# CS 245: Database System Principles

**Notes 5: Hashing and More**
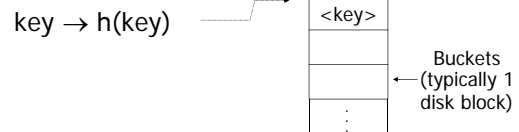
Hector Garcia-Molina

---

Hashing

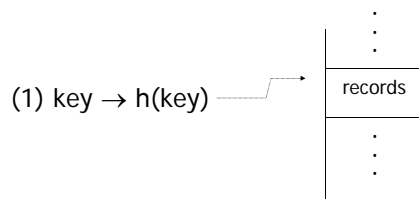key → h(key)  →  <key>

Buckets
← (typically 1 disk block)

---

Two alternatives

(1) key → h(key)  →  records

---

Two alternatives

(2) key → h(key)  →  key 1  →  record

Index

---

Two alternatives

(2) key → h(key)  →  key 1  →  record

Index

- Alt (2) for "secondary" search key

---

Example hash function

- Key = '$x_1 x_2 \ldots x_n$'  $n$ byte character string
- Have $b$ buckets
- h:  add $x_{1} + x_{2} + \ldots x_n$
  - compute sum modulo $b$

---

⊠ This may not be best function ...
⊠ Read Knuth Vol. 3 if you really
　　need to select a good function.

---

⊠ This may not be best function ...
⊠ Read Knuth Vol. 3 if you really
　　need to select a good function.

Good hash　　☞ Expected number of
function:　　　　keys/bucket is the
　　　　　　　　　same for all buckets

---

Within a bucket:
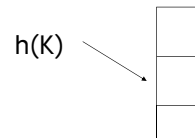
• Do we keep keys sorted?

• Yes, if CPU time critical
　　& Inserts/Deletes not too frequent

---

Next: example to illustrate
　　　　inserts, overflows, deletes

h(K)
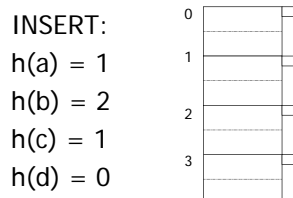
---

EXAMPLE  2 records/bucket

INSERT:
$h(a) = 1$
$h(b) = 2$
$h(c) = 1$
$h(d) = 0$

0
1
2
3

---

EXAMPLE  2 records/bucket

INSERT:
$h(a) = 1$
$h(b) = 2$
$h(c) = 1$
$h(d) = 0$
$h(e) = 1$

0　d
1　a
　　c
2　b
3

## EXAMPLE  2 records/bucket

INSERT:

$h(a) = 1$
$h(b) = 2$
$h(c) = 1$
$h(d) = 0$
$h(e) = 1$

| 0 | d |
| 1 | a → e |
| | c |
| 2 | b |
| 3 | |

---

## EXAMPLE:  deletion

Delete:
e
f

| 0 | a |
| 1 | b → d |
| | c |
| 2 | e |
| 3 | f |
| | g |

---

## EXAMPLE:  deletion

Delete:
e
f
c

| 0 | a |
| 1 | b → d |
| | c |
| 2 | e |
| 3 | f |
| | g |

maybe move "g" up

---

## EXAMPLE:  deletion

Delete:
e
f
c

| 0 | a |
| 1 | b → d |
| | c  d |
| 2 | e |
| 3 | f |
| | g |

maybe move "g" up

---

## Rule of thumb:

• Try to keep space utilization
  between 50% and 80%

$$\text{Utilization} = \frac{\text{\# keys used}}{\text{total \# keys that fit}}$$

---

## Rule of thumb:

• Try to keep space utilization
  between 50% and 80%

$$\text{Utilization} = \frac{\text{\# keys used}}{\text{total \# keys that fit}}$$

• If < 50%, wasting space
• If > 80%, overflows significant
   └─depends on how good hash
      function is & on # keys/bucket

3

## How do we cope with growth?

- Overflows and reorganizations
- Dynamic hashing

---

## How do we cope with growth?

- Overflows and reorganizations
- Dynamic hashing
  - Extensible
  - Linear

---

## Extensible hashing: two ideas

(a) Use $i$ of $b$ bits output by hash function

$$\xleftarrow{\quad} b \xrightarrow{\quad}$$

$$h(K) \rightarrow \boxed{00110101}$$

use $i \rightarrow$ grows over time....

