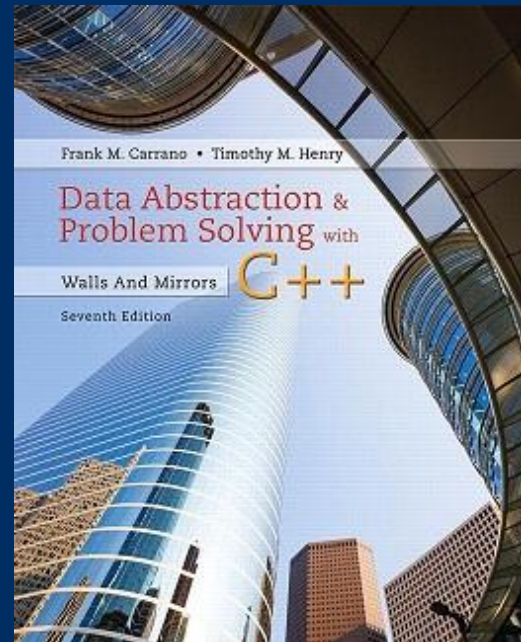


Processing Data in External Storage

CS 302 - Data Structures

M. Abdullah Canbaz

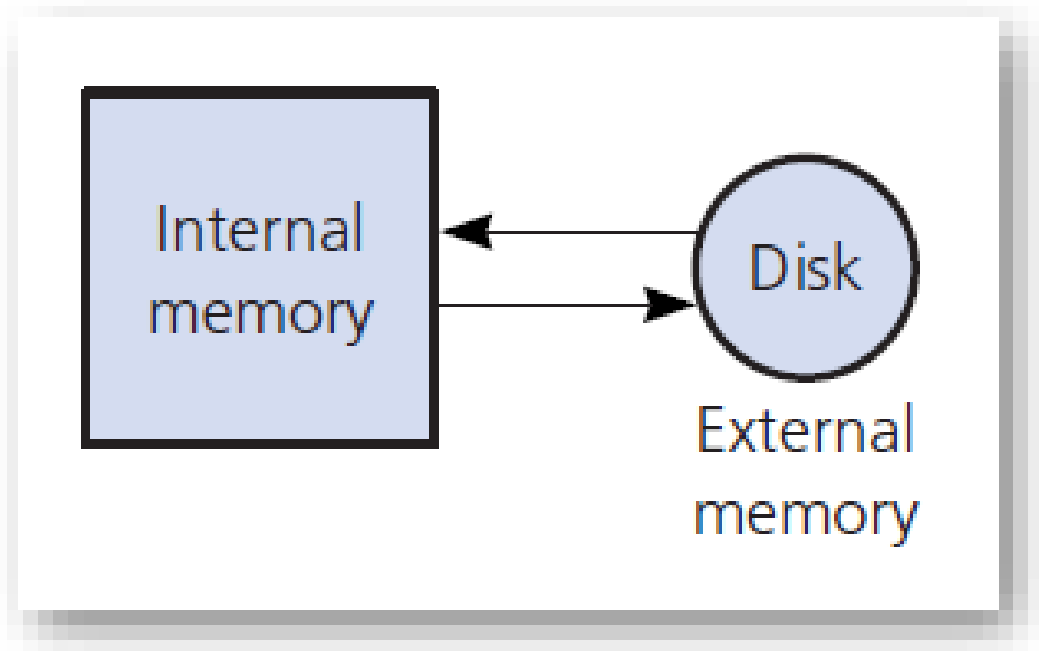




Look at External Storage

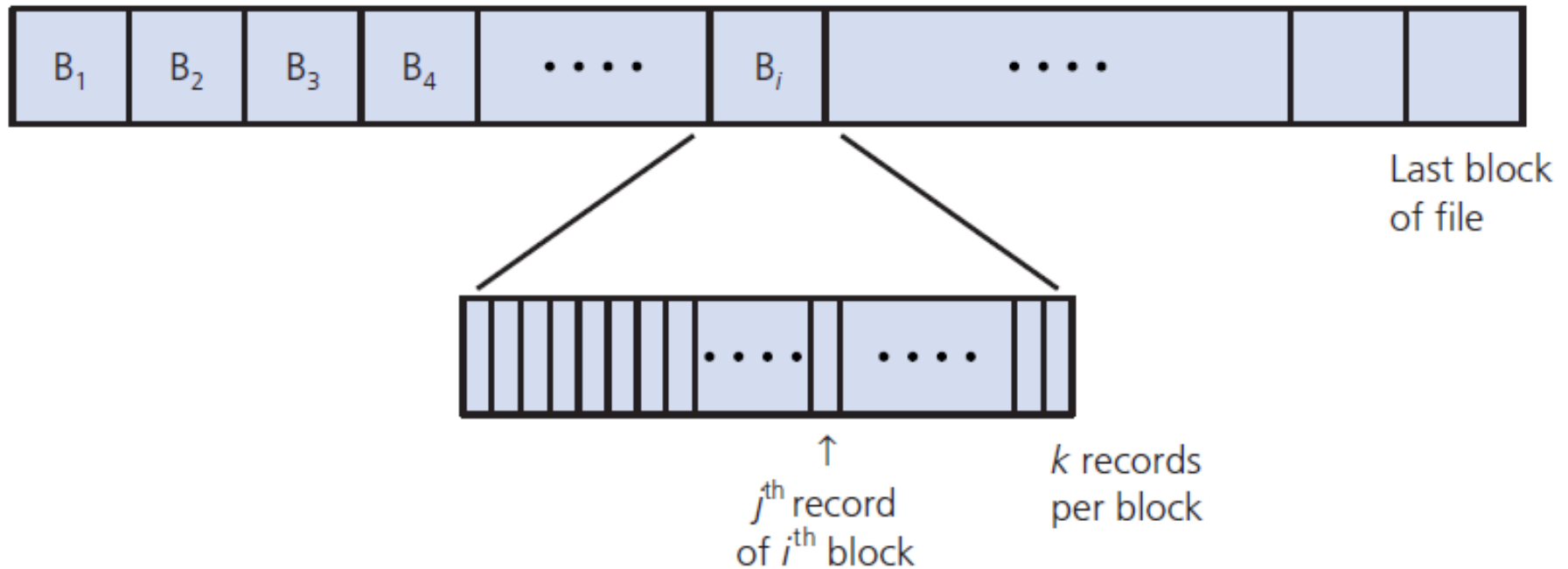
- External storage
 - Used when program reads/writes data to/from a C++ file
- Generally there is more external storage than internal memory
- Direct access files essential for external data collections

- Internal and external memory



Look at External Storage

- A file partitioned into blocks of records

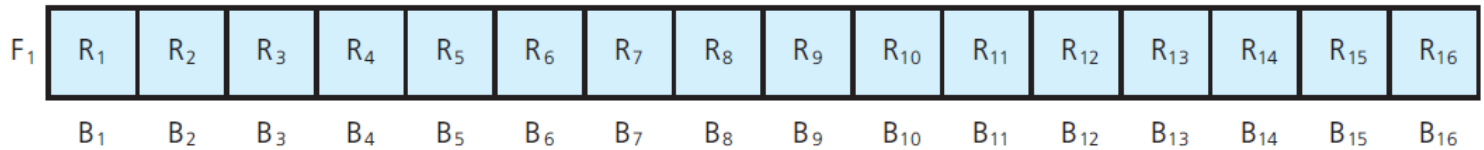




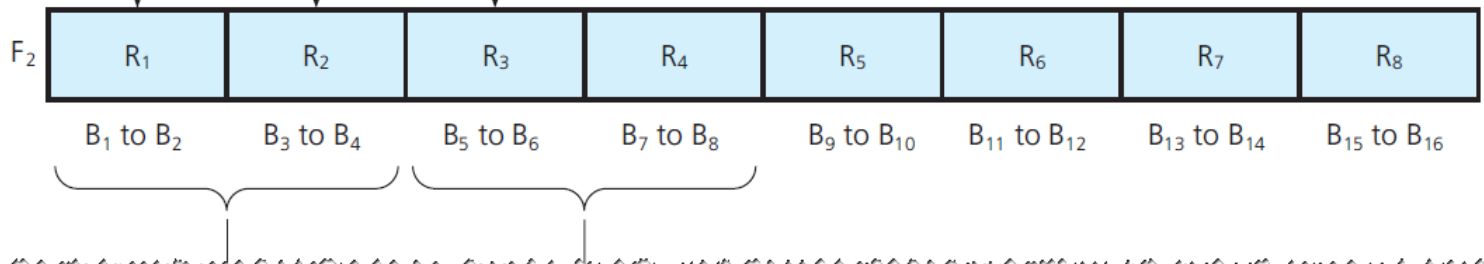
Look at External Storage

- Direct access I/O
 - Involves blocks instead of records
- Buffer stores data (blocks) temporarily
- Record updated within block (in buffer)
- Work to minimize block I/O
 - Takes more time for disk access

