

Lecture 10: External Merge Sort

Today's Lecture

1. External merge sort
2. External merge sort on larger files
3. Optimizations for sorting

1. External Merge Sort

Recap: External Merge Algorithm

- Suppose we want to merge two **sorted** files both much larger than main memory (i.e. the buffer)
- We can use the **external merge algorithm** to merge files of ***arbitrary length*** in **$2*(N+M)$ IO** operations with only **3 buffer pages!**

Our first example of an “IO aware”
algorithm / cost model

Why are Sort Algorithms Important?

- Data requested from DB in sorted order is **extremely common**
 - e.g., find students in increasing GPA order
- **Why not just use quicksort in main memory??**
 - What about if we need to sort 1TB of data with 1GB of RAM...

A classic problem in computer science!

More reasons to sort...

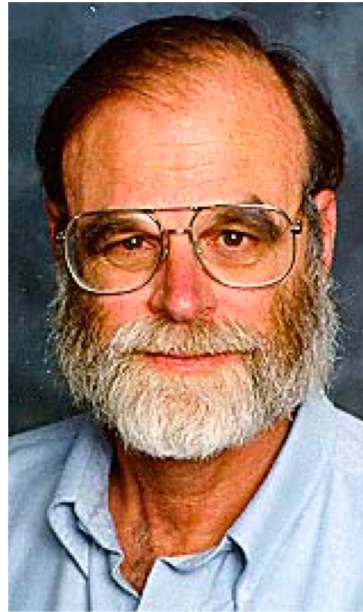
- Sorting useful for eliminating *duplicate copies* in a collection of records (Why?)
- Sorting is first step in *bulk loading* B+ tree index.
- *Sort-merge* join algorithm involves sorting

Coming up...

Next lecture

Do people care?

<http://sortbenchmark.org>



Sort benchmark bears his name

So how do we sort big files?

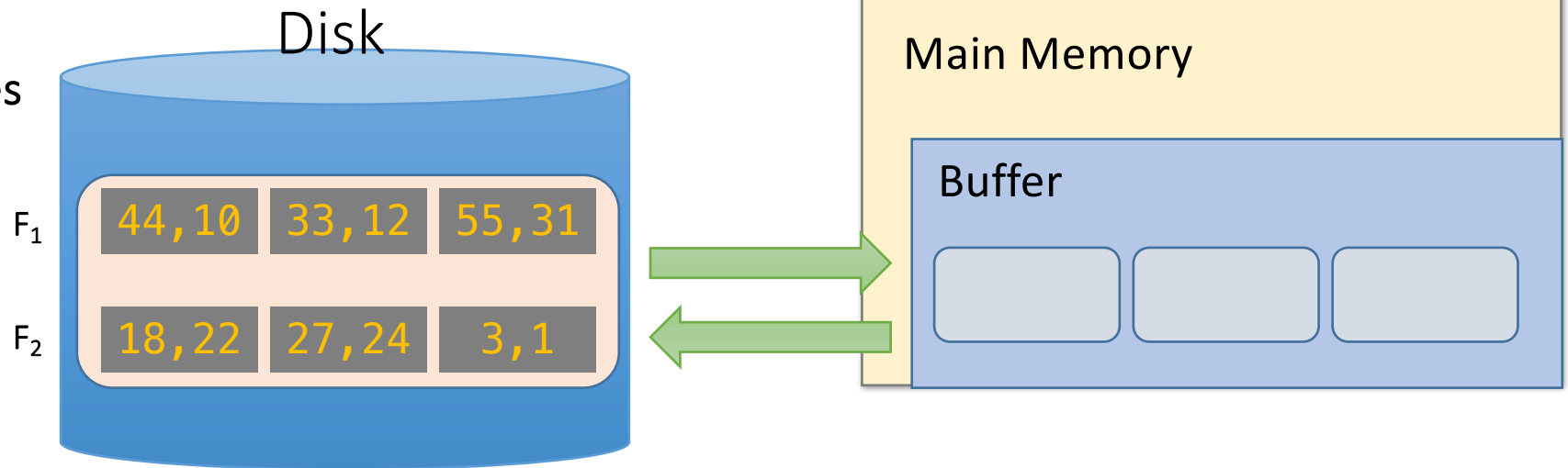
1. Split into chunks small enough to **sort in memory** (*“runs”*)
2. **Merge** pairs (or groups) of runs *using the external merge algorithm*
3. **Keep merging** the resulting runs (*each time = a “pass”*) until left with one sorted file!

