

Data Representation & External Sorting

DSCI 551

Wensheng Wu

Outline

- Representing data



- How are tables stored on storage devices?

- External Sorting

- How to sort 1TB data using 1GB of memory?

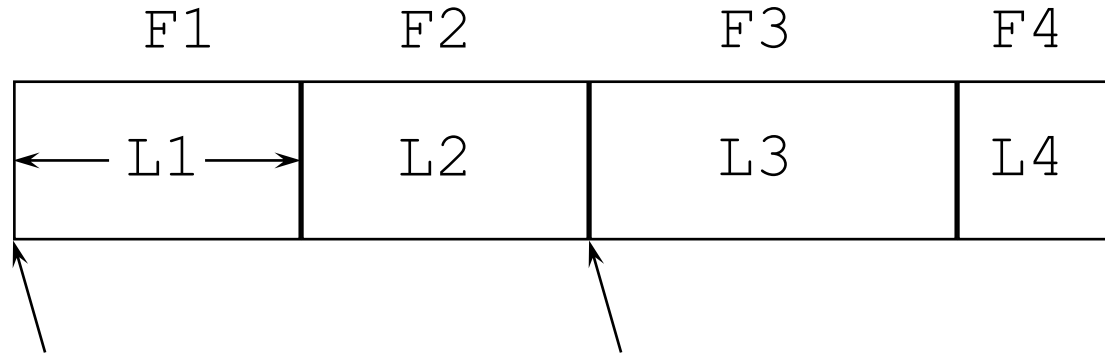
In-place sorting

Representing Data Elements

Base address (B)

- Relational database elements:
CREATE TABLE Product (
pid INT PRIMARY KEY, // tinyint, smallint, mediumint
name CHAR(20),
description VARCHAR(100),
? —→ maker CHAR(10) REFERENCES Company(name))
- A tuple/row is stored as a "record"

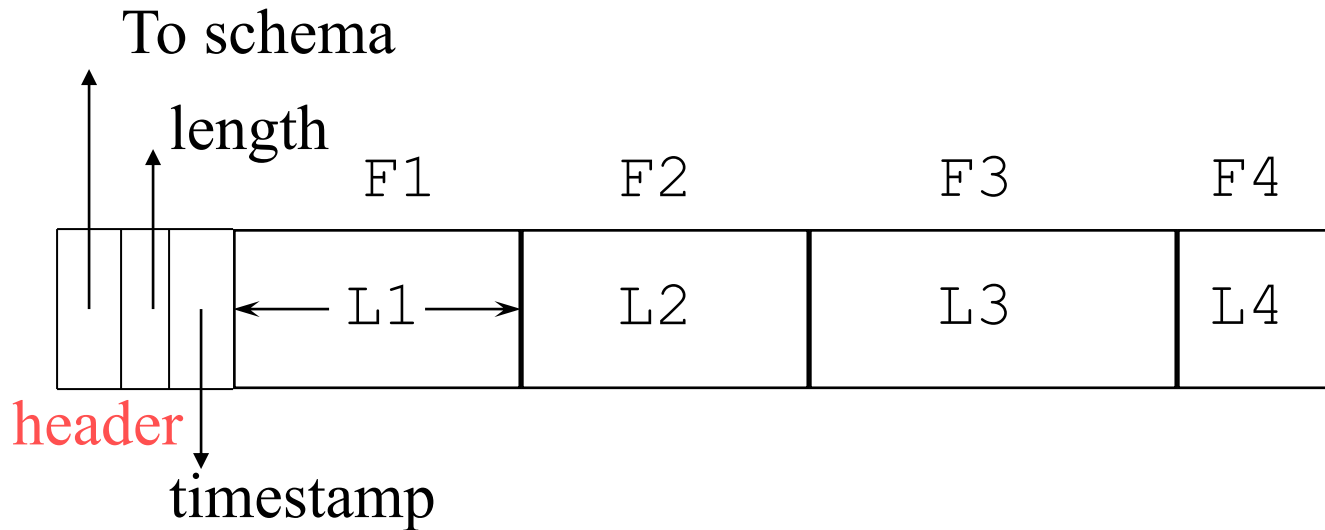
Record Formats: Fixed Length



Base address (B) Address = $B + L1 + L2$

- Information about field types is the same for all records in a file; stored in *system catalogs*.
- **Note the importance of schema information!**

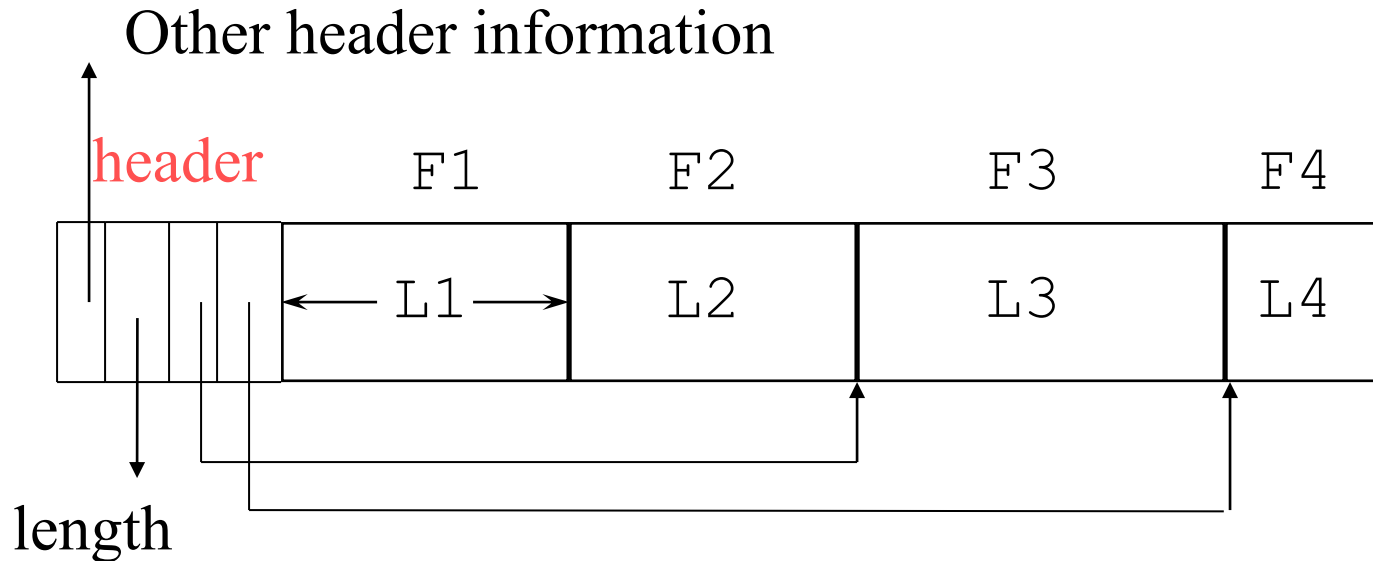
Record Header



Header:

- Pointer to schema: help finding fields
- Length: so we know where the record ends w/o consulting schema
- Timestamp: time when record last modified or read

Variable Length Records



Place the fixed fields first: F1, F2

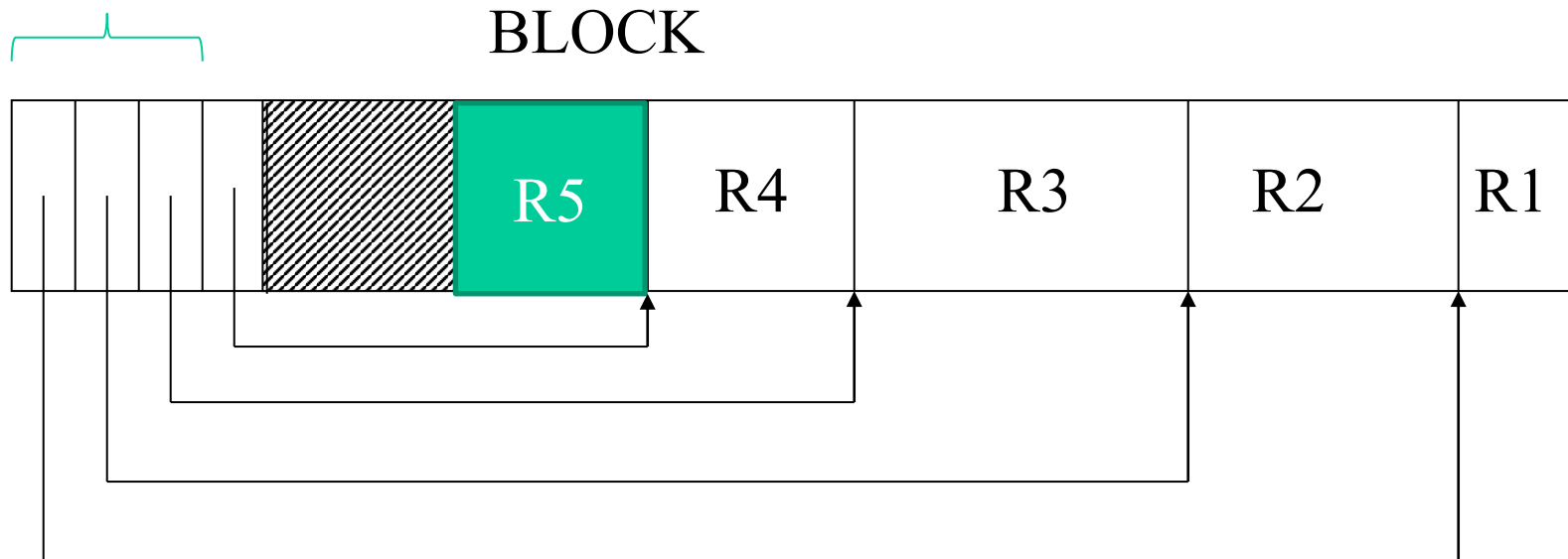
Then the variable length fields: F3, F4

Note: actually no need for pointer to F3, why?

Storing Records in Blocks

- Blocks have fixed size (typically 4KB)
 - But records may have variable-length

Offset table (slot directory)

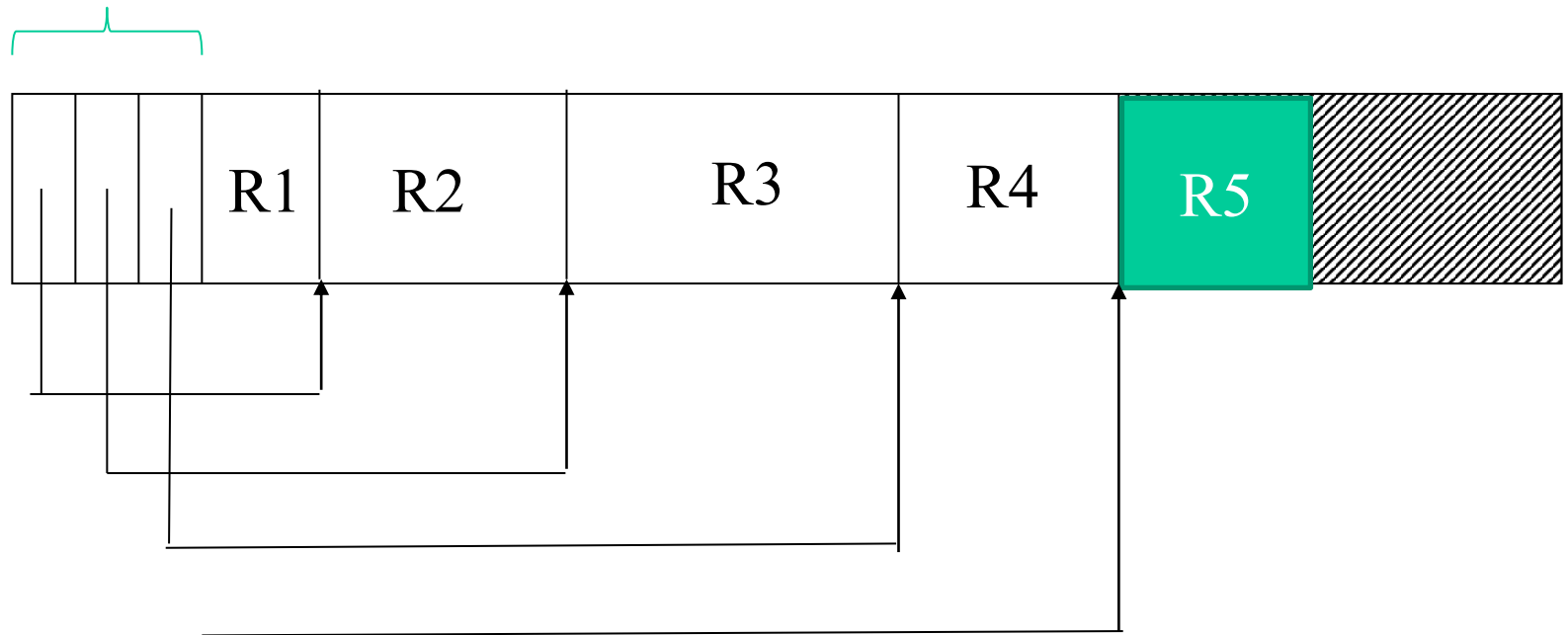


Why are records placed from the end?