---

(b) Use directory

$h(K)[i]$ ⟶ ⟶ to bucket

---

## Example: h(k) is 4 bits; 2 keys/bucket

$i = \boxed{1}$

⟶ 1 | 0001

⟶ 1 | 1001 | 1100

Insert 1010

---

## Example: h(k) is 4 bits; 2 keys/bucket

$i = \boxed{1}$

⟶ 1 | 0001

⟶ 1 | 1001 | 1010 1100

Insert 1010        1 | 1100

4

**Example:** h(k) is 4 bits; 2 keys/bucket

$i = 1$  $i = 2$

1
0001

1→2
1001
1010 1100

1→2
1100

**Insert 1010**

New directory
00
01
10
11

---

Example continued

i = 2
00
01
10
11

1
0001

2
1001
1010

2
1100

Insert:
0111
0000

---

Example continued

i = 2
00
01
10
11

0000
0001

1
0001  0111
0111

2
1001
1010

2
1100

Insert:
0111
0000

---

Example continued

i = 2
00
01
10
11

2
0000
0001

1→2
0001  0111
0111

2
1001
1010

2
1100

Insert:
0111
0000

---

Example continued

$i = 2$
00
01
10
11

0000  2
0001

0111  2

1001  2
1010

1100  2

Insert:
1001

---

Example continued

$i = 2$
00
01
10
11

0000  2
0001

0111  2

1001
1001

1010  1001  2
1010

1100  2

Insert:
1001

5

**Slide 31:**

Example continued

$i = 3$

| 0000 | 2 |
| 0001 | |

$i = 2$

00
01
10
11

| 0111 | 2 |

| 1001 | 3 |
| 1001 | |

1010 ~~1001~~ 2 3
~~1010~~

Insert:
1001

| 1100 | 2 |

000
001
010
011
100
101
110
111

CS 245          Notes 5          31

**Slide 32:**

Extensible hashing:  <u>deletion</u>

- No merging of blocks
- Merge blocks
   and cut directory if possible
   (Reverse insert procedure)

CS 245          Notes 5          32

**Slide 33:**

<u>Deletion example:</u>

- Run thru insert example in reverse!

CS 245          Notes 5          33

**Slide 34:**

<u>Note: Still need overflow chains</u>

- Example: many records with duplicate keys

insert 1100                    if we split:

| 1 |
| 1101 |
| 1100 |

| 2 |
| |

| 2 |
| 1100 |
| 1100 |

CS 245          Notes 5          34

**Slide 35:**

<u>Solution: overflow chains</u>

insert 1100                    add overflow block:

| 1 |
| 1101 |
| 1100 |

| 1 |
| 1101 | → | 1100 |
| 1101 |

CS 245          Notes 5          35

**Slide 36:**

| Summary |    Extensible hashing

⊕ Can handle growing files
   - with less wasted space
   - with no full reorganizations

CS 245          Notes 5          36

Extensible hashing

(+) Can handle growing files
     - with less wasted space
     - with no full reorganizations

(-) Indirection
      (Not bad if directory in memory)

(-) Directory doubles in size
      (Now it fits, now it does not)

CS 245                Notes 5               37

---

## Linear hashing

• Another dynamic hashing scheme

Two ideas:

(a) Use $i$ low order bits of hash

$$01110101$$

grows $\leftarrow \underbrace{\phantom{xx}}_{i}$

CS 245                Notes 5               38

---

## Linear hashing

• Another dynamic hashing scheme

Two ideas:

(a) Use $i$ low order bits of hash

$$01110101$$

grows $\leftarrow \underbrace{\phantom{xx}}_{i}$

(b) File grows linearly

CS 245                Notes 5               39

---

Example    $b$=4 bits,    $i$ =2,    2 keys/bucket

| 0000 | 0101 | | | ← Future growth buckets |
|------|------|------|------|
| 1010 | 1111 | | | |
| 00 | 01 | 10 | 11 | |

$m$ = 01 (max used block)

CS 245                Notes 5               40

---

Example    $b$=4 bits,    $i$ =2,    2 keys/bucket

| 0000 | 0101 | | | ← Future growth buckets |
|------|------|------|------|
| 1010 | 1111 | | | |
| 00 | 01 | 10 | 11 | |