External Merge Sort Algorithm

Example:

- 3 Buffer pages
- 6-page file

Orange file
= unsorted



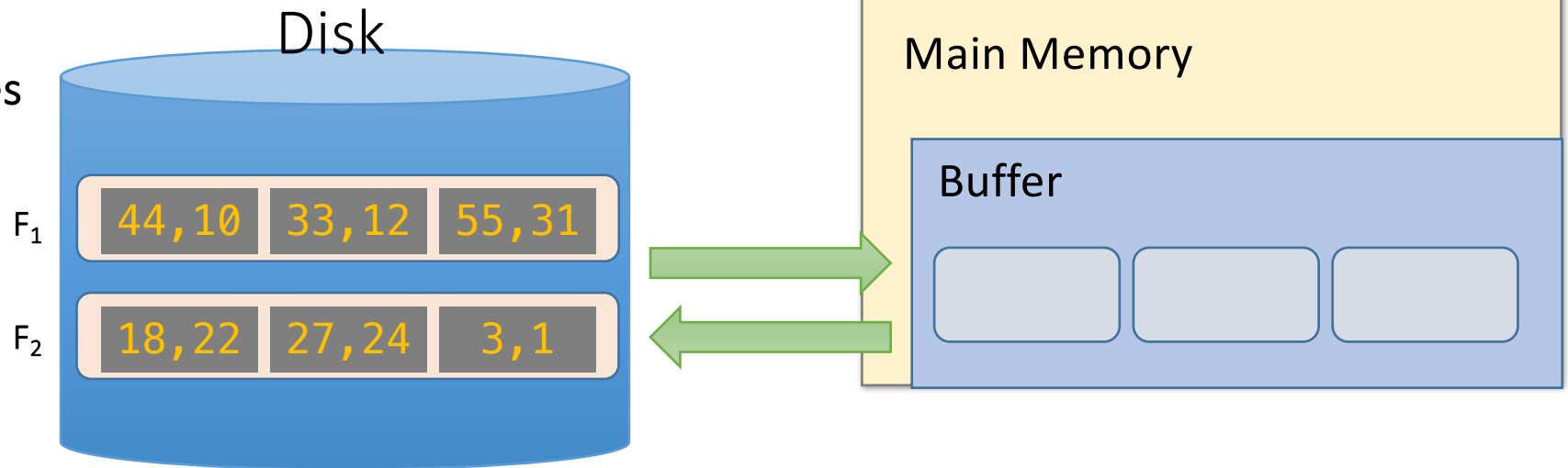
1. Split into chunks small enough to **sort in memory**

External Merge Sort Algorithm

Example:

- 3 Buffer pages
- 6-page file

Orange file
= unsorted



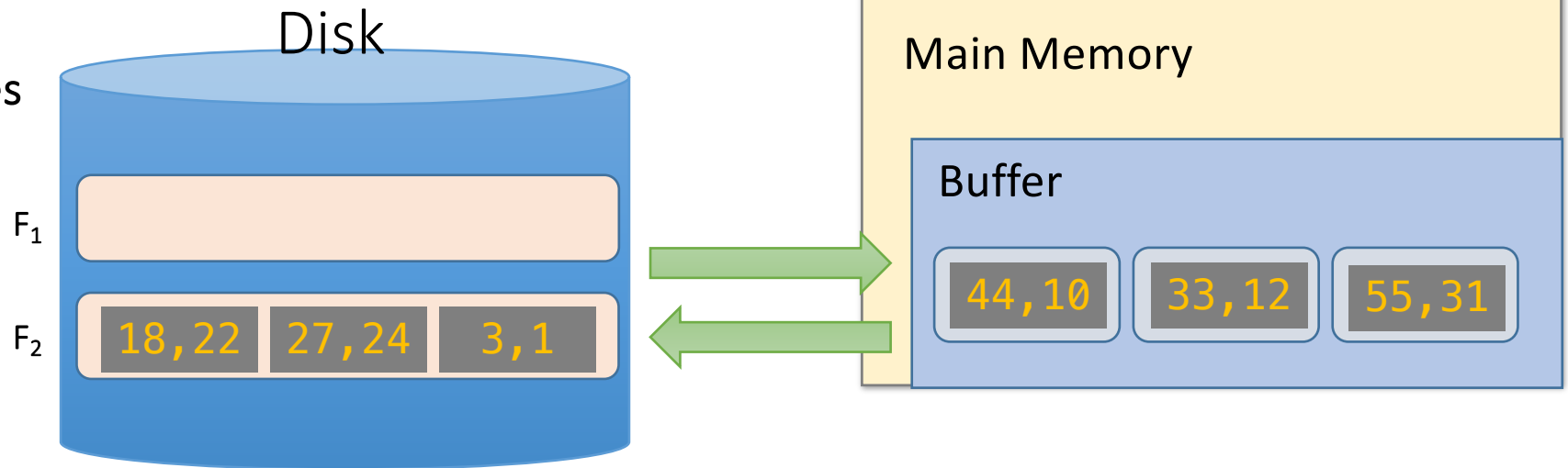
1. Split into chunks small enough to **sort in memory**

External Merge Sort Algorithm

Example:

- 3 Buffer pages
- 6-page file

Orange file
= unsorted



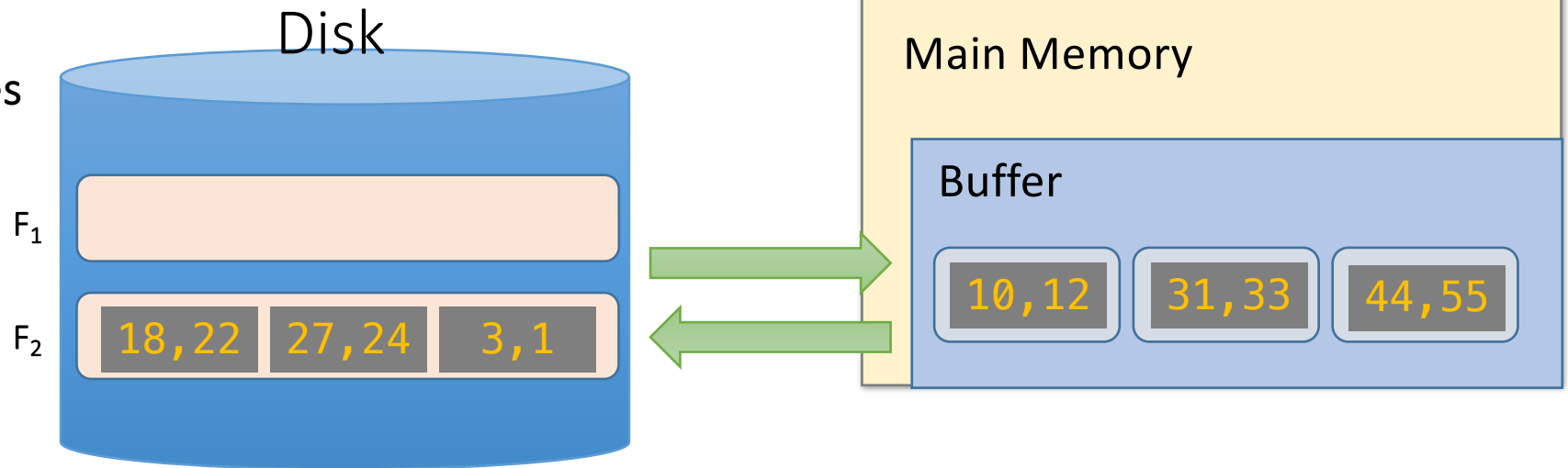
1. Split into chunks small enough to **sort in memory**

