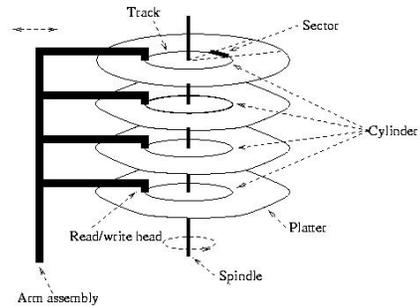


# Physical Data Organization and Indexing

## Chapter 11

1

## Physical Disk Structure



4

## Access Path

- Refers to the algorithm + data structure (*e.g.*, an index) used for retrieving and storing data in a table
- The choice of an access path to use in the execution of an SQL statement has no effect on the semantics of the statement
- This choice can have a major effect on the execution time of the statement

2

## Pages and Blocks

- Data files decomposed into *pages*
  - Fixed size piece of contiguous information in the file
  - Unit of exchange between disk and main memory
- Disk divided into page size *blocks* of storage
  - Page can be stored in any block
- Application's request for read item satisfied by:
  - Read page containing item to buffer in DBMS
  - Transfer item from buffer to application
- Application's request to change item satisfied by
  - Read page containing item to buffer in DBMS (if it is not already there)
  - Update item in DBMS (main memory) buffer
  - (Eventually) copy buffer page to page on disk

5

## Disks

- Capable of storing large quantities of data cheaply
- Non-volatile
- Extremely slow compared with cpu speed
- Performance of DBMS largely a function of the number of disk I/O operations that must be performed

3

## I/O Time to Access a Page

- *Seek latency* – time to position heads over cylinder containing page (avg = ~10 - 20 ms)
- *Rotational latency* – additional time for platters to rotate so that start of block containing page is under head (avg = ~5 - 10 ms)
- *Transfer time* – time for platter to rotate over block containing page (depends on size of block)
- *Latency* = seek latency + rotational latency
- Our goal – minimize average latency, reduce number of page transfers

6

## Reducing Latency

- Store pages containing related information close together on disk
  - *Justification*: If application accesses  $x$ , it will next access data related to  $x$  with high probability
- Page size tradeoff:
  - Large page size – data related to  $x$  stored in same page; hence additional page transfer can be avoided
  - Small page size – reduce transfer time, reduce buffer size in main memory
  - Typical page size – 4096 bytes

7

## Heap Files

- Rows appended to end of file as they are inserted
  - Hence the file is unordered
- Deleted rows create gaps in file
  - File must be periodically compacted to recover space

10

## Reducing Number of Page Transfers

- Keep cache of recently accessed pages in main memory
  - *Rationale*: request for page can be satisfied from cache instead of disk
  - Purge pages when cache is full
    - For example, use LRU algorithm
    - Record clean/dirty state of page (clean pages don't have to be written)

8

## Transcript Stored as a Heap File

666666	MGT123	F1994	4.0
123456	CS305	S1996	4.0
987654	CS305	F1995	2.0

page 0

717171	CS315	S1997	4.0
666666	EE101	S1998	3.0
765432	MAT123	S1996	2.0
515151	EE101	F1995	3.0

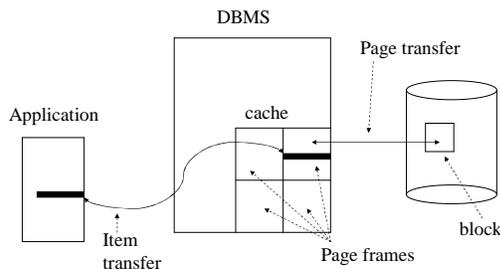
page 1

234567	CS305	S1999	4.0
878787	MGT123	S1996	3.0

page 2

11

## Accessing Data Through Cache



9

## Heap File - Performance

- Assume file contains  $F$  pages
- *Inserting a row*:
  - Access path is scan
  - Avg.  $F/2$  page transfers if row already exists
  - $F+1$  page transfers if row does not already exist
- *Deleting a row*:
  - Access path is scan
  - Avg.  $F/2+1$  page transfers if row exists
  - $F$  page transfers if row does not exist

