

Cost of Disk Access

When are Indexes Good? Bad?

Objective: Understand the implications of disk access in conjunction indexes

- Power of a dense index as a copy of data values
- Clustering
- Overhead

- many slides, Thanks Garcia-Molina,

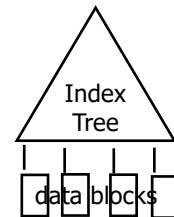
20: Index II

CS347

1

Indexes and Inserts

- All the data is in the data file, organized in blocks
 - sorted on primary key
 - (if possible) contiguously on disk



20: Index II

CS347

2

Simple I/O models

I/O models:

- Linear:
 - 1 average seek per block
 - n blocks, $f(n) = cn$, c = average seek time
- Affine:
 - seek first block,
 - weighted average of rotational latency + track to track seek time for each additional block
 - n block $f(n) = c + c'(n-1)$

5: Data on Disk

Database Management & Engineering

3

In this course

- Always the linear cost model, unless specified otherwise

20: Index II

CS347

4

Simple I/O models: How they fail

- buffering (hardware and software)
 - main memory
 - disk controller cache
- Prefetch – entire tracks into drive buffer
- blocks/pages aren't, necessarily, where *or* size as they are supposed to be (on board)
 - An indirection in case of bad pages/blocks

5: Data on Disk

Database Management & Engineering

5

What about SSDs?

20: Index II

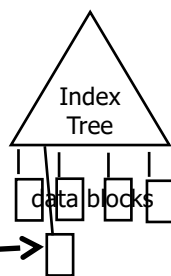
CS347

6

Data pages can overflow

If so

- fetch a new page from the heap
- new page is not contiguous



20: Index II

CS347

7

“Clustering”

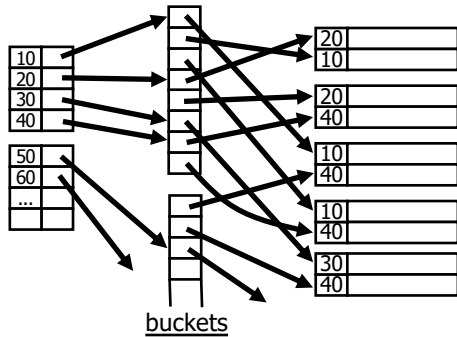
- Many uses/definitions of “clustering”.
- How well sequentially ordered blocks appear sequentially on disk.
 - 100% clustered → 100% sequentially accessible
 - 90% clustered → 9 out of 10 accesses is sequential
- Typically only the primary index is clustered - physical reality of disks.

// SQLserver, “clustered index” = primary index

15SmartIndexing

8

Secondary indexes are dense

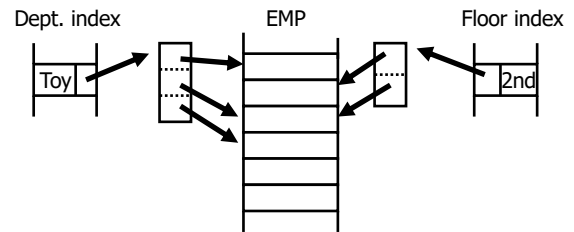


20: Index II

CS347

9

Query: Get employees in
(Toy Dept) \wedge (2nd floor)



→ Intersect toy bucket and 2nd Floor bucket to get set of matching EMP's

- Didn't have to touch the data file, except to find the output.

20: Index II

CS347

10

When do Secondary Indexes Make Sense?

Consider the example:

```
select *
from table
where city = 'austin'
```

table has one million rows

Number of cities =

- 100,000
- 100

Should we build/use a secondary index using a key called city?

20: Index II

CS347

11

Database Statistics are Stored in the Catalog/Dictionary

- 1 million records
- 10,000 data pages → avg 100 records per page
- fanout of a B+-tree = 100 using city as key

20: Index II

CS347

12

How Many Data Pages Do We Need to Read For the Query?

```
select *
from table
where city = 'austin'
```

Number of cities =

- 100,000

Repeat for 100 cities

20: Index II

CS347

13

What about primary index?

- Primary index guaranteed to be a [unique] key

- Suppose eid forms primary key

- Consider

```
select *
from table
where eid = 'miranker'    // how many results?
                        → 1 data page access
```

20: Index II

CS347

14

RDBMS Architecture

Storage Manager

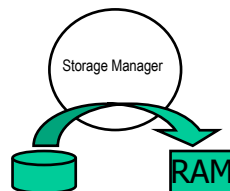
- Exploit memory hierarchy to compensate for slow disks.
 - working sets (from OS)
 - search algorithms

Specifics

- manage a *heap* of disk pages
- allocation of main memory (buffer management)

LRU-like replacement algorithm

- index methods, e.g. B+ tree (access paths)

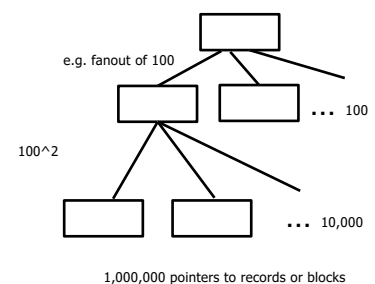


1: Introduction

Data Management

Locality of Reference at Top of the Tree

- *top of the index, usually pinned in memory*



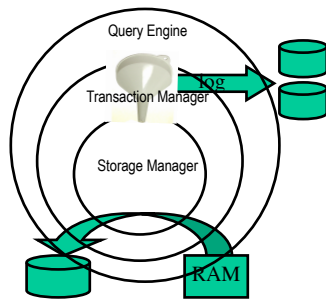
- number of leaves dominate the total number of nodes