External Merge Sort Algorithm

Example:

- 3 Buffer pages
- 6-page file

Orange file
= unsorted



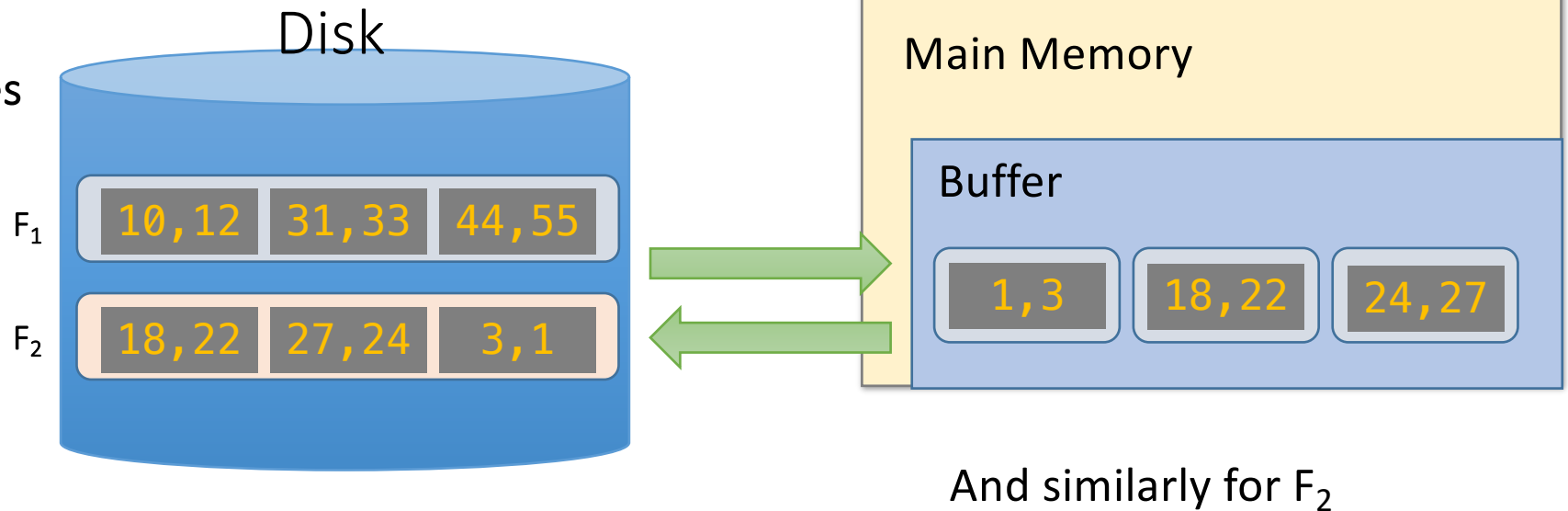
1. Split into chunks small enough to **sort in memory**

External Merge Sort Algorithm

Example:

- 3 Buffer pages
- 6-page file

Each sorted file is called a *run*

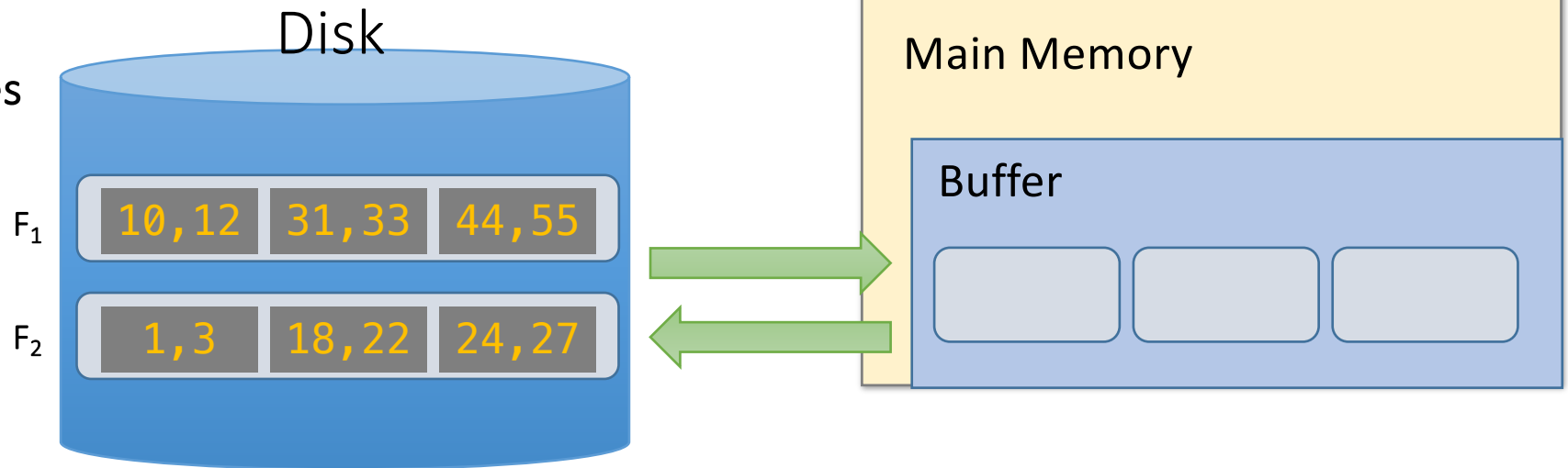


1. Split into chunks small enough to **sort in memory**

External Merge Sort Algorithm

Example:

- 3 Buffer pages
- 6-page file



2. Now just run the **external merge** algorithm & we're done!

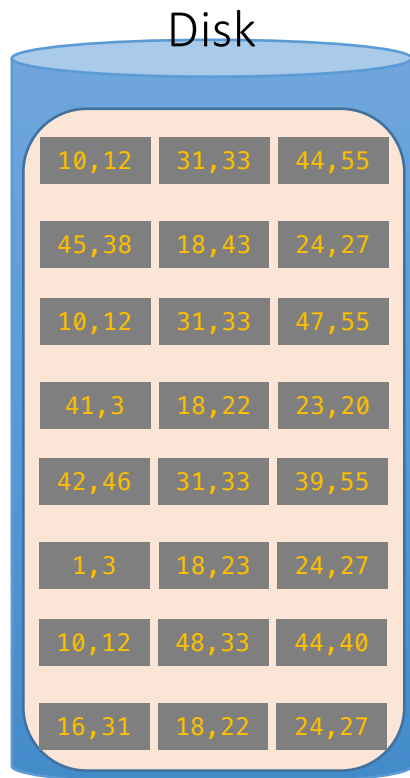
Calculating IO Cost

For 3 buffer pages, 6 page file:

1. Split into **two 3-page files** and **sort in memory**
 1. $= 1 R + 1 W$ for each file $= 2*(3 + 3) = 12$ IO operations
2. **Merge** each pair of sorted chunks *using the external merge algorithm*
 1. $= 2*(3 + 3) = 12$ IO operations
3. **Total cost = 24 IO**