Problem with this design?


- Records right after slot directory
- Free space after all records

Offset table (slot directory)



Outline

Mergesort(n)

- Representing data
 - How are tables stored on storage devices?
- External Sorting 
 - How to sort 1TB data using 1GB of memory?
 - 1GB \Rightarrow memory \Rightarrow 1GB run (sorted)
 - ...
 - 1GB \Rightarrow memory \Rightarrow 1GB run

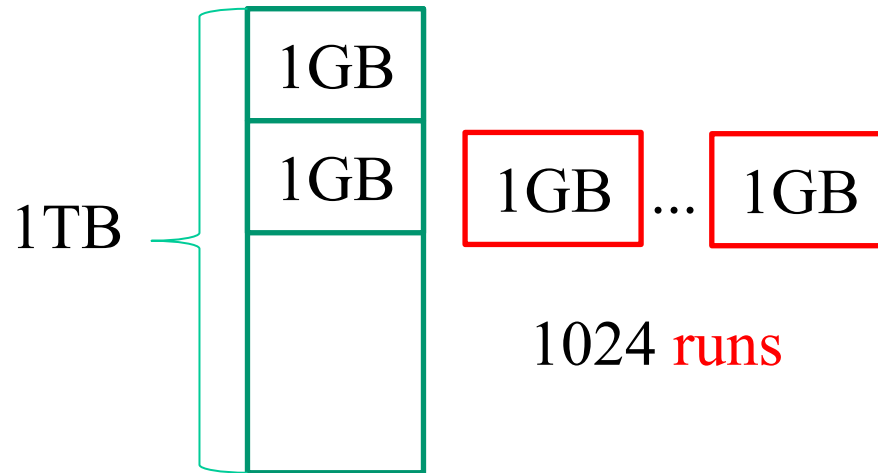
} Sorting (1024)

- merging: load one block from each run to memory

Notes

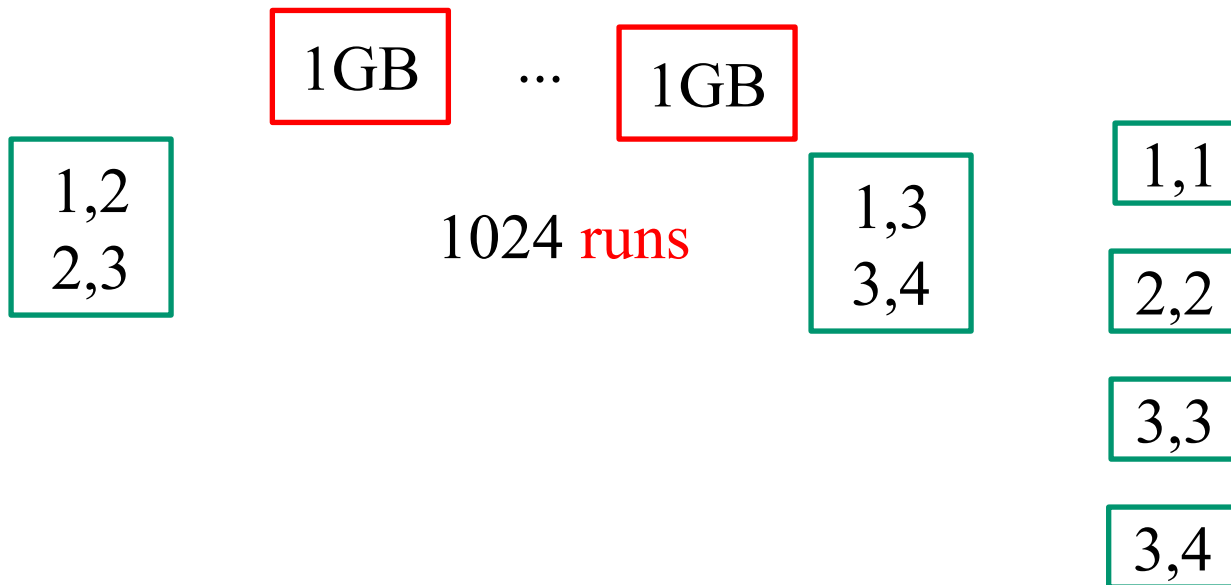
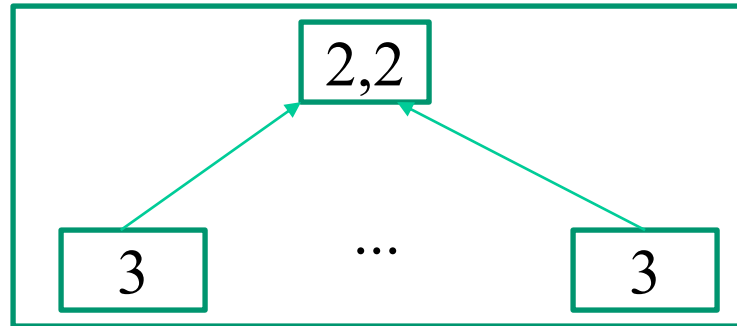
1GB

memory



Notes

Memory
1GB
4MB
(input buffer)



The I/O Model of Computation

- In main memory algorithms:
 - we care about CPU time
- In databases
 - time is dominated by I/O cost
- Assumption: cost is given only by I/O
- Consequence: need to redesign certain algorithms, e.g., sorting

Notes

- A block on storage devices loaded into a **page** in main memory
 - We sometimes interchange page with block
- Buffer pages
 - Often refer to pages in main memory used to store input, output, and intermediate data for an algorithm
- Run: a **sorted sublist** of input data

Notes

- Make a pass through data:
 - Loading the entire data from disk once

```
select bar, count(*)  
from Sells  
group by bar
```

Sorting

- Illustrates the difference in algorithm design when your data is not in main memory:

- Problem: sort 1TB of data with 1GB of RAM

R: 20 20 20 22 22 25 25 25

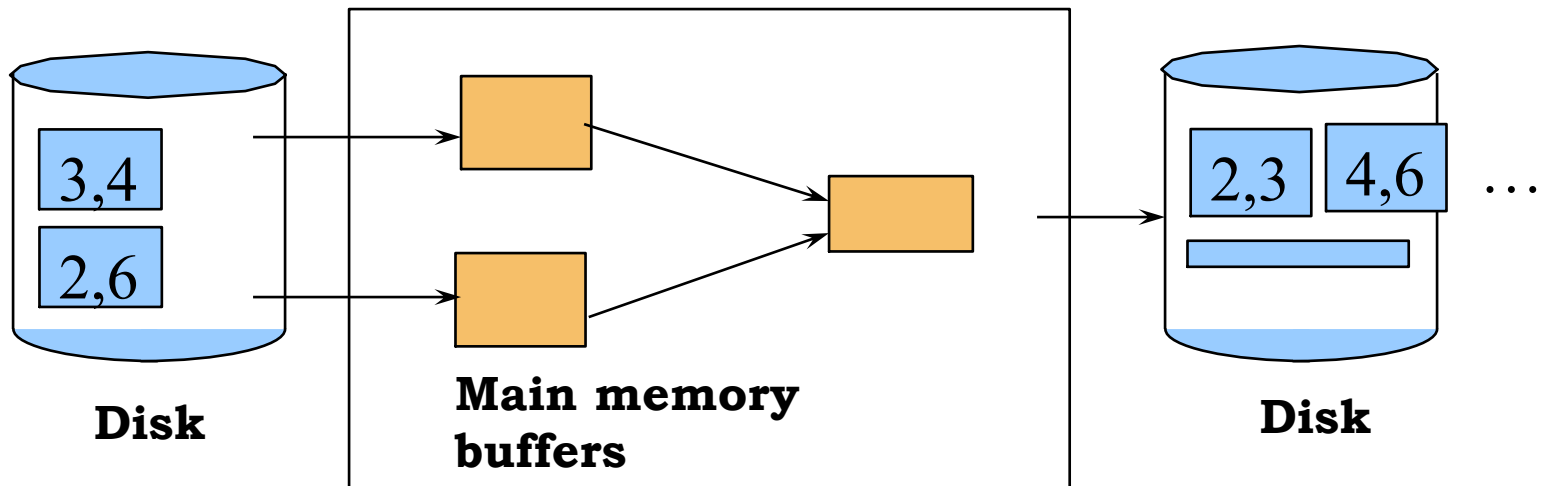
S: 20 20 23 24 25 26

- Arises in many places in database systems:

- Data requested in sorted order (ORDER BY)
 - Needed for grouping operations // group by age
 - First step in sort-merge join algorithm (R join_a S)
 - Duplicate removal
 - Bulk loading technique for creating B+-tree indexes

2-Way Merge-sort: Requires 3 Buffers

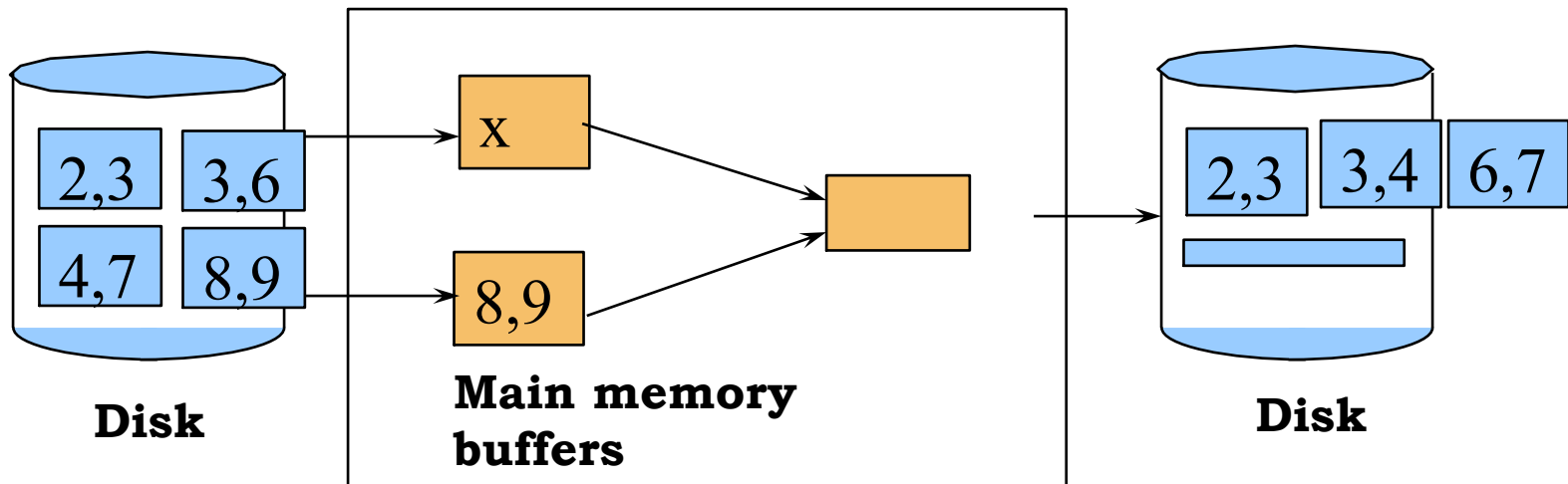
- Pass 0: Read a page, sort it, write it
 - only one buffer page is used
- Pass 1, 2, ..., etc.: merging two runs at a time
 - three buffer pages used.



M=3

2-Way Merge-sort: Requires 3 Buffers

- Pass 0: Read a page, sort it, write it
 - only one buffer page is used
- Pass 1, 2, ..., etc.: merging two runs at a time
 - three buffer pages used.



Two-Way External Merge Sort

- Each pass we read + write each page in file.

