

CS3223 Database Systems

Implementation

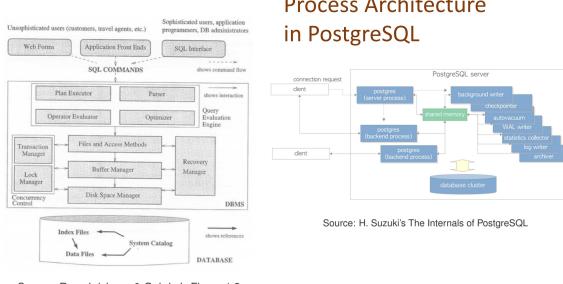
AY23/24 Sem 2, github.com/gerteck

Introduction

Course Details

- Prerequisite Knowledge: CS2040S, CS2102, CS2106 background (helpful).
- Reference Textbook: Raghu & Johannes Database M. Systems, 2002. Encouraged to read ahead based on schedule before the lecture.
- Course covers data structures, algorithms, different components making up database systems.

Architecture of DBMS



Source: Ramakrishnan & Gehrke's Figure 1.3

- OLTP:** Online Transaction Processing is a type of data processing that consists of executing a number of transactions occurring concurrently—online banking, shopping, order entry, or sending text messages, for example.
- OLAP:** Online Analytical Processing.
- Focusing on centralized database running on a single server.

1. Data Storage

References: R&G Chapt 8. (Storage & Indexing Overview), Chapt 9. (Storing Data: Disks and Files).

A DBMS stores

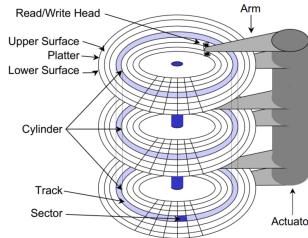
- Relations (Actual tables)
- System catalog (aka data dictionary) storing metadata about relations. (Relation schemas, structure of relations, constraints, triggers. View definitions, Indexes - derived info to speed up access to relations, Statistical information about relations for use by query optimizer.)
- Log files: Information maintained for data recovery.

DBMS Storage

Memory Hierarchy: Primary (registers, RAM), secondary (HDD, SSD), tertiary memory with capacity / cost / access speed / volatility tradeoffs.

- DBMS stores data on non-volatile disk for persistence.
- DBMS processes data in main memory (RAM).
- Disk access operations (I/O). Read: transfer data from disk to RAM. Write: transfer data from RAM to disk.
- Make use of index to speed up access, so that don't have to retrieve all the data when you run a query. Retrieve index and read only the block that contains specified data. Minimize I/O cost.

Magnetic Hard-Disk Drive HDD



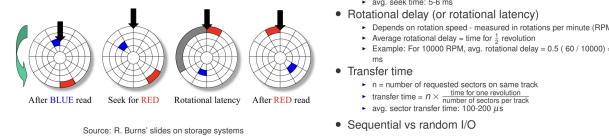
Source: R. Burns' slides on storage systems

- Cylinder, Track, Sector: Units of the HDD storage system. To read from different tracks, need to move the mechanical HDD arm.

Disk Access Time:

- command processing time: interpreting access command by disk controller.
- seek time: moving arms to position disk head on track.
- rotational delay: waiting for block to rotate under head.
- transfer time: actually moving data to/from disk surface.
- access time = seek time + rotational delay + transfer time. (CPT considered negligible).**

Disk Access Time Components



Concept of Sequential vs random I/O.

- Sequential:** Both sector on same track.
- Random:** Sectors on different track, require seeking (moving arm).
- Given a set of data, we hope to store the data contiguously, on the same track. (Minimize incurring random I/O). If data is too large, store on same track, but different surface (aka same cylinder).
- Complexity hidden to OS by disk controller. Shown as a sequence of memory locations.

Solid-State Drive: SSD

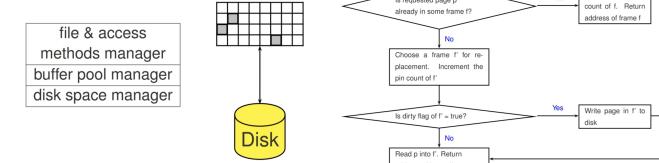
- Build with NAND flash memory without any mechanical moving parts. Lower power consumption.
- Random I/O:** 100x faster than HDD. (no moving parts)
- Sequential I/O:** slightly faster than HDD (2x)
- Disadvantages:** update to a page requires erasure of multiple pages (5ms) before overwriting page. Limited number of times a page can be erased ($10^5 - 10^6$)

Storage Manager Components

- Data is stored, retrieved in units called **disk blocks (or pages)**.
 - Each block = sequence of one or more contiguous sectors.
- Files & access methods layer (aka file layer)** - deals with organization and retrieval of data.
- Buffer Manager** - controls reading/writing of disk pages.

- Disk Space Manager** - keeps track of pages used by file layer.

Storage Manager Buffer Manager BM: Handling request for page p Components



Buffer Manager

- Buffer pool:** Main memory allocated for DBMS.
- Buffer pool is partitioned into block-sized pages called **frames**.
- Clients of buffer pool can request for disk page to be fetched into buffer pool, release a disk page in buffer pool.
- A page in the buffer is **dirty** if it has been modified & not updated on disk.
- Two variables** maintained for each frame in buffer pool:
 - pin count:** number of clients using page (initialized 0)
 - dirty flag:** whether page is dirty (initialized false)
- Free list: Keeps track of frames that are free / empty.
- Pin count:**
 - Incrementing pin count is **pinning** the requested page in its frame.
 - Decrementing is **unpinning** the page.

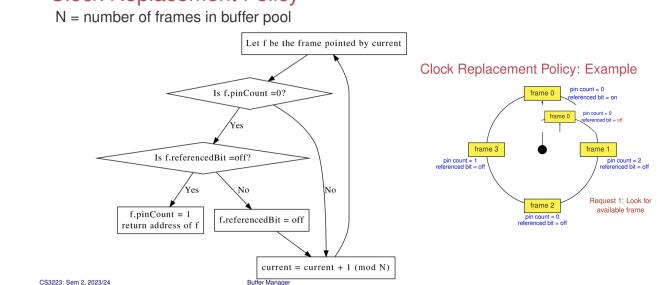
- Unpinning a page, dirty flag should be updated to true if page is dirty.
- A page in buffer can be replaced only when pin count is 0.
- Before replacing buffer page, needs to be written back to disk if its dirty flag is true.

- Buffer manager coordinates with transaction manager to ensure data correctness and recoverability.

Replacement Policies

- Replacement policy: Deciding which unpinned page to replace. (some examples:)
- Random, FIFO, Most Recently Used (MRU), Least Recently Used (LRU): (Use queue of pointers to frames with pin count = 0), most common, makes use of temporal locality.
- Clock:** cheaper popular variant of LRU
 - current** variable: points to some buffer frame.
 - Each frame has a **referenced bit**, turns on when its pin count turns 0.
 - Replace a page that has referenced bit off & pin count = 0.

Clock Replacement Policy



Files

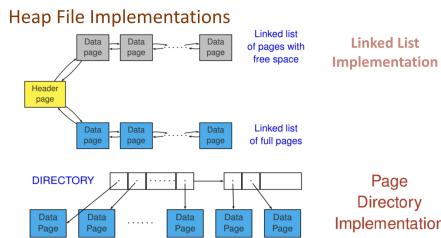
File Abstraction

- Each relation is a file of records.
- Each record has a unique record identifier called RID / TID.
- Common file operations: create/delete file, insert record, delete/get record with given RID, scan all records.

File Organization: Method of arranging data records in a file that is stored on disk.

- **Heap file:** Unordered file
- **Sorted file:** Records order on some search key.
- **Hashed file:** Records located in blocks via a hash function.

Heap File Implementations



- **Linked list implementation:** Two linked lists, one with pages with free space, other of completely full pages.
- **Page Directory Implementation:** Two leveled implementation. Each big block is a disk block with some metadata. Each disk block has a number of data pages.

Page Formats:

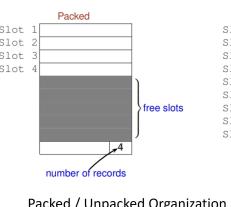
Records are organized within a page and referenced with the RID.

- **RID = (page id, slot number)**
- For **Fixed-Length Records**, Organization can be:

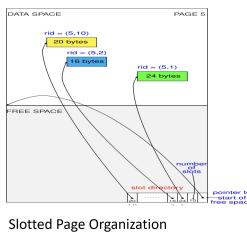
- **Packed Organization:** Store records in contiguous slots.
- For packed organization, memory organization is tough and costly when record in slot is deleted, need to move up a record. But as RID serves as a reference, but need to propagate change in RID.
- **Unpacked Organization:** Uses bit array to maintain free slots.
- For unpacked organization, more bookkeeping needed (use bitmap, 1 & 0 to check if occupied) to store records.

- For **Variable-Length Records**: We could assume some maximum size, then use packed organization. But wasteful. Instead, we can use **Slotted Page Organization**.

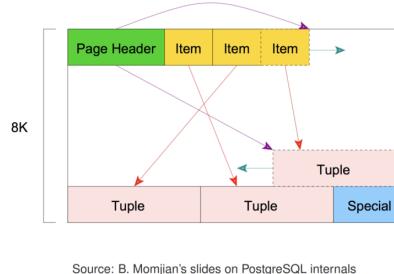
Fixed-Length Records:



Variable-Length Records:



PostgreSQL's Slotted Page Organization



Source: B. Momjian's slides on PostgreSQL internals

Record Formats: Organizing fields within a record.

- **Fixed-Length Records**
 - ▶ Fields are stored consecutively

F1	F2	F3	F4
----	----	----	----
- **Variable-Length Records**
 - ▶ Delimit fields with special symbols
 - ▶ Use an array of field offsets

F1	\$	F2	\$	F3	\$	F4
----	----	----	----	----	----	----

Each o_i is an offset to beginning of field F_i

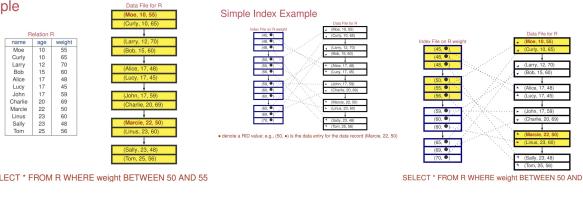
2. Indexing

Need some auxiliary data structure to make efficient queries.

Index

- An **index** is a data structure to speed up retrieval of data records based on some search key.
- A **search key** is a sequence of k data attributes, $k \geq 1$. (A search key is aka *composite search key* if $k > 1$, e.g. (state, city).)
- An index is a **unique index** if search key is a candidate key, otherwise it is **non-unique index**.
- An index is stored as a file, records in index file referred to as **data entries**.

Example



Index Types

Two main types of indexes

- **Tree-based Index:** Based on sorting of search key values (E.g. ISAM, B^+ -tree)
- **Hash-based Index:** Data entries accessed using hashing function (E.g. static/ extendible / linear hashing)
- Considerations when choosing an index:
 - Search Performance (Equality search: $k = v$, use hash-based.) (Range search, use tree)
 - Storage overhead
 - Update performance

Tree-based Indexing: B^+ -Tree

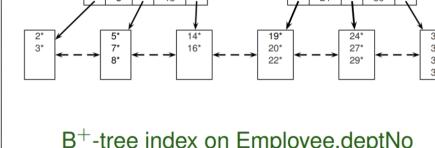
B^+ tree is a dynamic structure that adjusts to changes in the file gracefully, most widely used index structure as it adjusts well to changes and supports both equality and range queries.

- **Balanced tree:** Operations (insert, delete) on tree keep it balanced.
- **Internal nodes** direct the search.
- **Leaf nodes** contain the data entries. Leaf pages linked using page pointers for easy traversal of sequence of leaf pages in either direction.
- **Value d** is parameter of B^+ -tree, called order of the tree, is a measure of capacity of a tree node. Each node contains m entries, where $d \leq m \leq 2d$, except root node, where $1 \leq m \leq 2d$

B^+ -tree Index

B^+ -tree Index

Employee	
name	deptNo
Alice	5
Curly	19
Bob	39
Dave	38
Eve	14
Fred	33
Harry	2
John	34
Ken	6
Larry	27
Linus	24
Lucy	3
Marcie	22
Moe	29
Sally	20
Tom	7
	...