12

## Heap File - Performance

- Query
  - Access path is scan
  - Organization efficient if query returns all rows and order of access is not important
  - Organization inefficient if a *few* rows are requested
    - Average  $F/2$  pages read to get a single row

```
SELECT T.Grade
FROM Transcript T
WHERE T.StudId=12345 AND T.CrsCode = 'CS305'
      AND T.Semester = 'S2000'
```

13

## Transcript Stored as a Sorted File

111111	MGT123	F1994	4.0
111111	CS305	S1996	4.0
123456	CS305	F1995	2.0

page 0

123456	CS315	S1997	4.0
123456	EE101	S1998	3.0
232323	MAT123	S1996	2.0
234567	EE101	F1995	3.0

page 1

234567	CS305	S1999	4.0
313131	MGT123	S1996	3.0

page 2

16

## Heap File - Performance

- Organization inefficient when a subset of rows is requested:  $F$  pages must be read

```
SELECT T.Course, T.Grade
FROM Transcript T
WHERE T.StudId = 123456 -- equality search
```

```
SELECT T.StudId, T.CrsCode
FROM Transcript T
WHERE T.Grade BETWEEN 2.0 AND 4.0 -- range search
```

14

## Maintaining Sorted Order

- **Problem:** After the correct position for an insert has been determined, inserting the row requires (on average)  $F/2$  reads and  $F/2$  writes (because shifting is necessary to make space)
- **Partial Solution 1:** Leave empty space in each page: *fillfactor*
- **Partial Solution 2:** Use *overflow pages (chains)*.

– Disadvantages:

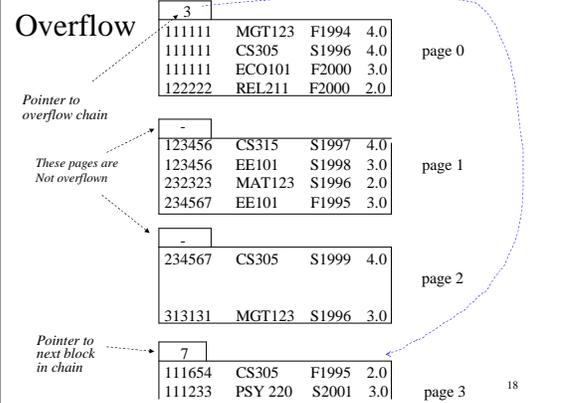
- Successive pages no longer stored contiguously
- Overflow chain not sorted, hence cost no longer  $\log_2 F$

17

## Sorted File

- Rows are sorted based on some attribute(s)
  - Access path is binary search
  - Equality or range query based on that attribute has cost  $\log_2 F$  to retrieve page containing first row
  - Successive rows are in same (or successive) page(s) and cache hits are likely
  - By storing all pages on the same track, seek time can be minimized
- Example – Transcript sorted on *StudId* :

```
SELECT T.Course, T.Grade FROM Transcript T WHERE T.StudId = 123456
SELECT T.Course, T.Grade FROM Transcript T WHERE T.StudId BETWEEN 111111 AND 199999
```



## Index

- Mechanism for efficiently locating row(s) without having to scan entire table
- Based on a **search key**: rows having a particular value for the search key attributes can be quickly located
- Don't confuse candidate key with search key:
  - Candidate key: *set* of attributes; *guarantees* uniqueness
  - Search key: *sequence* of attributes; *does not guarantee* uniqueness –just used for search

19

## Storage Structure

- Structure of file containing a table
  - Heap file (no index, not integrated)
  - Sorted file (no index, not integrated)
  - Integrated file containing index and rows (index entries contain rows in this case)
    - ISAM
    - B+ tree
    - Hash

22

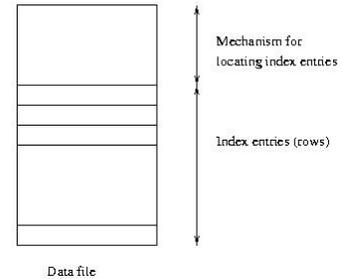
## Index Structure

- Contains:
  - *Index entries*
    - Can contain the data tuple itself (index and table are *integrated* in this case); or
    - Search key value and a pointer to a row having that value; table stored separately in this case – *unintegrated* index
  - *Location mechanism*
    - Algorithm + data structure for locating an index entry with a given search key value
  - Index entries are stored in accordance with the search key value
    - Entries with the same search key value are stored together (hash, B-tree)
    - Entries may be sorted on search key value (B-tree)