- N pages in the file => the number of passes

$$= \lceil \log_2 N \rceil + 1$$

- So total cost is:

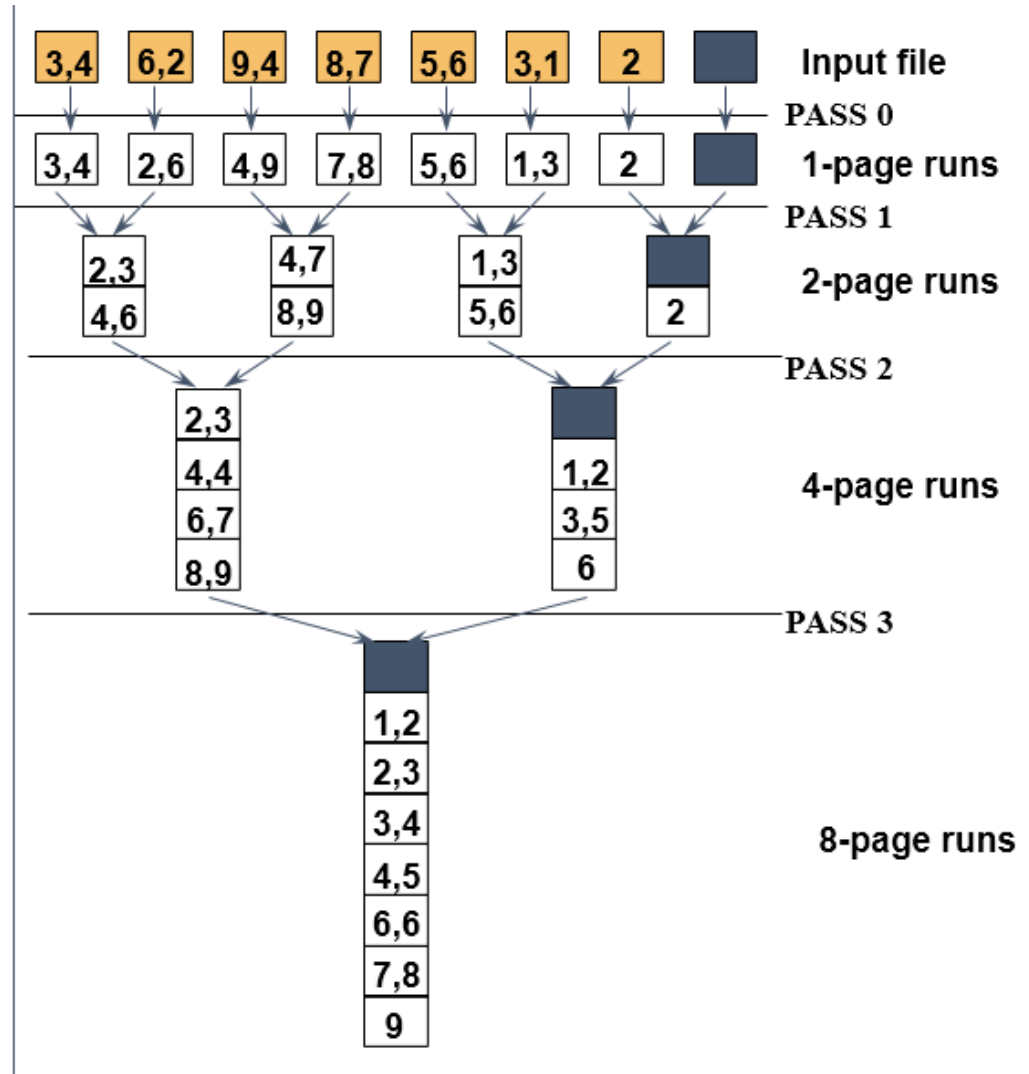
$$2N(\lceil \log_2 N \rceil + 1)$$

- Sort 4MB with buffer page size = 4KB: needs 11 passes

$$N=7, k=3$$

$$1 * 2 * 2 * \dots * 2 = 1 * 2^k \geq N \Rightarrow k \geq \log_2(N)$$

$$k = \text{ceil}[\log_2(N)] = \text{ceil}[2.8] = 3$$



Notes

$$\begin{aligned} N &= 4\text{MB}/4\text{KB} \\ &= 1024 \\ &= 2^{10} \end{aligned}$$

$$\log_2(2^{10}) = 10$$

Can We Do Better ?

- We have more main memory
 - Should use it to improve performance
-
- M: # of blocks (i.e., pages) in main memory
 - B(R): # of blocks of relation R
- M=5**, B(R)=108
- sorting: load 5 pages, sort them, write back as run
 $108/5 = \text{ceil}(21.6)$ runs
21 runs, 5 pages/run $\Rightarrow 21*5 = 105$ pages
1 run, 3 pages/run
- Merging: (M-1)-way merging \Rightarrow 4-way merging
take 4 runs, 5 pages/run \Rightarrow 20-page run
how many runs?

多路合并: (M-1)路合并

$M=5$, $B(R)=108$

sorting:

runs: $108/5 = 21.6 = 22$

size of run:

5 blocks for each of first 21 runs

3 blocks in last run

merging 1: (4-way)

1. take 4 runs, merge into a single run

size of run: $4*5 = 20$ blocks

2. take next 4 runs \Rightarrow 20 blocks

3. next 4 \Rightarrow 20 blocks

4. next 4 \Rightarrow 20 blocks

5. next 4 \Rightarrow 20 blocks

6. take last two runs $\Rightarrow 5+3 = 8$ blocks

$M=5$, $B(R)=108$

sorting:

runs: $108/5 = 21.6 = 22$

size of run:

5 blocks for each of first 21 runs

3 blocks in last run

merging 1: (4-way)

of runs: $6 = \text{ceiling}(22/4) = \text{ceiling}(5.5)$

size of run:

20 blocks: first 5

8 blocks: last run

merging 2: (4-way)

of runs: $2 = \text{ceiling}(6/4) = 2$

1. take first 4 runs $\Rightarrow 20 * 4 = 80$ blocks

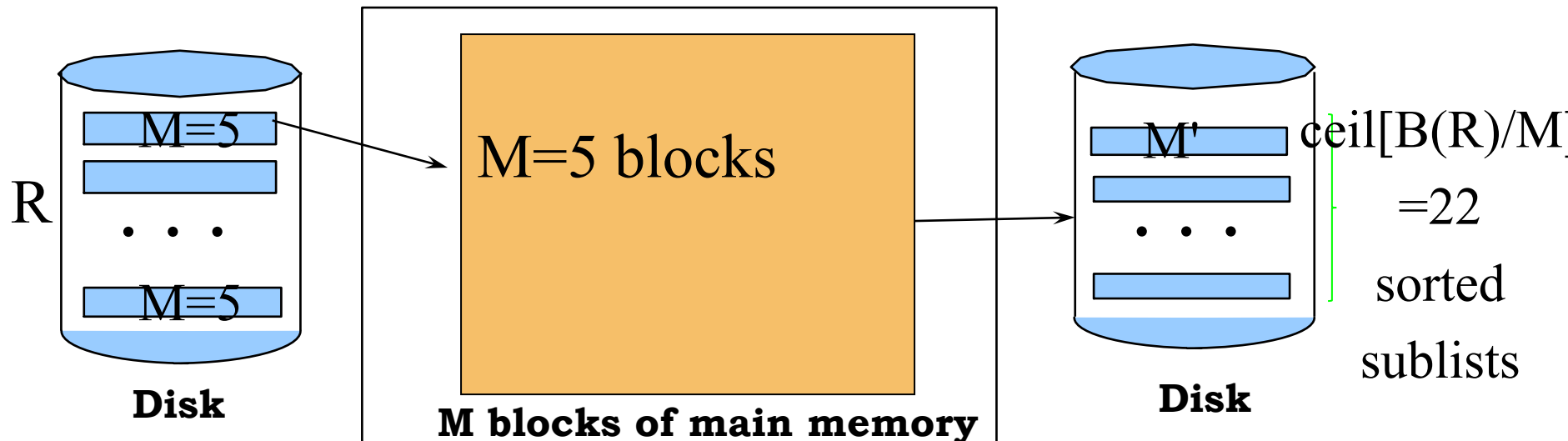
2. take last 2 runs $\Rightarrow 20 + 8 = 28$ blocks

