

CS3223 Database Sys Implementation

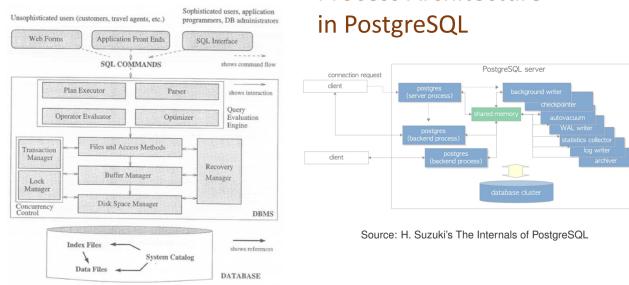
AY23/24 Sem 2, github.com/gertek

Introduction

Course Details

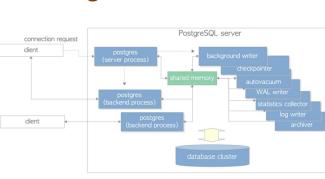
- Prerequisite knowledges: CS2040S, CS2102, CS2106 background (helpful).
- Reference Textbook: Raghu & Johannes Database M. Systems, 2002. Encouraged to read ahead based on schedule before the lecture.
- Course covers data structures, algorithms, different components making up database systems.

Architecture of DBMS



Source: Ramakrishnan & Gehrke's Figure 1.3

Process Architecture in PostgreSQL



Source: H. Suzuki's The Internals of PostgreSQL

- 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.

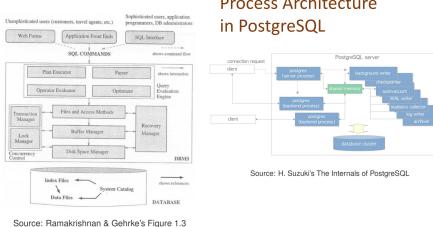
DBMS Storage

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.

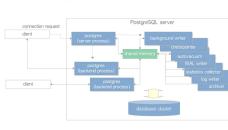
Magnetic Hard-Disk Drive HDD

Architecture of DBMS



Source: Ramakrishnan & Gehrke's Figure 1.3

Process Architecture in PostgreSQL



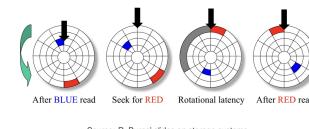
Source: H. Suzuki's The Internals of PostgreSQL

- 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



Source: R. Burns' slides on storage systems

- Seek time
 - avg. seek time: 5-6 ms
- Rotational delay (or rotational latency)
 - Depends on rotation speed - measured in rotations per minute (RPM)
 - Average rotational delay = time for $\frac{1}{2}$ revolution
 - Example: For 10000 RPM, avg. rotational delay = $0.5 \times (60 / 10000) = 3$ ms
- Transfer time
 - $n \times$ number of requested sectors on same track
 - transfer time = $n \times \frac{\text{time for one revolution}}{\text{number of sectors per track}}$
 - avg. sector transfer time: 100-200 µs
 - Sequential vs random I/O

- 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.

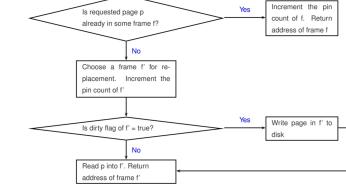
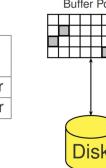
Solid-State Drive: SSD

- 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 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: Example**
 - N = number of frames in buffer pool
 - Diagram shows frames 0-4. Frame 0: pin count = 0, referenced bit = on. Frame 1: pin count = 0, referenced bit = off. Frame 2: pin count = 0, referenced bit = off. Frame 3: pin count = 0, referenced bit = off. Frame 4: pin count = 0, referenced bit = off.
 - Request 1: Look for available frame. Frame 0 is chosen because it has a referenced bit = on.
 - Request 2: Look for available frame. Frame 1 is chosen because it has a referenced bit = off and pin count = 0.
 - Request 3: Look for available frame. Frame 2 is chosen because it has a referenced bit = off and pin count = 0.
 - Request 4: Look for available frame. Frame 3 is chosen because it has a referenced bit = off and pin count = 0.
 - Request 5: Look for available frame. Frame 4 is chosen because it has a referenced bit = off and pin count = 0.

CS3223: Sem 2, 2023/24

Files

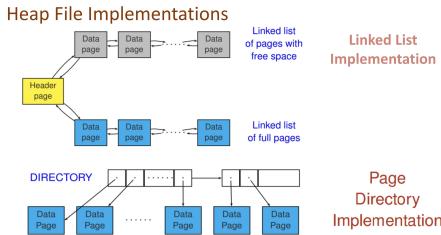
File Abstraction

- 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:

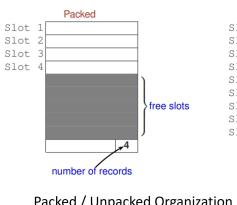
Records are organized within a page and referenced with the RID.

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

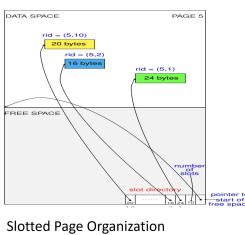
- **Packed Organization:** Store records in contiguous slots.
- For packed organization, memory organization is tough and costly when record in slot is deleted, need to move up a record. But as RID serves as a reference, but need to propagate change in RID.
- **Unpacked Organization:** Uses bit array 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**.

