Winter 2025 – Professor Ryan Rosario

## **Lecture 3/4: Relational Algebra**

#### Notes

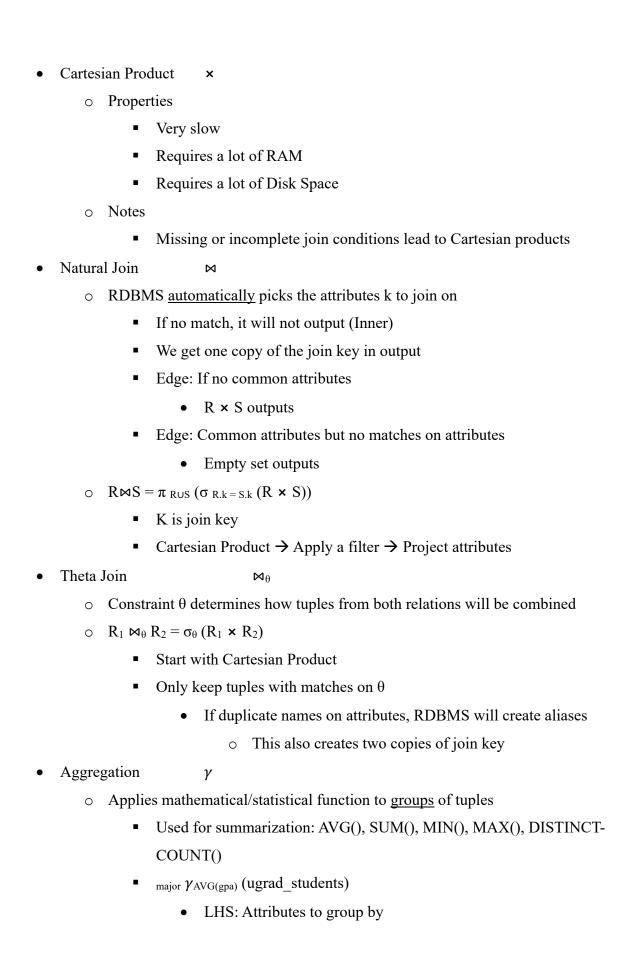
- Relational algebra works on unique sets of tuples
  - Using relational algebra operation removes duplicate tuples

## **Relational Operators**

• Selection

- σ
- Retrieves subset of tuples from a <u>single</u> relation that satisfies a particular constraint
  - $\sigma(R) \rightarrow$  Gets all tuples in the relation and all attributes
    - SQL WHERE
  - $\sigma_{\psi} \rightarrow$  Filters <u>tuples</u> (not attributes)
    - $\sigma_{\psi}(R) \rightarrow SELECT * FROM R WHERE \psi$ 
      - O  $\sigma$  likes > dislikes  $\wedge$  views > 1000000
      - ψ is a filtering predicate
- Projection

- π
- o Extracts attributes from a set of tuples and removes duplicates
  - Normally the last thing in a relational algebra expression
    - $\pi_{\text{title}}$  ( $\sigma_{\text{likes}} > \text{dislikes} \land \text{views} > 1000000$ )
- o Can create new attributes
  - $\blacksquare$   $\pi$  likes/(likes + dislikes)
  - $\pi$  likes/(likes + dislikes)  $\rightarrow$  favor
  - Does not create a new attribute within the tuple, but only the output
- o Can rename attributes using →
  - $\pi$  likes  $\rightarrow$  thumbs up
- Closest to SELECT DISTINCT



- RHS: Functions to apply
- Global Aggregate: γ<sub>AVG(gpa)→average</sub>(R)
- majory → SELECT DISTINCT major
- Closest to GROUP BY
- Rename

 $\circ$   $\rho_{\text{new}}(R)$  : renames relation R to S

 $\circ$   $\rho_{\text{new/old}}(R)$  : renames attribute old in relation R to new

- Set Union U
- Set Difference —
- Set Intersection ∩
  - o Division is not implemented in SQL

#### Join Notations

- Equijoin: Joins of equality
  - $\circ$  A == B
- Nonequijoin: Joins of not equality
  - $\circ$  A > 3
- <u>Inner</u>: If match, output match. If no match, do not output
- Outer: Potential to Output NULLS
  - <u>Left</u>: Always keep the tuples from the LHS. If match on RHS, output RHS. If no match, output special tuple/NULL.
  - Right: Always keep the tuples from the RHS. If match on LHS, output LHS. If no match, output special tuple/NULL.

