

SCIENCE PASSION TECHNOLOGY

Data Management 07 Physical Design & Tuning

Matthias Boehm

Graz University of Technology, Austria Computer Science and Biomedical Engineering Institute of Interactive Systems and Data Science BMK endowed chair for Data Management



Last update: Apr 30, 2022



Announcements/Org

- #1 Video Recording
 - Link in TeachCenter & TUbe (lectures will be public)
 - Hybrid: HSi13 / <u>https://tugraz.webex.com/meet/m.boehm</u>
 - Apr 25: no more COVID restrictions at TU Graz

#2 Course Evaluation and Exam

- Evaluation period: Jun 15 Jul 31
- Exams: Jun 27, 4pm (i13), Jul 07, 2.30pm (i12+i13), Jul 07, 5.30pm (i12+13), Jul 28, 5.30pm (i13)

#3 Exercises

- Exercise 1 in progress of being graded
- Exercise 2: May 03 + 7 late days in TeachCenter











Physical Design, and why should I care?

Performance Tuning via Physical Design

- Select physical data structures for relational schema and query workload
- #1: User-level, manual physical design by DBA (database administrator)
- #2: User/system-level automatic physical design via advisor tools







Agenda

- Compression Techniques
- Index Structures
- Table Partitioning
- Materialized Views







Compression Techniques



Background Storage System

Segments, Pages, Blocks

- Segments: storage unit of DB objects like relations (heap) and indexes
- Page: fixed-size memory region
- Block: smallest addressable unit on disk (e.g., POSIX block devices)

Buffer & Storage Management

- Buffer management at granularity of pages
- PostgreSQL default: 8KB
- Different table/page layouts (e.g., NSM, DSM, PAX, column)
- TID/RID Concept (pageID, slotID)
 → stable ID, even if records reorganized







ISDS



Background Storage System, cont.

- TID Concept (p, s)
 - TID := (page number, slot index)
 - Page slot directory holds tuple offsets (byte position) within page
- Example PostgreSQL
 - Recap: Papers(<u>PKey</u>, Title, Pages, CKey, JKey)
 - Hidden CTID system column (not shown on *, but usable)

SELECT CTID, PKey, Title, Pages FROM Papers		ctid tid	pkey integer	title character varying (512)	pages character va
	5681	(78,21)	731118	MV-IDX: Multi-Version Index in Action	671-674
	5682	(78,22)	731121	Hochperformante Analyse von Graph-Dat	311-330
	5683	(78,23)	731122	SPARQling Pig - Processing Linked Data wi	279-298
	5684	(78,24)	731123	RelaX: A Webbased Execution and Learnin	503-506
	5685	(78,25)	731129	Efficient In-Memory Indexing with General	227-246
	5686	(78,26)	731130	Datensicherheit in mandantenfähigen Clo	477-489
	5687	(78,27)	731131	In-Database Machine Learning: Gradient	247-266
	5688	(78,28)	731133	FlexY: Flexible; datengetriebene Prozessm	503-506
	5689	(78,29)	731134	Extending the MPSM Join	57-71
	5690	(78,30)	731137	Orthogonal key-value locking	237-256





Overview Database Compression

- Compression Overview
 - Fit larger datasets in memory, less I/O, better cache utilization
 - Some allow query processing directly on the compressed data
 - #1 Page-level compression (general-purpose GZIP, Snappy, LZ4)
 - #2 Row-level heavyweight/lightweight compression (e.g., Huffman)
 - #3 Column-level lightweight compression (NS, RLE, DICT, Delta, FOR → next slide)
 - #4 Specialized log and index compression

[Patrick Damme et al: Lightweight Data Compression Algorithms: An Experimental Survey. **EDBT 2017**]





Lightweight Database Compression Schemes

- Null Suppression
 - Compress integers by omitting leading zero bytes/bits (e.g., NS, gamma)
- Run-Length Encoding
 - Compress sequences of equal values by runs of (value, start, run length)

Dictionary Encoding

 Compress column w/ few distinct values as pos in dictionary (→ code size)

Delta Encoding

 Compress sequence w/ small changes by storing deltas to previous value

Frame-of-Reference Encoding

 Compress values by storing delta to reference value (outlier handling) 0000000 0000000 0000000 01101010 11 01101010



177	3 1 7 1 3 3 7 1 3 3 7 3
1,3,7	dictionary (code size 2 bit)
133	2 1 3 1 2 2 3 1 2 2 3 2





