

low c	ostl	y is thi	s oper	atio	n (naive solution)?			
<u>course</u>	per	<u>weekday</u>	nour	room	\			
TDA356	2	VR	Monday	13:15				
TDA356	2	VR	Thursday	08:00				
TDA356	4	HB1	Tuesday	08:00				
TDA356	4	HB1	Friday	13:15				
TIN090	1	HC1	Wednesday	08:00				
TIN090	1	HA3	Thursday	13:15				
SELECT * FROM Lectures WHERE course = 'TDA356' AND period = 2:				Go through all <i>n</i> rows, compare with the values for course and period = $2n$ comparisons				



### Index

- When relations are large, scanning all rows to find matching tuples becomes very expensive.
- An *index* on an attribute A of a relation is a data structure that makes it efficient to find those tuples that have a fixed value for attribute A.
  - Example: a hash table gives amortized O(1) lookups.

# Quiz!

# Asymptotic complexity (O(x) notation) is misleading here. Why?

The asymptotic complexity works for data structures in main memory. But when working with stored persistent data, the running time of the data structure, once in main memory, is negligible compared to the time it takes to read data from disk. What really matters to get fast lookups in a database is to minimize the number of disk blocks accessed (could use asymptotic complexity over disk block accessing though).

Indexes help here too though. If a relation is stored over a number of disk blocks, knowing in which of these to look is helpful.

# Typical costs

- Some typical costs of disk accessing for database operations on a relation stored over n blocks:
  - Query the full relation: n (disk operations)
  - Query with the help of index: k, where k is the number of blocks pointed to (1 for key).
  - Access index: 1
  - Insert new value: 2 (one read, one write)
  - Update index: 2 (one read, one write)

### Example:

```
SELECT *
FROM Lectures
WHERE course = 'TDA356'
AND period = 2;
```

Assume Lectures is stored in n disk blocks. With no index to help the lookup, we must look at all rows, which means looking in all n disk blocks for a total cost of n.

With an index, we find that there are 2 rows with the correct values for the course and period attributes. These are stored in two different blocks, so the total cost is 3 (2 blocks + reading index).

# Quiz!

### How costly is this operation?

#### SELECT \*

FROM Lectures, Courses WHERE course = code;

No index: Go through all *n* blocks in Lectures, compare the value for course from each row with the values for code in all rows of Courses, stored in all *m* blocks. The total cost is thus n \* maccessed disk blocks. Lectures: n disk blocks Courses: m disk blocks

Index on code in Courses: Go through all *n* blocks in Lectures, compare the value for course from each row with the index. Since course is a key, each value will exist at most once, so the cost is 2 \* n + 1accessed disk blocks (1 for fetching the index once).

# CREATE INDEX

 Most DBMS support the statement CREATE INDEX index name ON table (attributes);

– Example:

CREATE INDEX courseIndex ON Courses (code);

- Statement not in the SQL standard, but most DBMS support it anyway.
- Primary keys are given indexes implicitly (by the SQL standard).

### Important properties

- · Indexes are separate data stored by itself.
  - Can be created
    - ✓ on newly created relations
    - ✓ on existing relations
       will take a long time on large relations.
  - Can be dropped without deleting any table data.
- SQL statements do not have to be changed

- a DBMS automatically uses any indexes.

### Quiz!

# Why don't we have indexes on all attributes for faster lookups?

- Indexes require disk space.
- Modifications of tables are more expensive.
  - Need to update both table and index.
- Not always useful
  - The table is very small.
- We don't perform lookups over it (Note: lookups ≠ queries).
- Using an index costs extra disk block accesses.

## Rule of thumb

- Mostly queries on tables use indexes for key attributes.
- · Mostly updates be careful with indexes!



#### Example: Suppose that the Lectures relation is stored in 20 disk blocks, and that we typically perform three operations on this table: – insert new lectures (Ins) – list all lectures of a particular course (Q1)

list all lectures in a given room (Q2)

Let's assume that in an average week there are:

- 2 lectures for each course, and
- 10 lectures in each room.