## **Lecture 9: Functional Dependencies, Normal Forms**

## **Functional Dependencies**

• A, B  $\rightarrow$  C : Allowed

•  $A \rightarrow B$ , C: Not Allowed

#### 2NF Check

- Create Functional Dependencies
- Find Canonical Cover
  - o Simplify to Singletons
  - o Remove Composites comprised of Singletons
  - Remove Inferred Composites

#### Normal Form

- 1NF
  - o All attributes are flat (No arrays, maps)
  - o No columns logically reliant on each other (ex: Code1, Code2, Code3)
  - Unique Key exists (No duplicate tuples)
  - No null values
  - o Row ordering does not matter
- 2NF
  - o R is in 1NF
  - o Attribute appears in a CK
  - o Attributes depends on the **entire** key
    - If AB is CK, AB $\rightarrow$ C, and C $\rightarrow$ D, this is ok
  - o Not partially dependent on **any** composite candidate key
    - If AB is CK and B $\rightarrow$ E, this is not 2NF (partial dependency)
      - Decompose into 2NF: (B, E), (B, everything else)
- Codd 3NF
  - o R is in 2NF
  - o Non-key attributes depend only on an entire CK

- If all attributes are a part of a CK, R is in 3NF by definition
- o Non-key attributes must depend directly to an entire CK
  - If AB is CK, AB $\rightarrow$ C, and C $\rightarrow$ D, this is not 3NF
- Zaniolo 3NF (1 of following must be true)
  - $\circ$  A $\rightarrow$ B is trivial (LHS, RHS share an attribute)
  - o LHS is superkey for R
  - o RHS is contained in a CK
    - $A \rightarrow AB : B \text{ must be a part of a CK}$
    - If CK is AB,  $C \rightarrow B$  is OK as B is paritally in CK
- BCNF
  - o For every functional dependency:
    - A→B is trivial (LHS, RHS share an attribute)
    - LHS is superkey for R

Losslessness: Natural Join of subrelations returns original relation

- $R1 \bowtie R2 = R$
- Decomposition Tests
  - $\circ$  R1 U R2 = R
  - o R1  $\cap$  R2  $\neq$  {} (There exists a join key)
  - o R1  $\cap$  R2 is SK for R1, R2, or both
    - Closures: Determine if candidate is SK for R1 or R2
      - Attempt to derive all other attributes from SK, recheck for any changes in closure
      - Denoted with a + (ex: A+)
        - o If A+ == A, A is a superkey for R

## **Dependency Preservation**

- For each FD, check if you can use:
  - (Closure  $\cap$  R)+  $\cap$  R to derive new attributes
    - If no, then not Dependency Preserving

## **Lecture 10: OLAP, OLTP, Warehousing**

#### OLTP (Online Transaction Processing)

- Optimized for live traffic: random reads and writes, low latency, no patterned data
  - Databases optimize OLTP
- Random Access
- Normalization

## OLAP (Online Analytical Processing)

- Read-only, optimized to be fast with low latency
- Used for internal systems. Other systems write to system
  - o Incremental writes are very slow
- <u>Columnar/cubic/hypercube</u> data storage
  - o Length, height, depth

#### OLAP Cube Vocab

- Slice: 1D query, Removes one dimension (gender=M), WHERE
- Dice: 2D query, (gender=M, location=US), WHERE AND
  - Outputs a specific piece of data
- Rollup: Hierarchical attribute, aka ranges (space and time), Up the hierarchy
- Drill Down: Take a column from a slice, aka take a more precise data
  - Outputs a range of data, used in hierarchical
- Pivot: Long data to wide data
  - Wide: (One row per student)
    - Less redundancy on student data
    - Easy to understand
    - Can have NULLs
    - Not as flexible (you need to ALTER TABLE to add a column)
  - o Long: (Multiple rows per student)
    - More flexible
    - No NULLs

- Redundancy
- Harder to understand

## Data Warehouses

- Designed for analytics, not production
- Optimized for fast reads
- Optimized for OLAP, batch processing
  - Batch Processing: Low latency reads (SELECT) of aggregated or precomputed data
    - ETL Job: Data Extract, Transform, Load (into warehouse)
      - Atomic
      - Some data is lost if insert/update while processing
        - o Run ETLs on production database copy
        - o Dual Writes: Updates changes replica
      - App→ProductionDB→Replica→DW
      - App → ProductionDB, Replica →DW
- When data loads, groups are preaggregated
  - o By guideline (priori) decided by data engineer
  - o Automatically (2<sup>n</sup> dimensions)
- Stores normalized, denormalized/exploded tables for efficient reads
- Data is a little out of date

#### Data Warehouse Schemas

- Fact Table: Denormalized, contains the data to be analyzed
- Dimension Table: Data about attributes of each of facts
- Star Schema: Single fact table surrounded by several dimension tables
- <u>Snowflake Schema</u>: Central fact table where dimension tables may be joined to each other
- Constellation/Galaxy Schema: Multiple fact tables share dimension tables

# Lecture 11: DuckDB, Distributed Systems, MapReduce

## DuckDB

- Features
  - o Pattern matching with columns
  - o Pivot and Unpivot

<u>Data Lake</u>: Data storage where data is stored in different storage systems but can be queried from a single system

- Less preprocessing work, but more query minutia
  - o No schema, generally raw data
- <u>Hadoop</u>: Distributed, scalable file system, runs on Map-Reduce
  - o Better for disk operations

