Lecture 19:        Multikey Files

Last few days: Indexed Sequential Files

- Binary & Multiway Trees

- B trees

- B$^+$ trees

Today:   Multikey Files
- Secondary Keys

- Inverted Lists

Folk & Zoellick,   ch 6.5 - 6.7

# Multikey File Organization

- Direct files and Indexed Files support efficient data access by a <u>single</u> Field (the key).

- eg, Given an index on Employee Number, the query <u>"Retrieve employee #2317"</u> is fast,

- But, "Retrieve <u>all employees living in Toronto</u>" is slow, since we must scan the entire file.

- Multi key files support fast access by <u>several</u> different fields,

**Problem:** Given <u>non-key</u> field values, how can we <u>quickly</u> find all records having those field values? (ie, without scanning the entire file.)

**Solution:** Secondary Key Indices

- A secondary key is a field of a file that is indexed but is not the primary key.

- There are two main secondary index structur
  - Inverted Lists ⟵ today
  - Multi lists

# Terminology

- The _active values_ for a secondary key are the values that the field has in the data file.

## Inverted Lists have two parts:

- The _directory_ for the secondary key, which has an _entry_ for each active value,

- The _accession list_ for an active value is a list of pointers to the data file.
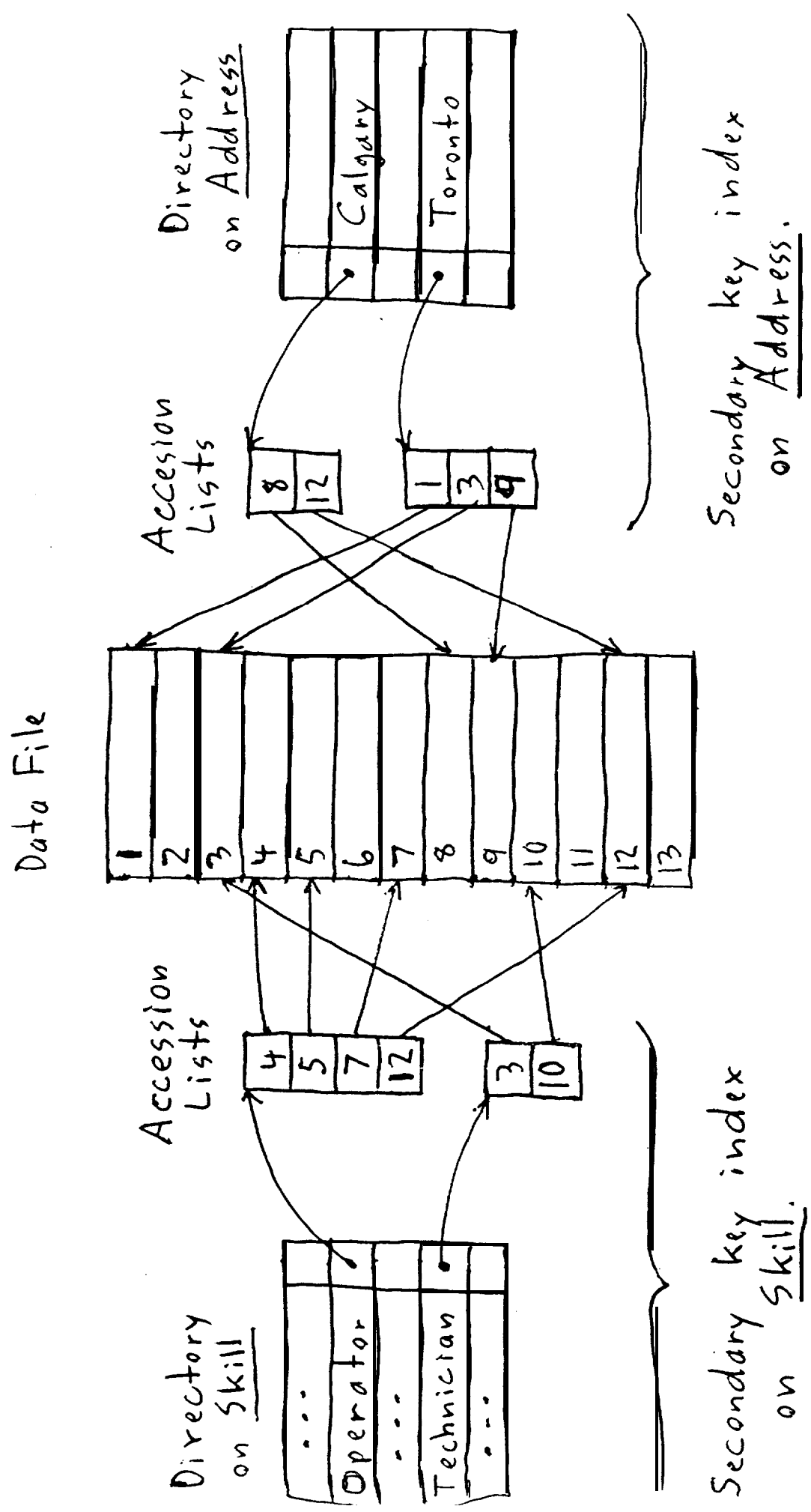
# Secondary Key Example