#### Let's also assume that

 each course has lectures stored in 2 blocks, and
 each room has lectures stored in 7 (some lectures are stored in the same block).

xamp	ole coi	ntinu	ed:	Insert new le List all lecture List all lecture	cture es of es in	s (Ins) a particula a given ro	ar course (Q1) om (Q2)
Indexes	No in	dex	Index for (course, period, weekday)			Index for room	Both indexe
Ins	2			4		4	6
Q1	20	)	3			20	3
Q2	20	J	20			8	8
cost	2 + 18p <sub>1</sub>	+ 18p <sub>2</sub>	4 - p <sub>1</sub> + 16p <sub>2</sub>		4	+ 16p <sub>1</sub> + 4	p <sub>2</sub> 6 - 3p <sub>1</sub> + 2p
The or	a artizad a	anat da	nondo or	the distributi	on of	the energy	tiona n ia
The ar propor propor differe	nortized of tion of op tion of op nt values	cost de peration peration of p <sub>1</sub> a	pends or is that are is that are nd p <sub>2</sub> we	the distribution e Q <sub>1</sub> queries, e Ins modifica get actual co	on of p <sub>2</sub> sir tions sts of	the opera nilarly for is 1 – p <sub>1</sub> - f:	tions. $p_1$ is $Q_2$ , and thus t – $p_2$ . For some
The ar propor propor differei	nortized of tion of op tion of op nt values	cost deperation peration of $p_1$ a 2 + 18p	pends on is that are is that are nd $p_2$ we $p_1 + 18p_2$	the distribution e $Q_1$ queries, e Ins modification get actual co $4 - p_1 + 16p_2$	on of $p_2$ sir tions sts of 4 + 1	the opera nilarly for is $1 - p_1 - f_1$ f: $16p_1 + 4p_2$	tions. $p_1$ is $Q_2$ , and thus t - $p_2$ . For some $6 - 3p_1 + 2p_2$
The an propor differen p <sub>1</sub> = p	nortized of tion of op tion of op nt values $\frac{1}{2} = 0.4$	cost deperation peration of $p_1$ a 2 + 18p 1	pends on is that are is that are nd $p_2$ we $p_1 + 18p_2$ 6.4	the distribution e $Q_1$ queries, e lns modifical get actual co $4 - p_1 + 16p_2$ 10	on of $p_2 siritionssts of4 + 1$	the operative for the operative for the operative formula $1 - p_1 - f_1$ for the formula $16p_1 + 4p_2$ of the formula $12$ of the operative formula $12$ oper	tions. $p_1$ is $Q_2$ , and thus t $-p_2$ . For some $6 - 3p_1 + 2p_2$ <b>5.6</b>
The an propor propor different $p_1 = p$ $p_1 = p$	mortized of tion of op tion of op nt values $\frac{1}{2} = 0.4$ $\frac{1}{2} = 0.1$	cost deperation peration of $p_1$ a 2 + 18p 1	pends on as that are nd $p_2$ we $p_1 + 18p_2$ 6.4 5.6	the distribution e $Q_1$ queries, e lns modifica get actual co $4 - p_1 + 16p_2$ 10 5.5	on of $p_2$ sir tions sts of 4 + 1	the operation of the o	tions. $p_1$ is $Q_2$ , and thus t $-p_2$ . For some $6 - 3p_1 + 2p_2$ <b>5.6</b> 5.9

KBB056	 •	KBB056	KC	Monday	08	1
KMB017	 •	KMB017	MVH12	Tuesday	08	
TDA357		KMB017	MVH12	Wednesday	15	-
TMS145	/	TDA357	HA4	Monday	10	-
UMF012		TDA357	HB1	Thursday	10	-
UMF018	$\langle \rangle$	TMS145	KC	Friday	08	-
	``	UMF012	MVF23	Friday	13	-
	$\langle \rangle$	UMF012	MVF23	Monday	13	-
		UMF018	MVF23	Tuesday	10	+







## Quiz!

- Indexes are incredibly useful (although they are not part of the SQL standard).
- Doing it wrong is costly.
- Requires knowledge about the internals of a DBMS.
- How is data stored? How large is a block?
- A DBMS should be able to decide better than the user what indexes are needed, from usage analysis.

So why don't they??

# Summary – indexes

- Indexes make certain lookups and joins more efficient.
  - Disk block access matters.
  - Multi-attribute indexes
- CREATE INDEX
- · Dense, sparse, multi-level and secondary
- · Usage analysis
  - What are the expected operations?
  - How much do they cost?
    - $\Sigma(\text{cost of operation})x(\text{proportion of operations of that kind})$