(a) 16 sorted runs,
1 block each,
in file F_1



(b) 8 sorted runs,
2 blocks each,
in file F_2

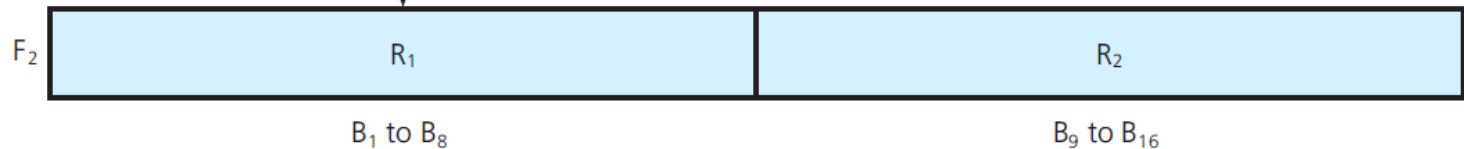


- Sorting a block of an external file F by merging the results of internal sorts and using two external work files F_1 and F_2

(c) 4 sorted runs,
4 blocks each,
in file F_1

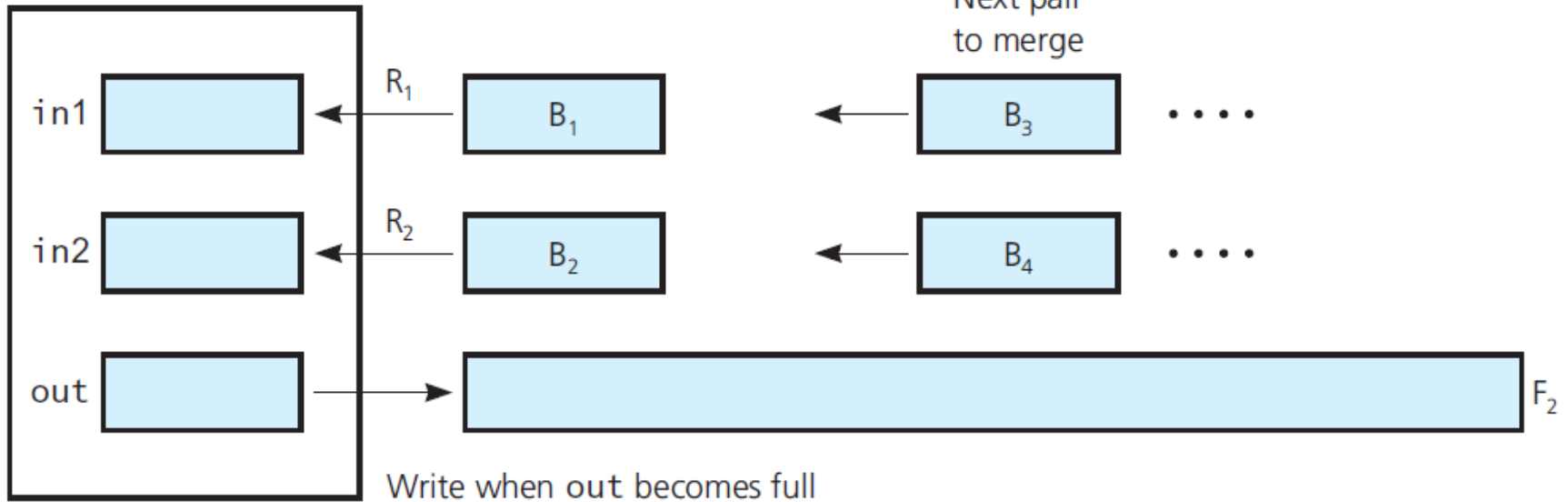


(d) 2 sorted runs,
8 blocks each,
in file F_2



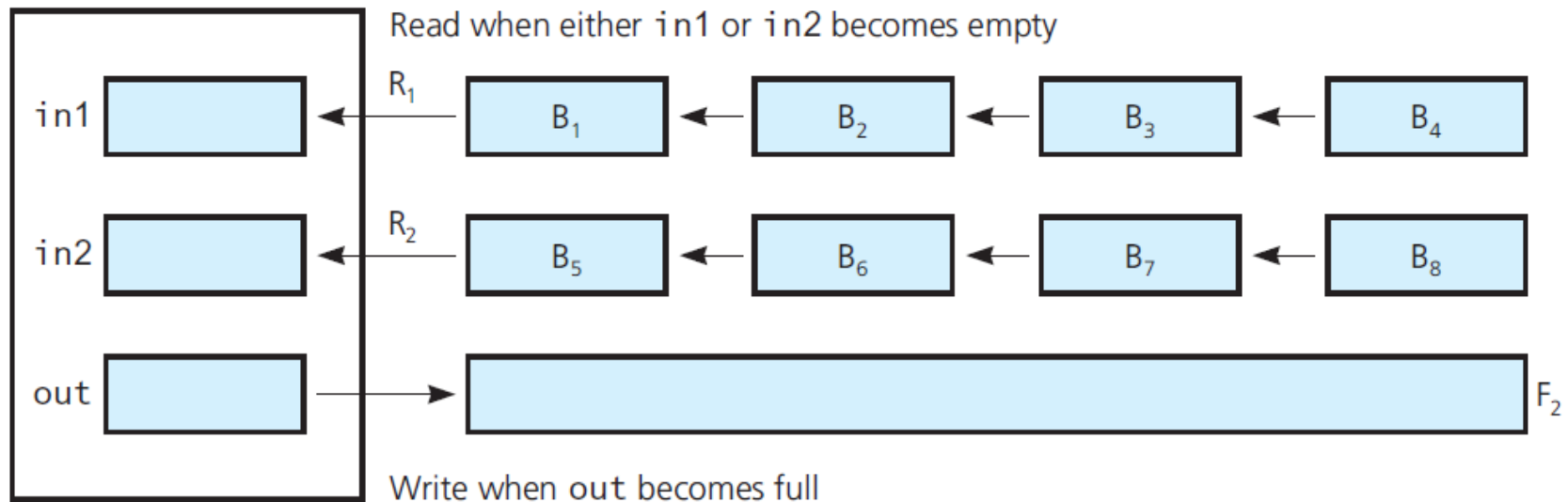
- Sorting a block of an external file F by merging the results of internal sorts and using two external work files F_1 and F_2

(a) Merging single blocks



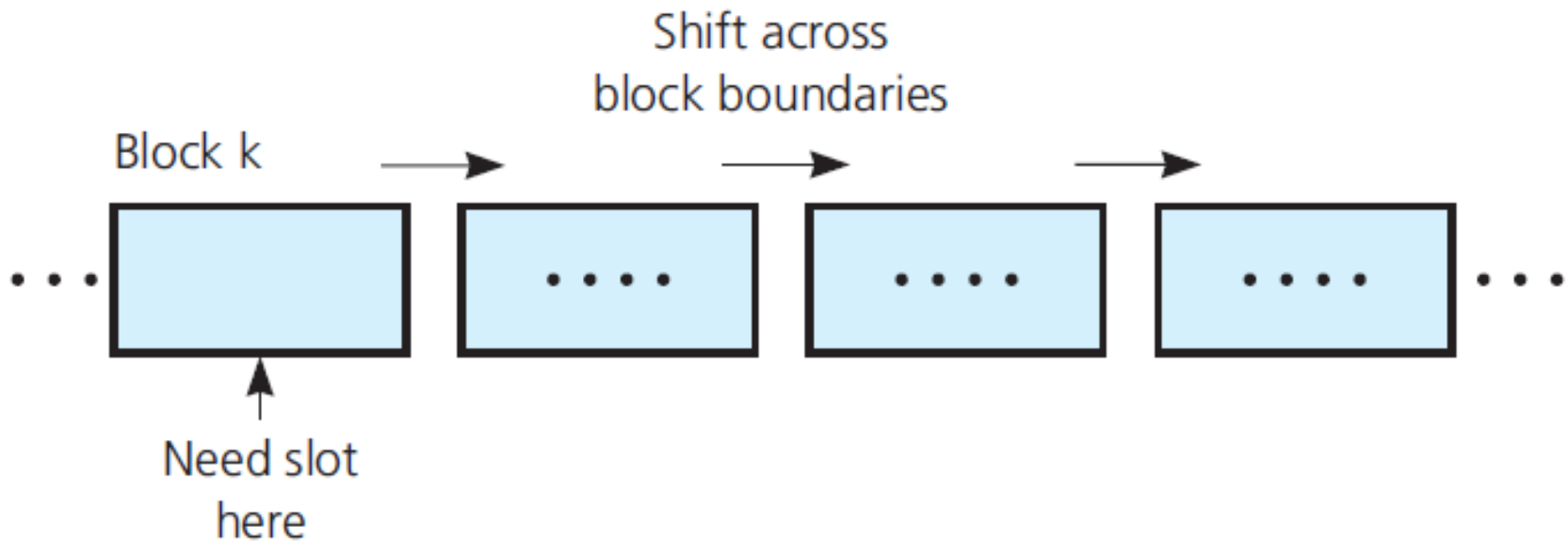
- Phase 2 of an external sort: Merging sorted runs