106



Index Structures





Index Scan

sorted

Table Scan

Overview Index Structures

- Table Scan vs Index Scan
 - For highly selective predicates, index scan asymptotically much better than table scan
 - Index scan higher per tuple overhead (break even ~5% output ratio)
 - Multi-column predicates: fetch/RID-list intersection





Additional Terminology

- Create Index
 - Create a secondary (nonclustered) index on a set of attributes
 - Clustered: tuples sorted by index
 - Non-clustered: sorted attribute with tuple references
 - Can specify uniqueness, order, and indexing method
 - PostgreSQL methods: <u>btree</u>, hash, gist, and gin

Binary Search

- pos = binarySearch(data,key=23)
- Given sorted data, find key position (insert position if non-existing)
- k-ary search for SIMD data-parallelism
- Interpolation search: probe expected pos in key range (e.g., search([1:10000], 9700))





CREATE INDEX ixStudLname

DROP INDEX ixStudLname;

ON Students **USING** btree

(Lname ASC NULLS FIRST);



Classification of Index Structures

ID Access Methods

[Theo Härder, Erhard Rahm: Datenbanksysteme: Konzepte und Techniken der Implementierung, **2001**]



ND Access Methods

- Linearization of ND key space + 1D indexing (Z order, Gray code, Hilbert curve)
- Multi-dimensional trees and hashing (e.g., UB tree, k-d tree, gridfile)
- Spatial index structures (e.g., R tree)



14

B-Tree Overview

[Rudolf Bayer, Edward M. McCreight: Organization and Maintenance of Large Ordered Indices. Acta Inf. (1) 1972]



- **History B-Tree**
 - Bayer and McCreight 1972, Block-based, Balanced, Boeing Labs
 - Multiway tree (node size = page size); designed for DBMS
 - Extensions: B+-Tree/B*-Tree (data only in leafs, double-linked leaf nodes)
- Definition B-Tree (k, h) $\left\lceil \log_{2k+1}(n+1) \right\rceil \le h \le \left| \log_{k+1}\left(\frac{n+1}{2}\right) \right| + 1$ All paths from root to leafs have equal length h All nodes (except root) have [k, 2k] key entries All nodes adhere to max constraints All nodes (except root, leafs) have [k+1, 2k+1] successors Data is a record or a reference to the record (RID) k=2





Index Structures





- Lookup Q_k within a node
 - Scan / binary search keys for Q_K, if K_i=Q_K, return D_i
 - If node does not contain key
 - If leaf node, abort search w/ NULL (not found), otherwise
 - Decent into subtree Pi with $K_i < Q_K \le K_{i+1}$
- Range Scan Q_{L<K<U}
 - Lookup Q_L and call next K while K<Q_U (keep current position and node stack)





B-Tree Insert

- Basic Insertion Approach
 - Always insert into leaf nodes!
 - Find position similar to lookup, insert and maintain sorted order
 - If node overflows (exceeds 2k entries) → node splitting

Node Splitting Approach

- Split the 2k+1 entries into two leaf nodes
- Left node: first k entries
- Right node: last k entries
- (k+1)th entry inserted into parent node
 → can cause recursive splitting
- Special case: root split (h++)
- B-Tree is self-balancing





Index Structures



B-Tree Insert, cont. (Example w/ k=1)





Insert 3 (2x split) 2 3

Insert 4



Note: Exercise 03? (B-tree insertion and deletion)





B-Tree Delete

- Basic Deletion Approach
 - Lookup deletion key, abort if non-existing
 - Case inner node: move entry from fullest successor node into position
 - Case leaf node: if underflows (<k entries) → merge w/ sibling





Index Structures

19



B-Tree Insert and Delete w/ k=2





Excursus: Prefix Trees (Radix Trees, Tries)

Generalized Prefix Tree

- Arbitrary data types (byte sequences)
- Configurable prefix length k'
- Node size: $s = 2^{k'}$ references
- Fixed maximum height h = k/k'
- Secondary index structure

Characteristics

- Partitioned data structure
- **Order-preserving** (for range scans)
- **Update-friendly**
- Properties
 - **Deterministic paths**
 - Worst-case complexity O(h)



reduced pointers

Level=1

012

P "value4"

P "value2"

Key

15

15



21

Excursus: Learned Index Structures

- A Case For Learned Index Structures
 - Sorted data array, predict position of key
 - Hierarchy of simple models (stages models)

[Tim Kraska, Alex Beutel, Ed H. Chi, Jeffrey Dean, Neoklis Polyzotis: The Case for Learned Index Structures. SIGMOD 2018]



Tries to approximate the CDF similar to interpolation search (uniform data)



 Follow-up Work on SageDBMS [Tim Kraska, Mohammad Alizadeh, Alex Beutel, Ed H. Chi, Ani Kristo, Guillaume Leclerc, Samuel Madden, Hongzi Mao, Vikram Nathan: SageDB: A Learned Database System. CIDR 2019]





BREAK (and Test Yourself)

Given B-tree below, insert key 9 and draw resulting B-tree (7/100 points)



Given B-tree below, delete key 27, and draw resulting B-tree (8/100 points)





23



BREAK (and Test Yourself), cont.