20: Index II

CS347

16

Transactions and Secondary Indexes



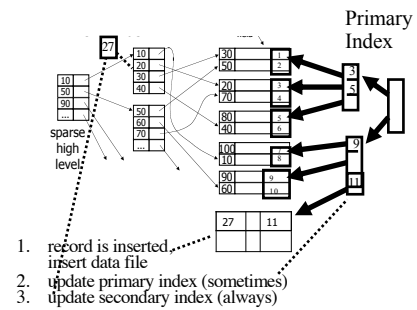
- recall, all writes/updates get written to a log
- atomicity - all or nothing
- *includes index records*

20: Index II

CS347

17

Transactions and Secondary Indexes What Happens On Insert?



20: Index II

CS347

18

Properties of Secondary Indexes

- They don't usually cluster well
- They must be dense - which can have it's own benefit
- Increase cost of dynamic operations
 - insert, delete, update
 - big impact on transaction performance
- ***Need to be careful about when to use them*** (text Ch. 8.4)
 - consider workload wrt:
 - transactions
 - database statistics (selectivity)
 - updates vs. reads

20: Index II

CS347

19

What About SSDs?

		9300 PRO (Read-Intensive, 1 Drive Write Per Day)			9300 MAX (Mixed-Use, 3 Drive Writes Per Day)		
Capacity ¹		3.04 TB	7.68 TB	15.36 TB	3.2 TB	6.4 TB	12.8 TB
Performance	Seq Read (MB/s) ²	3500	3500	3500	3500	3500	3500
	Seq Write (MB/s) ²	310	3500	3500	3100	3500	3500
	Rand Read (K IOPS) ³	835	850	850	835	850	850
	Rand Write (K IOPS) ³	105	145	150	210	310	310
Endurance (Terabytes Written in PB)		8.4	16.8	33.6	18.6	37.3	74.7

- Sequential access speed >> random access

20: Index II

CS347

20

SSD Internals

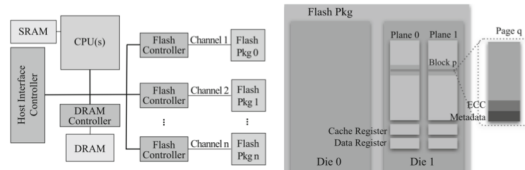


Fig. 2 Overview of the architecture of a solid-state drive

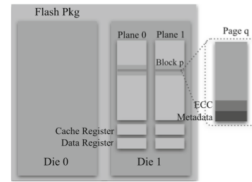


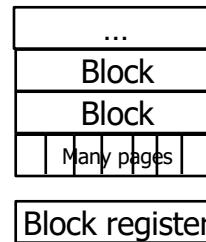
Fig. 1 The anatomy of a flash package

20: Index II

CS347

21

PLANE



20: Index II

Block read → many pages read in parallel

[like many word in RAM form a cache line]

CS347

22

But Also

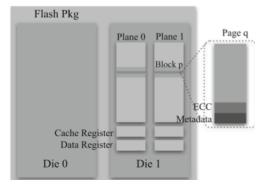


Fig. 1 The anatomy of a flash package

1. Many planes can be read in parallel
→ planes x pages/block = parallel page reads
2. Block is primitive unit for writing.
3. Block level wear.
4. Block level wear leveling
→ map logical blocks to physical blocks

20: Index II

CS347

23

DBA Moment: a word on transactions & real-world DBA heuristics

When developing database schema

- 1) Avoid secondary indexes on transactional data (e.g. point-of-sales)
 - why?
- 2) Separate transactional data from static data
 - e.g. (name,address) vs. balance
- 3) Minimize width of transactional rows
 - why?

20: Index II

CS347

24

A Practical (accepted) Dirty Trick

Account			
AID (pk)	name	address	balance
1234	dan	austin	1.00
6789	joe	houston	5.00

- start with good model
- then transform to denormalized model



- secondary indexes on "static" data,
- no impact on updates to dynamic data

AccountStatic		
AID	name	address
1234	dan	austin
6789	joe	houston

AccountDynamic	
AID	balance
1234	1.00
6789	5.00

20: Index II

CS347

25

A word on good design practice

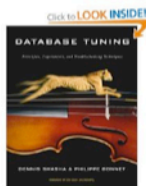
- do not consider the above at logical design time.
 - I.e. get the logic right, then optimize
 - good for getting system done on time
 - (separation of concerns)
- &
- long-term maintenance

20: Index II

CS347

26

If this is ever your job – money well spent:



Database Tuning: Principles, Experiments, and Troubleshooting Techniques (The Morgan Kaufmann Series in Data Management Systems) (Paperback)

~ Dennis Shasha (Author), Philippe Bonnet (Author) "Tuning rests on a foundation of informed common sense..." (more...)
 Key Phrases: [partitioning index](#), [multipoint queries](#), [summation query](#), [Morgan Kaufmann](#), [Tuning the Guts](#), [International Conference](#) (more...)
 (4 customer reviews)

List Price: ~~\$76.95~~

Price: **\$52.56** & eligible for free shipping with **Amazon Prime**

You Save: **\$23.39 (31%)**

20: Index II

CS347

27