(b) Merging long runs



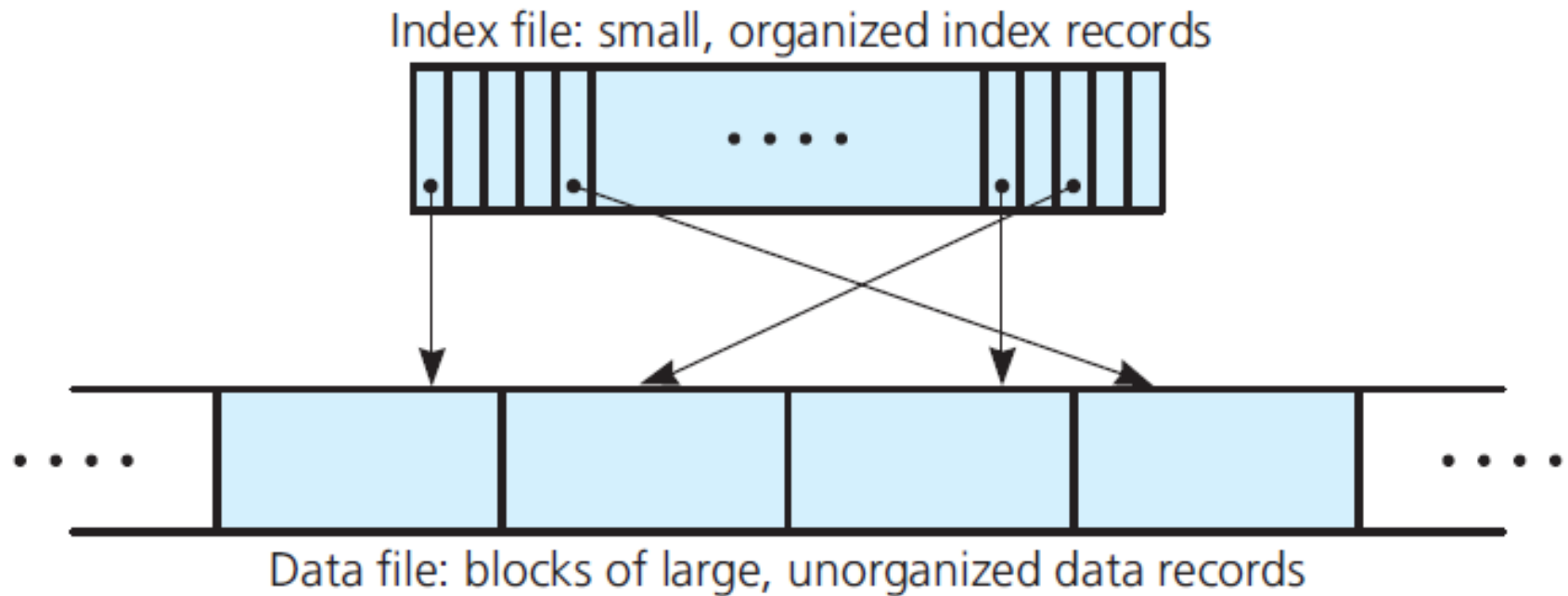
- Phase 2 of an external sort: Merging sorted runs