- Each node is either a **leaf node** (bottom-most level) or an internal node.
- Top-most internal node is the **root node** located at **level 0**.
- **Height of Tree** = number of level of internal nodes. (Leaf nodes are at level h where $h =$ height of tree).
- Nodes at same level are **sibling nodes** if they have the same parent node.
- **Leaf Nodes:**

- Leaf nodes store sorted data entries.
- $k*$ denote data entry of form (k, RID) , where k = search key value of corresponding data record, $RID =$ RID of data record.
- Lead nodes are doubly-linked to adjacent nodes.

• Internal Nodes:

- Internal nodes store index entries of the form $(p: pointer, k: separator)$ ($p_0, k_1, p_1, k_2, p_2, \dots, p_n$)
- $k_1 < k_2 < \dots < K_n$
- Each (k_i, p_i) is an **index entry**, k_i serves as **separator** between node contents pointed to by p_{i-1} & p_i
- p_i = disk page address (root node of an index subtree T_i)

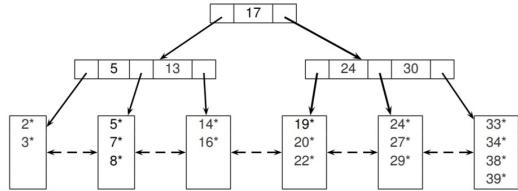
B^+ -tree Index Properties

Properties of B^+ -tree Index

- Dynamic index structure; adapts to data updates gracefully
- Height-balanced index structure
- Order of index tree, $d \in \mathbb{Z}^+$

1. Controls space utilization of index nodes
2. Each non-root node contains m entries, where $m \in [d, 2d]$
3. The root node contains m entries, where $m \in [1, 2d]$

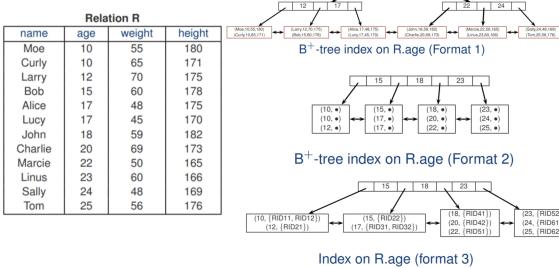
Example: B^+ -tree with order = 2



Formats of Data Entries in B-Tree

- Format 1: k^* is actual data record (with search key value k)
- Format 2: k^* is of form (k, rid) , where rid is record identifier of record with search key value k .
- Format 3: k^* is of form $(k, rid-list)$, where rid-list is list of record identifiers of data records with search key value k .
- Note, examples assume Format 2.

Formats of Data Entries: Example

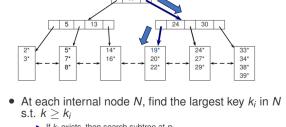


Index on R.age (format 3)

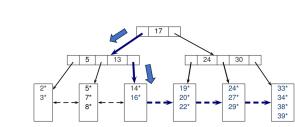
B^+ -tree Search Algorithms

- Search algorithm finds the leaf node a given data entry belongs to.
- We assume no duplicates, no data entries same key value. Note in practice, duplicates arise whenever search key does not contain candidate key, must be dealt with.

Equality Search ($k = 19$)



Range Search ($k \geq 15$)



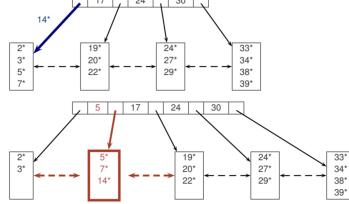
B^+ -Tree Insertion

- Algorithm for insertion takes an entry, finds the leaf node where it belongs, and inserts it there.
- Occasionally, a node is full and must be split. (More than $2d$ entries) When node is split, entry pointing to the node created by the split must be inserted into the parent.
- If the (old) root is split, a new root node is created and height of tree increases by 1.

Splitting of overflowed node

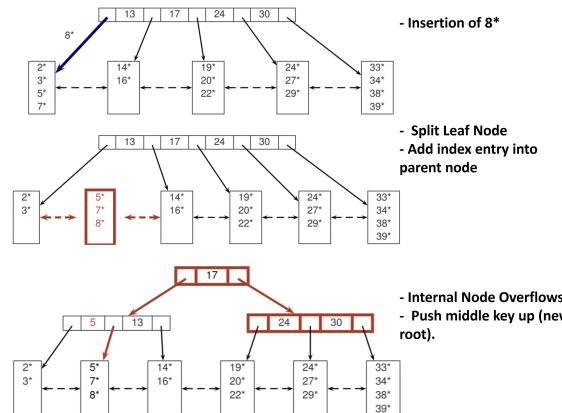
- Split overflowed leaf node by distributing $d + 1$ entries to new leaf node.
- Create a new entry index using smallest key in leaf node.
- Insert new index entry into parent node of overflowed node.

Inserting 14* (Splitting of overflowed node)

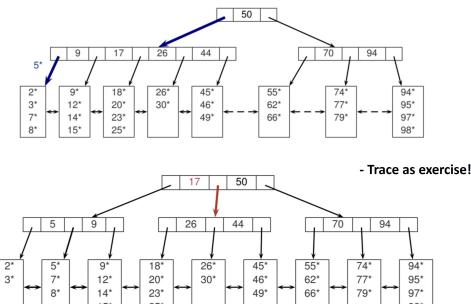


- Sometimes, node split is propagated upwards to ancestor internal nodes.
- When splitting an internal node, the middle key is pushed to parent node.

Inserting 8* (Propagation of node splits)



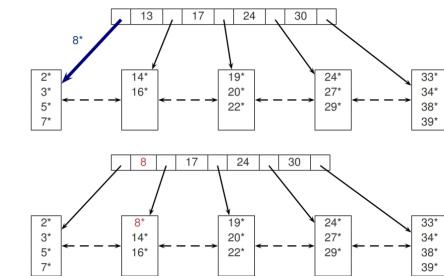
Inserting 5* (Propagation of node splits)



Redistributing of data entries in Overflow

- A node split can sometimes be avoided by distributing entries from overflowed node to a non-full adjacent sibling node.

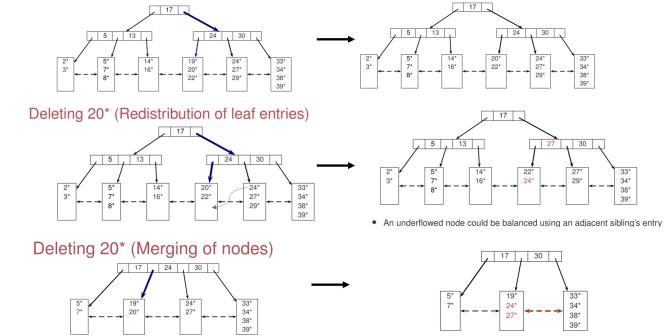
Inserting 8* (Redistribution of data entries)



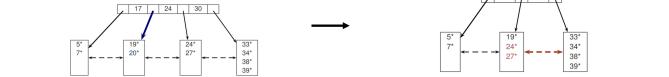
B^+ -Tree Deletion

- Algorithm for deletion takes an entry, finds leaf node it belongs to, and deletes it.
- Underflowed node: When node is at minimum occupancy before deletion, and goes below threshold, we must either redistribute entries from adjacent sibling, or merge node with sibling to maintain minimum occupancy.
- Merging: Underflowed node needs to be merged if each of adjacent sibling nodes has exactly d entries.

Deleting 19* (Simple Case)



Deleting 20* (Redistribution of leaf entries)



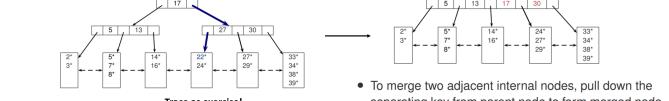
An underflowed node could be balanced using an adjacent sibling's entry

Deleting 20* (Merging of nodes)



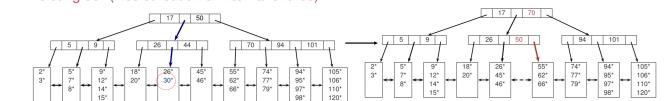
- Node mergers may propagate upwards.

Deleting 22* (Propagation of node merges)



To merge two adjacent internal nodes, pull down the separating key from parent node to form merged node

Deleting 30* (Redistribution of internal entries)



- Chances are high that redistribution is possible if node has two siblings, and unlike merging, redistribution is guaranteed to propagate no further than parent node. Also, pages have more space, reducing likelihood of split on subsequent insertions.

B⁺ Tree Bulk Loading

- Entries added to a B⁺ in two ways.

- Have existing collection of data records with B⁺ tree index on it.
When record added to collection, corresponding entry added to B⁺ tree. (Insert, Delete individually)
- Have collection of data records we want to create new B⁺ tree index on some key field(s). Start with an empty tree. Inserting one by one expensive due to overhead, systems provide **bulk loading** utility.

Bulk Loading:

- Sort data entries $k*$ to be inserted into B⁺ tree according to search key k . (Here, $d = 1$).
- Allocate empty page to serve as root. Insert a pointer to first page of (sorted entries into it).

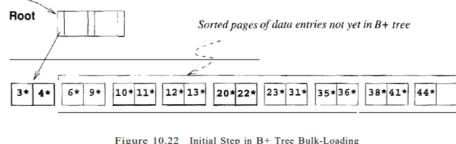


Figure 10.22 Initial Step in B+ Tree Bulk-Loading

- Add one entry to root page for each page of sorted data entries.
Proceed until root page is full. Here, we must split root and create a new root page.

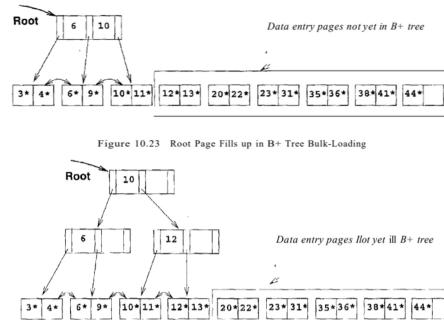


Figure 10.23 Root Page Fills up in B+ Tree Bulk-Loading

- To continue, entries for leaf pages **always inserted into right-most index page just above the leaf level**. When right-most page above leaf level fills up, it is split.

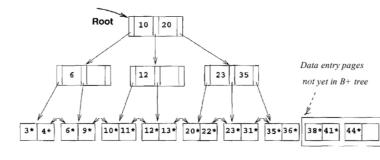


Figure 10.25 Before Adding Entry for Leaf Page Containing 38*

3. Hash-based Indexing

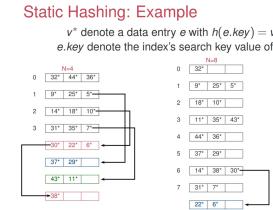
- Used for **equality queries**, not for range queries.
- Hashing techniques:** **Static hashing**, **dynamic hashing** (linear hashing, extendible hashing, etc).

Static Hashing

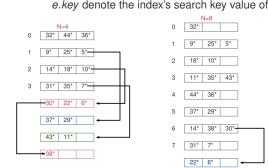
- Data stored in N buckets**, fixed at creation time. Bucket consists of **one primary data page & chain** of zero or more overflow data pages.
- v^* represents data entry e with $h(e.key) = v$, not search entry with RID.
- Problem with static hashing:** As data grows, longer overflow chain, efficiency drops. Need to periodically rehash and increase no. of buckets.

Static Hashing

- Data is stored in N buckets B_0, B_1, \dots, B_{N-1}
- N is fixed at creation time
- Hashing function $h(\cdot)$ is used to identify the bucket to store a record
 - A record with search key K is inserted into bucket B_i , where $i = h(K) \mod N$
 - $\{K\}$ maps the search key value into a bit string
- Each bucket consists of **one primary data page** & a chain of zero or more overflow data pages



Static Hashing: Example



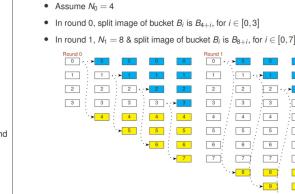
Dynamic Hashing: Linear Hashing

- Hash file grows / shrinks linearly, systematic **splitting of buckets**.
- Overflow pages needed as overflowed bucket may not be split immediately.
- Hashing function changes dynamically and at given instant, **at most two (successive) hashing functions** used by the scheme during search.
- Each bucket has primary data page & chain of zero+ overflow pages.
- Insert in bucket B_i overflows if all pages in B_i (primary + overflow) full.