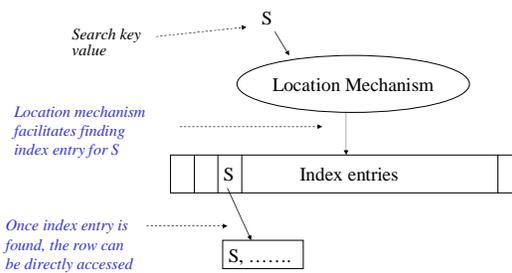
20

## Integrated Storage Structure

Contains table and (main) index

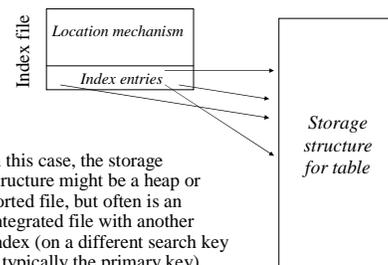


## Index Structure



21

## Index File With Separate Storage Structure



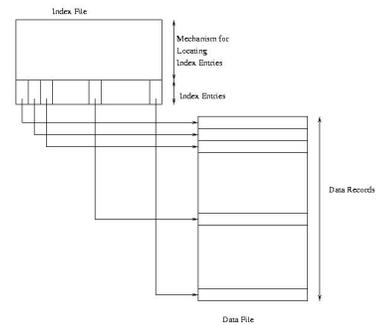
24

## Indices: The Down Side

- Additional I/O to access index pages (except if index is small enough to fit in main memory)
- Index must be updated when table is modified.
- SQL-92 does not provide for creation or deletion of indices
  - Index on primary key generally created automatically
  - Vendor specific statements:
    - CREATE INDEX ind ON Transcript (CrsCode)
    - DROP INDEX ind

25

## Clustered Secondary Index



28

## Clustered Index

- *Clustered index*: index entries and rows are ordered in the same way
  - An integrated storage structure is always clustered (since rows and index entries are the same)
  - The particular index structure (eg, hash, tree) dictates how the rows are organized in the storage structure
    - There can be at most one clustered index on a table
  - CREATE TABLE generally creates an integrated, clustered (main) index on primary key

26

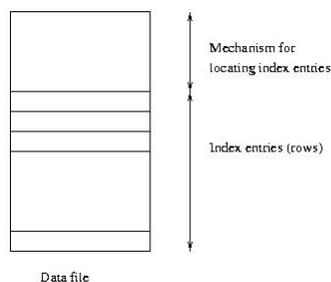
## Unclustered Index

- Unclustered (secondary) index: index entries and rows are not ordered in the same way
- An secondary index might be clustered or unclustered with respect to the storage structure it references
  - It is generally unclustered (since the organization of rows in the storage structure depends on main index)
  - There can be many secondary indices on a table
  - Index created by CREATE INDEX is generally an unclustered, secondary index

29

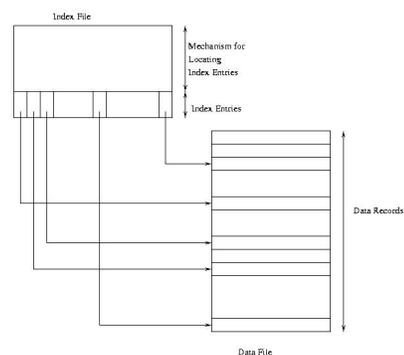
## Clustered Main Index

*Storage structure contains table and (main) index; rows are contained in index entries*



