#### **Sort-Merge Algorithm for External Sorting**

The algorithm has **two main phases**:

#### 1. Sorting Phase (Create Sorted Runs)

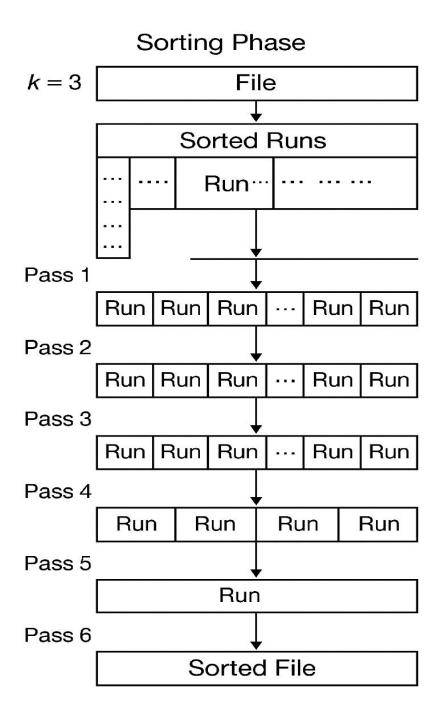
- Divide the file into chunks that fit into memory buffers.
- Load each chunk (say k blocks), sort it using an internal algorithm (like quicksort), and write it back to disk as a sorted subfile (run).
- Continue until the whole file is processed.

After this step, you have multiple **sorted runs** on disk.

## 2. Merging Phase (Combine Runs)

- Take k-1 subfiles at a time into memory (keeping 1 buffer for output).
- Merge them into a larger sorted subfile, writing back to disk block by block.
- Repeat in multiple passes, reducing the number of subfiles each time, until only 1 sorted file remains.

# **Example**



#### Suppose:

- File has 10,000 records.
- Each block stores 100 records.
- Buffer can hold 3 blocks (k = 3).

#### **Sorting Phase**

- Total blocks in file =  $10,000 \div 100 = 100$  blocks.
- Buffer can hold 3 blocks  $\rightarrow$  sort 3 blocks at a time.
- Number of runs = 100 ÷ 3 ≈ 34 runs (each run sorted internally).

So now we have 34 **sorted runs** stored on disk.

### **Merging Phase**

- Use k-1 = 2 runs per merge pass.
- Pass 1: Merge runs pair by pair  $\rightarrow$  34 runs  $\rightarrow$  17 runs.
- Pass 2: Merge again  $\rightarrow$  17 runs  $\rightarrow$  9 runs.
- Pass 3: Merge again  $\rightarrow$  9 runs  $\rightarrow$  5 runs.
- Pass 4: Merge again  $\rightarrow$  5 runs  $\rightarrow$  3 runs.
- Pass 5: Merge again  $\rightarrow$  3 runs  $\rightarrow$  2 runs.

Pass 6: Merge again → final single sorted file.

# (a) Sort-Merge Join ( $R \bowtie S$ on A = B)

#### **Steps:**

- 1. Sort relation R on attribute A and relation S on attribute B.
- 2. Use two pointers (i for R, j for S).
- 3. If R[i].A < S[j].B → move i.</li>
  If R[i].A > S[j].B → move j.
  If R[i].A = S[j].B → output the join pair and continue.

#### **Example:**

R:

A | Name

1 | Alex

3 | Ben

5 | Cathy

S:

B | Dept

1 | HR

3 | IT

```
5 | Finance
Output (R ⋈ S):
(1, Alex, HR)
(3, Ben, IT)
(5, Cathy, Finance)
```



# $\not$ (b) Sort–Merge Projection ( $\pi$ attributes (R))

### **Steps:**

- 1. From relation **R**, keep only required attributes into T'.
- 2. If projection key includes a primary key  $\rightarrow$  no duplicates possible.
- 3. Otherwise  $\rightarrow$  sort T' and remove duplicates.

### **Example:**

R:

(ID, Name, Dept)

- 1, Alex, HR
- 2, Ben, IT
- 3, Cathy, HR

Projection  $\pi$ Dept(R):

HR, IT, HR

→ After removing duplicates:

HR, IT



# 

### **Steps:**

- 1. Sort both R and S.
- 2. Scan both lists with two pointers.
- 3. If  $R[i] < S[j] \rightarrow \text{output } R[i]$ . If  $R[i] > S[j] \rightarrow \text{output } S[j]$ . If equal  $\rightarrow$  output one (skip duplicate).

# **Example:**

$$R = \{1, 3, 5\}$$

$$S = \{2, 3, 6\}$$

Output  $(R \cup S) = \{1, 2, 3, 5, 6\}$ 



# $\not$ (d) Sort–Merge Intersection (R $\cap$ S)

#### **Steps:**

1. Sort both R and S.

- 2. Scan both with two pointers.
- 3. If equal, output the tuple.
- 4. Otherwise advance the pointer of the smaller value.

#### **Example:**

$$R = \{1, 3, 5\}$$

$$S = \{2, 3, 5, 7\}$$

Output  $(R \cap S) = \{3, 5\}$ 



# (e) Sort-Merge Set Difference (R - S)

### **Steps:**

- 1. Sort both R and S.
- 2. Scan both.
- 3. If  $R[i] < S[j] \rightarrow \text{output } R[i]$ .
- 4. If  $R[i] = S[j] \rightarrow skip$  (remove common).
- 5. If  $R[i] > S[j] \rightarrow move j$ .

### **Example:**

$$R = \{1, 3, 5, 7\}$$

$$S = \{3, 6, 7\}$$

Output  $(R - S) = \{1, 5\}$