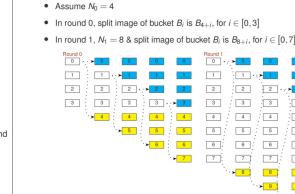
Linear Hashing (cont.)

- Assume initial file size of N_0 buckets
 - Buckets $B_0, B_1, \dots, B_{N_0-1}$
- File grows linearly by **splitting buckets** in rounds
 - At each round, buckets are split sequentially: B_i is split before B_{i+1}
- How to split a bucket B_i ?
 - Add a new bucket B_i' (known as split image of B_i)
 - Redistribute entries in B_i between B_i & B_i'
- File size increases by one bucket after each bucket split
- At the end of one round of splitting (i.e., every bucket at the start of the round has been split), file size is doubled.
- Let N_i denote the file size at the beginning of round i ($i = 0, 1, \dots$)
 - $N_0 = N_0$
- At the end of round i , N_i new buckets are added: $B_{N_0}, B_{N_0+1}, \dots, B_{N_i-1}$
 - In round i , the split image of B_i is B_{N_i+i} for $j \in [0, N_i - 1]$

Splitting Buckets



Redistributing Entries



Dynamicity of Linear Hashing

Dynamic Hashing

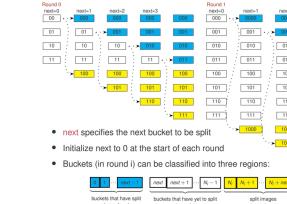
next = 1

- Recall that v^* denote a data entry e with $h(e.key) = v$ and e key denote the index's search key value of e
- Each round uses two hash functions: functions h_i and h_{i+1} for round i
 - $h_i(v) = h_{i+1}(v)$
 - B_i is the bucket for search key v if B_i had not been split, where $x = h_i(v)$
 - B_i' is the bucket for search key v if B_i had been split, where $y = h_{i+1}(v)$
- Keep track of which bucket to be split next using variable next

Linear (Hashing) Splitting of Buckets (Insertion)

- Number in buckets represent the rightmost bit values.
- Split Criteria:** variable, could be when some bucket overflow, space utilization of file above some threshold etc.
- Level:** We use level to denote splitting round number (use with next).

Splitting Buckets



When to split a bucket?

- The time to split the next bucket can be decided with various criteria:
 - Split whenever some bucket overflows
 - Split whenever space utilization of file is above some threshold
 - etc.
- Overflow pages are needed since an overflowed bucket might not be split immediately
- We shall assume that a bucket split is triggered whenever some bucket overflows
 - Bucket B_j overflows if $l(B_j) > l(B_{j-1})$ (i.e., primary + overflow pages) are full
- We use level to denote the splitting round number

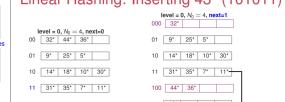
Examples: Linear Hashing Search & Insert

- Even though overflowed bucket split, not necessary mean enough space if bit value not right. Still require overflow page.
- Using **level** and **next**, we determine how many additional bits to consider in **second hashing function**.
- Second hashing function:** Is simply looking at the m rightmost bits of the hash.

Inserting data entry with search key k



Example: Linear Hashing: Inserting 43* (101011)



Linear Hashing Deletion

- Opposite of insertion.
- Two cases:** 1. (Next $\neq 0$), 2. (Next = 0, and level > 0)
 - Case 1: If $(next \neq 0)$ and $(level > 0)$
 - Locate bucket and delete entry
 - If the last bucket B_{next-1} becomes empty, it can be removed
 - Case 2: If $(next = 0)$ and $(level > 0)$
 - Update next to point to the last bucket in previous level $B_{next-1-1}$
 - Decrement level by one

Linear Hashing: Deletion

- Locate bucket and delete entry
- If the last bucket B_{next-1} becomes empty, it can be removed

Linear Hashing: Deletion (cont.)

- Case 2:** If $(next = 0)$ and $(level > 0)$
 - Update next to point to the last bucket in previous level $B_{next-1-1}$
 - Decrement level by one

Linear Hashing Performance, Summary

Linear Hashing: Performance

- One disk I/O unless the bucket has overflow pages
 - On average 1.2 disk I/O for uniform data distribution
 - Worst case: $O(I)$ cost is linear in the number of data entries
- Poor space utilization with skewed data distribution

Linear Hashing: Summary

- File grows dynamically by splitting buckets linearly in rounds
 - Each bucket split adds one new bucket to file
- The bucket B_i for a search key K is determined using two hash functions (h_i and h_{i+1})
 - Assume initial file has 2^n buckets & n bits per file
 - Round 0: $h_i(K) \mod N$ = value of last ($m-i$) bits of $h(K)$
 - Round 1: $h_{i+1}(K) \mod N$ = value of last ($m-i-1$) bits of $h(K)$
- In this course, we assume a bucket split is triggered whenever some bucket overflows due to an insertion
- Overflow pages are required for an overflown bucket if it is not the next bucket to be split

Dynamic Hashing: Extendible Hashing

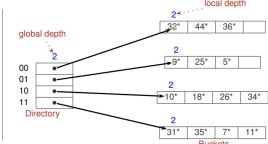
- Similar to Linear Hashing:** we want **bucket number to grow dynamically**, and use some **number of least significant bits of $h(k)$** to determine bucket address for search key k .
- Difference:** Add new bucket (as split image) when existing bucket overflows, No overflow pages (except when number of collisions exceed page capacity, two page entries collide if they have same $h(.)$ hash value).

Extendible Hashing

- Extendible hashing:** dynamically updatable disk-based index structure, implements hashing scheme utilizing a **directory of pointers to buckets**.
- Overflows handled by doubling the directory which logically doubles the number of buckets. **Physically, only the overflow bucket is split.**

Extendible Hashing

- Uses a directory of pointers to buckets
- Directory has 2^d entries
- Each entry has a unique d -bit address $b_1b_2 \dots b_db_d$
- Two directory entries are said to correspond if their addresses differ only in the i^{th} bit (i.e., b_i), such entries are called corresponding entries.
- Each bucket maintains a local depth denoted by $\ell \in [0, d]$
- All entries in a bucket with local depth ℓ have the same last ℓ bits in $h(.)$



Extendible Hashing Performance

- Performance:** At most 2 disk I/O for equality selection, at most 1 I/O if directory fits in main memory.
- Handling collision:** Two data entries **collide** if same hashed value, overflow pages need when number of collisions exceed page capacity.
- Compared with B+-tree index exact match queries (\log number of I/Os), E. Hashing better expected query cost $O(1)$ I/O.

Extendible Hashing: Handling Bucket Overflow

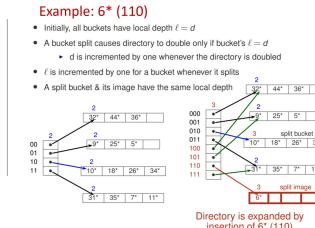
- Main idea: Determine if there is empty directory entry to point to new bucket. **2 Cases:** decision to split, or use empty directory entry.

Case 1: Split bucket local depth = global depth

Extendible Hashing

Handling Bucket Overflow (Case 1)

- When a bucket overflows, it is split
 - Allocate a new bucket called its **split image**
 - Redistribute entries (including new entry) between split bucket & its split image
- Case 1:** Split bucket's local depth is equal to global depth
 - When the directory is doubled,
 - Each new directory entry (except for the entry for the split image) points to the same bucket as its corresponding entry
- Number of directory entries pointing to a bucket = $2^{d-\ell}$

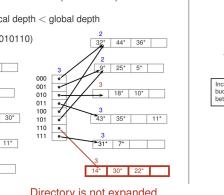


Case 2: Split bucket local depth < global depth.

Extendible Hashing

Handling Bucket Overflow (Case 2)

- Case 2:** Split bucket's local depth < global depth
- Example: Inserting 22* (010110)



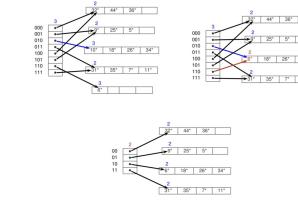
Extendible Hashing Deletion

- To delete entry, simply locate Bucket and delete.
- Merging:** Mergeable if entries can fit within a bucket, and same local depth, j differs on in l^{th} bit.

Extendible Hashing: Deletion

- Locate bucket B_j containing entry & delete entry
- If B_j becomes empty, B_j can be merged with the bucket B_i where both buckets have the same local depth ℓ and $i \& j$ differs only in the l^{th} bit
 - B_i is deallocated
 - B_i 's local depth is decremented by one
 - Directory entries that point to B_i are updated to point to B_j
- More generally, B_i & B_j (with same local depth ℓ and $i \& j$ differs only in the l^{th} bit) can be merged if their entries can fit within a bucket
- If each pair of corresponding entries point to the same bucket, directory can be halved
 - d is decremented by one

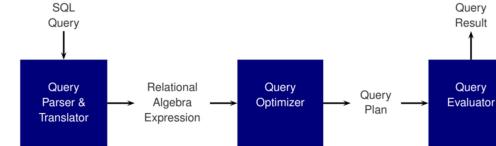
Example: Deleting 10* (1010)



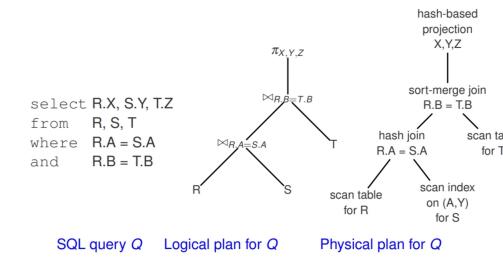
4. Query Evaluation: Sorting & Selection

DBMS describes data it manages using tables, indexes (metadata), which is stored in special tables (system catalogs). This data used to find best way to evaluate a query.

Query Processing



Query Plans



- SQL queries translate into extended form of relational algebra
- Query evaluation plans represented as tree of relational operators, with labels identifying algorithm to use at each node.
- R. operators building blocks for evaluating queries, implementation of operators optimized for good perf.

External Sorting

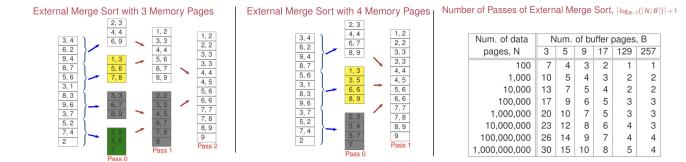
Sorting collection of records on some (search) key is useful and required in variety of situations, including

- Some sorted table of results: `SELECT * FROM student ORDER BY age .`
- Bulk loading a B^+ -tree index
- Implementation of other relational algebra operators (e.g. projection, join), which require some sorting step.

When data to be sorted too large to fit into available main memory. Need some **external sorting algorithm**. Algos seek to minimize cost of disk accesses.

External Merge Sort

- Main Idea:** Pass 0: Creating initial sorted runs (each of X memory pages), then continue during merging passes till you get final sorted pass.
- Sort entire file by breaking into smaller subfiles, sorting subfiles and merging using minimal amount of main memory at given time.
- Each sorted subfile is referred to as a run.**
- Sorting 11-page data R using 3 vs 4 memory pages:



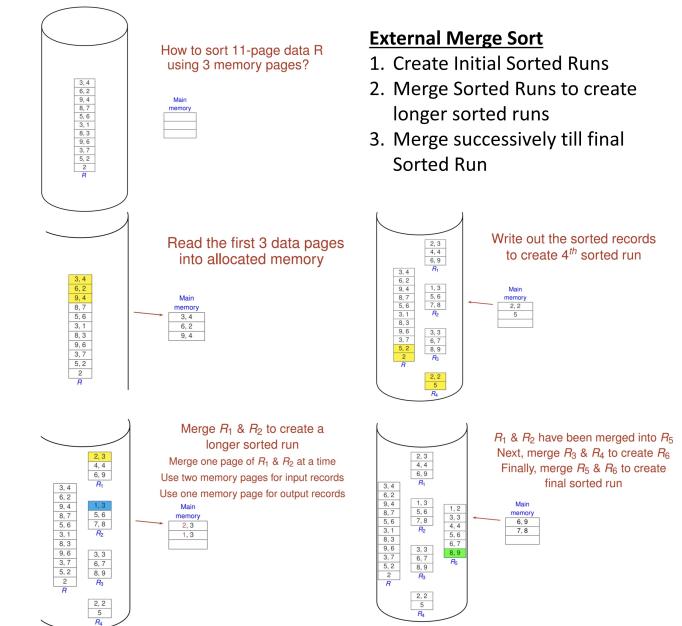
External Merge Sort Analysis

- Note:** We consider only I/O costs, which approx by counting no. of pages read/written as per cost model. (Simple cost model to convey main idea).

