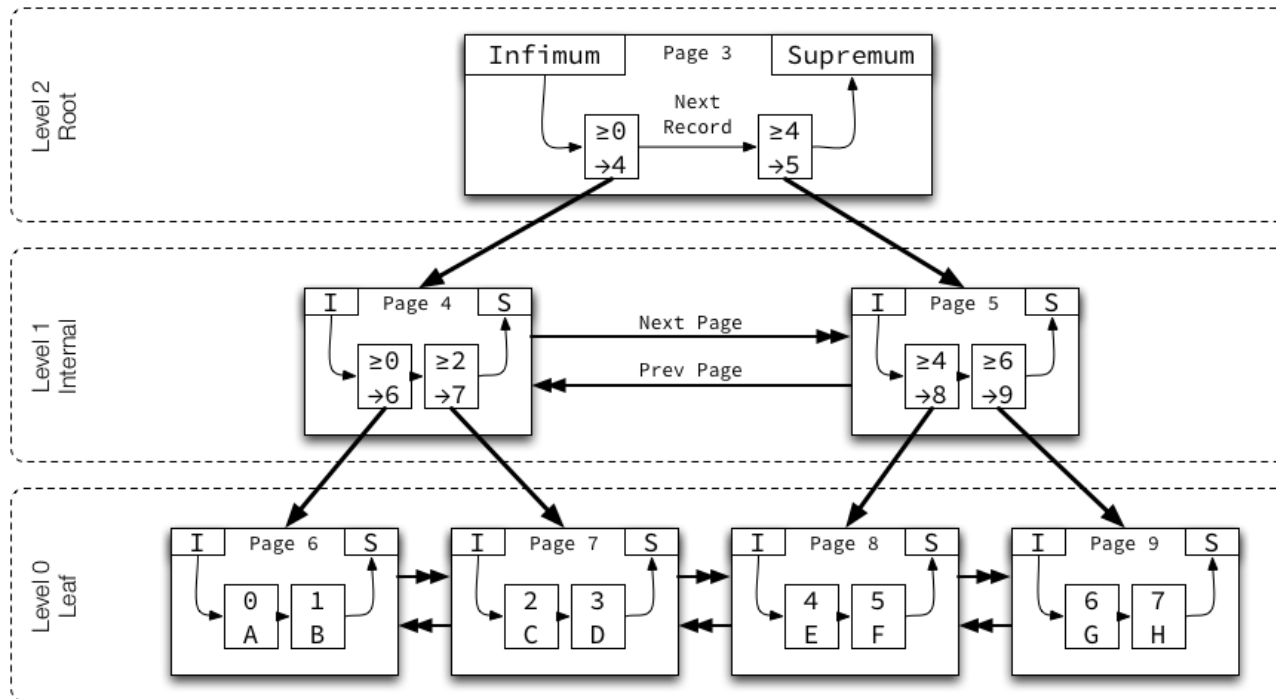
- Every InnoDB table has a special index called the clustered index (by default built on the primary key)
- Record locks always lock index records
 - Even if a table is defined with no indexes
 - InnoDB creates a hidden clustered index and uses this index for record locking
- Accessing a row through the clustered index is fast because the index search leads directly to the page with all the row data.
- All InnoDB indexes are B+ trees where the index records are stored in the leaf pages of the tree
- You can configure the page size for all InnoDB tablespaces in a MySQL instance with the variable `innodb_page_size`
 - default size of an index page is 16KB

InnoDB Index structure

- Root page is allocated when the INDEX is created and is stored in the data dictionary
 - It can never be relocated or removed
- All pages at each level are double-linked to each other
- All pages have anchors for the beginning and the end of the linked list of records
 - Statically defined: Infimum – lowest key, supremum – highest key
- Within a page, records are singly-linked in ascending order
 - Records not stored in ascending order
- Non-leaf pages contain page addresses to a child node
- Leaf pages contain the actual data record (non-key data)

Tree Levels in InnoDB

B+Tree Structure



Levels are numbered starting from 0 at the leaf pages, incrementing up the tree.

Pages on each level are doubly-linked with previous and next pointers in ascending order by key.

Records within a page are singly-linked with a next pointer in ascending order by key.

Infimum represents a value lower than any key on the page, and is always the first record in the singly-linked list of records.

Supremum represents a value higher than any key on the page, and is always the last record in the singly-linked list of records.

Non-leaf pages contain the minimum key of the child page and the child page number, called a "node pointer".

InnoDB: Secondary Index

- You can create multiple indexes on a table
 - These additional indexes that are not on the primary key are secondary indexes
- Each record in a secondary index contains the primary key columns for the row, as well as the columns specified for the secondary index
- InnoDB uses this primary key value to search for the row in the clustered index

Index Optimizations in InnoDB

- The change buffer is a special data structure that caches changes to secondary index pages when affected pages are not in the buffer pool
- The buffered changes are merged later when the pages are loaded into the buffer pool by other read operations.
- secondary indexes are usually non-unique, and inserts into secondary indexes happen in a relatively random order.
- Merging cached changes at a later time, when affected pages are read into the buffer pool by other operations, avoids substantial random access I/O that would be required to read-in secondary index pages from disk
- Periodically, the purge operation that runs when the system is mostly idle, writes the updated index pages to disk
- The purge operation can write disk blocks for a series of index values more efficiently than if each value were written to disk immediately

InnoDB locks

- MySQL sets record locks on every index record that is scanned in the processing of a SQL statement
- Types of object locks
 - Record lock: This is a lock on an index record.
 - Gap lock: This is a lock on a gap between index records, or a lock on the gap before the first or after the last index record.
 - Next-key lock: This is a combination of a record lock on the index record and a gap lock on the gap before the index record.
- InnoDB uses next-key locks for searches and index scans
- If one session has a shared or exclusive lock on record R in an index, another session cannot insert a new index record in the gap immediately before R in the index order

InnoDB hash indexes

- Based on the observed pattern of searches, MySQL builds a hash index using a prefix of the index key.
 - Hash indexes are built on demand for those pages of the index that are often accessed.
- The prefix of the key can be any length, and it may be that only some of the values in the B+tree appear in the hash index.
- If a table fits almost entirely in main memory, a hash index can speed up queries by enabling direct lookup of any element, turning the index value into a sort of in memory pointer.

Summary: Tree-based Index

- Tree-structured indexes are ideal for range-searches, also good for equality searches.
- ISAM is a static structure.
 - Only leaf pages modified; overflow pages needed.
 - Overflow chains can degrade performance unless size of data set and data distribution stay constant.
- B+ tree is a dynamic structure.
 - Inserts/deletes leave tree height-balanced; $\log_F N$ cost.
 - High fanout (**F**) means depth rarely more than 3 or 4.
 - Almost always better than maintaining a sorted file.
- InnoDB provides many optimizations to speed up the access to a record