- Shifting across block boundaries



Indexing an External File

- A data file with an index





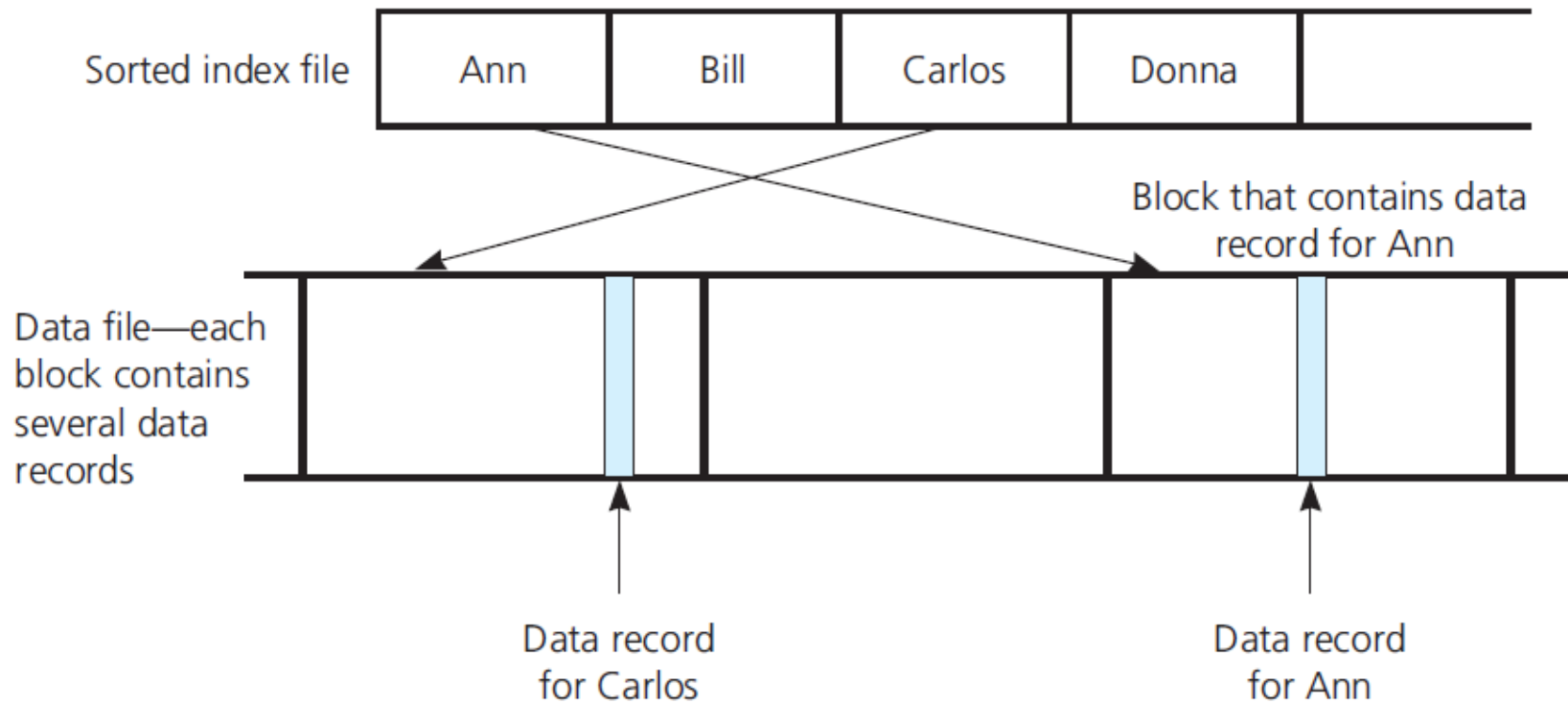
Indexing an External File

Advantages of an index file

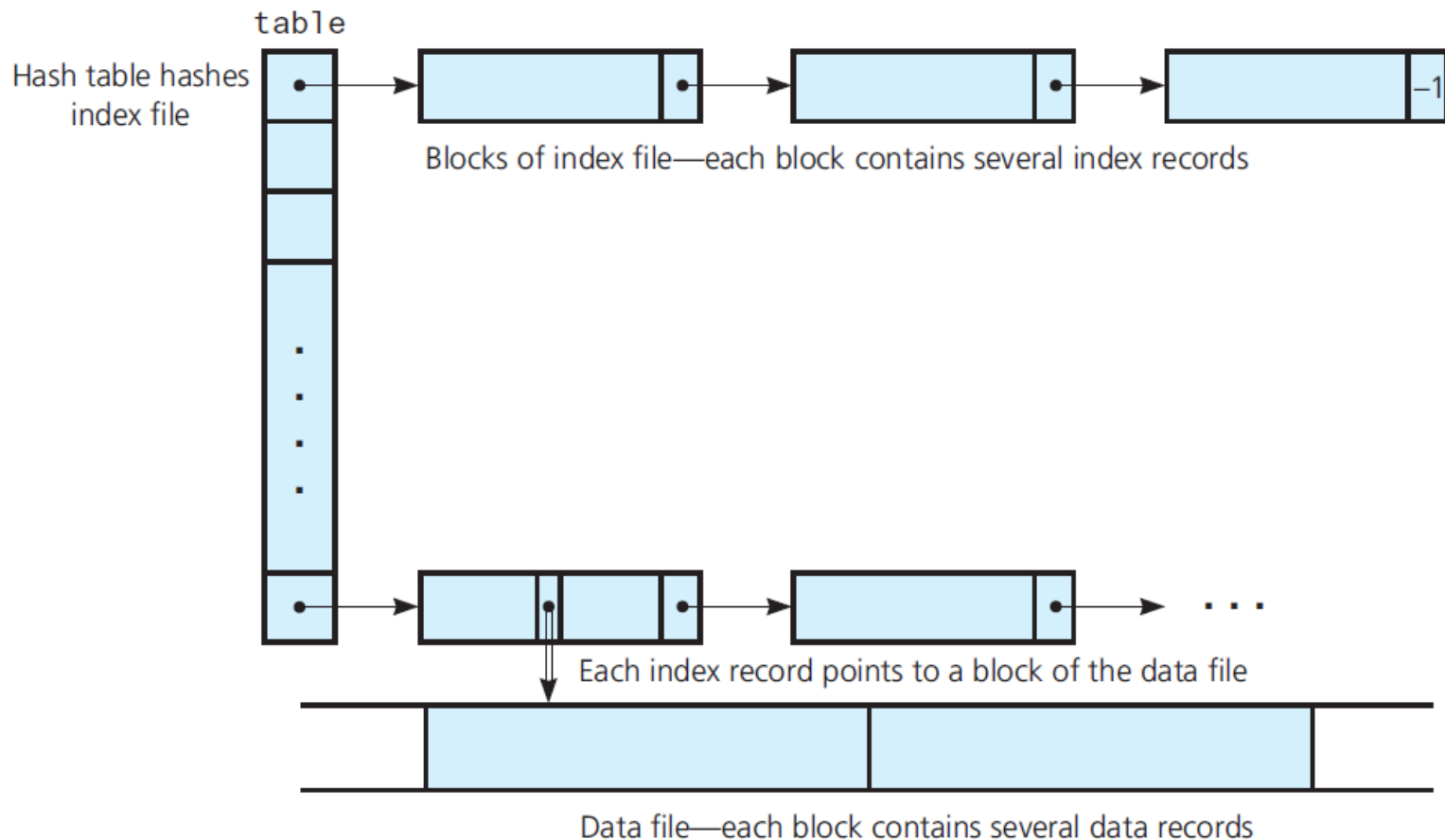
- Index record smaller than a data record
- Data file need not be kept in any particular order
- Possible to maintain several indexes simultaneously
- Index file reduces number of block accesses
- Shift index records, not data records

Indexing an External File

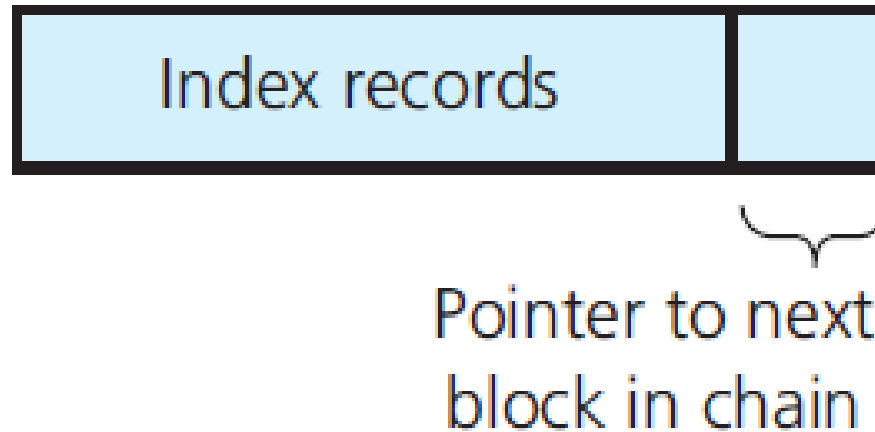
- A data file with a sorted index file



- A hashed index file



- A single block with a pointer





External Hashing

Addition when an index file uses external hashing

1. Add data record into data file.
2. Add corresponding index record into index file

$$i = h(\textit{searchKey})$$



External Hashing

Removal when an index file uses external hashing

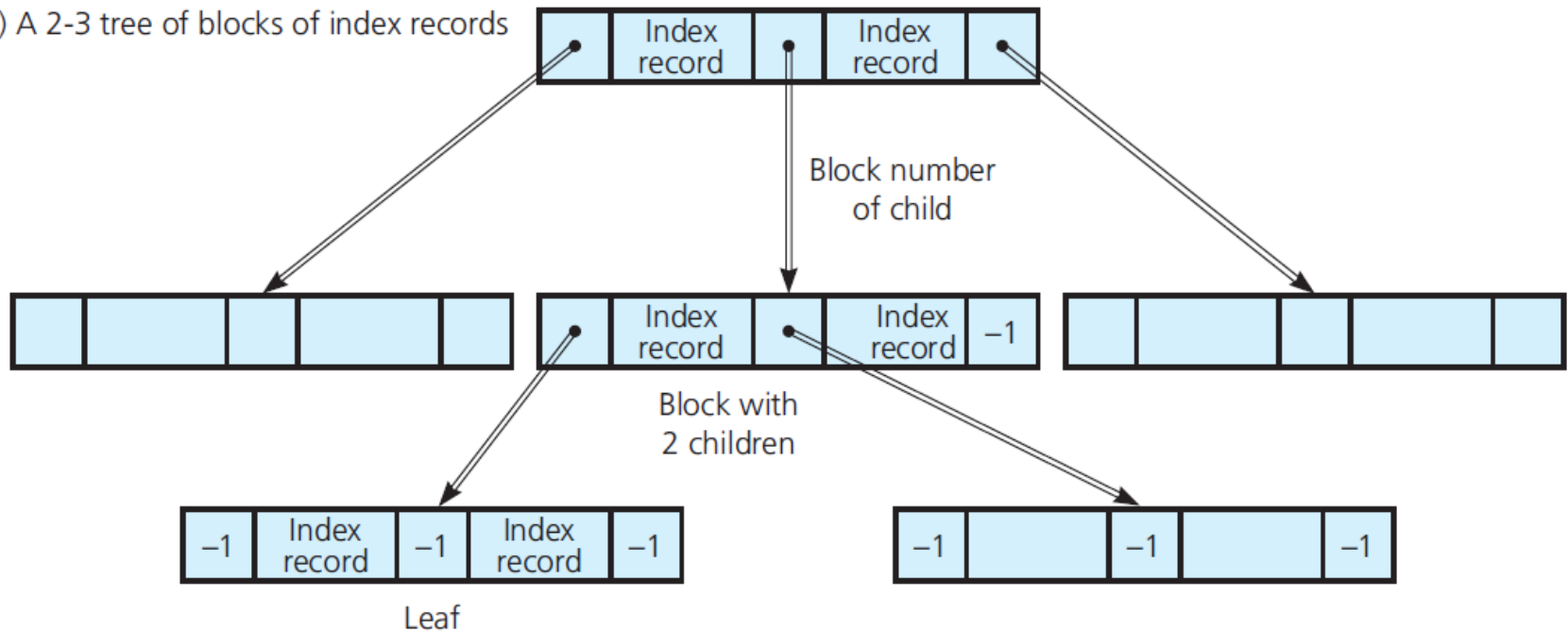
1. Search index file for corresponding index record