$m$ = 01 (max used block)

Rule   If $h(k)[i] \leq m$, then
         look at bucket $h(k)[i]$
         else, look at bucket $h(k)[i] - 2^{i-1}$

CS 245                Notes 5               41

---

Example    $b$=4 bits,    $i$ =2,    2 keys/bucket

• insert 0101

| 0000 | 0101 | | | ← Future growth buckets |
|------|------|------|------|
| 1010 | 1111 | | | |
| 00 | 01 | 10 | 11 | |

$m$ = 01 (max used block)

Rule   If $h(k)[i] \leq m$, then
         look at bucket $h(k)[i]$
         else, look at bucket $h(k)[i] - 2^{i-1}$

CS 245                Notes 5               42

## Slide 43

Example   $b$=4 bits,   $i$ =2,   2 keys/bucket

0101

- insert 0101
- can have overflow chains!

← Future growth buckets

| 0000 | 0101 | | |
| 1010 | 1111 | | |
| 00 | 01 | 10 | 11 |

$m$ = 01 (max used block)

Rule  If h(k)[$i$] ≤ $m$, then
       look at bucket h(k)[i ]
       else, look at bucket h(k)[$i$] - $2^{i-1}$

CS 245          Notes 5          43

## Slide 44

### Note

- In textbook, n is used instead of m
- n=m+1

n=10

← Future growth buckets

| 0000 | 0101 | | |
| 1010 | 1111 | | |
| 00 | 01 | 10 | 11 |

$m$ = 01 (max used block)

CS 245          Notes 5          44

## Slide 45
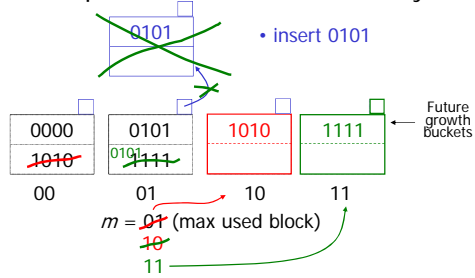
Example   $b$=4 bits,   $i$ =2,   2 keys/bucket

← Future growth buckets

| 0000 | 0101 | | |
| 1010 | 1111 | | |
| 00 | 01 | 10 | 11 |

$m$ = 01 (max used block)

CS 245          Notes 5          45

## Slide 46

Example   $b$=4 bits,   $i$ =2,   2 keys/bucket

← Future growth buckets

| 0000 | 0101 | 1010 | |
| ~~1010~~ | 1111 | | |
| 00 | 01 | 10 | 11 |

$m$ = ~~01~~ (max used block)
10

CS 245          Notes 5          46

## Slide 47

Example   $b$=4 bits,   $i$ =2,   2 keys/bucket

0101

- insert 0101

← Future growth buckets

| 0000 | 0101 | 1010 | |
| ~~1010~~ | 1111 | | |
| 00 | 01 | 10 | 11 |

$m$ = ~~01~~ (max used block)
10

CS 245          Notes 5          47

## Slide 48

Example   $b$=4 bits,   $i$ =2,   2 keys/bucket

0101

- insert 0101

← Future growth buckets

| 0000 | 0101 | 1010 | |
| ~~1010~~ | 1111 | | |
| 00 | 01 | 10 | 11 |

$m$ = ~~01~~ (max used block)
~~10~~
11

CS 245          Notes 5          48

### Slide 49

Example  $b=4$ bits,  $i=2$,  2 keys/bucket

0101

• insert 0101

| 0000 | 0101 | 1010 | 1111 |
| ~~1010~~ | 0101 0101 1111 | | |
| 00 | 01 | 10 | 11 |

Future growth buckets

$m =$ ~~01~~ (max used block)
~~10~~
11

CS 245    Notes 5    49

### Slide 50

Example Continued: How to grow beyond this?

$i = 2$

| 0000 | 0101 | 1010 | 1111 |
| | 0101 | | |
| 00 | 01 | 10 | 11 |

. . .

$m = 11$ (max used block)

CS 245    Notes 5    50

### Slide 51

Example Continued: How to grow beyond this?

$i = $ ~~2~~ 3

| 0000 | 0101 | 1010 | 1111 | | |
| | 0101 | | | | |
| 000 | 001 | 010 | 011 | | |
| 100 | 101 | 110 | 111 | | |

. . .