 Which of the following trees are valid – i.e., satisfy the constraints of – B-trees with k=1? Mark each tree as valid or invalid and name the violations (4/100 points)







Table Partitioning



25



Overview Partitioning Strategies

- Horizontal Partitioning
 - Relation partitioning into disjoint subsets

Vertical Partitioning

 Partitioning of attributes with similar access pattern

Hybrid Partitioning

- Combination of horizontal and vertical fragmentation (hierarchical partitioning)
- Derived Horizontal Partitioning









ISD



Correctness Properties

- #1 Completeness
 - $R \rightarrow R_1, R_2, ..., R_n$ (Relation R is partitioned into *n* fragments)
 - Each item from R must be included in at least one fragment
- #2 Reconstruction
 - $R \rightarrow R_1, R_2, ..., R_n$ (Relation R is partitioned into *n* fragments)
 - Exact reconstruction of fragments must be possible

#3 Disjointness

- $R \rightarrow R_1, R_2, ..., R_n$ (Relation R is partitioned into *n* fragments)
- $\mathbf{R}_i \cap \mathbf{R}_j = \emptyset \ (1 \le i, j \le n; i \ne j)$



Horizontal Partitioning

- Row Partitioning into n Fragments R_i
 - Complete, disjoint, reconstructable
 - Schema of fragments is equivalent to schema of base relation

Partitioning

- Split table by n selection predicates P_i (partitioning predicate) on attributes of R
- Beware of attribute domain and skew



 $\begin{aligned} \mathbf{R}_i &= \mathbf{\sigma}_{\mathbf{P}_i}(\mathbf{R}) \\ (1 \leq i \leq n) \end{aligned}$

Reconstruction Union of all fragments Bag semantics, but no duplicates across partitions R1 R2

 $\bigcup_{\substack{U \\ U \\ R3}} R = \bigcup_{1 \le i \le n} R$



Table Partitioning

Vertical Fragmentation

- Column Partitioning into n Fragments Ri
 - **Complete**, **reconstructable**, but not disjoint (primary key for reconstruction via join)
 - Completeness: each attribute must be included in at least one fragment

Partitioning

- Partitioning via projection
- Redundancy of primary key
- Reconstruction
 - Natural join over primary key

 $R_i = \pi_{PK,A_i}(R)$

A1 A2 PK

РК	A1

РК	A2

Hybrid horizontal/vertical partitioning

 $R = R_1 \bowtie R_i \bowtie R_n \bowtie / R_i = \bigcup R_{ii}$ $\rightarrow R = \cup R_i \lor / R_i = R_{1i} \bowtie R_{ii} \bowtie R_{ni}$



- $(1 \leq i \leq n)$

 $R = R_1 \bowtie R_i \bowtie R_n$ $(1 \leq i \leq n)$

Derived Horizontal Fragmentation

- Row Partitioning R into n fragements
 R_i, with partitioning predicate on S
 - Potentially complete (not guaranteed), restructable, disjoint
 - Foreign key / primary key relationship determines correctness

Partitioning

- Selection on independent relation S
- Semi-join with dependent relation R to select partition R_i

Reconstruction

- Equivalent to horizontal partitioning
- Union of all fragments













Exploiting Table Partitioning

- Partitioning and query rewriting
 - #1 Manual partitioning and rewriting
 - #2 Automatic rewriting (spec. partitioning)
 - #3 Automatic partitioning and rewriting
- Example PostgreSQL (#2)
 - CREATE TABLE Squad(JNum INT PRIMARY KEY, Pos CHAR(2) NOT NULL, Name VARCHAR(256)) PARTITION BY RANGE(JNum);
 - CREATE TABLE Squad10 PARTITION OF Squad FOR VALUES FROM (1) TO (10);
 - CREATE TABLE Squad20 PARTITION OF Squad FOR VALUES FROM (10) TO (20);
 - CREATE TABLE Squad24 PARTITION OF Squad FOR VALUES FROM (20) TO (24);