$$i = h(\textit{searchKey})$$

2. Remove data record from data file

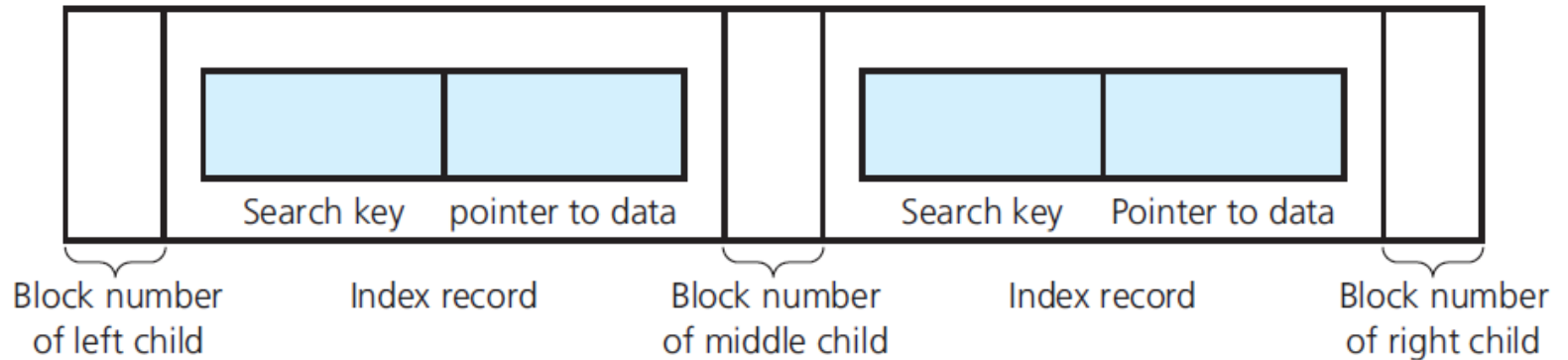
- An index file organized as a 2-3 tree

(a) A 2-3 tree of blocks of index records

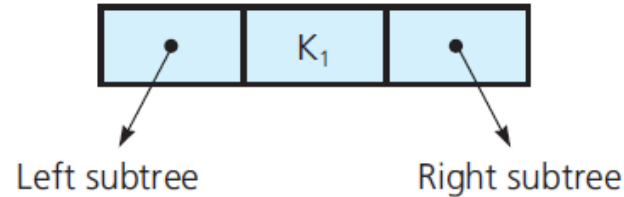


- An index file organized as a 2-3 tree

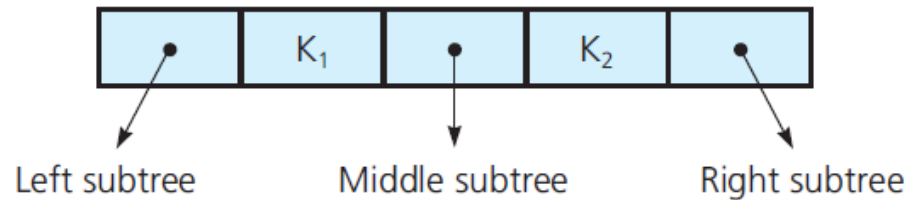
(b) The detail of a single block



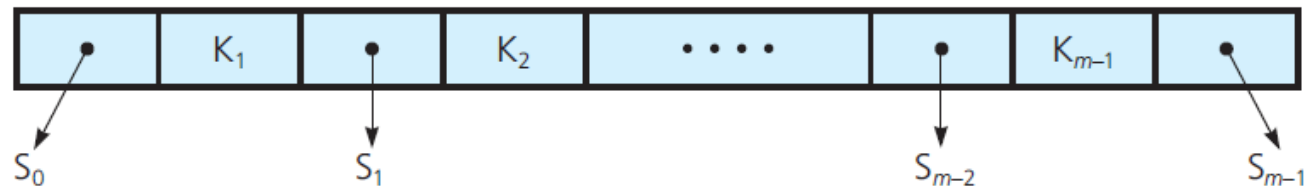
- (a) A node with two children
has one search key



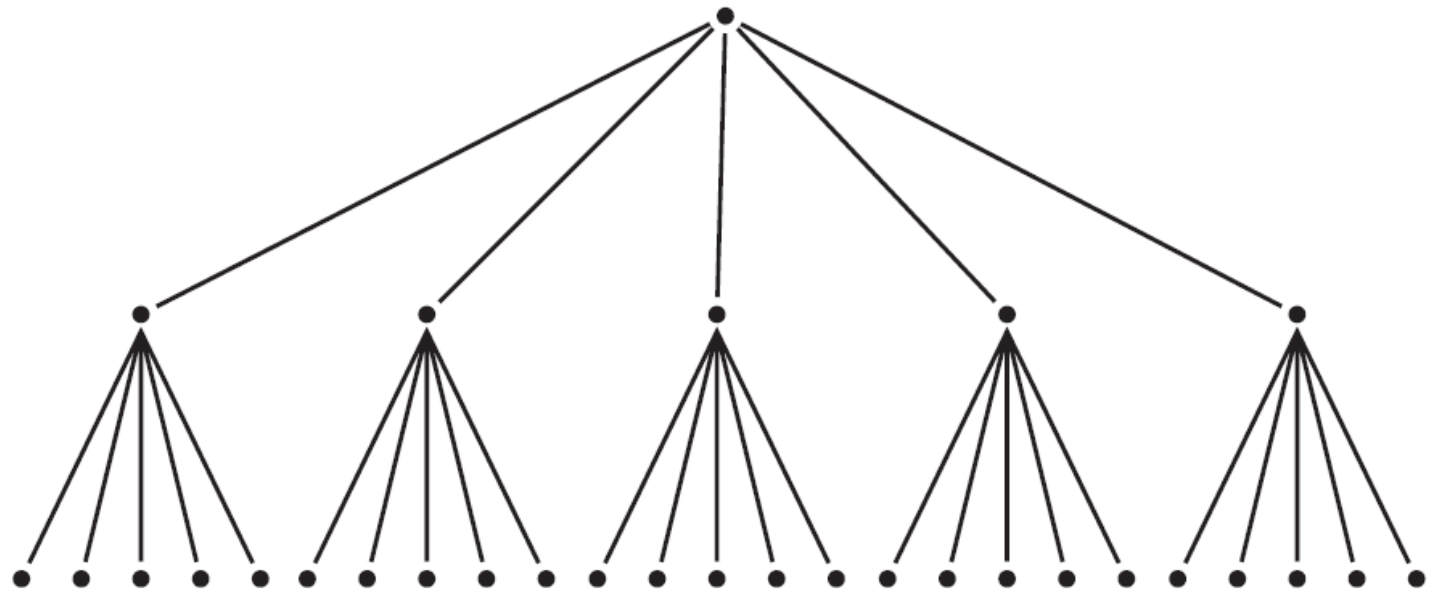
- (b) A node with three children
has two search keys



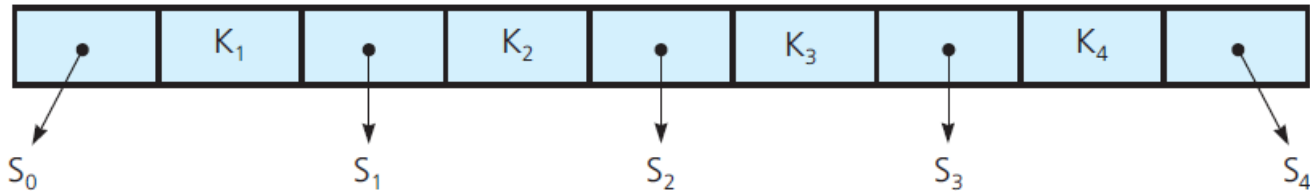
- (c) A node with m children
has $m - 1$ search keys