External Merge Sort

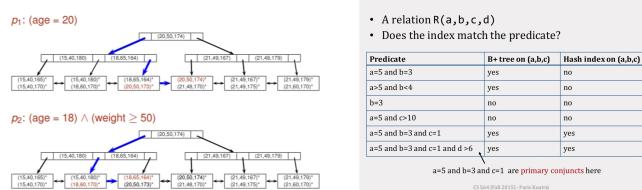
- Pass i , $i \geq 1$: Merging of sorted runs
 - Use $B-1$ buffer pages for input & one buffer page for output
 - Performs $(B-1)$ -way merge
- Analysis:**
 - N_0 = number of sorted runs created in pass 0 = $[N/B]$
 - Total number of passes = $\lceil \log_{B-1}(N_0) \rceil + 1$
 - Size of each sorted run = B pages (except possibly for last run)
 - * Each pass reads N pages & writes N pages

External Merge Sort Steps



Examples of Index matching CNF Selection

Example: B⁺-tree on (age,weight,height)



Primary and Covered Conjuncts

- Primary Conjuncts:** Subset of conjuncts in selection predicate p that index I matches.
- In general, only subset of conjuncts of predicate matches index.
- Covered Conjunct:** Conjunct C in predicate p covered if all attributes in C appear in the key, or *include column(s)* of index I .
- Primary conjuncts subset of covered conjuncts.

Cost of Evaluation of Selection Predicate p

Notation	Meaning
r	relational algebra expression
$\ r\ $	number of tuples in output of r
$ r $	number of pages in output of r
b_d	number of data records that can fit on a page
b_i	number of data entries that can fit on a page
F	average fanout of B ⁺ -tree index (i.e., number of pointers to child nodes)
h	height of B ⁺ -tree index (i.e., number of levels of internal nodes)
$h = \lceil \log_F(\lceil \frac{\ r\ }{b_i} \rceil) \rceil$ if format-2 index on table R	
B	number of available buffer pages

Cost of B⁺-tree Index Evaluation of p

Let p' = primary conjuncts of p , p_c = covered conjuncts of p

1. Navigate internal nodes to locate first leaf page

$$Cost_{internal} = \begin{cases} \lceil \log_F(\lceil \frac{\|R\|}{b_i} \rceil) \rceil & \text{if } I \text{ is a format-1 index,} \\ \lceil \log_F(\lceil \frac{\|R\|}{b_i} \rceil) \rceil & \text{otherwise.} \end{cases}$$

2. Scan leaf pages to access all qualifying data entries

$$Cost_{leaf} = \begin{cases} \lceil \frac{\|\sigma_{p'}(R)\|}{b_d} \rceil & \text{if } I \text{ is a format-1 index,} \\ \lceil \frac{\|\sigma_{p'}(R)\|}{b_i} \rceil & \text{otherwise.} \end{cases}$$

3. Retrieve qualified data records via RID lookups

$$Cost_{rid} = \begin{cases} 0 & \text{if } I \text{ is covering or format-1 index,} \\ \|\sigma_{p_c}(R)\| & \text{otherwise.} \end{cases}$$

Cost of RID lookups could be reduced by first sorting the RIDs

$$\lceil \frac{\|\sigma_{p_c}(R)\|}{b_d} \rceil \leq Cost_{rid} \leq \min\{\|\sigma_{p_c}(R)\|, |R|\}$$

Example

- B⁺-tree index I = (age, weight, height), Format 2
- Query: `select * from R where p`
 - $p : (age = 18) \wedge (weight = 60) \wedge (height = 3)$
 - $p : (age = 18) \wedge (height > 18) \wedge (weight = 60)$
- $\|R\| = 12$, $\|\sigma_p(R)\| = 9$, $\|\sigma_{p_c}(R)\| = 2$
- $b_d = b_i = 2$, Height of $I = 2$
- Evaluation cost of p using p using $2 + \lceil \frac{9}{2} \rceil + 2 = 9$

Cost of Hash Index Evaluation of p

- Let p' = primary conjuncts of p

For format-1 index

Cost to retrieve data records: at least $\lceil \frac{\|\sigma_{p'}(R)\|}{b_d} \rceil$

For format-2 index

Cost to retrieve data entries: at least $\lceil \frac{\|\sigma_{p'}(R)\|}{b_i} \rceil$

$$\text{Cost to retrieve data records} = \begin{cases} 0 & \text{if } I \text{ is a covering index,} \\ \lceil \frac{\|\sigma_{p'}(R)\|}{b_i} \rceil & \text{otherwise.} \end{cases}$$

Evaluating Non-Disjunctive / Disjunctive Conjuncts

Evaluating Non-Disjunctive Conjuncts

- Consider the query $\sigma_p(R)$, where $p = (age = 21) \wedge (weight \geq 70) \wedge (height = 180)$
- Suppose the available unclustered indexes on R are
 - a hash index H_{age} on (age), and
 - a B⁺-tree index T_{weight} on (weight)
- What are the possible strategies to evaluate the following predicates?

Evaluating Disjunctive Conjuncts

- Suppose the available unclustered indexes are
 - a hash index H_{age} on (age), and
 - a B⁺-tree index T_{weight} on (weight)
 - What are the possible strategies to evaluate the following predicates?
- $p_1 = (age = 21) \vee (height = 180)$
- $p_2 = ((age = 21) \vee (height = 180)) \wedge (weight \geq 70)$
- $p_3 = (age = 21) \vee (weight \geq 70)$

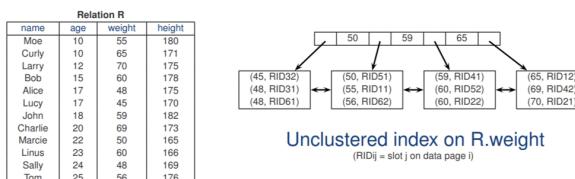
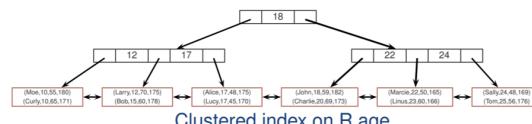
Possible strategies to evaluate Disjunctive / Non-Disjunctive predicates

- File scan, Use both (fetch RID, take Union), Use B+ tree etc.

Clustered vs. Unclustered Index (B+Tree)

- Clustered Index:** Order of its data entries is the same or ‘close to’ order of the data records (in pages).
- Layman Terms: If clustered, if we do a file scan, records will be in order with respect to the attribute.
- An index using Format 1 for data entries is a clustered index.
- Logically, at most one clustered index for each relation.
- Implication:** Tutorial 3 Q4: When doing index scan with RID lookup, for clustered index, RID page I/O incurred will be the number of leaf pages.
- Unclustered Index:** Order of data entries not same as actual order of data records. To retrieve each tuple / entry, need to do separate RID lookup / page retrieval.
- Implication:** Tutorial 3 Q4: When doing index scan with RID lookup, for unclustered index, RID page I/O incurred will be number of pages of tuples (> no. of leaf pages).

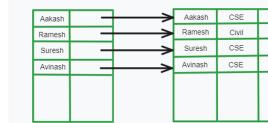
Clustered vs Unclustered Index: Example



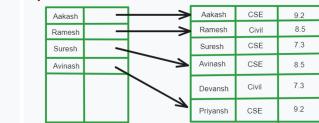
Dense vs. Sparse Index (B+Tree)

- Dense index:** there is an index record for every search key value in the data.
- The total number of records in the index table and main table are the same.
- Gives quick access to records, effective for range searches as each key value has an entry, but takes more storage, and insertion and deletion higher overhead.
- For *unclustered index*, must be dense.
- Sparse index:** Some search key value has no have index record. (Main table index points records in specific gap (range of space where index resides in).
- Uses less storage space, lessen effect of insert/delete on index maintenance operations. Time to locate data in index table more, and sparse index records need to be clustered.
- For sparse index, records need to be clustered (in order).

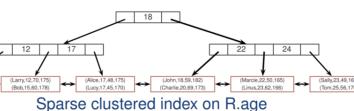
Dense Index



Sparse Index



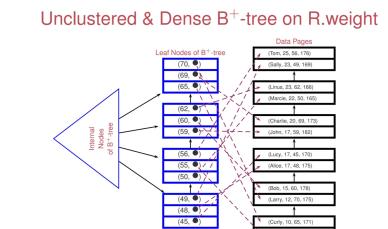
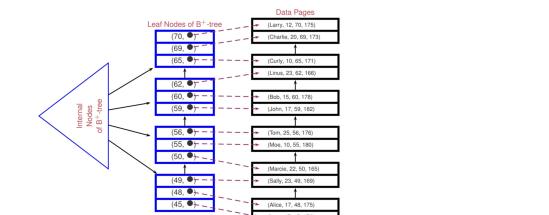
Dense vs Sparse Index: Example



name	age	weight	height
Moe	10	55	180
Curly	10	65	171
Larry	12	70	175
Bob	15	60	178
Alice	17	48	175
Lucy	17	45	170
John	18	59	182
Charlie	20	69	173
Marcie	22	50	165
Linus	23	62	166
Sally	24	48	169
Tom	25	56	176

Dense clustered index on R.age
(RID = slot # on data page i)

Clustered & Dense B⁺-tree on R.weight



PostgreSQL: Buffer Replacement Policy

PostgreSQL: Open source relational database management system.

- **Port Number:** By default, server listens on port number 5432 for client connections.
- PostgreSQL uses variant of Clock policy as default buffer replacement policy.

Overview of PostgreSQL

Shared-memory Data Structures

- As multiple backend server processes may access database at same time, access to shared-memory structures (buffer pool) controlled to ensure consistent access and updates. Use **locks** (spin locks, light-weight locks) to control access.
- **Locking Protocol:** Before accessing shared-memory structure, process acquire lock, upon completion of access, process release lock.

Buffer Manager

- Two main types of Buffers used: Shared Buffer, Local Buffer.
- **Shared Buffer:** Used for holding page from globally accessible relation.
- **Local Buffer:** Used for holding page from temporary relation locally accessible to specific process.

Management of Shared Buffers

- Initially, all shared buffers maintained in free list.
- **Free list:** Buffer in free list if contents invalid. When new buffer needed, check if buffer available in free list, returned to satisfy buffer request.
- If no available buffer, use **buffer replacement policy** to select victim buffer for eviction to make room for new request. (Use variant of Clock algorithm)
- **Record Deleted/Modified:** Not immediately removed/changed. Multiple versions of record maintained to support **multiversion concurrency control**.
- **Vacuuming Process:** Periodically, vacuuming process runs to remove obsolete versions of records that can be safely deleted from relations. If entire page of records removed, buffer holding page becomes invalid, returned to free list.
- **bgwriter:** background writer process that writes out dirty shared buffers to partly help speed up buffer replacement.

PostgreSQL Buffer Pool

- **Buffer Pool:** Implemented as array of disk blocks, index to each array entry referred to a `buffer_id`, each disk block location identified by buffer tag.
- Pin count for each buffer frame known as reference count `refcount`.
- Each buffer frame associated with a buffer descriptor, stores metadata about contents.
- Given a buffer tag, hash-based buffer table is used to efficiently locate buffer id of buffer frame that stores the disk block (corresponding to given tag) if disk block resident in buffer pool.

Implementing LRU Buffer Replacement Policy

- Simple approach to implement LRU Policy: **Stack LRU Method**.
- Use doubly Linked List to link up buffer pages. Page closer to the front is more recently used, than page closer to tail of list.
- Whenever buffer page referenced ("Used"), moved to front of list.
- When replacement page sought from list, **unpinned** buffer page closest to tail (LRU) selected for eviction.
- **Whenever buffer accessed, position needs to be adjusted:**

1. If accessed page in buffer pool already, containing buffer needs to be moved to top of stack.
2. If accessed page not in buffer pool, **free buffer available** to hold page, selected buffer from free list needs to be inserted onto top of stack.
3. If accessed page not in buffer pool, **free list empty**, selected victim buffer moved from current stack position to top of the stack.
4. If buffer in buffer pool **returned to the free list**, buffer removed from stack.