Data file

## Unclustered Secondary Index



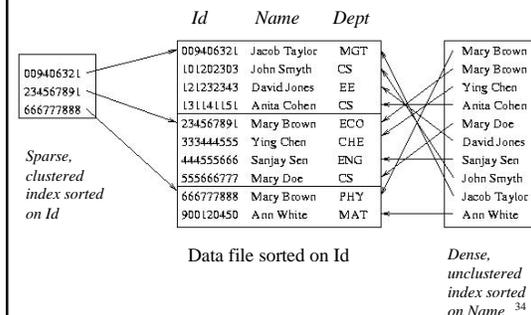
Data File

## Clustered Index

- Good for range searches when a range of search key values is requested
  - Use location mechanism to locate index entry at start of range
    - This locates first row.
  - Subsequent rows are stored in successive locations if index is clustered (not so if unclustered)
  - Minimizes page transfers and maximizes likelihood of cache hits

31

## Sparse Vs. Dense Index



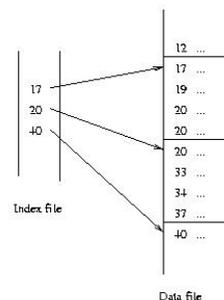
## Example – Cost of Range Search

- Data file has 10,000 pages, 100 rows in search range
- Page transfers for table rows (assume 20 rows/page):
  - Heap: 10,000 (entire file must be scanned)
  - File sorted on search key:  $\log_2 10000 + (5 \text{ or } 6) \approx 19$
  - Unclustered index:  $\leq 100$
  - Clustered index: 5 or 6
- Page transfers for index entries (assume 200 entries/page)
  - Heap and sorted: 0
  - Unclustered secondary index: 1 or 2 (all index entries for the rows in the range must be read)
  - Clustered secondary index: 1 (only first entry must be read)

32

## Sparse Index

*Search key should be candidate key of data file (else additional measures required)*



## Sparse vs. Dense Index

- *Dense index*: has index entry for each data record
  - Unclustered index *must* be dense
  - Clustered index need not be dense
- *Sparse index*: has index entry for each page of data file

33

## Multiple Attribute Search Key

- CREATE INDEX Inx ON Tbl (Att1, Att2)
- Search key is a *sequence* of attributes; index entries are lexically ordered
- Supports finer granularity equality search:
  - “Find row with value (A1, A2) ”
- Supports range search (tree index only):
  - “Find rows with values between (A1, A2) and (A1', A2' ) ”
- Supports partial key searches (tree index only):
  - Find rows with values of Att1 between A1 and A1'
  - But not “Find rows with values of Att2 between A2 and A2' ”

36

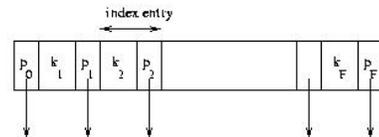
## Locating an Index Entry

- Use binary search (index entries sorted)
  - If  $Q$  pages of index entries, then  $\log_2 Q$  page transfers (which is a big improvement over binary search of the data pages of a  $F$  page data file since  $F \gg Q$ )
- Use multilevel index: Sparse index on sorted list of index entries

37

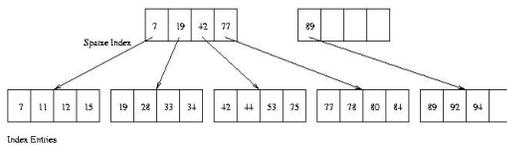
## Index Sequential Access Method (ISAM)

- Generally an integrated storage structure
  - Clustered, index entries contain rows
- Separator entry =  $(k_i, p_i)$ ;  $k_i$  is a search key value;  $p_i$  is a pointer to a lower level page
- $k_i$  separates set of search key values in the two subtrees pointed at by  $p_{i-1}$  and  $p_i$



40

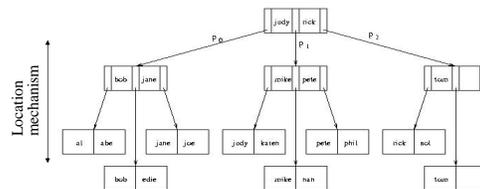
## Two-Level Index



- Separator level is a sparse index over pages of index entries
- Leaf level contains index entries
- Cost of searching the separator level  $\ll$  cost of searching index level since separator level is sparse
- Cost or retrieving row once index entry is found is 0 (if integrated) or 1 (if not)