- Nodes with two, three, and m children and their search keys



The format of each node in the above tree



- A full tree whose internal nodes have five children

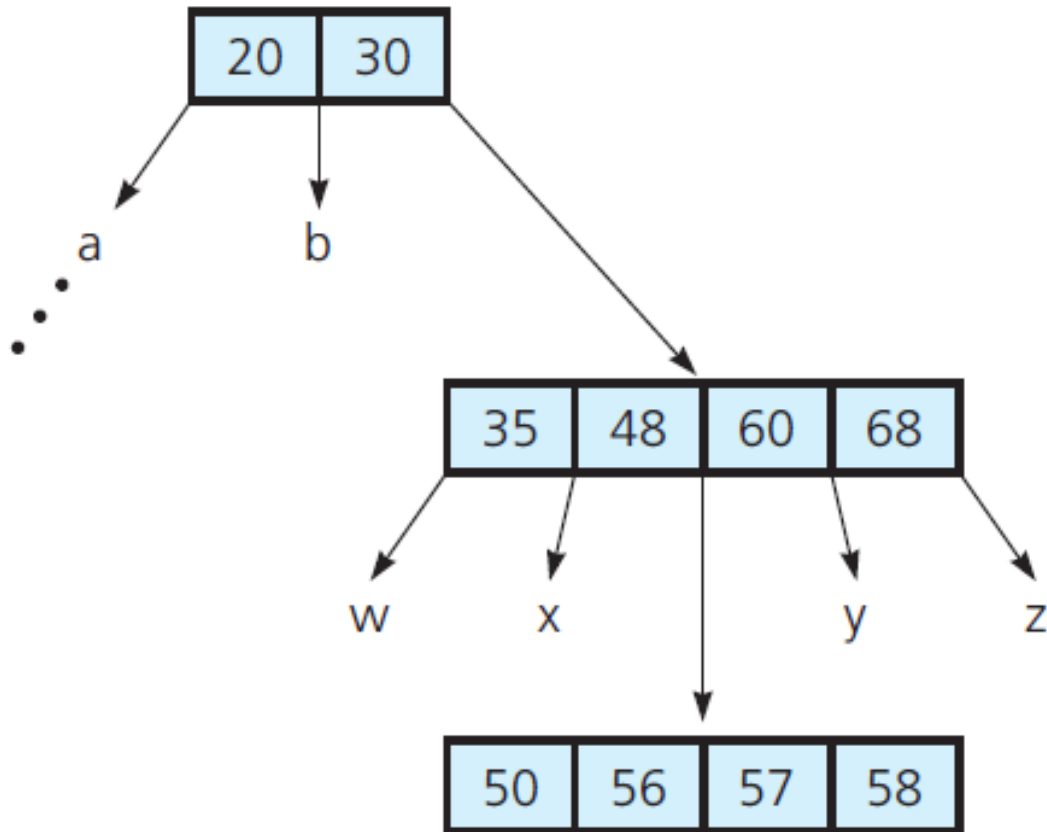


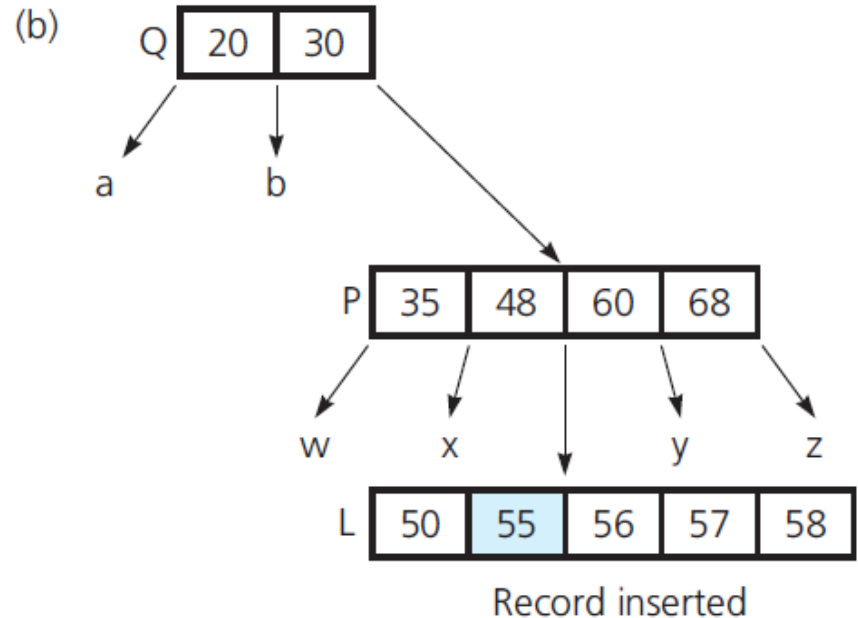
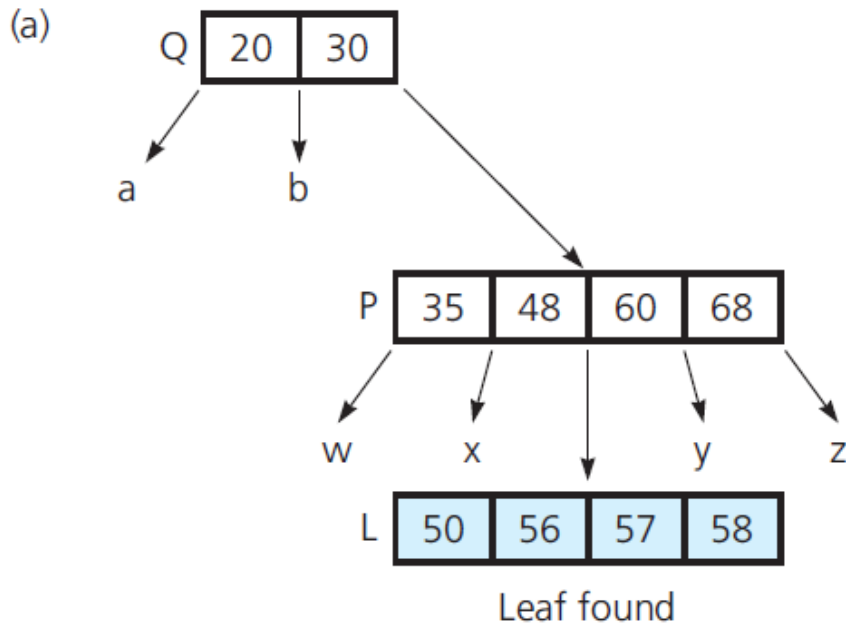
B-Trees

Adding a record to a B-tree

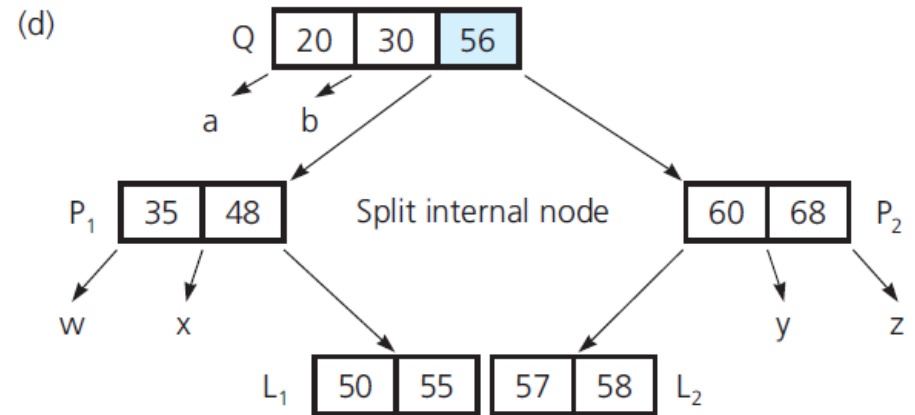
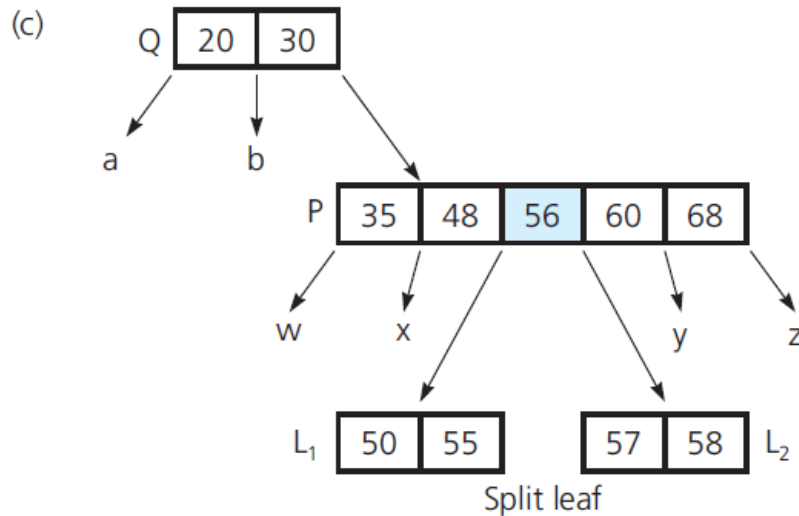
1. Add data record to data file
2. Add a corresponding index record to index file