$m = 11$ (max used block)

CS 245    Notes 5    51

### Slide 52

Example Continued: How to grow beyond this?

$i = $ ~~2~~ 3

| 0000 | 0101 | 1010 | 1111 | | |
| | 0101 | | | 100 | |
| 000 | 001 | 010 | 011 | | |
| ~~100~~ | 101 | 110 | 111 | | |

. . .

$m = $ ~~11~~ (max used block)
100

CS 245    Notes 5    52

### Slide 53

Example Continued: How to grow beyond this?

$i = $ ~~2~~ 3

| 0000 | ~~0101~~ | 1010 | 1111 | | 0101 |
| | ~~0101~~ | | | | 0101 |
| 000 | 001 | 010 | 011 | 100 | 101 |
| ~~100~~ | ~~101~~ | 110 | 111 | | |

. . .

$m = $ ~~11~~ (max used block)
~~100~~
101

CS 245    Notes 5    53

### Slide 54

✉ When do we expand file?

• Keep track of:   $\dfrac{\text{# used slots}}{\text{total # of slots}} = U$

CS 245    Notes 5    54

⊠ When do we expand file?

- Keep track of: $\dfrac{\text{\# used slots}}{\text{total \# of slots}} = U$

- If U > threshold then increase $m$
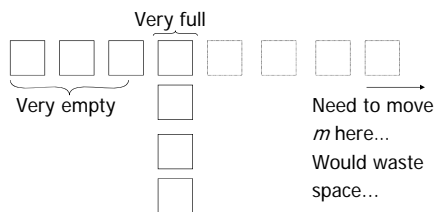  (and maybe $i$)

---

| Summary | Linear Hashing

⊕ Can handle growing files
  - with less wasted space
  - with no full reorganizations

⊕ No indirection like extensible hashing

⊖ Can still have overflow chains

---

Example: BAD CASE



Very full

Very empty

Need to move
$m$ here...
Would waste
space...

---

Summary

Hashing
  - How it works
  - Dynamic hashing
      - Extensible
      - Linear

---

Next:

- Indexing vs Hashing
- Index definition in SQL
- Multiple key access

---

Indexing vs Hashing

- Hashing good for probes given key
    e.g.,       SELECT ...
                FROM R
                WHERE R.A = 5

## Indexing vs Hashing

- INDEXING (Including B Trees) good for Range Searches:

  e.g.,    SELECT
           FROM R
           WHERE R.A > 5

## Index definition in SQL

- Create index name on rel (attr)
- Create unique index name on rel (attr)
            └──→ defines candidate key

- Drop INDEX name

---

Note  CANNOT SPECIFY TYPE OF INDEX
          (e.g. B-tree, Hashing, ...)
      OR PARAMETERS
          (e.g. Load Factor, Size of Hash,...)

      ... at least in SQL...

---

Note  ATTRIBUTE LIST $\Rightarrow$ MULTIKEY INDEX
                             (next)
      e.g., CREATE INDEX foo ON R(A,B,C)

---

## Multi-key Index

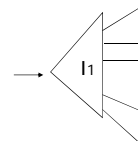Motivation: Find records where
            DEPT = "Toy" AND SAL > 50k

---

## Strategy I:

- Use one index, say Dept.
- Get all Dept = "Toy" records
         and check their salary

## Strategy II:

• Use 2 Indexes; Manipulate Pointers

Toy → [ | | | | ]    [ | | | | | | ] ← Sal
                                        > 50k

---

## Strategy III:

• Multiple Key Index

One idea:

---

## Example



10k
15k
17k
21k

Art
Sales
Toy

Dept
Index

12k
15k
15k
19k

Salary
Index

Example
Record

Name=Joe
DEPT=Sales
SAL=15k

---

For which queries is this index good?

☐ Find RECs Dept = "Sales" $\wedge$ SAL=20k
☐ Find RECs Dept = "Sales" $\wedge$ SAL $\geq$ 20k
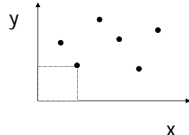☐ Find RECs Dept = "Sales"
☐ Find RECs SAL = 20k

---

## Interesting application:

• Geographic Data



DATA:
⟨X$_1$,Y$_1$, Attributes⟩
⟨X$_2$,Y$_2$, Attributes⟩
⋮

---

## Queries:

• What city is at ⟨Xi,Yi⟩?
• What is within 5 miles from ⟨Xi,Yi⟩?
• Which is closest point to ⟨Xi,Yi⟩?