merging 3:

take 2 runs \Rightarrow a single run

External Merge-Sort

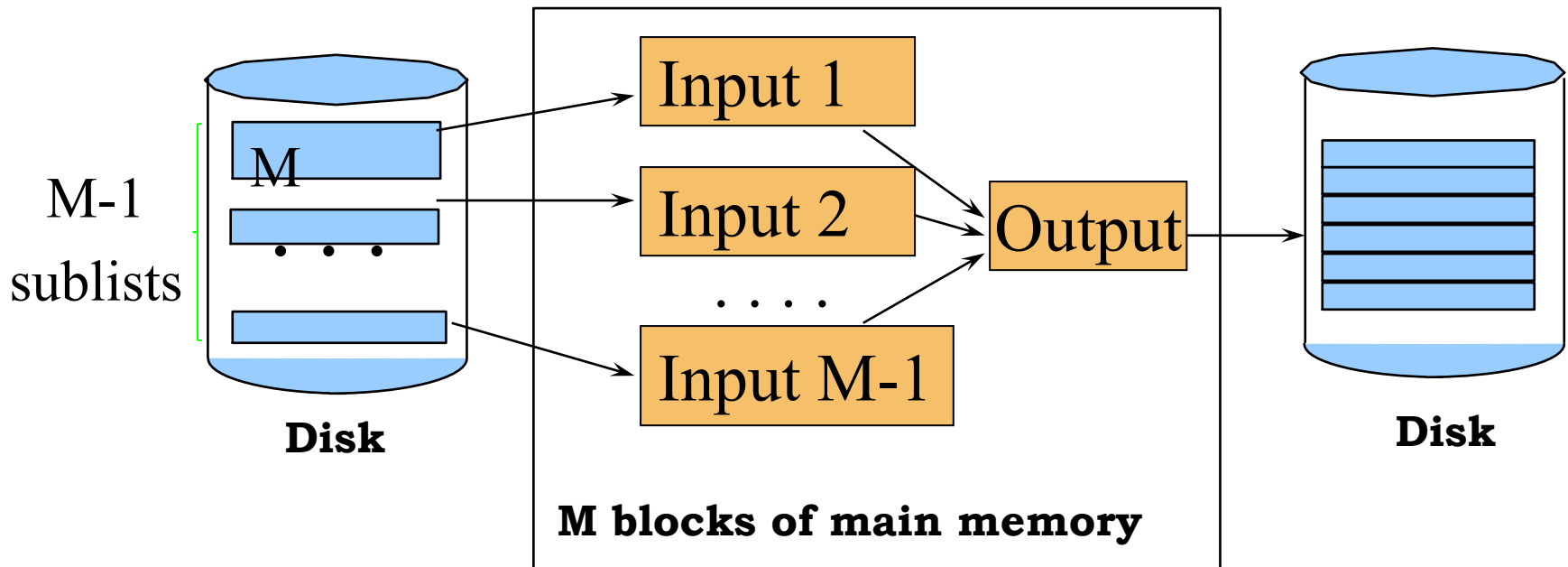
- Pass 0: load M blocks in memory, sort
 - Result: $\text{ceil}(B(R)/M)$ sorted sublists of size M
 - Each sorted sublist is a run



$$B(R)=110$$

Pass One (merging)

- Merge $M - 1$ runs into a new run
- Result: each run has now $M (M - 1)$ blocks