| loc | E# | NAME | ADDRESS | AGE | SEX | SKILL |
|-----|-----|------|---------|-----|-----|-------|
| 1 | 001 | Hicks | Toronto | 36 | M | Programmer |
| 2 | 020 | McLeod | Montreal | 51 | M | Analyst |
| 3 | 023 | Lucas | Toronto | 25 | F | Technician |
| 4 | 025 | Bradley | Ottawa | 35 | F | Operator |
| 5 | 030 | Date | Montreal | 45 | M | Operator |
| 6 | 045 | Loomis | Vancouver | 45 | F | Analyst |
| 7 | 046 | Mader | Edmonton | 38 | M | Operator |
| 8 | 048 | Wu | Calgary | 50 | F | Programmer |
| 9 | 055 | Bair | Toronto | 28 | M | Analyst |
| 10 | 060 | Uhlig | Vancouver | 24 | M | Technician |
| 11 | 062 | Orilia | Montreal | 21 | M | Designer |
| 12 | 070 | Fry | Calgary | 34 | F | Operator |
| 13 | 075 | Riley | Ottawa | 40 | F | Designer |

| ADDRESS | loc |
|---------|-----|
| Calgary | 8, 12 |
| Edmonton | 7 |
| Montreal | 2, 5, 11 |
| Ottawa | 4, 13 |
| Toronto | 1, 3, 9 |
| Vancouver | 6, 10 |

| SKILL | loc |
|-------|-----|
| Analyst | 2, 6, 9 |
| Designer | 11, 13 |
| Operator | 4, 5, 7, 12 |
| Programmer | 1, 8 |
| Technician | 3, 10 |

Key is E#

Secondary keys are ADDRESS, SKILL

# Inverted Lists

14-6

## Directory on Address

| | |
|---|---|
| • | Calgary |
| • | Toronto |

## Accession Lists

| 8 | 12 |
|---|---|

| 1 | 3 | 9 |
|---|---|---|

## Data File

| 1 |
|---|
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |
| 7 |
| 8 |
| 9 |
| 10 |
| 11 |
| 12 |
| 13 |

## Accession Lists

| 4 | 5 | 7 | 12 |
|---|---|---|---|

| 3 | 10 |
|---|---|

## Directory on Skill

| ... | |
|---|---|
| • | Operator |
| ... | |
| • | Technician |
| ... | |

Secondary key index on Address.

Secondary key index on Skill.

Note:

- A secondary key of the data file is a primary key of its directory.

- Directories may be implemented in many ways. eg, as arrays in main memory (if small enough), as flat files, as B trees as hash files, etc.

# Retrieval : Example

— Retrieve all records of operators:

(1) Find Operator entry in the Skill directory

(2) Follow the pointer to the accession list.

(3) Follow each pointer in the accession list to find all operator records in the data file.

— Retrieve all records of employees living in Toronto. (exercise)

# Insertion: Example

Insert a record with

$$E\# = 14$$

Name = Marvin

Address = Toronto

Age = 25

Sex = M

Skill = Technician

Note: All secondary indexes must be updated.

<u>Note:</u> How the record is inserted into the data file depends on how the data file is organized.

- If the data file is serial & unsorted, then just append the new record to the file.

- If the data file is sorted (by E#) then add the record to a differential file

- If there is a B-tree on E#, then insert the record into the B-tree.

After inserting the record into the data file, the secondary indexes must be updated.

# Insertion Algorithm

First, insert the new record into the data file.