	RDBMS/Database	Data Warehouse	Data Lake
Data	Structured	Structured	Unstructured
	Tabular	Columnar/cube	
	Clean	Fairly Clean	
Schema	Highly normalized	Not normalized	No schema
	Fixed Schema	Fixed Schema	
Performance	Fast for Transactions	Fast for SELECTS	Mixed
(Price may vary for all)	Fast for Reads, Writes		
Data Quality	Very high	Very high, but comes	Messy
		from mixed sources	
Users	SWEs, SRE, DS	Analysts, DS, BA,	Anyone
		SWE	
Analytics	Rowwise for	Denormalized	Anything
	Transactions	Prejoined data	
	Joins are painful	Preaggregates	

Distributed Systems

- Architectures
  - o Replication: Decentralized, all nodes expected to have same data
    - All writes populate to all other nodes
  - o Sharding: Centralized, different data on each machine
    - Typically, a hybrid model is used
- File Systems
  - o Distributed
    - Distributed File Systems are logical
      - Systems are clusters, not individual
      - Files are divided into blocks
    - Loss of relational processing
  - Traditional
    - Storage devices are divided into partitions
    - Partitions are formatted with <u>file systems</u>
      - File systems are physical

Big Data: Datasets too big for RAM where parallel computing is required to handle data

• Maximizes usage of compute resources, increases amount of work per unit time, increase throughput (not speed)

MapReduce: Data Processing Platform

- Ideas
  - MapReduce functions on each split independently
  - o Operates under Shared-Nothing Model → Extreme parallelism
- Map (Process and Transform)
  - Map: Outputs key-value pairs
    - Represented by relational algebra select  $\sigma$
  - Combine (Preaggregate, optional)
    - Not used when there is not enough RAM
    - Allows for fewer disk writes, fewer network transfers
    - Combiner is normally the same as the reducer function

- o Partition: Merge results from split
  - Merge based on record
- Shuffle
  - o Partitions send data to reducers using hash(key) % n
    - n is number of reducers
  - o Reducers receive all key-value pairs with same hash
  - Represented by relational algebra aggregation Υ, otherwise it is represented by
    Map and/or Reduce phase
- Reduce (Aggregate)
  - Mergesort
    - Merge based on split
  - o Reduce (Aggregation): Call reducer function on output to aggregate
    - Each reduce task outputs one output file containing all values associated with processed keys
      - If no reducers: Number of files = Number of map tasks = Number of input splits
    - Aggregation function represented by relational algebra reduce
    - Key Skew: Reducing of outlier data takes significant more time than a regular case
      - Solution: Custom partitioner
        - Hash uniformly by appending a salt (such that the same word maps to different buckets)
          - Outputs a partial count
        - o Remove uniform number
        - o Aggregate (and do a 2<sup>nd</sup> Map Reduce Job)

## Lecture 12: MapReduce, Spark, Streaming Systems

Pipeline: Processes input via series of algebraic steps

**Spark** (Resilient Distributed Datasets)

- Operations completed in RAM
  - o Changes only persist after commits
- Lazy Evaluation
  - More efficient for iterative problems
  - o head(), first(), take(), count(), collect(), collectAsList(), reduce(), show()
- Joins have poor performance
- File System is treated as truth

## RDDs (Resilient Distributed Datasets)

- Processes datasets rather than files
- Immutable: You cannot modify relations → Increases performance
- Resilient: If one node fails, another replicate exists somewhere else in system
- Distributed: parts of data are in different nodes of system

## **Streaming Systems:**

- Incomplete, Unbounded, continuous data
  - o Streaming System is ground truth
  - o Process data as it arrives without storing it
    - Lambda architecture: Streaming layer processes unstoreable data,
      Complete layer stores batched data
  - o Events: Individual records
  - o Messages: Data about event
- Streaming Models
  - o Producer/Consumer
    - Producer: Creates/sends data to send into system
      - INSERT
    - Consumer: Reads data and processes it

- SELECT → This is extremely inefficient
- o Push Model
  - Bipartite Graph
    - Direct connection between producer, consumer
    - Issues
      - o Not Flexible: producers, consumers are hardcoded
      - If consumer goes down, entire system needs to be reconfigured
      - o Producers are independent → Decentralized
- Data Loss
  - O Use Cases:
    - 1. Replicate a write to the Database
    - 2. Computing Aggregates
  - o Cases:
    - Is data loss a big deal in the streaming system?
      - 1. Yes
      - 2. No, assuming the messages are missing at random
    - Is receiving duplicate messages a big deal?
      - 1. Not really as long as a Primary Key exists
      - 2. Usually no, duplicates may wash out
    - Is receiving data in the wrong order a big deal?
      - 1. No. Relation is not ordered!
      - 2. No. Relation is aggregated!