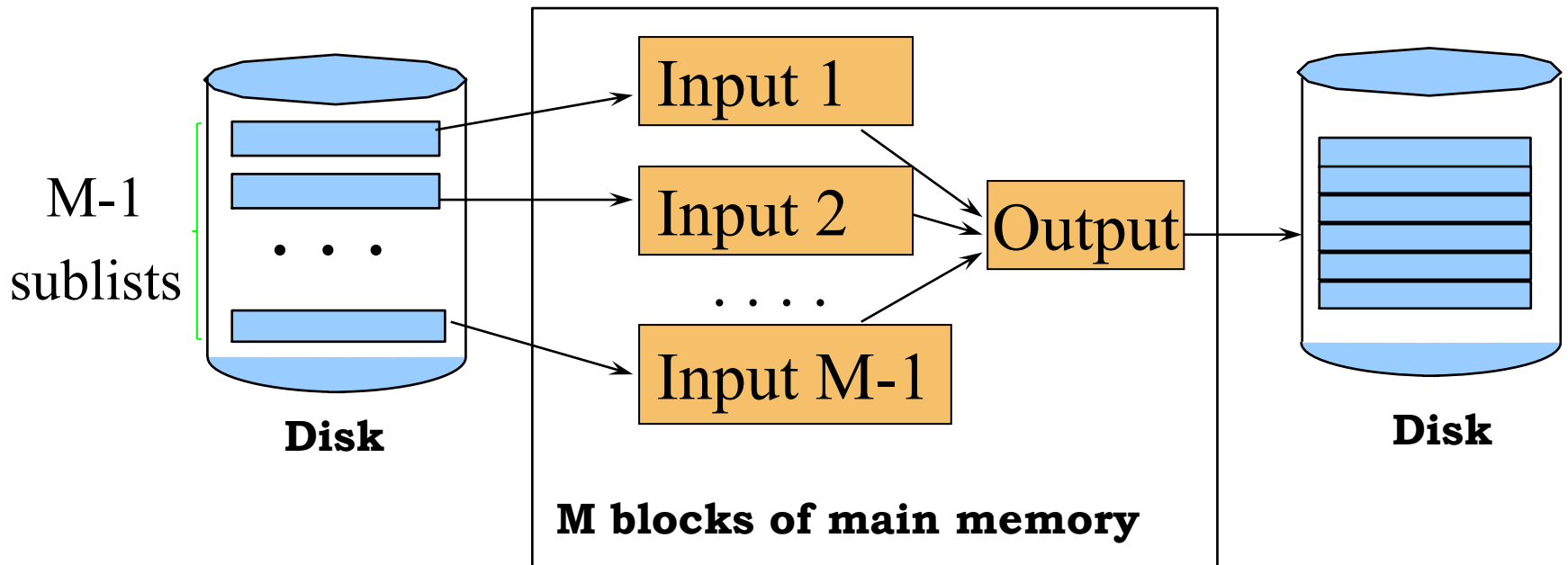


Cost of Two-Pass, Multiway Merge Sort

- Pass 0: sort B/M sublists of size M , write
 - Cost: $2B(R)$ $2B(R)$ 是怎么来的? 是因为 not in-place吗
- Pass 1: merge B/M sublists, write
 - Cost: $2B(R)$
- Total cost: $4B(R)$
- Assumption: $B(R) \leq M^2$
 - $B/M \leq M - 1$ or 保证是2个pass 就能排序完成
 - $B \leq M(M-1) \sim M^2$

Generalized to k Passes


- Merge every $M - 1$ runs into a new run
- Result: each run has now $M (M - 1)^k$ blocks

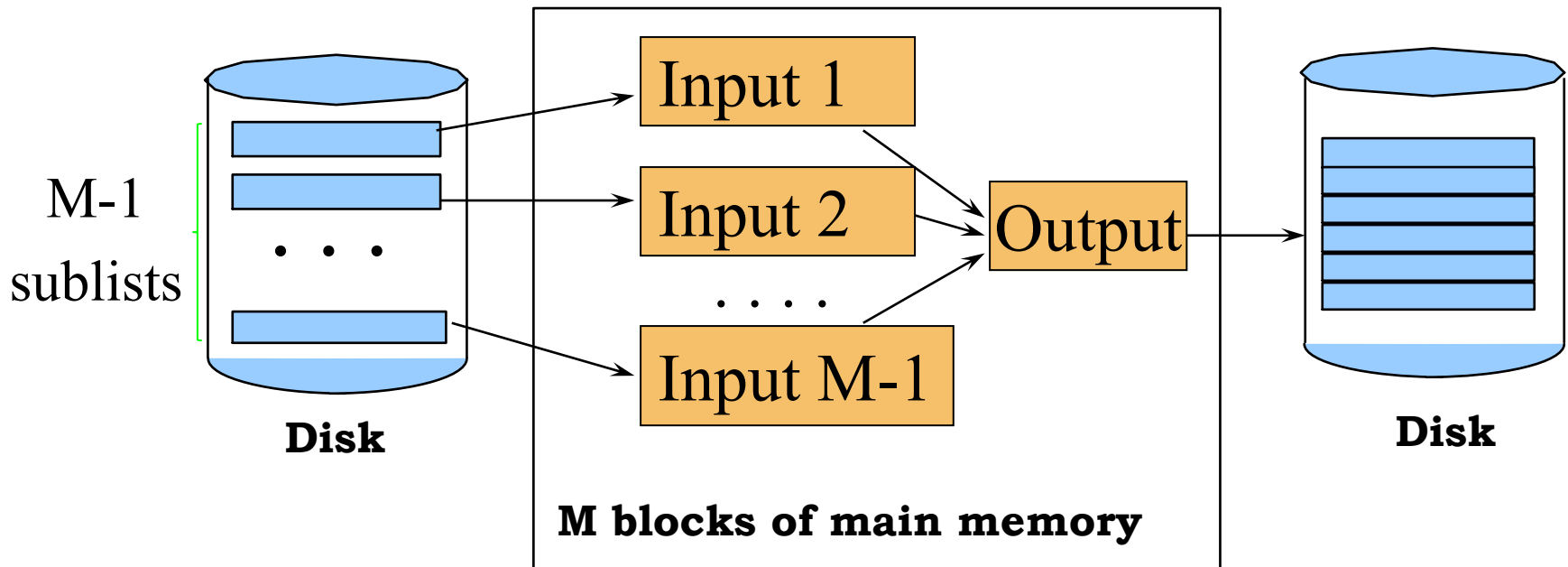


If k is the last pass

$$1 * 2^k \geq N$$

- Merge $M - 1$ runs into a single run

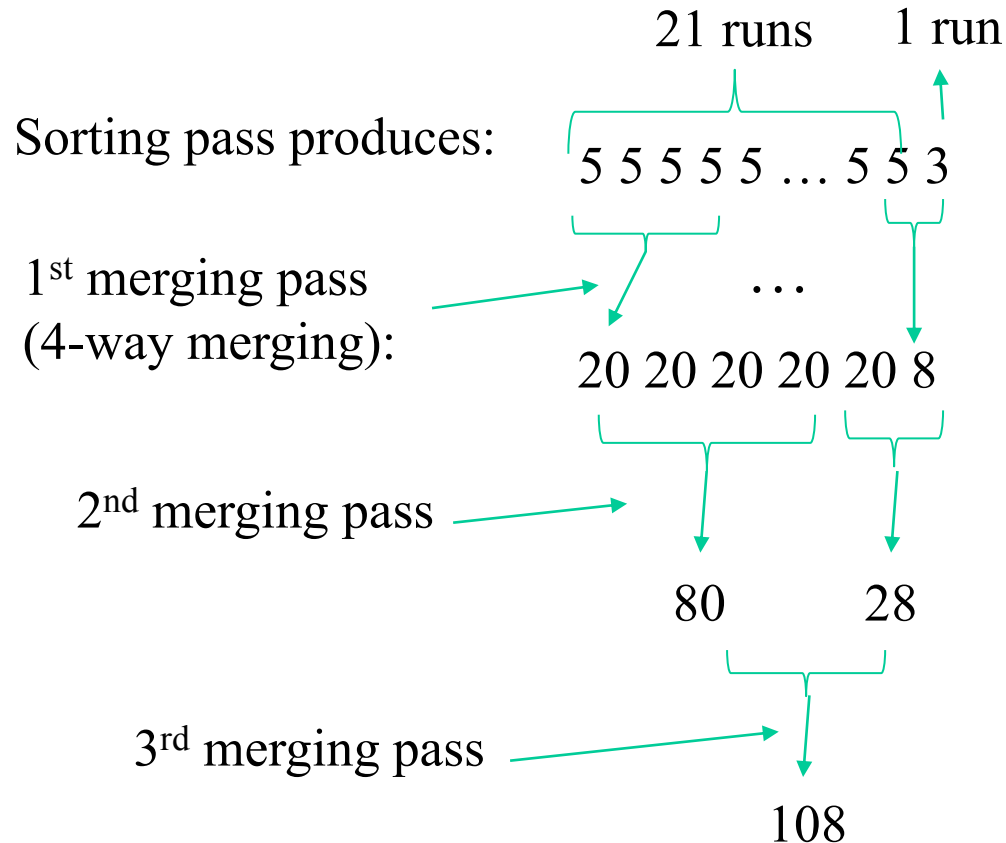
- We must have $M (M - 1)^k \geq B(R)$  $k = \lceil \log_{M-1} \lceil B / M \rceil \rceil$