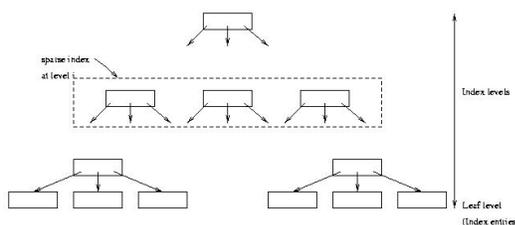
38

## Index Sequential Access Method



41

## Multilevel Index



- Search cost = number of levels in tree
- If  $\Phi$  is the fanout of a separator page, cost is  $\log_{\Phi} Q + 1$
- Example: if  $\Phi = 100$  and  $Q = 10,000$ , cost = 3 (reduced to 2 if root is kept in main memory)

39

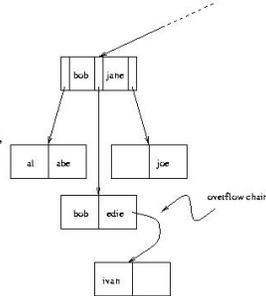
## Index Sequential Access Method

- The index is static:
  - Once the separator levels have been constructed, they never change
  - Number and position of leaf pages in file stays fixed
- Good for equality and range searches
  - Leaf pages stored sequentially in file when storage structure is created to support range searches
    - if, in addition, pages are positioned on disk to support a scan, a range search can be very fast (static nature of index makes this possible)
- Supports multiple attribute search keys and partial key searches

42

## Overflow Chains

- Contents of leaf pages change
- Row deletion yields empty slot in leaf page
- Row insertion can result in overflow leaf page and ultimately overflow chain
  - Chains can be long, unsorted, scattered on disk
  - Thus ISAM can be inefficient if table is dynamic



43

## Insertion and Deletion in B+ Tree

- Structure of tree changes to handle row insertion and deletion – *no* overflow chains
- Tree remains *balanced*: all paths from root to index entries have same length
- Algorithm guarantees that the number of separator entries in an index page is between  $\Phi/2$  and  $\Phi$ 
  - Hence the maximum search cost is  $\log_{\Phi/2} Q + 1$  (with ISAM search cost depends on length of overflow chain)

46

## B+ Tree

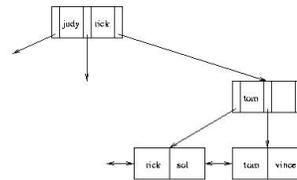
- Supports equality and range searches, multiple attribute keys and partial key searches
- Either a secondary index (in a separate file) or the basis for an integrated storage structure

$\emptyset$  Responds to dynamic changes in the table

44

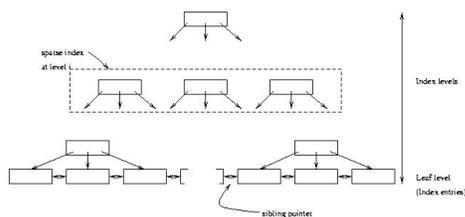
## Handling Insertions - Example

- Insert "vince"



47

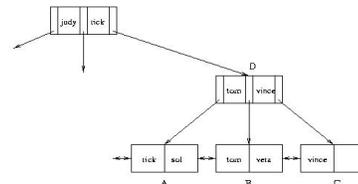
## B+ Tree Structure



- Leaf level is a (sorted) linked list of index entries
- Sibling pointers support range searches in spite of allocation and deallocation of leaf pages (but leaf pages might not be physically contiguous on disk) <sup>45</sup>

## Handling Insertions (cont'd)

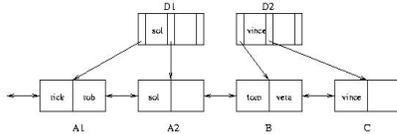
- Insert "vera": Since there is no room in leaf page:
  1. Create new leaf page, C
  2. Split index entries between B and C (but maintain sorted order)
  3. Add separator entry at parent level



48

## Handling Insertions (con't)

- Insert “rob”. Since there is no room in leaf page A:
  1. Split A into A1 and A2 and divide index entries between the two (but maintain sorted order)
  2. Split D into D1 and D2 to make room for additional pointer
  3. Three separators are needed: “sol”, “tom” and “vince”