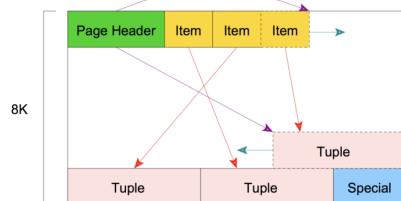
Fixed-Length Records:



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

O ₁	O ₂	O ₃	O ₄	F ₁	F ₂	F ₃	F ₄
----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------

Each O_i is an offset to beginning of field F_i

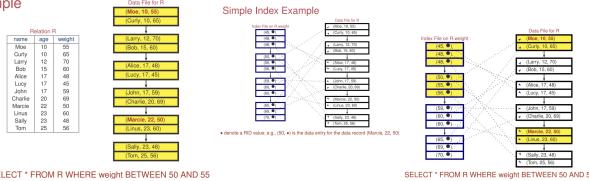
2. Indexing

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



w/o index, need to retrieve all pages/records

w. simple index, only require retrieval of select pages

Index Types

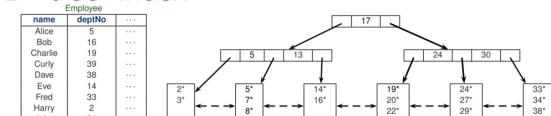
Two main types of indexes

- **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

B⁺-tree Index

B⁺-tree Index

Employee	
Alice	5
Bob	16
Charlie	19
Carl	39
Dave	38
Eve	14
Fred	33
Harry	24
John	34
Kate	8
Larry	27
Linus	3
Lily	22
Marcie	29
Moe	20
Sally	7
Tom	7



B⁺-tree index on Employee.deptNo

- 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 $(p: \text{pointer}, k: \text{separator})$ ($p_0, k_1, p_1, k_2, p_2, \dots, p_n$)
- $k_1 < k_2 < \dots < K_n$
- 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)

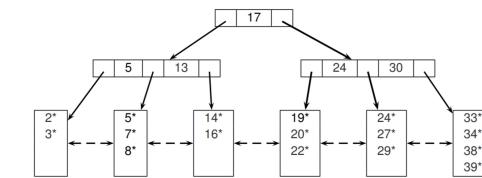
Properties of B⁺-tree Index

- Dynamic index structure; adapts to data updates gracefully
- Height-balanced index structure

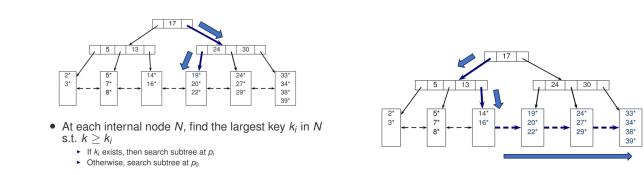
Order of index tree, $d \in \mathbb{Z}^+$

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

Example: B⁺-tree with order = 2



Equality Search ($k = 19$)



Range Search ($k \geq 15$)

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

Relation R			
name	age	weight	height
Moe	10	55	180
Curly	10	65	171
Larry	12	70	175
Bob	15	60	178
Alice	17	49	175
Lucy	17	45	170
John	18	59	182
Charlie	20	69	173
Marie	22	50	165
Linus	23	60	166
Sally	24	48	169
Tom	25	56	176

B⁺-tree index on R.age (Format 1)

```

graph TD
    Root[18] --> Node1[12]
    Root --> Node2[17]
    Root --> Node3[20]
    Node1 --> Leaf1["Moe(10, 55, 180)"]
    Node1 --> Leaf2["Curly(10, 65, 171)"]
    Node2 --> Leaf3["Larry(12, 70, 175)"]
    Node3 --> Leaf4["Bob(15, 60, 178)"]
    Node3 --> Leaf5["Alice(17, 49, 175)"]
    Node3 --> Leaf6["Lucy(17, 45, 170)"]
    Node3 --> Leaf7["John(18, 59, 182)"]
  
```

B⁺-tree index on R.age (Format 2)

```

graph TD
    Root[18] --> Node1[15]
    Root --> Node2[18]
    Root --> Node3[23]
    Node1 --> Leaf1["(10, Moe)"]
    Node1 --> Leaf2["(10, Curly)"]
    Node2 --> Leaf3["(12, Larry)"]
    Node3 --> Leaf4["(15, Bob)"]
    Node3 --> Leaf5["(17, Alice)"]
    Node3 --> Leaf6["(17, Lucy)"]
    Node3 --> Leaf7["(18, John)"]
  
```

Index on R.age (format 3)

```

graph TD
    Root[18] --> Node1[15]
    Root --> Node2[18]
    Root --> Node3[23]
    Node1 --> Leaf1["(10, [Moe, Curly])"]
    Node1 --> Leaf2["(12, [Larry])"]
    Node2 --> Leaf3["(15, [Bob])"]
    Node3 --> Leaf4["(18, [Alice, Lucy, John])"]
    Node3 --> Leaf5["(23, [Bob, Alice, Lucy, John, Moe, Curly])"]
  
```

B⁺-Tree Insertion

B⁺-Tree Deletion

B⁺-Tree Bulk Loading