## Lecture 13: Stream Processing, Probabilistic Data Structures

## Streaming Use Cases

- 1. Replicating writes to a backup database
  - Database handles everything
    - Performant but less flexible
  - o Separate multiplexer sends requests to both production and replica databases
  - Streaming Systems
- 2. Processing aggregates from a stream
- 3. Making real-time decisions

## **Streaming Test Cases**

- Can we tolerate data loss?
  - o 1. No! We need to maintain data integrity
  - o 2. Yes, as long as messages are missing at random
    - If not at random, bias can occur over long periods of time in a global aggregate
  - o 3. Yes, but only when an event consists of many messages at a high velocity
    - No if we receive the message once every N seconds
- Can we tolerate message delay?
  - o 1. To an extent yes
    - Data is complete, but not up-to-date
      - Client must adapt, but is ok for data warehouses
  - o 2. Yes. Data will still be aggregated unless it arrives extremely late
  - o 3. Yes. As long as it's not extremely late
- Can we tolerate out of order messages?
  - o 1. Yes! There's no order in relations anyway, and we can add a time stamp
  - o 2. Yes except for rolling/cumulative statistics or rolling windows
  - o 3. Yes
    - No if we are dealing with a sequential/complex event

Message Broker: Lives in message bus, communicates with producers and consumers

- Maintains state of streaming system
  - o Messages stored in broker buffer until consumer ACKs it
  - o However, client and broker buffers may become full, losing packets
    - Delete oldest message in queue
    - Delete random message in queue
- Determining which client gets a message
  - Load balancing: deliver message randomly or based on shard key
  - o Fan out: deliver message to all clients of a group
  - Hybrid is commonly used

# Computing aggregates in streams

- Compute over window w
  - Issue: Stragglers (events that arrive in window but are not included in computation)
    - Ignore stragglers
    - Open window w for a little longer
      - We may store more data in buffer than needed
    - Update aggregation afterwards
      - Aggregation function needs to be updatable
        - o Easy: count(), min(), max(), avg(), sd()
        - Hard: mode(), median()
- Compute over all time
- Update Rule (Can use with or without windows)
  - Compute aggregates across all messages of a given window, even if all messages are evicted from the window

## Bloom Filter (Filters Duplicate messages)

- Returns false positives with small probability
  - o Bit string of length m associated with k hash functions

HyperLogLog (outputs amount of distinct values in a set)

- Treat bitstring as a binary number, provide an estimate of number of distinct objects
  - o Log2log2(2^32) can be stored in 5 bits
    - Average error of 2%

## Lecture 14: CAP, NoSQL (MongoDB, Neo4j)

#### **NoSQL**

- "Not only SQL" → DML or DDL is something other than SQL
- "No SQL" → No use of relational model

#### CAP Theorem

- It is impossible for a <u>distributed</u> system to satisfy all 3 of the following:
  - o Consistency: Every read receives the *most recent* write, or an error is thrown
  - Availability: Every request receives some kind of response (not an error), without guaranteeing the most recent write. *It always returns a response*.
  - Partition Tolerance: The system continues to function even if messages are lost or are delayed by the network between or among nodes
- MongoDB <u>supports</u> CP/AP

## NoSQL Data Models

- Key-value store (Redis, Memcached)
  - o Redis is single node, so CAP does not apply  $\rightarrow$  AP
  - o Redis commands take a key and a value, or just a key
  - Everything is stored in RAM (for indexing)
- Columnar (Cassandra, HBase, DynamoDB)
  - o Fixed number of columns but each row has a subset of columns
    - Key-value pair where the value is a collection of key-value pairs
  - o Cassandra
    - Strictly ordered, not relational
- Document store (MongoDB, CouchDB)
  - o Documents are schemaless → Stored in JSON/BSON
    - No joins
    - Documents ⊆ collections ⊆ databases
  - o CP
- Consistency from allowing writes on only primary replica

- Availability from using secondary replica sets, one or more with the exact same data
- Horizontal Sharding → Data is distributed across nodes of network based on shard key
- ODMs (like Mongoose) support CRUD
  - Create, Read, Update, Delete
- o Indexed by B-Tree → Simple and compound fields
- Graph (Neo4j, OrientDB, Giraph)
  - Vocabulary
    - <u>Nodes</u> represent entities
    - <u>Edges</u> represent foreign keys
      - Nodes and edges have attributes called <u>properties</u> and also <u>labels</u>
  - o Neo4j does not have CAP → Graphs cannot be distributed
  - o Indexed with B-Tree → indexes properties of nodes and relationships

## Lecture 15: Data System Architecture, Disk Storage, B+ Trees

## Data System Architecture

Data System:	Use:
DBMS	Production Data
Key-value store/cache	Fast access to cache, frequently accessed data
Data Warehouse	Dashboarding, Reporting, Fast access to aggregated data
Search Engine	Quickly locate records
Big Data System, Relational	ETL Jobs
Layer	

## Moving Data

- ETL Jobs
  - o Runs nightly, transfers data from one system to another
- Replication
  - o Simply copies data to multiple locations of same type
- Streaming / Message Bus
  - o Publisher/Subscriber model: Subscribers read data from bus and push to databases

#### **HDDs**

- Vocabulary
  - o <u>Tracks</u> are divided into <u>sectors</u>
    - Sectors can achieve locality if data is on same cylinder
  - O Data is addressed by block number and transferred via blocks
    - Pages are the same as blocks
    - We want to minimize seeks and block transfers
- Disk Performance
  - o Seek Time: Time it takes to move head from parked state to a track
  - o Rotational Latency: Once head is on a track, the time to find a sector
  - o Access Time: Time from request to start of transfer.