Enhanced LRU (ELRU) Buffer Replacement Policy

- A disadvantage of LRU is that a page that is accessed only once could evict a frequently accessed page from the buffer.
- To address this limitation, ELRU, which is a variant of LRU, keeps track of additional page access information to manage the eviction of buffer pages in two separate groups.
- In contrast to LRU which uses only the last access time of buffer pages, ELRU maintains the last two access times of buffer pages.

Specifications of ELRU

- Specifically, let $B = B_1 \cup B_2$ denote the set of unpinned buffer pages that can be selected for replacement.
- B_1 is the set of buffer pages that have been accessed only once.
- B_2 is the set of buffer pages that have been accessed at least twice.
- The ELRU policy selects a replacement page from B for replacement as follows. If B_1 is non-empty, ELRU selects the least recently accessed page in B_1 as the replacement page; i.e., ELRU applies the conventional LRU policy to select the replacement page from B_1 . Otherwise, if B_1 is empty, ELRU selects the page in B_2 with the smallest second-last access time as the replacement page.

PostgreSQL: Normal Buffer Replacement Strategy (Clock-Sweep)

PostgreSQL standard Buffer Replacement Strategy

- There is a "free list" of buffers that are prime candidates for replacement. In particular, buffers that are completely free (contain no valid page) are always in this list. We could also throw buffers into this list if we consider their pages unlikely to be needed soon; however, the current algorithm never does that. The list is singly-linked using fields in the buffer headers; we maintain head and tail pointers in global variables. (Note: although the list links are in the buffer headers, they are considered to be protected by the `buffer_strategy_lock`, not the buffer-header spinlocks.)
- To choose a victim buffer to recycle when there are no free buffers available, we use a simple clock-sweep algorithm, which avoids the need to take system-wide locks during common operations. It works like this:
- Each buffer header contains a usage counter, which is incremented (up to a small limit value) whenever the buffer is pinned. (This requires only the buffer header spinlock, which would have to be taken anyway to increment the buffer reference count, so it's nearly free.)
- The "clock hand" is a buffer index, `nextVictimBuffer`, that moves circularly through all the available buffers. `nextVictimBuffer` is protected by the `buffer_strategy_lock`.

Clock-Sweep Buffer Replacement Algorithm

- The **algorithm** for a process that needs to obtain a victim buffer is:
 1. Obtain `buffer_strategy_lock`.
 2. If buffer free list is nonempty, remove its head buffer. Release `buffer_strategy_lock`. If the buffer is pinned or has a nonzero usage count, it cannot be used; ignore it go back to step 1. Otherwise, pin the buffer, and return it.
 3. Otherwise, the buffer free list is empty. Select the buffer pointed to by `nextVictimBuffer`, and circularly advance `nextVictimBuffer` for next time. Release `buffer_strategy_lock`.
 4. If the selected buffer is pinned or has a nonzero usage count, it cannot be used. Decrement its usage count (if nonzero), reacquire `buffer_strategy_lock`, and return to step 3 to examine the next buffer.
 5. Pin the selected buffer, and return.
- (Note that if the selected buffer is dirty, we will have to write it out before we can recycle it; if someone else pins the buffer meanwhile we will have to give up and try another buffer. This however is not a concern of the basic select-a-victim-buffer algorithm.)

5. Query Evaluation: Projection & Join

5.1 Projection: $\pi_{A_1, \dots, A_m}(R)$

- $\pi_L(R)$ projects columns given by list L from relation R.
- $\pi_L^*(R)$ same as $\pi_L(R)$ but preserves duplicates.
- Example: select distinct age from R

Projection Operation $\pi_{A_1, \dots, A_m}(R)$

- Projection involves two tasks:

1. Remove unwanted attributes. (from tuples)
2. Eliminate any duplicate tuples produced.

- Two approaches to Project:

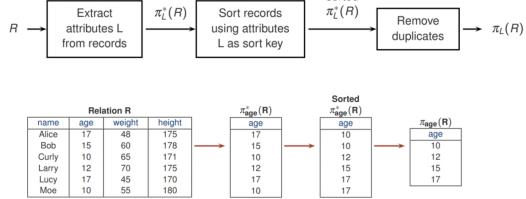
1. Projection based on sorting.
2. Projection based on hashing.

5.1.1 Sort-based Approach

Simple Sort-based Approach

- Treat each step as a black box, push tuples through pipeline to sort.

Consider $\pi_L(R)$ where L denote some sequence of attributes of R



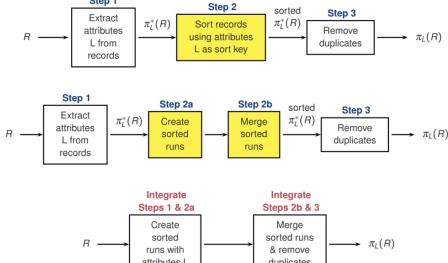
Cost Analysis

- Step 1:**
 - Cost to scan records = $|R|$
 - Cost to output temporary result = $|\pi_L^*(R)|$
- Step 2:**
 - Cost to sort records = $2|\pi_L^*(R)|(\log_m(N_0) + 1)$
 - N_0 = number of initial sorted runs, m = merge factor
- Step 3:**
 - Cost to scan records = $|\pi_L(R)|$

Optimized Sort-based Approach

- By opening black box and examining the sort algorithm (external merge sort), we can optimize the sort-based approach.

Optimized Sort-based Approach



5.1.2 Hash-based Approach

- For , we build a **main-memory hash table T** to detect and remove duplicates.
- Cost = $|R|$ if T fits in main memory.
- Two Phases:** Partitioning phase and Duplicate Elimination phase.

1. Partitioning Phase

Partition R into R_1, R_2, \dots, R_{B-1}

- Hash on $\pi_L(t)$ for each tuple $t \in R$
- $R = R_1 \cup R_2 \cup \dots \cup R_{B-1}$
- $\pi_L^*(R_1) \cap \pi_L^*(R_2) = \emptyset$
- Basically, each hash table does not overlap, and all tables make up R.

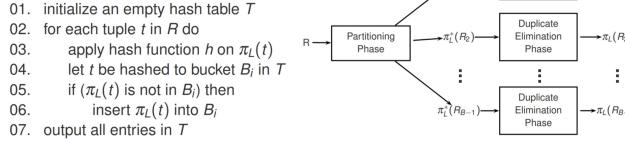
2. Duplicate Elimination Phase

Eliminates duplicates for each partition $\pi_L^*(R_i)$.

- For each tuple t , hash t into bucket B_j with diff. hash function, insert if not already in.
- Output all tuples in all buckets once done.
- Partition with * means contain duplicate.
- Partition with no * means no duplicate.
- $\pi_L(R) = \text{duplicate-free union of } \pi_L(R_1), \dots, \pi_L(R_{B-1})$

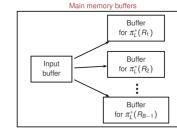
Overview of Hash-based Approach

- Build a main-memory hash table T to detect & remove duplicates



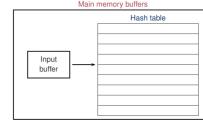
Partitioning Phase

- Use one buffer for input & $(B - 1)$ buffers for output
- Read R one page at a time into input buffer
- For each tuple t in input buffer,
 - project out unwanted attributes from t to form t'
 - apply a hash function h on t' to distribute t' into one output buffer
 - flush output buffer to disk whenever buffer is full



Duplicate Elimination Phase

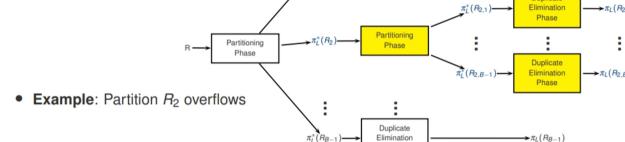
- For each partition R_i ,
 - Initialize an in-memory hash table
 - Read $\pi_L^*(R_i)$ one page at a time; for each tuple t read,
 - Hash t into bucket B_j with hash function $h' \neq h$
 - Insert t into B_j if $t \notin B_j$
 - Output tuples in hash table



Hash-based Approach Partition Overflow

- Partition Overflow Problem:** When hash table for $\pi_L^*(R)$ (Partitioned table) is larger than available memory buffers.
- Recursively apply hash-based partitioning to the overflowed partition.

Hash-based Approach: Partition Overflow



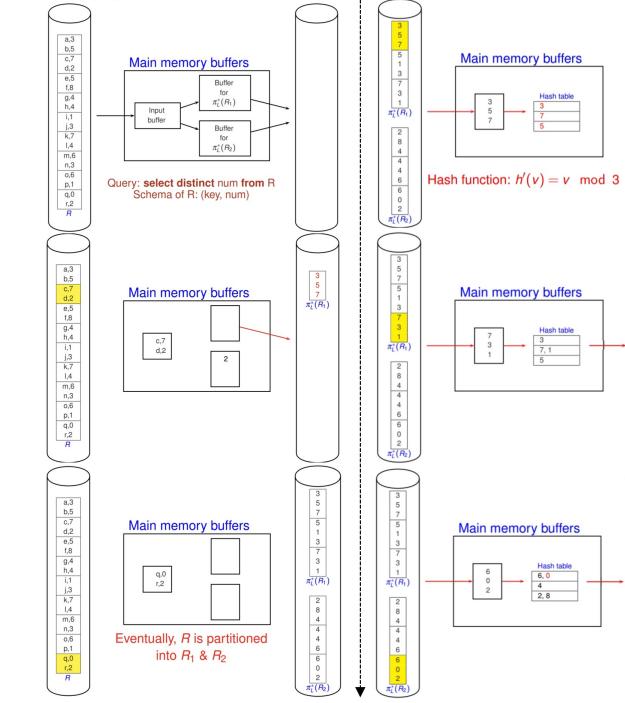
- Example: Partition R_2 overflows

Notation

Notation	Meaning
r	relational algebra expression
$ r $	number of tuples in output of r
$ r $	number of pages in output of r
b_d	number of data records that can fit on a page
b_s	number of data entries that can fit on a page
F	average fanout of B+-tree index (i.e., number of pointers to child nodes)
h	height of B+-tree index (i.e., number of levels of internal nodes)
$h = \lceil \log_F(\lceil r \rceil / b_s) \rceil$	if format-2 index on table R
B	number of available buffer pages

Illustration of Partitioning, Duplicate Elimination Phase

Partitioning Phase



5.2 Join: $R_{\bowtie_\theta} S$

- Join is where there is match between values of specified columns.
- Two-table joins, multiple-table joins, self-joins etc.
- Example: `SELECT * FROM customer c, orders o WHERE c.name = o.name;`

Database Schema:	
• Employee (eid, ename, city, did)	
• Department (did, dname, city, managerId)	
Example 1: Find (eid, managerId) pairs where managerId is manager of eid.	
<code>SELECT eid, managerId FROM Employee WHERE managerId = managerId</code>	
Example 2: Find (eid, did) pairs where eid and did co-located in same city.	
<code>SELECT eid, did FROM Employee WHERE city = Department.city</code>	

General Join Conditions

- Multiple equality-join conditions**
 - Example: (R.A = S.A) and (R.B = S.B)
 - Algorithms:
 - Index Nested Loop Join: use index on all or some of join attributes
 - Sort-Merge Join: need to sort on combination of attributes
 - Other algorithms essentially unchanged
- Inequality-join conditions**
 - Example: (R.A < S.A)
 - Algorithms:
 - Index Nested Loop Join: requires a B+-tree index
 - Sort-Merge Join: not applicable
 - Hash-based Joins: not applicable
 - Other algorithms essentially unchanged

Join Algorithms

- Iteration-Based:** Block nested loop.
- Index-Based:** Index nested loop.
- Partition-Based:** Sort-merge join, hash join.

Join Algorithms Analysis

- Factors to Consider:**
 - Type of join predicates (equality predicates e.g. $R.A_i = S.B_j$), (inequality predicates e.g. $R.A_i < S.B_j$)
 - Size of join operands
 - Available buffer space, available access methods.
- Given a join $R_{\bowtie_\theta} S$, R is **outer relation**, S is **inner relation**

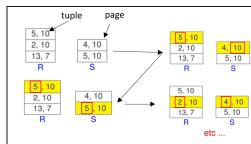
Tuple/Page-based Nested Loop Join (Naive)

- Simplest join algo is **tuple-at-a-time** nested loop evaluation. We scan outer relation R, and for each tuple in R, scan inner relation S.
- Simple refinement is **join page-at-a-time**. For each page of R, retrieve each page of S, and write out matching tuples.
- Importance of page-orientated operations for minimizing disk I/O.
- Observation:** Choose outer relation R to be smaller of two relations.