- A B-tree of degree 5

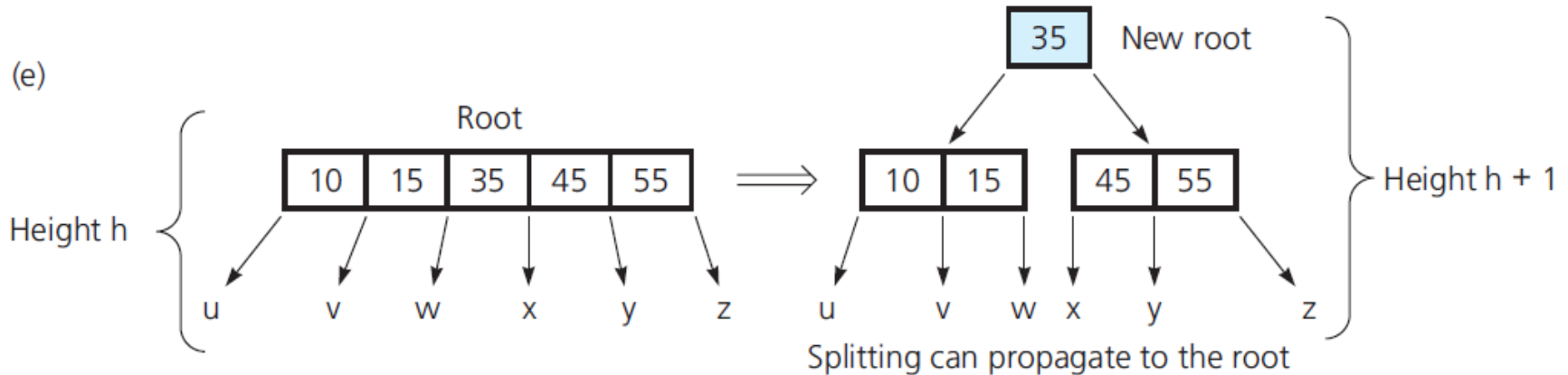




- (a through d) The steps for adding 55 to a B-tree;
(e) splitting the root



- (a through d) The steps for adding 55 to a B-tree;
(e) splitting the root

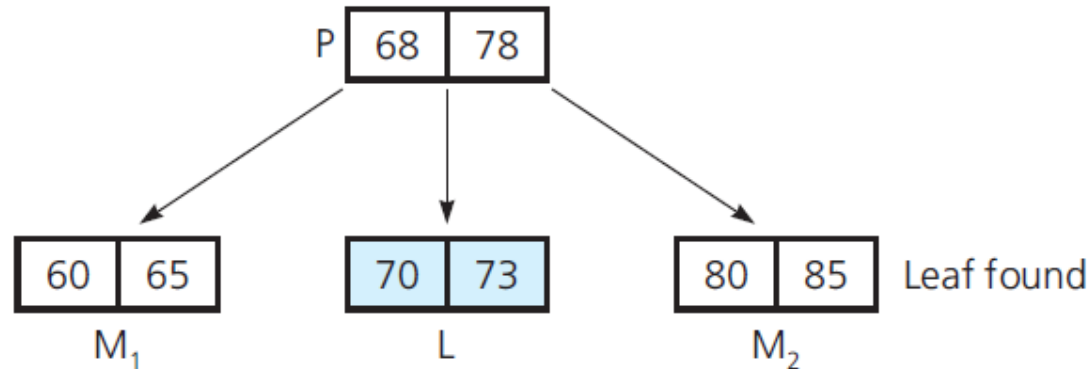


- (a through d) The steps for adding 55 to a B-tree;
(e) splitting the root

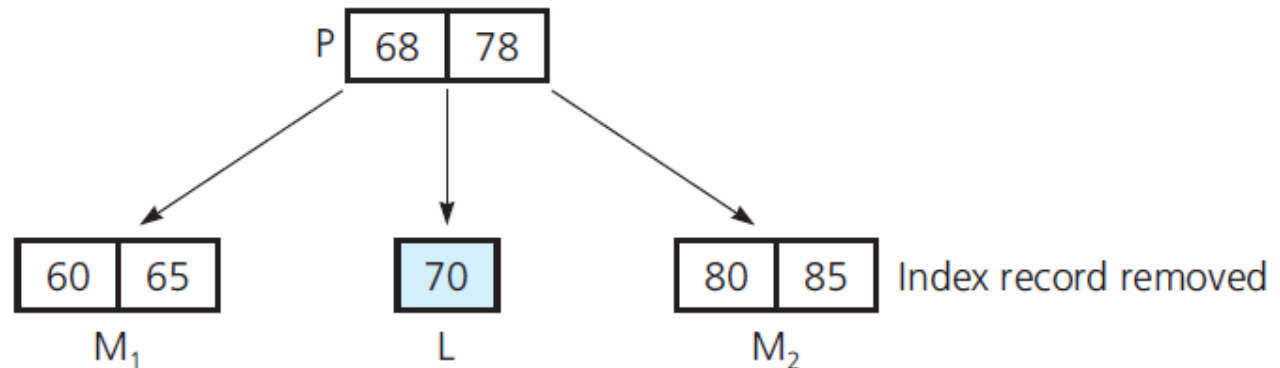
Removing a record from a B-tree

1. Locate index record in index file
2. Remove data record from data file

(a)

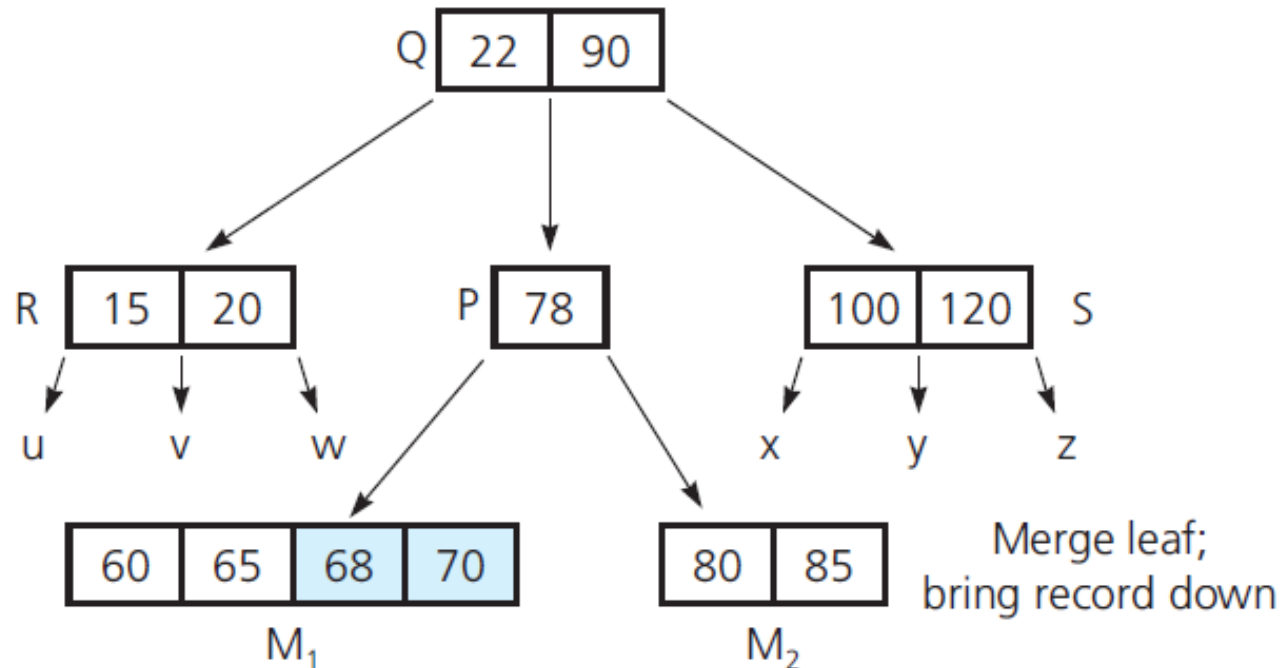


(b)



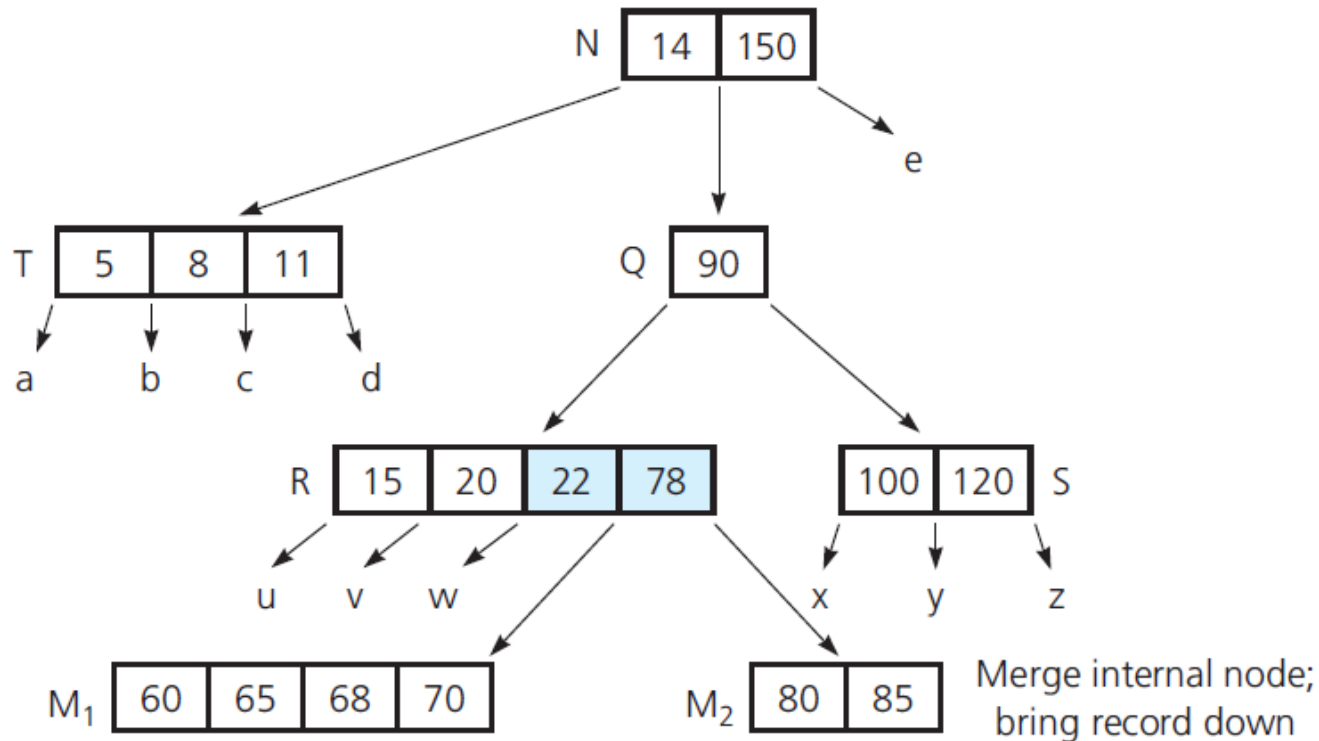
- (a through e) The steps for removing 73 ;
(f) removing the root