- Access Time = Seek Time + Rotational Latency
- Data Transfer Rate: Time to transfer one sector from disk to RAM
- O <u>IOPS</u>: Cloud computing, burstiness, Number of Read/Writes per second
- Data Access
  - Sequential (HDD)
    - Access blocks in a pre-determined order
      - Very fast if data is contiguous
  - Random Access (SSD)
    - Blocks fetched from arbitrary locations
      - Bad on HDD
- Speeding up Disk Reads/Seeks
  - o <u>Buffering</u>: Read blocks into memory, when out of RAM, evict based on policy
    - Evict using Most Recently Used (MRU)
  - o Read-Ahead: Capture the next few blocks when retrieving a block from disk
  - Scheduling: Group requests by cylinder
  - o File Organization: Organize blocks and records into files
    - Records  $\subseteq$  Blocks  $\subseteq$  Files  $\subseteq$  Relations
  - o Non-Volatile Write Buffers: Queue Read/Write requests in NVRAM on HDD

## Organizing Records into Files

- Fixed types
  - o Pointer to first record
  - o Pointer to first empty record (Free List)
    - On deletion, update previous record's pointer to next record
    - Update the free list to store new space
  - Overflow Block
    - Create a new block, store the overflow, update the pointers
- Variable-length types
  - Table of Contents
    - Variable Sized Attributes, Data for fixed sized attributes, null bitmap, data for variable sized attributes

• We store a pointer to the first record in the file header

## Organizing Records into Blocks

- Blocks contain headers
  - Number of records in the block
  - Location of end of free space in the block
  - Location and size of each record
- Modifying the file means you need to rebuild the entire file

## The Index Data Structure (for indexing files)

- Index: Key → Block Number, Byte Offset
  - o Ordered Indices: Based on sorted ordering of sort-keyvalues
  - o Hash Indices: Based on uniform distribution of values across a range of buckets
- Types of Ordered Indices
  - o Primary (Clustering)
    - Both values of keys in index and records in corresponding data file are sorted by search key
      - Dictates ordering of underlying data file
      - Index and Records are separate
    - Does not need to be a Primary Key or Unique, but generally is
      - There is at most one primary index in a table
    - If there are duplicates, the first value is returned
    - Good for sequential access
  - Secondary (Non-Clustering)
    - Index and keys in records are sorted differently
      - Generally Random Access performance
  - Dense: Every value of the search key from the records is also in the index
  - O Sparse: Only some of the possible key values are in the index
    - To search for a key, we find the nearest key that is less than the one we are trying to find
    - We walk the pointers until we find what we need

- Keys in index and data are still ordered in the same way
- Index may fit in RAM
- More block transfers
  - To minimize block transfers, pick keys associated with records at the beginning of blocks

#### B+ Trees

## Concept

- o Good for range queries (non-equi), ok for equality (equi)
- Root and interior nodes form sparse index
- o Leaf nodes for a dense index on top of the records

#### Guarantees

- o Requires all paths have the same length from the root
- Balanced and self balancing
- Has O(n) fan out to reduce number of lock reads
- o Represents 1 disk seek and Block Transfer at each level of the tree
- o Left is  $\leq$ , Right is  $\geq$
- $\circ$  Holds up to n pointers but must contain roof (n/2) pointers
- o Some values may be duplicated within the tree

## Complexities

- $\circ$  Insertions and deletions:  $O(\log(\text{roof n/2}) k)$
- Additional space overhead
- Wasted space since nodes only need to be half full
- Reading each level of tree + sibling leaves = 1 seek + 1 block transfer

## **Lecture 16: Hash Indices, Nested Loop Join Algorithms**

#### Hash Indices

- Concepts
  - o Good for equi joins (equality)
  - Each key index maps into a bucket via a hash function h
    - Bucket contains a pointer to a linked list of records or blocks of records in RAM
    - H is hash mod function  $\rightarrow$  h: K  $\rightarrow$  B
  - Insertion
    - Compute Bucket via hash
    - Add the record
      - If does not fit, add an overflow bucket
  - Deletion
    - Compute h(K)
    - Find all records that match key K, delete them
      - If block is empty, delete the block
  - Hash Collisions/Bucket Overflow
    - Too little buckets causes too many collisions
      - Static Hashing: We can preallocate extra buckets beforehand
    - Too many buckets wastes resources
    - Bucket skew: Some buckets have a lot of records and overflow buckets while others don't
  - Ideal Buckets = [(number of records in relation R) / ratio of records per bucket] \*
    1.2
  - o Avoiding Bucket Overflow
    - Choose hash function based on anticipated size of file
      - Performance degradation lessened, space is wasted
    - Reorganize the hash structure due to file growth
      - Slow
    - Dynamically Resize Buckets