Tuple-based Nested Loop Join

Tuple-based Algorithm:

for each tuple $r \in R$ do
 for each tuple $s \in S$ do
 if (r matches s) then
 output (r, s) to result



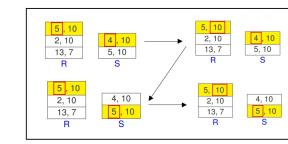
Tuple-based join scans S once for every tuple in R.

I/O Cost Analysis:
 $|R| + |R| \times |S|$
 Single line $|R|$ is number of pages in R
 Double line $|R|$ is number of tuples in R

Page-based Nested Loop Join

Page-based Algorithm:

for each page P_R of R do
 for each page P_S of S do
 for each tuple $r \in P_R$ do
 for each tuple $s \in P_S$ do
 if (r matches s) then
 output (r, s) to result



Page-based optimization of tuple-based nested loop join.
 • Scan S once for every page in R.
 • As a result, it prevents S from being scanned as many times, thus decreasing the number of disk reads.

I/O Cost Analysis:
 $|R| + |R| \times |S|$
 scan R scan S

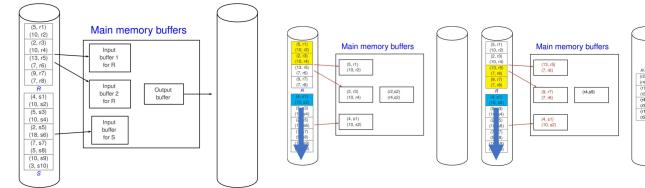
5.2.1 Block Nested Loop Join

- Simple nested loops join algo does not effectively utilize buffer pages.
- If enough memory, read whole of R (smaller relation), use one of extra buffer page to scan larger relation S. Last buffer page used as output buffer. Each relation scanned just once, optimal!
- Generalization:** Break relation R into *blocks* that can fit into available buffer pages, and scan all of S for each block of R.
- R is outer relation**, as it is scanned only once. **S is inner relation**, scanned multiple times.

Block Nested Loop Join

- Motivation:** How to better exploit buffer space to minimize number of I/Os?
- Assume $|R| \leq |S|$. So choose R as outer & S as inner
- Buffer space allocation:** Allocate one page for S, for each tuple r in R, if r matches s then output (r, s) to result

Block Nested Loop Join: Example



Algorithm (using B buffer pages):
 while (scan of R is not done) do
 read next $\lceil \frac{|R|}{B} \rceil$ pages of R into buffer
 for each page P_S of S do
 read P_S into buffer
 for each tuple r in buffer and each tuple $s \in P_S$ do
 if (r matches s) then output (r, s) to result

I/O Cost: $|R| + (\lceil \frac{|R|}{B} \rceil \times |S|)$

5.2.2 Index Nested Loop Join

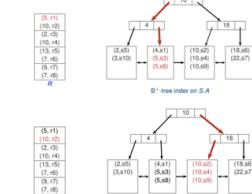
- If exists index on one of the relations on the join attribute(s), take advantage by making **indexed relation to be inner relation (S)**.
- For each tuple in R, use index to retrieve matching tuples of S.
- Cost of scanning R is M, and cost of retrieving matching S tuples depends on kind of index and number of matching tuples.

Index Nested Loop Join

Precondition: There is an index on the join attribute(s) of inner relation S.
Idea: for each tuple $r \in R$, use r to probe S's index to find matching tuples.

Analysis:

- Let $R.A_i = S.B_j$ be the join condition
- Uniform distribution assumption:
 each R-tuple joins with $\lceil \frac{|S|}{|\pi_B(S)|} \rceil$ number of S-tuples
- For a format-1 B+-tree index on S,
 - I/O Cost = $|R| + \lceil \frac{|R|}{B} \rceil \times J$
 - $J = \log_F(\lceil \frac{|S|}{Bd} \rceil)$



Cost of retrieving matching S tuple for each R tuple:

- If the index on S is a B+-tree index, the cost to find the appropriate leaf is typically 2-4 I/Os. If the index is a hash index, the cost to find the appropriate bucket is 1-2 I/Os.
- Once we find the appropriate leaf or bucket, the cost of retrieving matching S tuples depends on whether the index is clustered. If it is, the cost per outer tuple $r \in R$ is typically just one more I/O. If it is not clustered, the cost could be one I/O per matching S-tuple (since each of these could be on a different page in the worst case).

5.2.3 Sort-Merge Join

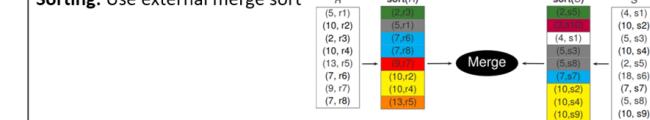
- Basic idea: Sort both relations on join attribute and look for qualifying tuples by essentially merging the two relations.
- Exploit partitioning, compare R tuples with only S tuples in same partition (rather than all S tuples), avoid enumeration of cross-product of R & S. (Works only for equality join conditions.)

Sort-Merge Join

Idea: Sort both relations based on join attributes & merge them.

Partition: Sorted relation R consists of partitions R_i of records, where records have same values for the join attribute(s)

Sorting: Use external merge sort



Merging:

- Assume R is outer relation & S is inner relation
- Each tuple in R-partition merges with all tuples in matching S-partition
- A pointer is maintained for each sorted join operand
- Each pointer is initialized to the first tuple in sorted operand
- Search for matching partitions by advancing the pointer that is pointing to a "smaller" tuple
- Need to remember position of first tuple in matching S-partition to enable rewinding of S-pointer

Example: $R: 2 \ 5 \ 7 \ 10 \ 13$ $S: 4 \ 5 \ 5 \ 10 \ 10 \ 13 \ 22$ $R \bowtie S: 5.5 \ 5.5 \ 10.10 \ 10.10 \ 10.10 \ 10.10$

Analysis: I/O cost = Cost to sort R + Cost to sort S + Merging Cost

- Cost to sort R** = $2|R|(\log_m(N_R) + 1)$ if using external merge sort
 - N_R = number of initial sorted runs of R, m = merge factor
- Cost to sort S** = $2|S|(\log_m(N_S) + 1)$ if using external merge sort
 - N_S = number of initial sorted runs of S, m = merge factor
- If each S partition is scanned at most once during merging,
 - Merging cost** = $|R| + |S|$
- Worst case occurs when each tuple of R requires scanning entire S!
 - Merging cost** = $|R| + ||R|| \times |S|$

Conventional Sort-Merge Join



Conventional Sort-Merge Join:

- Sort R:** Create sorted runs of R; Merge sorted runs of R
- Sort S:** Create sorted runs of S; Merge sorted runs of S
- Join R & S:** Merge sorted R and sorted S.

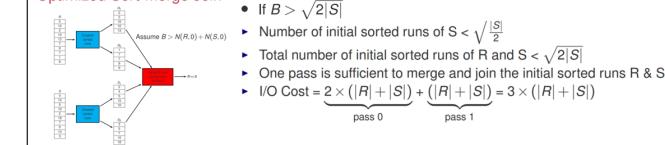
Optimization: Combine merge phase of sorting & merge phase of join.

- It's not necessary to merge sorted runs into a single run before performing join
- If $B > N(R, i) + N(S, j)$ for some i & j, sorting of R and S can stop
 - $N(R, i)$ = total number of sorted runs of R at the end of pass i of sorting R
- Sort R:** Create sorted runs of R; Merge sorted runs of R partially
- Sort S:** Create sorted runs of S; Merge sorted runs of S partially
- Join R & S:** Merge remaining sorted runs of R & S and join at same time.

Optimized Analysis:

Optimized Sort-Merge Join

- Assume $|R| \leq |S|$
- If $B > \sqrt{2|S|}$
- Number of initial sorted runs of $R < \sqrt{\frac{|S|}{2}}$
- Total number of initial sorted runs of R and S $< \sqrt{2|S|}$
- One pass is sufficient to merge and join the initial sorted runs R & S
- I/O Cost** = $2 \times (\lceil \frac{|R|}{B} \rceil + |S|) + 3 \times (\lceil \frac{|R|}{B} \rceil + |S|)$

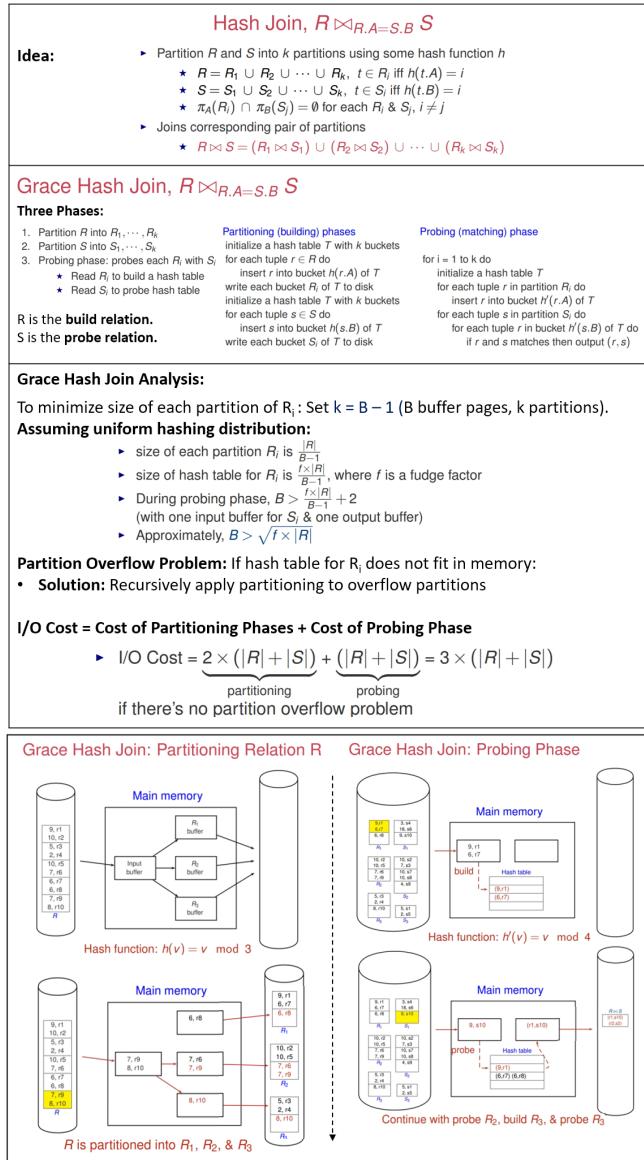


5.2.4 Hash Join

- Hash join algorithm, like sort-merge join, identifies partitions in R & S in partitioning phase, and compares tuples in corresponding partitions.
- Idea:** Unlike sort-merge join, use **hashing to identify partitions** over sorting. Hash both relations on join attribute using **same hash function h**.
- Partitioning (building) phase of hash join similar to partitioning in hash-based projection.

Grace Hash Join

- University of Tokyo, (1981), GRACE parallel relational database machine. Enables handling of large data stream quite efficiently in parallel.



6. Query Evaluation & Optimization

6.1 Set Operations

Set operations include **cross product** ($R \times S$), **intersection** ($R \cap S$), **union** ($R \cup S$), **difference** ($R - S$).

Note that **intersection & cross-product** are just special cases of join! (with equality on all fields as join condition for intersection, and with no join condition for cross-product).

Union: Main point to address is elimination of duplicates.

Difference: Variation of technique for duplicate elimination. Both can use sorting-based approach or hashing-based approach.

• Sorting approach for $R \cup S$:

- Sort R using all attributes
- Sort S using all attributes
- Merge the sorted operands to combine them and discard duplicates

• Hashing approach for $R \cup S$:

- Partition R into $\{R_1, \dots, R_k\}$ using hash function h on all attributes
- Partition S into $\{S_1, \dots, S_k\}$ using hash function h on all attributes
- For $i = 1$ to k
 - Build an in-memory hash table T_i (using hash function h') for S_i
 - Scan R_i : for each tuple $t \in R_i$, probe T_i . Discard t if it is in T_i ; otherwise, insert t into T_i
 - Write T_i to disk

- Algorithms for $R - S$ are similar to those for $R \cup S$

6.2 Aggregation Operations

Simple Aggregation

- E.g. `select count(*) from Movies`
- Aggregate Operators:** SUM, COUNT, AVG, MIN, MAX
- Algorithm:** Maintain some running information while scanning table.

