

Slotted Pages (5)

Problem:

- 1 transaction T_1 updates data item i_1 on page P_1 to a very small size (or deletes i_1)
- 2 transaction T_2 inserts a new item i_2 on page P_1 , filling P_1 up
- 3 transaction T_2 commits
- 4 transaction T_1 aborts (or T_3 updates i_1 again to a larger size)

TID concept \Rightarrow create an indirection

but where to put it? Would have to move i_1 and i_2 .

Slotted Pages (7)

One possible slot implementation:

T	S	O	O	O	L	L	L
---	---	---	---	---	---	---	---

- 1 if $T \neq 11111111_b$, the slot points to another record
- 2 otherwise the record is on the current page
 - 1 if $S = 0$, the item is at offset O , with length L
 - 2 otherwise, the item was moved from another page
 - ★ it is also placed at offset O , with length L
 - ★ but the first 8 bytes contain the original TID

The original TID is important for scanning.

Record Layout

The tuples have to be materialized somehow.

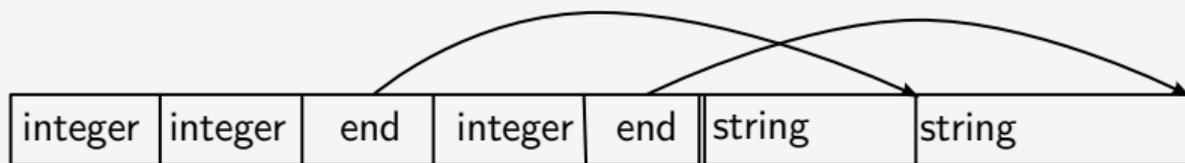
One possibility: serialize the attributes

integer	integer	length	string	integer	length	string	
---------	---------	--------	--------	---------	--------	--------	--

Problem: accessing an attribute is $O(n)$ in worst case.

Record Layout (2)

It is better to store offset instead of lengths



- splits tuple into two parts
- fixed size header and variable size tail
- header contains pointers into the tail
- allows for accessing any attribute in $O(1)$

Record Layout (3)

For performance reasons one should even reorder the attributes

- split strings into length and data
- re-order attributes by decreasing alignment
- place variable-length data at the end
- variable length has alignment 1

Gives better performance without wasting any space on padding.

NULL Values

What about NULL values?

- represent an unknown/unspecified value
- is a special value outside the regular domain

Multiple ways to store it

- either pick an invalid value (not always possible)
- or use a separate NULL bit

NULL bits allow for omitting NULL values from the tuple

- complicates the access logic
- but saves space
- useful if NULL values are common.

Compression

Some DBMS apply compression techniques to the tuples

- most of the time, compression is **not** added to save space!
- disk is cheap after all
- compression is used to **improve performance!**
- reducing the size reduces the bandwidth consumption

Some people really care about space consumption, of course.
But outside embedded DBMSs it is usually an afterthought.

Compression (2)

What to compress?

- the larger data compressed chunk, the better the compression
- but: DBMS has to handle updates
- usually rules out page-wise compression
- individual tuples can be compressed more easily

How to compress?

- general purpose compression like LZ77 too expensive
- compression is about performance, after all
- most system use special-purpose compression
- byte-wise to keep performance reasonable

Compression (3)

A useful technique for integer: variable length encoding

length (2 bits)	data (0-4 bytes)
-----------------	------------------

	Variant A	Variant B
00	1 byte value	NULL, 0 bytes value
01	2 bytes value	1 byte value
10	3 bytes value	2 bytes value
11	4 bytes value	4 bytes value

Compression (4)

The length is fixed length, the compressed data is variable length



Problem: locating compressed attributes

- depends on preceding compression
- would require decompressing all previous entries
- not too bad, but can be sped up
- use a lookup tuples per length byte

Compression (5)

Another popular technique: dictionary compression

Dictionary:

1	Berlin
2	München
3	Passauerstraße
...	...

Tuples:

city	street	number
1	3	5
2	3	7
...

- stores strings in a dictionary
- stores only the string id in the tuple
- factors out common strings
- can greatly reduce the data size
- can be combined with integer compression

Long Records

Data is organized in pages

- many reasons for this, including recovery, buffer management, etc.
- a tuple must fit on a single page
- limits the maximum size of a tuple

What about large tuples?

- sometimes the user wants to store something large
- e.g., embed a document
- SQL supports this via BLOB/CLOB (variable-length character data)

Requires some mechanism so handle these large records.

Long Records (2)

Simply spanning pages is not a good idea:

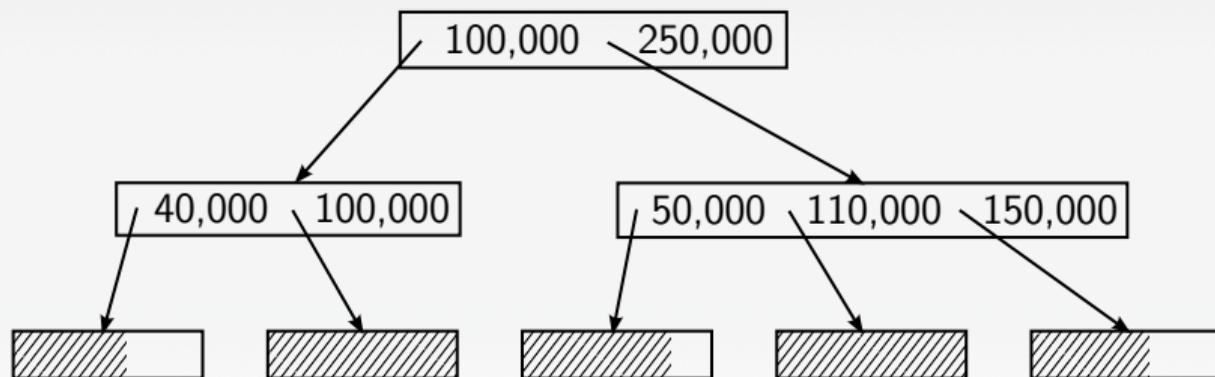
- must read an unbounded number of pages to access a tuple
- greatly complicates buffering
- a tuple might not even fit into main memory!
- updates that change the size are complicated
- intermediate results during query processing

Instead, keep the main tuple size down

- BLOBS/CLOBS are stored separate from the tuple
- tuple only contains a pointer
- increases the costs of accessing the BLOB, but simplifies tuple processing

Long Records (3)

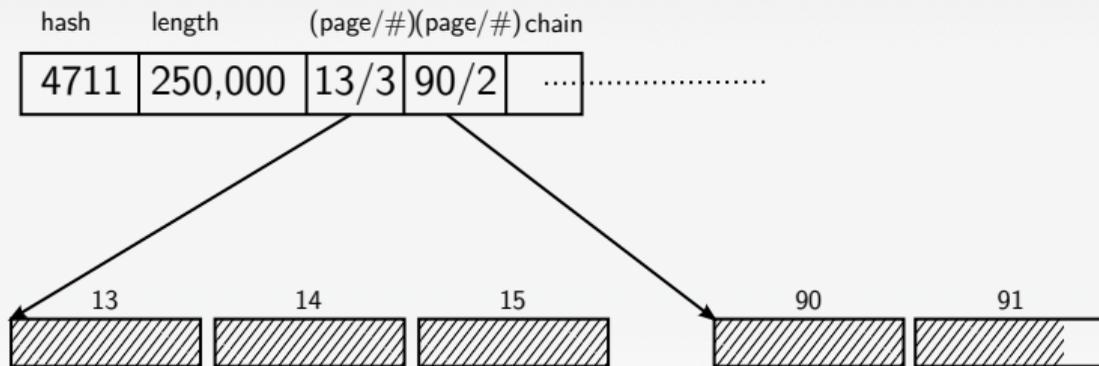
BLOBs can be stored in a B-Tree like fashion



- (relative) offset is search key
- allows for accessing and updating arbitrary parts
- very flexible and powerful
- but might be over-sophisticated
- SQL does not offer this interface anyway

Long Records (4)

Using an extent list is simpler



- real tuple points to BLOB tuple
- BLOB tuple contains a header and an extent list
- in worst case the extent list is chained, but should rarely happen
- extent list only allows for manipulating the BLOB in one piece
- but this is usually good enough
- hash and length to speed up comparisons

Long Records (5)

It makes sense to optimize for short BLOBs/CLOBs

- users misuse BLOBs/CLOBs
- they use CLOB to avoid specifying a maximum length
- but most CLOBs are short in reality
- on the other hand some BLOBs are really huge
- the DBMS cannot know
- so BLOBs can be arbitrary large, but short BLOBs should be more efficient

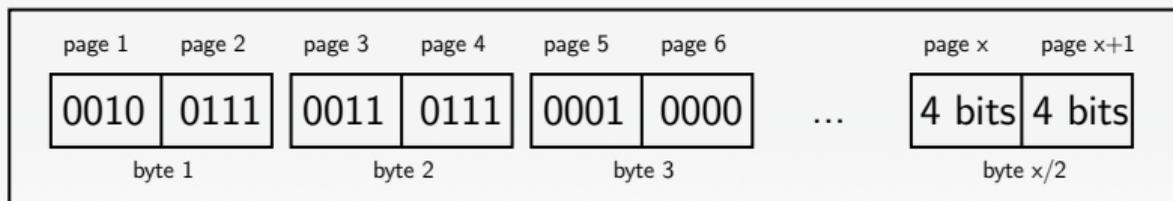
Approach:

- 1 BLOBs smaller than TID are encoded in BLOB TID
- 2 BLOBs smaller than page size are stored in BLOB record
- 3 only larger BLOBs use the full mechanism

Free Space Inventory

Problem: Where do we have space for incoming data?

Traditional solution: free space bitmap



Each nibble indicates the fill status of a given page.

Free Space Inventory (2)

Encode the fill status in 4 bits (some system use only 1 or 2):

- must approximate the status
- one possibility: $\text{data size} / \frac{\text{page size}}{2^{\text{bits}}}$
- loss of accuracy in the lower range
- logarithmic scale is often better
- $\lceil \log_2(\text{text size}) \rceil$
- or a combination (logarithmic for lower range, linear for upper range)

Encodes the free space (alternative: the used space) in a few bits.

Allocation

Allocating pages (or parts of a page) benefits from application knowledge

- often larger pieces are inserted soon after each other
- e.g. a set of tuples
- or one very large data item
- should be allocated close to each other

Allocation interface is usually

allocate(min, max)

- *max* is a hint to improve data layout
- some interfaces (e.g., segment growth) even implement over-allocation
- reduces fragmentation