- Wastes space, minimizes collisions
- <u>Dynamic Hashing</u>: Grow or shrink amount of buckets
  - Linear Hashing, Extendible Hashing

## **Query Processing**

- Optimizer: Optimizes query into query/execution plan
- Process
  - o Parsing and Translation
    - Checks for table existence, converts to relational algebra
    - Optimizer selects most optimal plan
  - Optimization
    - Criterium:
      - Block Transfers and Seeks
      - Number of tuples
      - Current state of CPU/RAM
      - Data transfer speed
      - Disk Space
      - Time
  - Evaluation
    - Execution Time
    - Block reads and disk seeks

#### **Formulas**

- t<sub>T</sub>: Time to transfer 1 block from disk to memory
- t<sub>s</sub>: Average block access time (including rotational latency)
- Average time to transfer B blocks with S random accesses: Bt<sub>T</sub> + St<sub>s</sub>
- SELECT  $\sigma$ 
  - o b<sub>r</sub> is number of blocks in the file for R
  - Execution time:
    - Sequential File
      - $t_s + b_r * t_T$

- Non-Contiguous File then we must seek for each block
  - $b_r * t_s + b_r * t_T = b_r (t_T + t_s)$
- Search on key using ISAM, there is only one match
  - Best Case: Key/record is in first block
    - $\circ$   $t_s + b_r * t_T$
  - Worst Case: Key/record is in last block
    - $\circ$   $t_s + t_T$
  - Average Complexity for Uniform Distribution
    - $o t_s + (b_r / 2) * t_T$
- Indexed using clustering B+ tree, search on equality and unique key
  - h: height of the B+ tree, can be depth
  - Traverse from root to leaf node
    - $\circ$  h(t<sub>s</sub> + t<sub>T</sub>)
  - Search for record pointer in block
    - $0 \rightarrow \text{it's already in RAM}$
  - Retrieve record associated with the key
    - $\circ$   $t_s + t_T$
  - Total Time doing disk I/O
    - $h(t_s + t_T) + t_s + t_T = (h + 1)(t_s + t_T)$
- Indexed using a primary/clustering B+ tree, search on a key v such that K
  v or K > v
  - Find v and go to leaf node
    - $\circ$  h(t<sub>s</sub> + t<sub>T</sub>)
  - Fetch the first record
    - $o t_s + t_T$
  - Because it's primary, we can walk to pointers
  - Total:
    - $o h(t_s + t_T) + t_s + b * t_T$
    - o b: number of record blocks we need to read to find a match

- Three for loops: Iterate through each pairing of outer, inner relations
  - Does not require an index
  - o No restrictions on join condition  $\theta$
  - o Checks per tuple
  - O Super Complex: O(n<sub>r</sub> \* n<sub>s</sub>)
    - $n_r$ : number of tuples in R
    - $n_s$ : number of tuples in S
  - o Assuming maximum of 1 block of R, S can fit in RAM and R, S are sequential
    - Note: head can only be in one place at a time
    - b<sub>r</sub>: number of blocks in R
    - b<sub>s</sub>: number of blocks in S
  - Seeks and Block Transfers
    - Outer Relation Inner Relation
    - $n_r + b_r$  seeks
    - $n_r * b_s + b_r$  block transfers

#### Block Nested Loop Join

- Idea: Process outer relation by block and not tuple
  - o Each block in inner relation is read only once for each block in outer relation
- Complexity: Outer Relation Inner Relation
  - $\circ$   $b_r + b_r$  seeks
  - o  $b_r + b_s * b_r$  block transfers
  - o Total Disk I/O:
    - $(2 b_r)t_s + (br * b_s + b_r)t_T = 2 b_r * t_s + (t_T * b_r)(b_s + 1)$

#### Indexed Nested Join Loop

- Idea: Place an index on the inner relation to replace file scans on lookups
  - For each tuple in R, we lookup a key from the tuple is S and we can retrieve all matching tuples
- Complexity (Assuming index already exists)
  - o  $b_r (t_T + t_s)$  time to scan

- $\circ$  c \*  $n_r$  lookups
  - c is disk cost of a select from the index

$$\bullet \quad c = h(t_T + t_s) + (t_T + t_s)$$

- Find a key, fetch a record
- o Total Disk I/O:

$$b_r (t_T + t_s) + c * n_r$$

## Lecture 17: Merge Join, Hash Join, Spark Joins, Indexing Strategies

Merge Join/Interleaved Linear Scan

- Idea: Sort both relations before we join them, so we only need to scan each table once
  - o Sort based on common attributes in R  $\cap$  S
- Complexity:
  - $\circ$   $b_r + b_s$  block transfers
  - o 2 seeks (1 for S, 1 for R)

