# CS 525: Advanced Database Organization

## 05: Hashing and More

Boris Glavic

Slides: adapted from a course taught by
Hector Garcia-Molina, Stanford InfoLab

---

Hashing
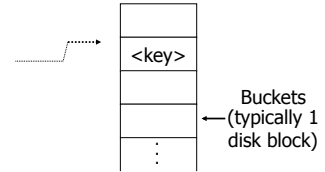
key → h(key)                <key>

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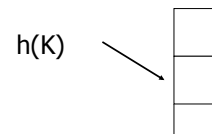
➼ This may not be best function …
➼ Read Knuth Vol. 3 if you really
  need to select a good function.

Good hash          ☞ Expected number of
function:              keys/bucket is the
                       same for all buckets
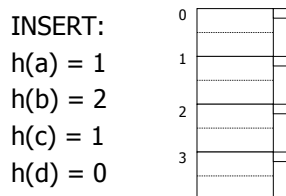
<u>Within a bucket:</u>

• Do we keep keys sorted?

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

<u>Next:</u> example to illustrate
                    inserts,
overflows, deletes

h(K)

<u>EXAMPLE</u>  2 records/bucket

INSERT:       0
h(a) = 1      1
h(b) = 2      2
h(c) = 1
h(d) = 0      3

<u>EXAMPLE</u>  2 records/bucket

INSERT:       0    d
h(a) = 1      1    a
h(b) = 2           c
h(c) = 1      2    b
h(d) = 0      3
h(e) = 1

## EXAMPLE  2 records/bucket

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

## EXAMPLE:  deletion

Delete:
e
f

## EXAMPLE:  deletion

Delete:
e
f
c



maybe move "g" up

## EXAMPLE:  deletion

Delete:
e
f
c



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
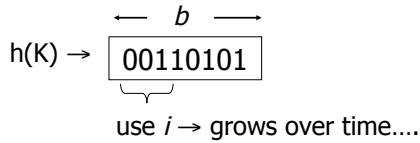  function is & on # keys/bucket

## 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

$$\longleftarrow b \longrightarrow$$

$h(K) \rightarrow$ | 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 = $ | 1 |

| 1 |
| 0001 |

| 1 |
| 1001 |
| 1100 |

Insert 1010

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

$i = $ | 1 |

| 1 |
| 0001 |

| 1 |
| 1001 |
| 1010 | 1100 |

| 1 |
| 1100 |

Insert 1010

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

*i* = 1

1
0001

1 2
1001
1010 1100

Insert 1010

1 2
1100

*i* = 2

00
01
10
11

New directory

CS 525    COMPUTER SCIENCE    Notes 5 - Hashing    25    IIT College of Science and Letters    ILLINOIS INSTITUTE OF TECHNOLOGY

---

## Example continued

i = 2
00
01
10
11

1
0001

2
1001
1010

2
1100

Insert:
0111
0000

CS 525    COMPUTER SCIENCE    Notes 5 - Hashing    26    IIT College of Science and Letters    ILLINOIS INSTITUTE OF TECHNOLOGY

---

## Example continued

0000
0001

i = 2
00
01
10
11

1
0001 0111
0111

2
1001
1010

2
1100

Insert:
0111
0000

CS 525    COMPUTER SCIENCE    Notes 5 - Hashing    27    IIT College of Science and Letters    ILLINOIS INSTITUTE OF TECHNOLOGY

---

## Example continued

2
0000
0001

i = 2
00
01
10
11

1 2
0001 0111
0111

2
1001
1010

2
1100

Insert:
0111
0000

CS 525    COMPUTER SCIENCE    Notes 5 - Hashing    28    IIT College of Science and Letters    ILLINOIS INSTITUTE OF TECHNOLOGY

---

## Example continued

0000  2
0001

*i* = 2
00
01
10
11

0111  2

1001  2
1010

Insert:
1001

1100  2

CS 525    COMPUTER SCIENCE    Notes 5 - Hashing    29    IIT College of Science and Letters    ILLINOIS INSTITUTE OF TECHNOLOGY

---

## Example continued

0000  2
0001

*i* = 2
00
01
10
11

0111  2

1001
1001

1010 1001  2
1010

Insert:
1001

1100  2

CS 525    COMPUTER SCIENCE    Notes 5 - Hashing    30    IIT College of Science and Letters    ILLINOIS INSTITUTE OF TECHNOLOGY

5

$i = \boxed{2}$

| 00 |
| 01 |
| 10 |
| 11 |

$i = \boxed{3}$

```
0000  2
0001

0111  2

1001  3
1001

1010  1001  2→3
      1010

1100  2
```

000
001
010
011
100
101
110
111

Insert:
1001

---

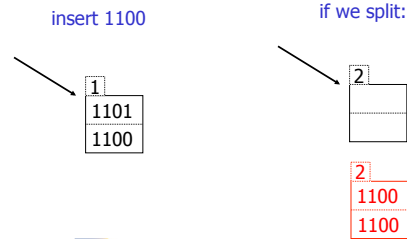## Extensible hashing: __deletion__

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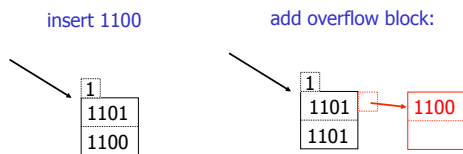
---

## Deletion example:

- Run thru insert example in reverse!