49

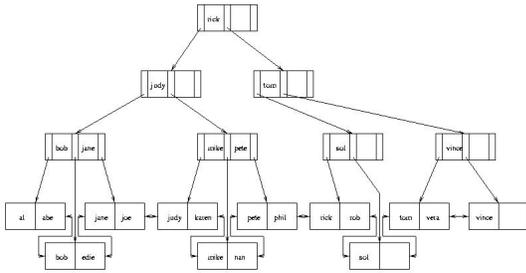
## Hash Index

- Index entries partitioned into *buckets* in accordance with a *hash function*,  $h(v)$ , where  $v$  ranges over search key values
  - Each bucket is identified by an address,  $a$
  - Bucket at address  $a$  contains all index entries with search key  $v$  such that  $h(v) = a$
- Each bucket is stored in a page (with possible overflow chain)
- If index entries contain rows, set of buckets forms an integrated storage structure; else set of buckets forms an (unclustered) secondary index

52

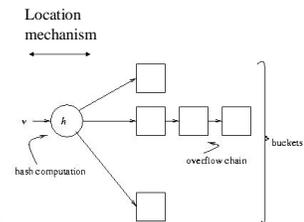
## Handling Insertions (cont'd)

- When splitting a separator page, push a separator up
- Repeat process at next level
- Height of tree increases by one



## Equality Search with Hash Index

- Given  $v$ :
1. Compute  $h(v)$
  2. Fetch bucket at  $h(v)$
  3. Search bucket



Cost = number of pages in bucket (cheaper than B<sup>+</sup> tree, if no overflow chains)

53

## Handling Deletions

- Deletion can cause page to have fewer than  $\Phi/2$  entries
  - Entries can be redistributed over adjacent pages to maintain minimum occupancy requirement
  - Ultimately, adjacent pages must be merged, and if merge propagates up the tree, height might be reduced
  - See book
- In practice, tables generally grow, and merge algorithm is often not implemented
  - Reconstruct tree to compact it

51

## Choosing a Hash Function

- Goal of  $h$ : map search key values randomly
  - Occupancy of each bucket roughly same for an average instance of indexed table
- Example:  $h(v) = (c_1 * v + c_2) \bmod M$ 
  - $M$  must be large enough to minimize the occurrence of overflow chains
  - $M$  must not be so large that bucket occupancy is small and too much space is wasted

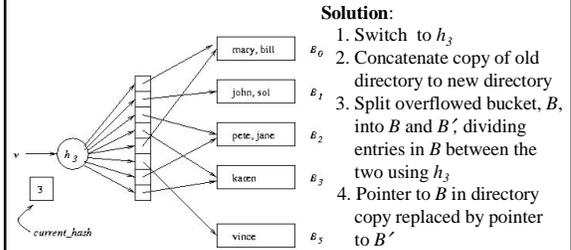
54

## Hash Indices – Problems

- Does not support range search
  - Since adjacent elements in range might hash to different buckets, there is no efficient way to scan buckets to locate all search key values  $v$  between  $v_1$  and  $v_2$
- Although it supports multi-attribute keys, it does not support partial key search
  - Entire value of  $v$  must be provided to  $h$
- Dynamically growing files produce overflow chains, which negate the efficiency of the algorithm

55

## Example (cont'd)



### Solution:

- Switch to  $h_3$
- Concatenate copy of old directory to new directory
- Split overflowed bucket,  $B_2$ , into  $B$  and  $B'$ , dividing entries in  $B$  between the two using  $h_3$
- Pointer to  $B$  in directory copy replaced by pointer to  $B'$

Note: Except for  $B'$ , pointers in directory copy refer to original buckets.

$current\_hash$  identifies current hash function.