(c)

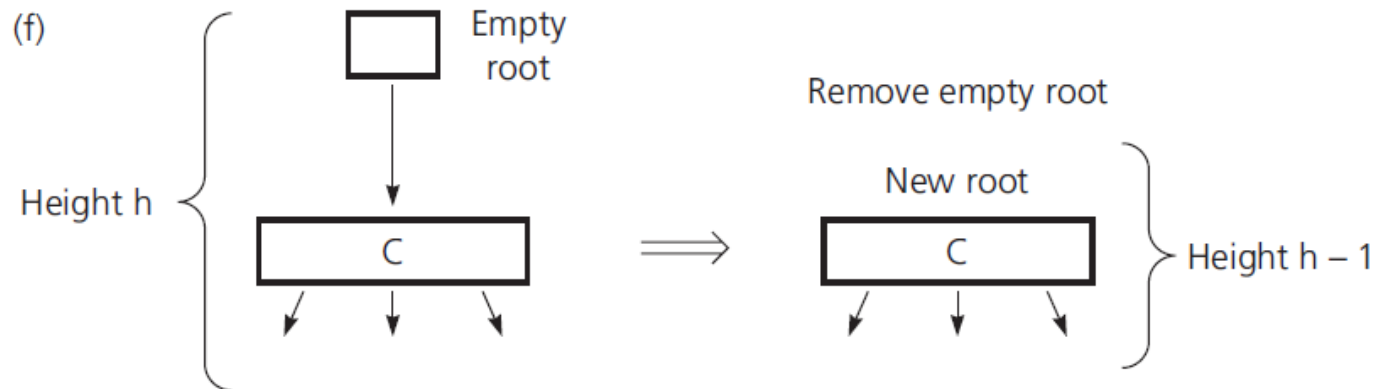
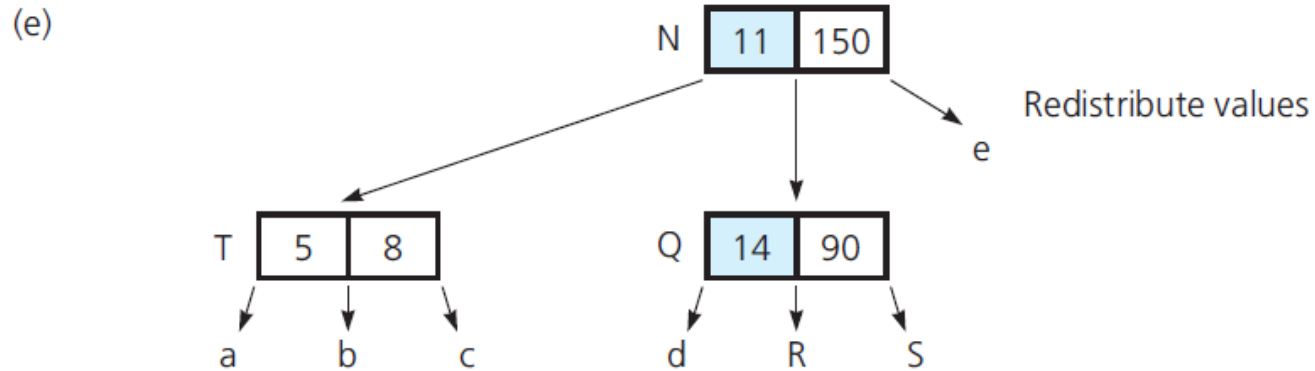


- (a through e) The steps for removing 73 ;
(f) removing the root

(d)

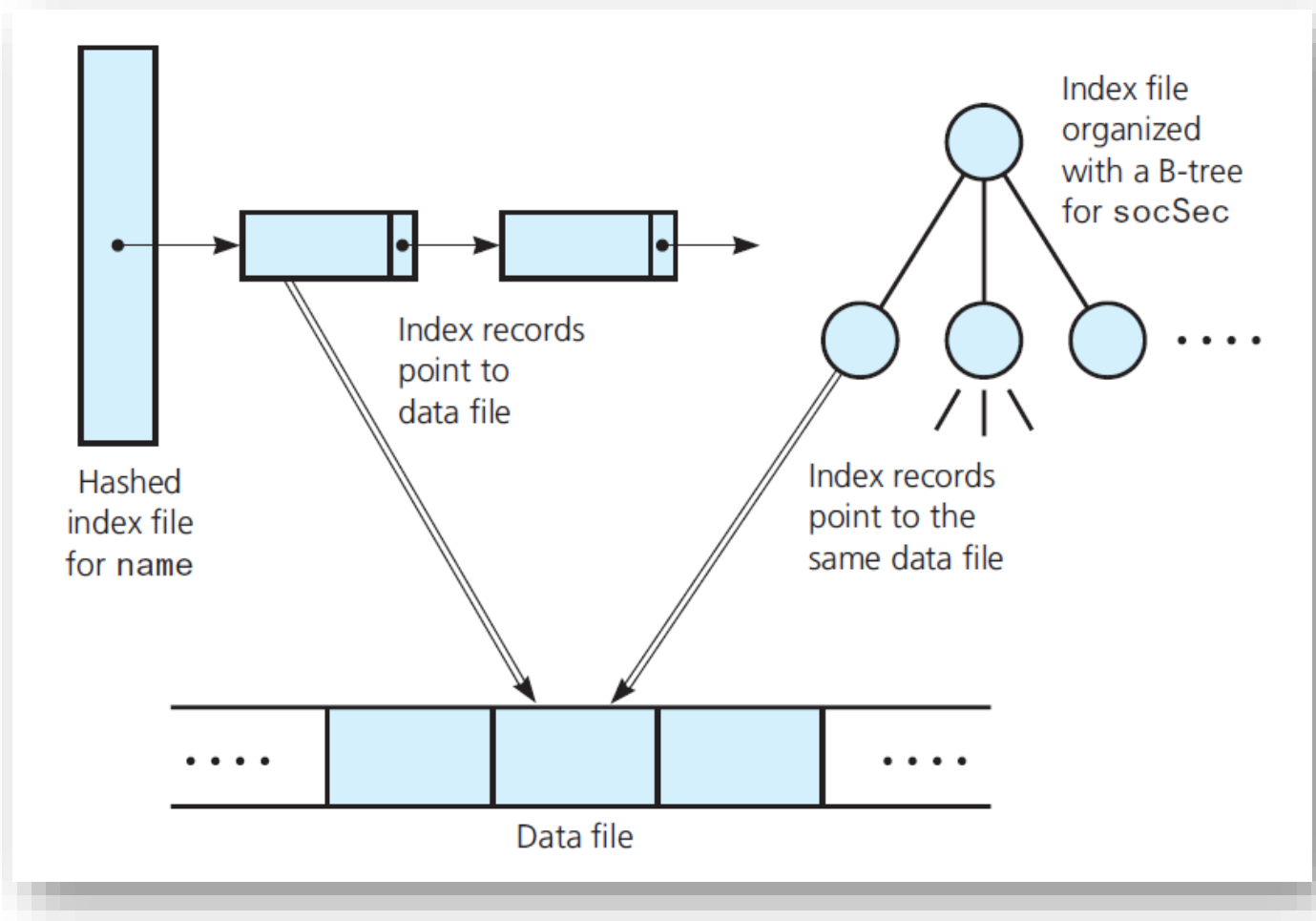


- (a through e) The steps for removing 73 ;
(f) removing the root



- (a through e) The steps for removing 73 ;
- (f) removing the root

Multiple Indexing



- Multiple index files



Multiple Indexing

A removal by name must update both indexes

1. Search **name** index file for **jones** and remove index record.
2. Remove appropriate data record from data file, noting **socSec** value **ssn** of this record.
3. Search **socSec** index file for **ssn** and remove this index record.