Cost of External Merge Sort

- Number of passes: $1 + \lceil \log_{M-1} \lceil B / M \rceil \rceil$
- **Cost = 2B * (# of passes)** M=5, B(R)=108
- E.g., with 5 buffer pages, to sort 108-page file:
 - Pass 0: produces $\lceil 108/5 \rceil = 22$ runs (21 sorted runs of 5 pages each + last run of only 3 pages)
 - Pass 1: $\lceil 22/4 \rceil = 6$ (5 sorted runs of 20 pages each + last run of only 8 pages)
 - Pass 2: 2 sorted runs, 80 pages and 28 pages
 - Pass 3: Sorted file of 108 pages

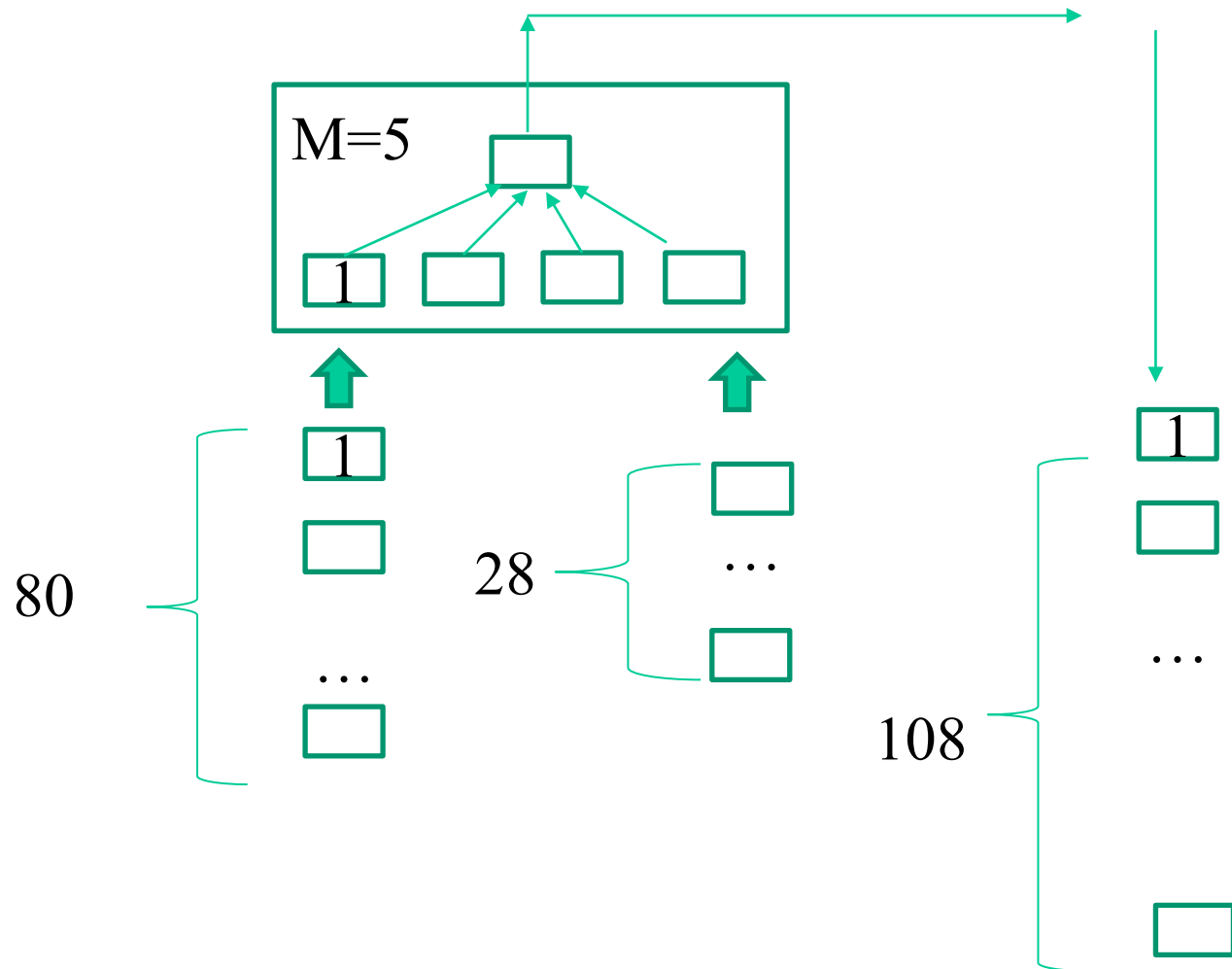
Example Illustrated



of passes:
 $1 + \text{ceil}(\log_4(N))$

$N = \# \text{ of runs by sorting}$
 $= 22$
 $= \text{ceil}(108/5)$
 $= \text{ceil}(B(R)/M)$

Example



Sorting 1TB using 1GB Memory

- $B(R) = 1\text{TB}/4\text{KB} = 256\text{M}$ (blocks), $M = 1\text{GB}/4\text{KB} = 256\text{K}$ (pages)
- Sorting phase produces 1024 runs = 1K runs
- Merging:
 - Can do: $1\text{GB}/4\text{KB}-1 = 256\text{K}-1$ ways of merging
 - Can we finish merging in one merging pass?

of passes

- 2-way: $1 + \log_2(N)$
 - $m = 3$
 - $m' = 1$
 - $B = N$
- $(m-1)$ -way of merging:
 - $1 + \log_{(m-1)} (B/m')$