Then, for each field, F, with a secondary index,

- Let $v$ be the value of field F in the new record

- Insert $v$ into the directory for F (if it is not already there)

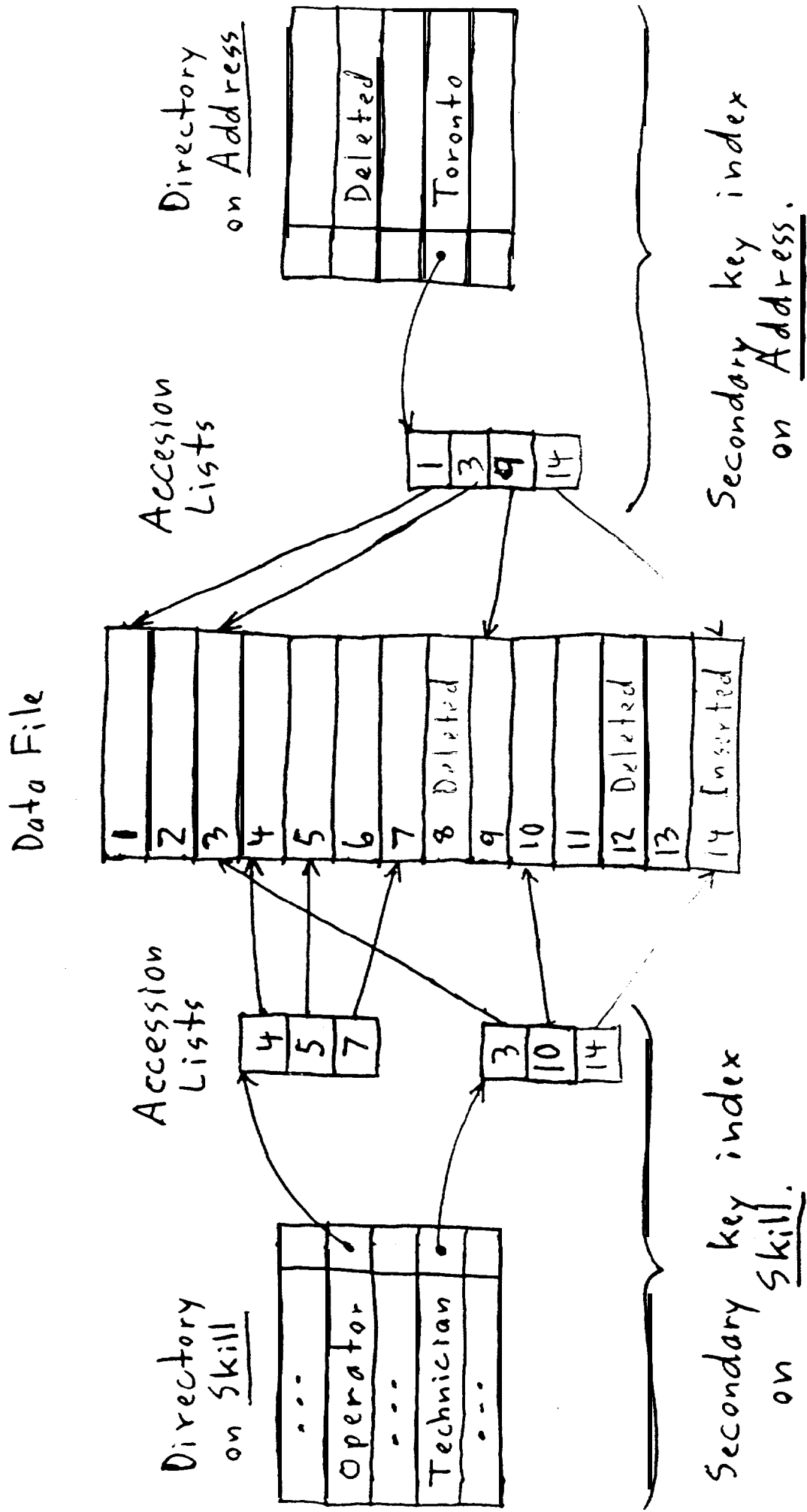- Insert a pointer to the new record into the accession list for $v$.

# Deletion

**Examples:**

- Delete record 12

- Delete record 9 ( Accession list for Calgary becomes empty).

Note:

- Again, all secondary indexes must be updated.

- Again, how a record is deleted from the data file depends on how the data file is organized.

14-13

Inverted Lists

Directory on Address

| | Deleted | | Toronto | |

Accession Lists

| 1 | 3 | 9 | 14 |

Data File

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 Deleted | 9 | 10 | 11 | 12 Deleted | 13 | 14 Inserted |

Accession Lists

| 4 | 5 | 7 |

| 3 | 10 | 14 |

Directory on Skill

| ... | Operator | ... | Technician | ... |

Secondary key index on Skill.

Secondary key index on Address.

Final File Organization

# Deletion Algorithm

First, delete the record from the data file.

Then, for each field, F, with a secondary index

- Let $v$ be the value of field F in the deleted record.

- Remove the pointer to the deleted record from the accession list for $v$.

- If the accession list for $v$ is now empty, delete $v$ from the directory for F.