#### Hash Join

- Idea: Build and Probe Sides
  - Build Side S: Builds a hash table from build side relation to variety of things in memory
    - Must fit in memory
    - Value can be: Linked list of records, Linked list of pointers to records on disk, or linked list on disk
  - Probe Side R: Checks keys if it is in hash table on build side
    - If match, there might be a match on join key
      - Hash collisions → Lookup hash on join side to find records with proper key
- Procedure
  - o Construct a hash table using the values of the join key from S (and hash them)
  - O Hash table value can be one of three things:
    - Linked list, in memory, of pointers to records on disk (common)
    - Linked list, on disk, of blocks of records (disk-based)
    - Linked list, on memory, of records (memory-based)
  - o Do a full-table scan of R and search for records that match the join key
    - Due to hash collisions, we must walk the linked list to find all the records that match the key (and not just the hash)
- Complexity
  - Common Approach

- Build Phase (S):
  - b<sub>s</sub> Block Transfers + 1 Seek
- Probe Phase (R):
  - b<sub>r</sub> seeks + m seeks
- Total:
  - m is number of matches on the join key
  - Block Transfers:  $b_r + b_s + m$
  - Seeks:  $b_r + m + 1$
- o Disk-Based
  - Build Phase (S):
    - 1 seek + b<sub>s</sub> block transfers
  - Probe Phase (R):
    - 1 seek + b<sub>r</sub> block transfers
  - Total
    - Seeks: 2
    - Block Transfers:  $b_r + b_s$

## Misc questions

- Difference between hash table and hash index
  - o In hash join, RDBMS creates an ephemeral hash table
    - Lives as long as the join does, some part is in memory
    - Keys are in memory → Less disk I/O
  - A hash index is persisted on disk
    - Lives in the database schema
    - Lots of disk I/O

## Spark's Selection of Best Join condition

- Size of tables involved in the join
- Join condition
- Configuration values (defaults or spec of cluster)
- Hints given by the user

## Spark Join Algorithms

- Broadcast based joins: Table is broadcasted to all executors in the cluster
  - o Automatically chosen if the table is less than 10MB
  - Broadcast Hash Join
    - Small dataframe R that fits in memory on all executors and driver
    - R is divided across worker nodes
    - Only supports equijoins
  - Broadcast Nested-Loop Join
    - Process cartesian join in each executor as a nested loop join after broadcasting R across the cluster
    - Least performant, used as a last resort
    - Can handle all join types
    - Complexity
      - $n_r * b_s + b_r$  Block Transfers
      - $n_r + b_r$  Seeks
- Shuffle based joins
  - o Called when neither relation fits in RAM on every executor and the driver
    - Heavily impacted by key shew
    - Lots of Disk I/O and Network I/O
  - Shuffle Hash Join
    - Each row of R and S is shuffled across the network to a partition such that all rows with the same value on the join key ends up in the same partition.
       Local joins afterwards
      - Doesn't require keys to be sortable
      - Only works for equijoins
      - A lot more expensive than a broadcast
      - Uses additional RAM and time
  - o Shuffle Sort-Merge Join
    - After R, S shuffled to executed partitioned by key-value, each executor may contain rows of R and S with multiple different key-values

- All rows of the same key-value must be co-located in the same executor
- Notes
  - Keys must be sortable
  - Works on any equijoin

# **Indexing Strategies**

- Spatial Data and Temporal Data
  - o Nearest Neighbor Query: Circle within certain distance around user
  - o If 1D, divide and conquer via subintervals
  - o If 2D, Use a k-d tree
    - O(logN) time on average, O(N) in sequential files
  - O Data visualization: Random Forest, kNN, Voronoi tesselation

## **Lecture 18: Transactions, Isolation Levels**

<u>Concurrency</u>: Two transactions start, execute, and complete in overlapping time periods to share a resource

• Increasing concurrency/performance decreases Isolation and Consistency

<u>Parallelism</u>: Two operations in a transaction literally execute at the same time

<u>Serial</u>: Transaction runs entirely in full before another one begins

#### **Transactions Commits**

- Changes to database does not occur until commit is called
  - o Begin transaction, begin work, begin
  - o Atomic → Either all execute correctly, or nothing happens
    - Failures
      - Constraint (FK, PK, UK)
      - Hardware/Software Failure
- Lifecycle
  - o Committed: Transaction completes successfully and entirely
  - Partially Committed: After all steps run successfully, but before the result has been written to the database
  - o Failed: All changes are rolled back and transaction becomes aborted
  - Active: Transaction is still processing data
- Errors
  - External Writes: RDBMS commits a transaction but the application does not complete its work
    - App must ask the database system to issue a compensating transaction to undo the actions of the first transaction

#### **ACID**

- Atomicity
  - No partial results
- Consistency

- We do not gain or lose data
- <u>Isolation</u>
  - o Concurrent transactions operate independently and do not intersect
    - Maximizes consistency at the cost of performance
- <u>Durability</u>:
  - Once a transaction is committed, you cannot roll back and the data persists on disk