## Queries

- Find points with $Y_i > 20$
- Find points with $X_i < 5$
- Find points "close" to $i = <12,38>$
- Find points "close" to $b = <7,24>$

---

- Many types of geographic index structures have been suggested
  - kd-Trees (very similar to what we described here)
  - Quad Trees
  - R Trees
  - ...

---

## Two more types of multi key indexes

- Grid
- Partitioned hash

---

## Grid Index



To records with key1=$V_3$, key2=$X_2$

---

## CLAIM

- Can quickly find records with
  - key 1 = $V_i$ ∧ Key 2 = $X_j$
  - key 1 = $V_i$
  - key 2 = $X_j$

---

## CLAIM

- Can quickly find records with
  - key 1 = $V_i$ ∧ Key 2 = $X_j$
  - key 1 = $V_i$
  - key 2 = $X_j$

- And also ranges....
  - E.g.,  key 1 ≥ $V_i$ ∧ key 2 < $X_j$

14

## Slide 85

• How do we find entry i,j in linear structure?

max number of i values N=4

i, j

| 0, 0 | ← position S+0 |
| 0, 1 | ← position S+1 |
| 0, 2 | ← position S+2 |
| 0, 3 | ← position S+3 |
| 1, 0 | ← position S+4 |
| 1, 1 | |
| 1, 2 | |
| 1, 3 | |
| 2, 0 | |
| 2, 1 | ← position S+9 |
| 2, 2 | |
| 2, 3 | |
| 3, 0 | |

$pos(i, j) =$

## Slide 86

• How do we find entry i,j in linear structure?
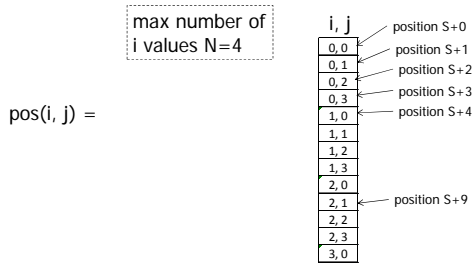
max number of i values N=4

i, j

| 0, 0 | ← position S+0 |
| 0, 1 | ← position S+1 |
| 0, 2 | ← position S+2 |
| 0, 3 | ← position S+3 |
| 1, 0 | ← position S+4 |
| 1, 1 | |
| 1, 2 | |
| 1, 3 | |
| 2, 0 | |
| 2, 1 | ← position S+9 |
| 2, 2 | |
| 2, 3 | |
| 3, 0 | |

$pos(i, j) = S + iN + j$

Issue: Cells must be same size, and N must be constant!

⬇

Issue: Some cells may overflow, some may be sparse…

## Slide 87

Solution: Use Indirection

Buckets

X1  X2  X3

$V_1$
$V_2$
$V_3$
$V_4$

*Grid only contains pointers to buckets

Buckets

## Slide 88

With indirection:

• Grid can be regular without wasting space
• We do have price of indirection

## Slide 89

Can also index grid on value ranges

Salary

| 0-20K | 1 |
| 20K-50K | 2 |
| 50K- ∞ | 3 |

Grid

Linear Scale

| 1 | 2 | 3 |
| Toy | Sales | Personnel |

## Slide 90

Grid files

⊕ Good for multiple-key search
⊝ Space, management overhead
    (nothing is free)
⊝ Need partitioning ranges that evenly split keys

15

## Slide 91

Partitioned hash function

Idea:

010110 1110010

Key1 → (h1)   (h2) ← Key2

CS 245                Notes 5                91

## Slide 92

EX:

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | |
| h1(sales) | =1 | 001 | |
| h1(art) | =1 | 010 | |
| . | | 011 | |
| h2(10k) | =01 | 100 | |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | |
| h2(40k) | =00 | 111 | |
| . | | | |

Insert ⇒ <Fred,toy,10k>,<Joe,sales,10k>
<Sally,art,30k>

CS 245                Notes 5                92

## Slide 93

EX:

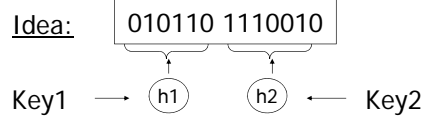| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | |
| h1(sales) | =1 | 001 | <Fred> |
| h1(art) | =1 | 010 | |
| . | | 011 | |
| h2(10k) | =01 | 100 | |
| h2(20k) | =11 | 101 | <Joe><Sally> |
| h2(30k) | =01 | 110 | |
| h2(40k) | =00 | 111 | |
| . | | | |

Insert ⇒ <Fred,toy,10k>,<Joe,sales,10k>
<Sally,art,30k>

CS 245                Notes 5                93

## Slide 94

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | <Fred> |
| h1(sales) | =1 | 001 | <Joe><Jan> |
| h1(art) | =1 | 010 | <Mary> |
| . | | 011 | |
| h2(10k) | =01 | 100 | <Sally> |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | <Tom><Bill> |
| h2(40k) | =00 | 111 | <Andy> |
| . | | | |

• Find Emp. with Dept. = Sales ∧ Sal=40k

CS 245                Notes 5                94

## Slide 95

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | <Fred> |
| h1(sales) | =1 | 001 | <Joe><Jan> |
| h1(art) | =1 | 010 | <Mary> |
| . | | 011 | |
| h2(10k) | =01 | 100 | <Sally> |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | <Tom><Bill> |
| h2(40k) | =00 | 111 | <Andy> |
| . | | | |

• Find Emp. with Dept. = Sales ∧ Sal=40k

CS 245                Notes 5                95

## Slide 96

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | <Fred> |
| h1(sales) | =1 | 001 | <Joe><Jan> |
| h1(art) | =1 | 010 | <Mary> |
| . | | 011 | |
| h2(10k) | =01 | 100 | <Sally> |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | <Tom><Bill> |
| h2(40k) | =00 | 111 | <Andy> |
| . | | | |

• Find Emp. with Sal=30k

CS 245                Notes 5                96

## Slide 97

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | <Fred> |
| h1(sales) | =1 | 001 | <Joe><Jan> |
| h1(art) | =1 | 010 | <Mary> |
| . | | 011 | |
| h2(10k) | =01 | 100 | <Sally> |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | <Tom><Bill> |
| h2(40k) | =00 | 111 | <Andy> |
| . | | | |

- Find Emp. with Sal=30k

*look here*

## Slide 98

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | <Fred> |
| h1(sales) | =1 | 001 | <Joe><Jan> |
| h1(art) | =1 | 010 | <Mary> |
| . | | 011 | |
| h2(10k) | =01 | 100 | <Sally> |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | <Tom><Bill> |
| h2(40k) | =00 | 111 | <Andy> |
| . | | | |

- Find Emp. with Dept. = Sales

## Slide 99

| | | | |
|---|---|---|---|
| h1(toy) | =0 | 000 | <Fred> |
| h1(sales) | =1 | 001 | <Joe><Jan> |
| h1(art) | =1 | 010 | <Mary> |
| . | | 011 | |
| h2(10k) | =01 | 100 | <Sally> |
| h2(20k) | =11 | 101 | |
| h2(30k) | =01 | 110 | <Tom><Bill> |
| h2(40k) | =00 | 111 | <Andy> |
| . | | | |

- Find Emp. with Dept. = Sales    *look here*

## Slide 100

Summary

### Post hashing discussion:

- Indexing vs. Hashing
- SQL Index Definition
- Multiple Key Access
  - Multi Key Index
    - Variations: Grid, Geo Data
  - Partitioned Hash

## Slide 101

### Reading Chapter 5

- Skim the following sections:
  - Sections 14.3.6, 14.3.7, 14.3.8
    [Second Ed: 14.6.6, 14.6.7, 14.6.8]
  - Sections 14.4.2, 14.4.3, 14.4.4
    [Second Ed: 14.7.2, 14.7.3, 14.7.4]
- Read the rest

## Slide 102

The BIG picture....

- Chapters 11 & 12 [13]: Storage, records, blocks...
- Chapters 13 & 14 [14]: Access Mechanisms
  - Indexes
  - B trees
  - Hashing
  - Multi key
- Chapters 15 & 16 [15, 16]: Query Processing

NEXT