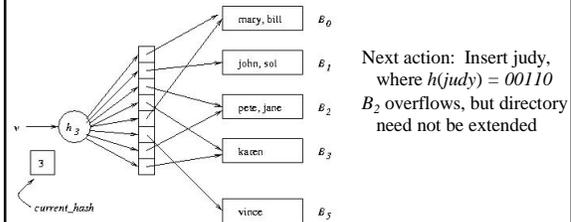
58

## Extendable Hashing

- Eliminates overflow chains by splitting a bucket when it overflows
- Range of hash function has to be extended to accommodate additional buckets
- Example:** family of hash functions based on  $h$ :
  - $h_k(v) = h(v) \bmod 2^k$  (use the last  $k$  bits of  $h(v)$ )
  - At any given time a unique hash,  $h_k$ , is used depending on the number of times buckets have been split

56

## Example (cont'd)



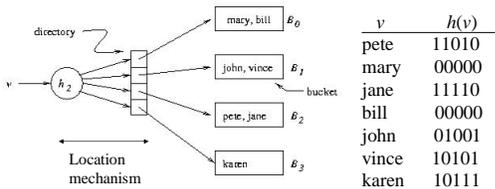
Next action: Insert judy, where  $h(judy) = 00110$ .  $B_2$  overflows, but directory need not be extended

**Problem:** When  $B_i$  overflows, we need a mechanism for deciding whether the directory has to be doubled

**Solution:**  $bucket\_level[i]$  records the number of times  $B_i$  has been split. If  $current\_hash > bucket\_level[i]$ , do not enlarge directory

59

## Extendable Hashing – Example



$v$	$h(v)$
pete	11010
mary	00000
jane	11110
bill	00000
john	01001
vince	10101
karen	10111

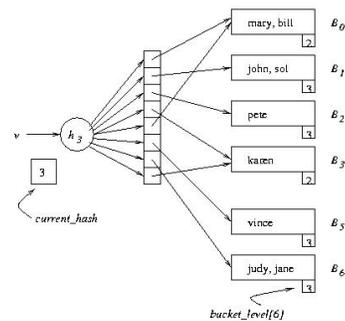
Extendable hashing uses a directory (level of indirection) to accommodate family of hash functions

Suppose next action is to insert sol, where  $h(sol) = 10001$ .

**Problem:** This causes overflow in  $B_1$

57

## Example (cont'd)



$bucket\_level[6]$

60

## Extendable Hashing

- Deficiencies:
  - Extra space for directory
  - Cost of added level of indirection:
    - If directory cannot be accommodated in main memory, an additional page transfer is necessary.

61

## Choosing an Index (cont'd)

Example 3:

```
SELECT T.CrsCode, T.Grade
FROM Transcript T
WHERE T.StudId = :id AND T.Semester = 'F2000'
```

- Equality search on *StudId* and *Semester*.
- If the primary key is (*StudId*, *Semester*, *CrsCode*) it is likely that there is a main, clustered index on this sequence of attributes.
- If the main index is a B<sup>+</sup> tree it can be used for this search.
- If the main index is a hash it cannot be used for this search. Choose B<sup>+</sup> tree or hash with search key *StudId* (since *Semester* is not as selective as *StudId*) or (*StudId*, *Semester*)

64

## Choosing An Index

- An index should support a query of the application that has a significant impact on performance
  - Choice based on frequency of invocation, execution time, acquired locks, table size

Example 1: 

```
SELECT E.Id
FROM Employee E
WHERE E.Salary < :upper AND E.Salary > :lower
```

- This is a range search on *Salary*.
- Since the primary key is *Id*, it is likely that there is a clustered, main index on that attribute that is of no use for this query.
- Choose a secondary, B<sup>+</sup> tree index with search key *Salary*

## Choosing An Index (cont'd)

Example 3 (cont'd):

```
SELECT T.CrsCode, T.Grade
FROM Transcript T
WHERE T.StudId = :id AND T.Semester = 'F2000'
```

- Suppose Transcript has primary key (*CrsCode*, *StudId*, *Semester*). Then the main index is of no use (independent of whether it is a hash or B<sup>+</sup> tree).

65

## Choosing An Index (cont'd)

Example 2: 

```
SELECT T.StudId
FROM Transcript T
WHERE T.Grade = :grade
```

- This is an equality search on *Grade*.
- Since the primary key is (*StudId*, *Semester*, *CrsCode*) it is likely that there is a main, clustered index on these attributes that is of no use for this query.
- Choose a secondary, B<sup>+</sup> tree or hash index with search key *Grade*

63