Group-by Aggregation

- `select year, count(*) from Movies group by year`
- Sorting Approach:** Sort relation on grouping attribute(s), scan sorted relation to compute aggregate for each group.
- Hashing Approach:** Scan relation to build a hash table on grouping attribute(s). For each group, maintain (grouping-value, running-information).
- Using Index:** If there is covering index for query, avoid table scan. Aggregation operation can be computed from index's data entries instead of data records.
- Ordered scan using index:** For group-by aggregation, if set of group-by attributes forms a prefix of a B +tree index's search key, the data entries (and data records if necessary) for each group can be retrieved without an explicit sorting.,

6.3 Query Evaluation

6.3.1 Query Evaluation Approaches

Materialized Evaluation

- An operator is evaluated only when each of its operands has been completely evaluated or materialized.
- Intermediate results are materialized to disk

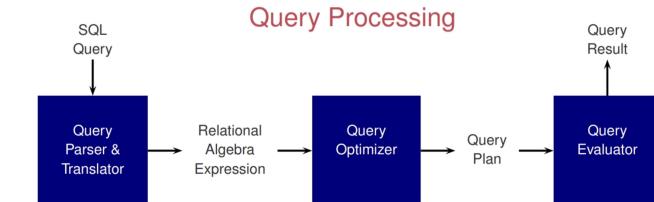
Pipelined Evaluation

- The output produced by a operator is passed directly to its parent operator.
- Execution of operators is interleaved. Top-down, demand-drive approach.
- An operator O is a **blocking operator** if O may not be able to produce any output until it has received all the input tuples from its child operator(s). Examples of blocking operator: external merge sort, sort-merge join, Grace hash join.
- Iterator Interface** for operator evaluation. "open, getNext, close" functions to initialize state of iterator, generate next output tuple and deallocate state information.
- For examples of Iterator on Hash-based Projection, Nested Loop Join, Table Scan, see slides.

Pipelined Evaluation with Partial Materialization

Also possible to have a mix of both approaches.

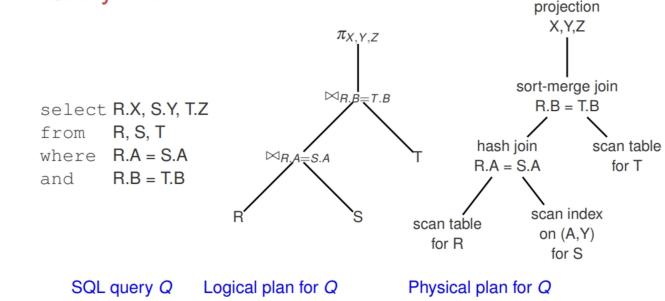
6.4 Query Processing



Query Plans

SQL Query $Q \rightarrow$ Logical Plan for $Q \rightarrow$ Physical Plan for Q .

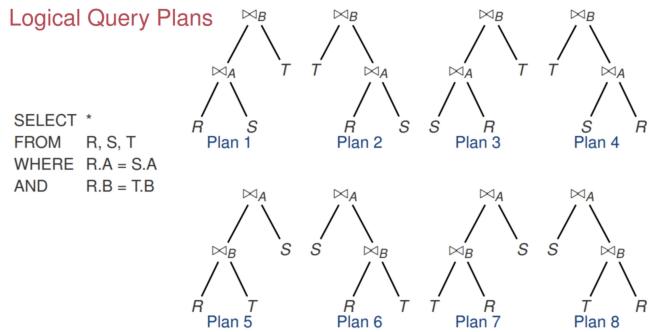
Query Plans



Logical Query Plans

- A query generally has many equivalent *logical query plans*.
- E.g. for join, the number of logical query plans grows exponentially.

Logical Query Plans

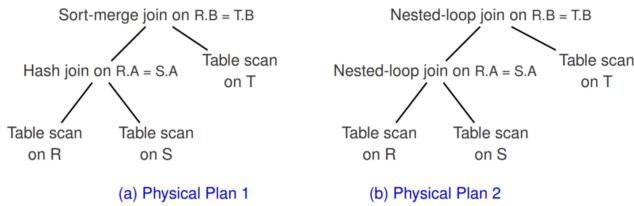


Physical Query Plans

- Subsequently, each logical plan can be implemented by many physical query plans.

Physical Query Plans

- Example:** Two possible physical query plans for $(R \bowtie_A S) \bowtie_B T$



Note: Join Plan Notation



6.5 Query Optimization

The process of finding a good query evaluation plan is called query optimization. Optimizing a relational algebra expression involves (three) basic steps:

- Search Space:** What is space (subset) of query plans being considered? (Since number of possible plans too large).
 - Plan Enumeration:** How to enumerate space of query plans.
 - Cost Model:** Estimate cost of each enumerated plan, choose plan with lowest estimated cost.
- SQL queries optimized by decomposing into blocks, optimize single block at a time. A query block can be expressed as a **relational algebra expression**.

Relational Algebra Equivalence Rules

Relational Algebra Equivalence Rules

$\text{attributes}(R)$ = Set of attributes in schema of relation R
 $\text{attributes}(p)$ = Set of attributes in predicate p

- Commutativity of binary operators**
 - $R \times S \equiv S \times R$
 - $R \bowtie S \equiv S \bowtie R$
- Associativity of binary operators**
 - $(R \times S) \times T \equiv R \times (S \times T)$
 - $(R \bowtie S) \bowtie T \equiv R \bowtie (S \bowtie T)$
- Idempotence of unary operators**
 - $\pi_L(\pi_L(R)) \equiv \pi_L(R)$
if $L \subseteq \text{attributes}(R)$
 - $\sigma_p(\sigma_p(R)) \equiv \sigma_p(R)$
if $\text{attributes}(p) \subseteq \text{attributes}(R)$
- Commutating selection with projection**
 - $\pi_L(\sigma_p(R)) \equiv \sigma_p(\pi_L(R))$
if $\text{attributes}(p) \subseteq \text{attributes}(R)$
 - $\sigma_p(R \bowtie S) \equiv \sigma_p(R) \bowtie_{p'} S$
if $\text{attributes}(p) \cap \text{attributes}(R) \subseteq L_R$ and $\text{attributes}(p) \cap \text{attributes}(S) \subseteq L_S$
 - $\sigma_p(R \cup S) \equiv \sigma_p(R) \cup \sigma_p(S)$
 - $\sigma_{p_1}(\sigma_{p_2}(R)) \equiv \sigma_{p_1 \cup p_2}(R)$
- Commuting projection with binary operators**

Let $L = L_R \cup L_S$, where $L_R \subseteq \text{attributes}(R)$ and $L_S \subseteq \text{attributes}(S)$

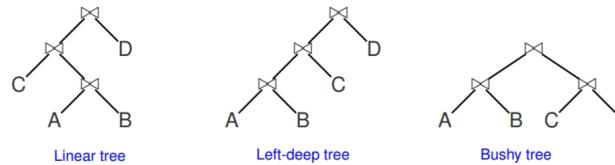
 - $\pi_L(R \times S) \equiv \pi_{L_R}(R) \times \pi_{L_S}(S)$
 - $\pi_L(R \bowtie_{p'} S) \equiv \pi_{L_R}(R) \bowtie_{p'} \pi_{L_S}(S)$
if $\text{attributes}(p) \cap \text{attributes}(R) \subseteq L_R$ and $\text{attributes}(p) \cap \text{attributes}(S) \subseteq L_S$
 - $\pi_L(R \cup S) \equiv \pi_L(R) \cup \pi_L(S)$

- Using equivalence, we can find equivalent logical query plans.
- Generally, we want to push down selection, projection operations to reduce size of table for subsequent joining.

Types of Query Plan Trees

- Linear:** Query plan is linear if at least one operand for each join operation is a base relation; otherwise, the plan is **bushy**.
- A linear query plan is **left-deep** if every right join operand is a base relation, and **right-deep** if every left join operand is a base relation.

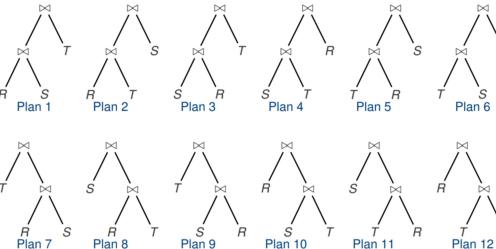
- Example:** Consider the query $A \bowtie B \bowtie C \bowtie D$



6.5.1 Search Space (Pruning)

Consider search space for simple two join operations: (12. Growth of search space is exponential!)

Query Plan Search Space for $R \bowtie S \bowtie T$



- Not all algebraically equivalent plans considered, make cost of optimization prohibitively expensive.
- Use heuristics to prune search space:** (System R Optimizer) Enumerates only left-deep query plans, avoids cross-product query plans (only join with condition), Considers early selections & projections.

6.5.2 Query Plan Enumeration

Dynamic programming (using solution of subset to solve superset problem).

Straightforward DP formulation

For each relation, consider best access plan. Then, for each subset of relations, consider best plan to join.

Dynamic programming formulation

```

Input: A SPJ query q on relations  $R_1, R_2, \dots, R_n$ 
Output: An optimal query plan for q
01. for i = 1 to n do
02.   optPlan({ $R_i$ }) = best access plan for  $R_i$ 
03. for i = 2 to n do
04.   for each  $S \subseteq \{R_1, \dots, R_n\}$ ,  $|S| = i$  do
05.     bestPlan = dummy plan with cost(bestPlan) =  $\infty$ 
06.     for each  $S_j, S_k, |S_j| \in [1, i]$ ,  $S = S_j \cup S_k$  do
07.       p = best way to join optPlan( $S_j$ ) and optPlan( $S_k$ )
08.       if ( $\text{cost}(p) \leq \text{cost}(\text{bestPlan})$ ) then
09.         bestPlan = p
10.   optPlan( $S$ ) = bestPlan
11. return optPlan({ $R_1, \dots, R_n$ })

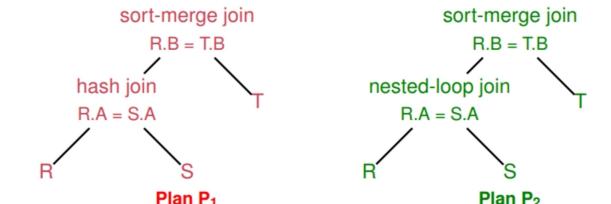
```

Enhanced dynamic programming approach (System R Optimizer)

Consider **sort order** of query plan's output. Also use heuristics to prune search space.

- Maintains $\text{optPlan}(S_i, o_i)$ instead of $\text{optPlan}(S_i)$.
- o_i captures the sort order of output produced by query plan wrt S_i , null if output is unordered or a sequence of attributes.
- $\text{optPlan}(S_i, o_i)$ = cheapest query plan for relation with output ordered by i if not null

System R Optimizer: Example



- Suppose that $\text{Cost}(\text{hash join of } R \text{ & } S) < \text{Cost}(\text{nested loop join of } R \text{ & } S)$
- In **basic dynamic programming approach**
 - $\text{optPlan}(\{R, S\}) = \text{hash join of } R \text{ and } S$
 - Therefore, plan P_2 will not be considered
- Assume that the output of nested-loop join is sorted on $R.B$
- In **enhanced dynamic programming approach**
 - $\text{optPlan}(\{R, S\}, \text{null}) = \text{hash join of } R \text{ and } S$
 - $\text{optPlan}(\{R, S\}, (R.B)) = \text{nested-loop join of } R \text{ and } S$
 - Therefore, plan P_2 will be considered during plan enumeration for $\{R, S, T\}$

6.5.3 Cost Estimation of Query Plans

Cost estimation involves the following:

- What is the **evaluation cost** of each operation? Cost model depends on: size of input operands, available buffer pages, available indexes, etc.
- What is the **output size** of each operation?

Use a **cost model** for each operator's algorithms.

Estimation Assumptions

- Uniformity assumption:** uniform distribution of attribute values
- Independence assumption:** independent distribution of values in different attributes
- Inclusion assumption:** For Join R(R.A=S.B)S, if no. tuples in A \leq B, consider $\pi_A(R) \subseteq \pi_B(S)$.

Database Statistics: Relation cardinality (no. of tuples), no. of distinct values / highest / lowest / frequent values in each column, column group statistics, histograms etc.