#### Relational databases follow ACID

• High consistency

## NoSQL transactions follow BASE

- Basically Available
  - o Ensure availability of data by spreading and replicating across nodes of cluster
- Soft-state
  - o Because of lack of consistency, data values may change over time
- Eventually Consistent
  - o Over time, data may become consistent (immediate consistency is not guaranteed
- High availability with eventual consistency

#### Thread Chart notes

- Consistent: Sum of initial registers is same as sum of threads after serialization
  - Conflict Serializable: If a schedule can be rearranged to become serial/matches behavior of a serial schedule
  - o <u>Serializability</u>
    - Improved throughput
    - Reduced waiting time
- Writes are global after commits
  - o Rollbacks no longer affect registers after commit
- Swaps move one instruction up and another down across two threads
  - You cannot swap incompatible instructions

- Errors
  - o WAR, WAW are not compatible
  - o RAR is ok
- Problems
  - <u>Dirty Writes</u>: T<sub>i</sub> writes/modifies a register that T<sub>j</sub> has already written/modified while T<sub>i</sub> is running
  - <u>Dirty Reads</u>: When transaction can read uncommitted changes from other transactions
    - Only read uncommitted allows dirty reads
    - Causes data inconsistencies
  - Fuzzy Read/Non-Repeatable Read: T<sub>i</sub> reads a value and then another transaction overwrites the value and commits it
    - T<sub>i</sub> reads the data twice and receives different values
  - One transaction inserts/deletes rows on a table in between two fetches (SELECTs) in another transaction
    - Occurs on a range of records
    - Occurs in Read Uncommitted, Read Committed, and Repeatable Read
    - The existing individual records do not change in T<sub>2</sub>
    - Preventable using serialization or predicate locking

#### Precedence Graph

- Draw edge if:
  - o WAR, RAW, WAW
  - Denote duplicate edges for correction
- Properties
  - o If no cycles  $\rightarrow$  S is conflict serializable and guarantees consistency
  - o If  $T_i \rightarrow T_j$ , exists in the graph, then any serial schedule S' == S if  $T_i$  finishes before  $T_j$  begins executing
  - Use topological sort to find proper ordering

## **Isolation Layers**

# Most Consistency, Isolation. Least Concurrency

- Serializable
  - o Maximum Consistency, Poor performance
  - o Little concurrency, equivalent to a serial schedule
- Repeatable Read
  - $\circ$  If  $T_i$  reads X, no other transaction  $T_j$  can update it until  $T_i$  commits
    - Lots of locking!
    - MySQL
- Read Committed
  - o T<sub>i</sub> can only read data from T<sub>j</sub> that is committed
    - DB2, SQLServer, Postgres
- Read Uncommitted
  - o T<sub>i</sub> can read data written from any other transaction even if not committed

Least Consistency, Isolation. Most Concurrency

## **Lecture 19: Transactions, Locking**

## Locking

- Shared Locks
  - o Allows T<sub>i</sub> to read a value of X
  - o Transactions can hold multiple shared/read locks on X
- Exclusive Locks
  - o Allows read/write access to T<sub>i</sub>
    - But we must use an exclusive lock to write/update X
  - o If T<sub>i</sub> holds an exclusive lock on X, no other transaction can update it
    - Other transaction will block

## Two-Phase Locking (2PL)

- Properties
  - o Ensures serializability
    - As soon as a transaction releases one lock, it can no longer acquire any new locks
  - No Anomalies
    - Each transaction pulls on the data it needs and no other transaction can modify these values
  - No guarantee of consistency
  - Not schedule specific
- Growing Phase
  - o T<sub>i</sub> can acquire all needed locks in the phase, but it cannot release any of them
- Shrinking Phase
  - o Begins as soon as a transaction releases a lock
  - o T<sub>i</sub> can release locks but cannot acquire any more

## Locking Issues

- <u>Starvation</u>: A given thread will never get the lock because other transactions keep requesting shared locks
  - o Solution: Lock Manager: Grants a lock on data item A in a particular mode M if:

- No other transaction holds a lock on A that conflicts with M
- There is no other transaction that is waiting for a lock on A and that made its lock request before T<sub>i</sub>
- Queue processes lock requests via a Lock Table
- Timing
  - Unlocking Too Early
    - Lose Serializability
    - Dirty Reads, Non-Repeatable Reads
  - Unlocking too late
    - Deadlock
- <u>Deadlock</u>: T<sub>1</sub> wants T<sub>2</sub>'s lock, T<sub>2</sub> wants T<sub>1</sub>'s lock, and nothing happens
  - o Prevention:
    - Lock Manager or roll backs
      - Wait-for graphs computed to detect deadlock
    - Acquire all required locks at once
    - Acquire locks in non-cyclical order
    - Preemption: Pause the other transaction until current transaction finishes
      with lock. Roll back
    - Lock timeouts
  - Detection and Recovery
    - Choose a victim transaction to abort and roll back
    - Rollback
    - Repeat until no deadlock
      - Oftentimes Round Robin works better