---

## Note: Still need overflow chains

- Example: many records with duplicate keys

insert 1100

```
1
1101
1100
```

if we split:

```
2
```

```
2
1100
1100
```

---

## Solution: overflow chains

insert 1100

```
1
1101
1100
```

add overflow block:

```
1
1101  →  1100
1101
```

---

| Summary | Extensible hashing |

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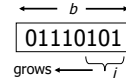  - with no full reorganizations

Summary    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)

---

## Linear hashing

• Another dynamic hashing scheme

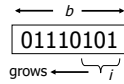Two ideas:

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

$\longleftarrow b \longrightarrow$

01110101

grows $\longleftarrow$ $i$

---

## Linear hashing

• Another dynamic hashing scheme

Two ideas:

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

$\longleftarrow b \longrightarrow$

01110101

grows $\longleftarrow$ $i$

(b) File grows linearly

$\longrightarrow$

---

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)

---

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

---

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

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

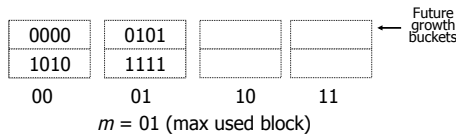CS 525  Notes 5 - Hashing  43  IIT College of Science and Letters · ILLINOIS INSTITUTE OF TECHNOLOGY

---

## Note

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

n=10

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

← Future growth buckets

$m$ = 01 (max used block)

CS 525  Notes 5 - Hashing  44  IIT College of Science and Letters · ILLINOIS INSTITUTE OF TECHNOLOGY

---

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

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

← Future growth buckets

$m$ = 01 (max used block)

CS 525  Notes 5 - Hashing  45  IIT College of Science and Letters · ILLINOIS INSTITUTE OF TECHNOLOGY

---

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

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

← Future growth buckets

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

CS 525  Notes 5 - Hashing  46  IIT College of Science and Letters · ILLINOIS INSTITUTE OF TECHNOLOGY

---

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

| 0101 |
|------|

• insert 0101

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

← Future growth buckets

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

CS 525  Notes 5 - Hashing  47  IIT College of Science and Letters · ILLINOIS INSTITUTE OF TECHNOLOGY

---

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

| 0101 |
|------|

• insert 0101

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

← Future growth buckets

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

CS 525  Notes 5 - Hashing  48  IIT College of Science and Letters · ILLINOIS INSTITUTE OF TECHNOLOGY

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

0101

• insert 0101

Future growth buckets

| 0000 | 0101 | 1010 | 1111 |
| 1010 | 0101 1111 | | |

00      01      10      11

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

## Example Continued: How to grow beyond this?

$i = 2$

| 0000 | 0101 | 1010 | 1111 |
| | 0101 | | |

00      01      10      11       . . .

$m$ = 11 (max used block)

## 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)

## Example Continued: How to grow beyond this?

$i = 2$ 3

| 0000 | 0101 | 1010 | 1111 | | |
| | 0101 | | | | |

000    001    010    011    100
100    101    110    111      . . .

$m$ = 11 (max used block)
100

## 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

☛ When do we expand file?

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

☛ When do we expand file?

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

• 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
-> **Point Queries**

## Indexing vs Hashing

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

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

-> **Range Queries**

## 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 standard SQL...

    Vendor specific extensions allow
that

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



| | |
|---|---|
| Art | |
| Sales | |
| Toy | |

Dept Index

10k
15k
17k
21k

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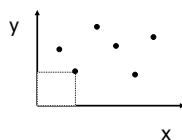☐ Find RECs Dept = "Sales"
☐ Find RECs SAL = 20k

## Interesting application:

• Geographic Data



DATA:
$\langle X_1, Y_1, \text{Attributes} \rangle$
$\langle X_2, Y_2, \text{Attributes} \rangle$
⋮

## Queries:

• What city is at $\langle X_i, Y_i \rangle$?
• What is within 5 miles from $\langle X_i, Y_i \rangle$?
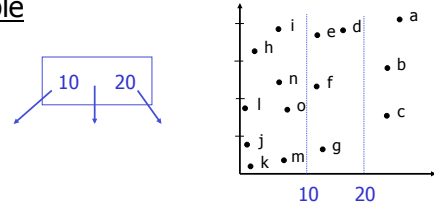• Which is closest point to $\langle X_i, Y_i \rangle$?

## Example

## Example

## Example

## Example

## Example

## Example



- Search points near f
- Search points near b

## Queries

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

## Next

- Even more index structures ☺