Size Estimation

For selection operations across joins, to estimate $\|q\|$, each predicate term potentially filters out some tuples.

- Reduction factor:** denoted $rf(t_i)$, is fraction of tuples in input relation that satisfy t_i .

- Reduction factor** also known as **Selectivity Factor**.

Size Estimation

$$\text{query } q = \sigma_p(e) \quad \|e\| = \prod_{i=1}^m \|R_i\| = \|R_1\| \times \|R_2\| \times \cdots \times \|R_m\|$$

$$p = t_1 \wedge t_2 \wedge \cdots \wedge t_n$$

$$e = R_1 \times R_2 \times \cdots \times R_m$$

$$rf(t_i) = \frac{\|\sigma_{t_i}(e)\|}{\|e\|}$$

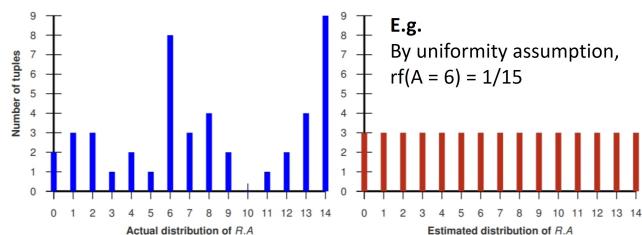
- Assuming the terms in p are statistically independent,

$$\|q\| \approx \|e\| \times \prod_{i=1}^n rf(t_i)$$

A. Selectivity Factor Estimation (within Relation)

- Simple Uniformity Assumption:** Using number of unique tuple values, find rf:

$$\text{e.g. } \|R\| = 45, \|\pi_A(R)\| = 15, rf(A = c) \approx \frac{1}{\|\pi_A(R)\|} = \frac{1}{15}$$



Size Estimation: Example

- Consider a relation $R(A, B, C)$
 - output size of query $Q: \sigma_{B=10 \wedge C=23}(R)$
 - By uniformity assumption,
 - $rf(B=10) \approx \frac{1}{10}$
 - $rf(C=23) \approx \frac{1}{50}$
 - By independence assumption,
 - $rf((B=10) \wedge (C=23)) \approx \frac{1}{10} \times \frac{1}{50}$
 - $\|Q\| \approx 1000 \times \frac{1}{10} \times \frac{1}{50} = 2$

B. Join Selectivity (inter-Relation selectivity)

- Join selectivity factor:** Selectivity factor for join predicates
- Inclusion assumption:** Assume every (smaller) relation tuple joins with some (larger) relation tuple.

Join Selectivity

Consider query Q: $\text{SELECT * FROM } R \text{ JOIN } S \text{ ON } R.A = S.B$

$$rf(R.A = S.B) = \frac{\|R \bowtie_{R.A=S.B} S\|}{\|R\| \times \|S\|}$$

Inclusion assumption: Consider $R \bowtie_{R.A=S.B} S$

$$\text{If } \|\pi_A(R)\| \leq \|\pi_B(S)\|, \text{ then } \pi_A(R) \subseteq \pi_B(S)$$

Join selectivity estimation:

- Assume $\|\pi_A(R)\| \leq \|\pi_B(S)\|$
- By inclusion assumption, every R-tuple joins with some S-tuple
- By uniformity assumption, there are $\frac{\|S\|}{\|\pi_B(S)\|}$ S-tuples corresponding to each S.B value
- Therefore, each R-tuple joins with $\frac{\|S\|}{\|\pi_B(S)\|}$ S-tuples
- Thus, $\|Q\| \approx \|R\| \times \frac{\|S\|}{\|\pi_B(S)\|}$

$$rf(R.A = S.B) \approx \frac{1}{\max\{\|\pi_A(R)\|, \|\pi_B(S)\|\}}$$

Join Selectivity: Example

Consider query Q: $R \bowtie_{dept} S$

S		
name	dept	course
Alice	CS	CS101
Bob	CS	CS111
Carol	CS	CS302
Dave	CS	MA105
Eve	CS	MA203
Fred	CS	MU108
George	EE	PH113
Henry	EE	PH203
Ivy	EE	
Jane	EE	

$$\begin{aligned} \|R\| &= 10, \|\pi_{dept}(R)\| = 2 \\ \|S\| &= 8, \|\pi_{dept}(S)\| = 4 \\ rf(R.dept = S.dept) &\approx \frac{1}{\|\pi_{dept}(S)\|} = \frac{1}{4} \\ \|Q\| &\approx \|R\| \times \frac{\|S\|}{\|\pi_{dept}(S)\|} = 10 \times \frac{8}{4} = 20 \end{aligned}$$

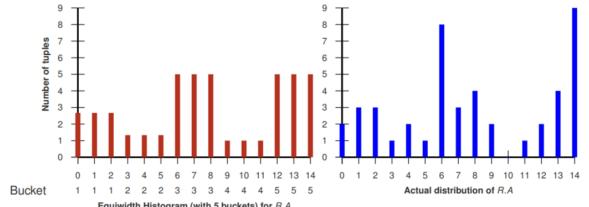
Estimation using Histograms

- Histogram:** stat. info maintained by DBMS to estimate data distribution.
- Main idea using histogram:** Partition attribute's domain into sub-ranges called buckets, assume value distribution within each bucket is uniform.
- Equiwidth Histograms:** Each bucket has (almost) equal number of values
- Equidepth Histograms:** Each bucket has (almost) equal number of tuples, sub-ranges of adjacent buckets might overlap.
- MCV: Most Common Values:** Separately keep track of the frequencies of the top-k most common values and exclude MCV from histogram's buckets.
- Accuracy of histograms:** Equiwidth < Equidepth < w. MCV excluded.

Type of Histograms & Estimation Accuracy

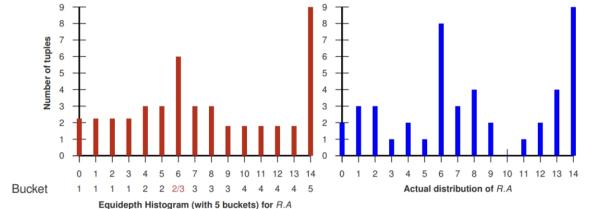
Equiwidth Histograms (with 5 buckets)

Bucket No	Value Range	No. of Tuples	Estimated Number of Tuples per Bucket Value
1	[0, 2]	2+3=5	8/3 = 2.67
2	[3, 5]	1+2=3	4/3 = 1.33
3	[6, 8]	8+3=11	15/3 = 5
4	[9, 11]	2+0=2	3/3 = 1
5	[12, 14]	2+4=6	15/3 = 5



Equidepth Histograms (with 5 buckets)

Bucket No	Value Range	No. of Tuples
1	[0, 3]	2+3=5
2	[4, 6]	2+1=3
3	[6, 8]	2+3=5
4	[9, 13]	2+0=2
5	[14, 14]	2+4=6

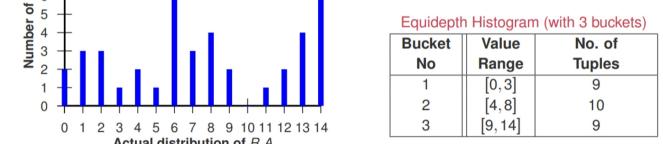


Estimation with Histograms: Example

- Query $Q_1 : \sigma_{A=6}(R), \|Q_1\| = 8$**
 - Without histogram: $\|Q_1\| \approx 45/15 = 3$
 - Equiwidth histogram: $\|Q_1\| \approx 15/3 = 5$
 - Equidepth histogram: $\|Q_1\| \approx (1/3 \times 9) + (1/3 \times 9) = 6$
- Query $Q_2 : \sigma_{A \in [7,12]}(R), \|Q_2\| = 12$**
 - Without histogram: $\|Q_2\| \approx 45/15 \times 6 = 18$
 - Equiwidth histogram: $\|Q_2\| \approx (2/3 \times 15) + 3 + (1/3 \times 15) = 18$
 - Equidepth histogram: $\|Q_2\| \approx (2/3 \times 9) + (4/5 \times 9) = 13.2$

Improved Histogram Estimation with MCV

MCV (k=2)	Value	No. of Tuples
6	8	9
14	9	14



- Query $Q_1 : \sigma_{A=6}(R), \|Q_1\| = 8$**
 - Equidepth histogram: $\|Q_1\| \approx (1/3 \times 9) + (1/3 \times 9) = 6$
 - Equidepth histogram with MCV: $\|Q_1\| \approx 8$
- Query $Q_2 : \sigma_{A \in [7,12]}(R), \|Q_2\| = 12$**
 - Equidepth histogram: $\|Q_2\| \approx (2/3 \times 9) + (4/5 \times 9) = 13.2$
 - Equidepth histogram with MCV: $\|Q_2\| \approx (2/5 \times 10) + (4/6 \times 9) = 12.2$

7. Transaction Management

7.1 Transactions

A transaction is an abstraction representing a logical unit of work. Ensure 4 properties of Xacts / Txn. (**ACID**).

Concurrency control manager component ensures isolate, recovery
manager component ensures durability.

- **Atomicity:** Either all or none of the actions in Xact happen.
- **Consistency:** If each Xact is consistent, DB starts consistent, ends up consistent.
- **Isolation:** Execution of one Xact is isolated from other Xacts.
- **Durability:** If a Xact commits, its effects persist.

Notation

- **Transaction Schedule:** a list of actions from a set of Xacts, where the order of the actions within each Xact is preserved.
- **Serial Schedule:** Actions of Xacts are not interleaved.
- **Read From:** T_j reads object O from T_i in schedule if last write action on O is by I , and J reads after.
- **Final Write:** T_i performs final write on O in a schedule S if it does the last write action on O .
- **Correctness of Interleaved Xact Executions:** Correct if “equivalent” to some serial schedule over the same set of Xacts.

View Equivalent / Serializable

View Equivalent: Two schedules over same set of Xacts are view equivalent if they have **same read froms** and the **same final write**.

View Serializable: A schedule is a **view serializable schedule (VSS)** if it is view equivalent to some serial schedule over the same set of Xacts.

Testing for View Serializability: Given a schedule S , construct a directed graph (denoted by $VSG(S)$) to capture the read-from and final-write relations among the transactions in S . If $VSG(S)$ is cyclic, then S is not VSS. If acyclic, it is VSS iff there exists serial schedule produced from topological ordering of $VSG(S)$ that is view equivalent to S .

Conflicting Actions

Conflict: Two actions on same object conflict if at least 1 write, different Xacts. (E.g. W + R, W + W, R + W).

Anomalies with Interleaved Xact Executions

- **Dirty Read Problem (WR Conflict):** T_2 reads object produced by uncommitted Xact (which then aborts).
- **Unrepeatable Read Problem (RW Conflict):** T_2 updates object, reading yields different values before and after while T_1 in progress.
- **Lost Update Problem (WW Conflict):** T_2 overwrites value of object while T_1 still in progress.
- **Phantom Read:** Re-executing query on search condition gives different results (bc new added Object) (prevent by predicate / index locking).

Conflict Equivalent / Serializable

Conflict Equivalent: Two schedules over same set of Xacts are *conflict equivalent* if they order every pair of conflicting actions of two committed Xacts in the same way.

Conflict Serializable: A schedule is a *conflict serializable schedule (CSS)* if it is conflict equivalent to a serial schedule over same set of Xacts.

Testing for conflict serializability

1. A schedule is conflict serializable iff its conflict serializability graph is acyclic. ($CSG(S)$: contains node for each committed Xact in S , each edge is a conflict).
2. A schedule that is conflict serializable \rightarrow view serializable.
3. If S is view serializable and S has no blind writes, then S is also conflict serializable.

Blind write: Did not read before write.

7.2 Recovery

- **Cascading Abort:** Recursive aborting process: If T_1 reads from T_2 , former must abort when latter aborts for correctness.
- **Recoverable:** If T_2 reads from T_1 , T_2 commits only after T_1 . Guarantees committed Xacts will not be aborted.
- **Cascadeless:** Whenever T_i reads from T_j in S , $Commit_j$ must precede this read action. Only can read from committed Xact.
- Cascadeless schedule is recoverable.
- **Before-Images:** (for recovery) log before action & restore (must be strict).
- **Strict:** for every $W_i(O)$ in S , O is not read/written by another Xact until T_i aborts or commits.
- Strict schedule is cascadeless.

8. Concurrency Control