J#	Pos	Name
1	GK	Manuel Neuer
12	GK	Ron-Robert Zieler
22	GK	Roman Weidenfeller
2	DF	Kevin Großkreutz
4	DF	Benedikt Höwedes
5	DF	Mats Hummels
15	DF	Erik Durm
16	DF	Philipp Lahm
17	DF	Per Mertesacker
20	DF	Jérôme Boateng
3	MF	Matthias Ginter
6	MF	Sami Khedira
7	MF	Bastian Schweinsteiger
8	MF	Mesut Özil
9	MF	André Schürrle
13	MF	Thomas Müller
14	MF	Julian Draxler
18	MF	Toni Kroos
19	MF	Mario Götze
21	MF	Marco Reus
23	MF	Christoph Kramer
10	FW	Lukas Podolski
11	FW	Miroslav Klose

Table Partitioning



Exploiting Table Partitioning, cont.

Example, cont.
 SELECT * FROM Squad
 WHERE JNum > 11 AND JNum < 20





Table Partitioning

Excursus: Database Cracking

 Core Idea: Queries trigger physical reorganization (partitioning and indexing)



[Stratos Idreos, Martin L. Kersten, Stefan Manegold: Database Cracking. **CIDR 2007**]







Materialized Views





Overview Materialized Views

- Core Idea of Materialized Views
 - Identification of frequently re-occuring queries (views)
 - Precompute subquery results once, store and reuse many times





35



View Selection and Usage

- Motivation
 - Shared subexpressions very common in analytical workloads
 - Ex. Microsoft's Analytics Clusters (typical daily use -> 40% CSE saving)
- #1 View Selection
 - Exact view selection (query containment) is NP-hard
 - Heuristics, greedy and approximate algorithms

#2 View Usage

- Given query and set of materialized view, decide which views to use and rewrite the query for produce correct results
- Generation of compensation plans





[Alekh Jindal, Konstantinos Karanasos, Sriram Rao, Hiren Patel: Selecting Subexpressions to Materialize at Datacenter Scale. **PVLDB 2018**]



[Leonardo Weiss Ferreira Chaves, Erik Buchmann, Fabian Hueske, Klemens Boehm: Towards materialized view selection for distributed databases. **EDBT 2009**] 36



View Maintenance – When?

- Materialized view creates redundancy → Need for #3 View Maintenance
- Eager Maintenance (writer pays)
 - Immediate refresh: updates are directly handled (consistent view)
 - On Commit refresh: updates are forwarded at end of successful TXs

Deferred Maintenance (reader pays)

- Maintenance on explicit user request
- Potentially inconsistent base tables and views

Lazy Maintenance (async/reader pays)

- Same guarantees as eager maintenance
- Defer maintenance until free cycles or view required (invisible for updates and queries)



[Jingren Zhou, Per-Åke Larson, Hicham G. Elmongui: Lazy Maintenance of Materialized Views. **VLDB 2007**]





Materialized Views



ISDS



INF.01017UF Data Management / 706.010 Databases – 07 Physical Design and Tuning Matthias Boehm, Graz University of Technology, SS 2022



Materialized Views in PostgreSQL

- View Selection
 - Manual definition of materialized view only
 - With or without data

CREATE MATERIALIZED VIEW TopScorer AS

SELECT P.Name, Count(*)
FROM Players P, Goals G
WHERE P.Pid=G.Pid AND G.GOwn=FALSE
GROUP BY P.Name
ORDER BY Count(*) DESC
WITH DATA;

REFRESH MATERIALIZED VIEW TopScorer;

- View Usage
 - Manual use of view
 - No automatic query rewriting

View Maintenance

- Manual (deferred) refresh
- Complete, no incremental maintenance
- Note: Community work on IVM

[Yugo Nagata: Implementing Incremental View Maintenance on PostgreSQL, **PGConf 2018**], patch in 2019

[Yugo Nagata: The Way for Updating Materialized Views Rapidly, **PGConf 2020**, <u>https://www.pgcon.org/events/pgcon_2020/sessions/session/56/</u>

slides/47/pgcon2020_nagata_the_way_to_update_materialized_views_rapidly.pdf]

NameCountJames Rodríguez6Thomas Müller5Robin van Persie4Neymar4Lionel Messi4Arjen Robben3



Conclusions and Q&A

- Compression Techniques
- Index Structures
- Table Partitioning
- Materialized Views
- Next Lectures (Part A)
 - 08 Query Processing [May 09]
 - 09 Transaction Processing and Concurrency [May 16]
- Next Lectures (Part B)
 - 10 NoSQL (key-value, document, graph) [May 23]
 - 11 Distributed Storage and Data Analysis [May 30]
 - 12 Data Stream Processing Systems [Jun 13, Patrick]