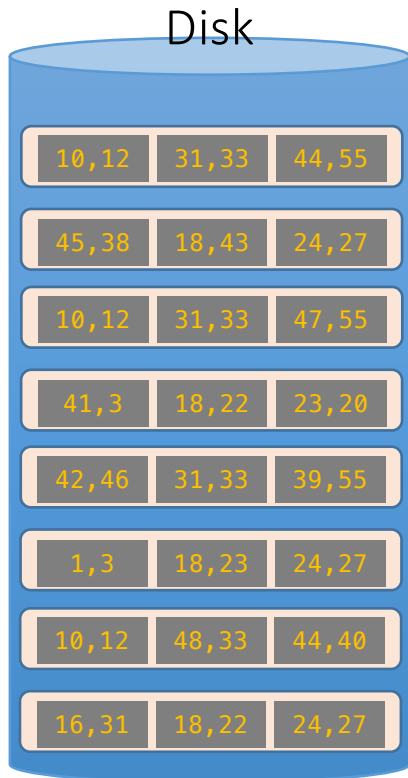
2. External Merge Sort on larger files

Running External Merge Sort on Larger Files



Assume we still only have 3 buffer pages (*Buffer not pictured*)

Running External Merge Sort on Larger Files



1. Split into files small enough to sort in buffer...

Assume we still only have 3 buffer pages (*Buffer not pictured*)

Running External Merge Sort on Larger Files

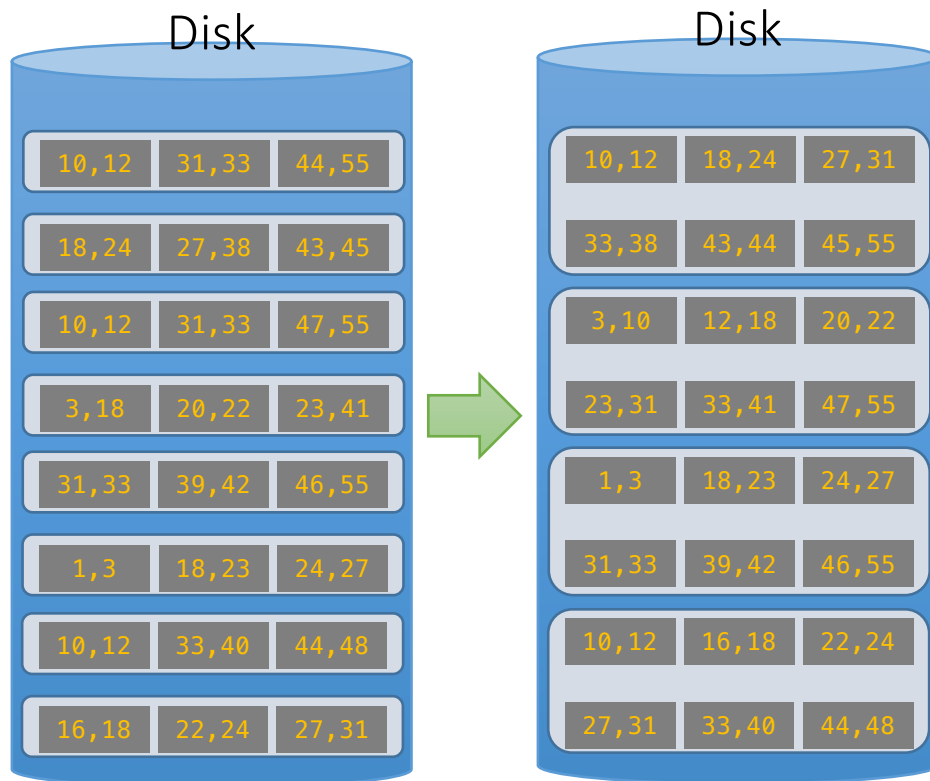


1. Split into files small enough to sort in buffer... and sort

Assume we still only have 3 buffer pages (*Buffer not pictured*)

Call each of these sorted files a *run*

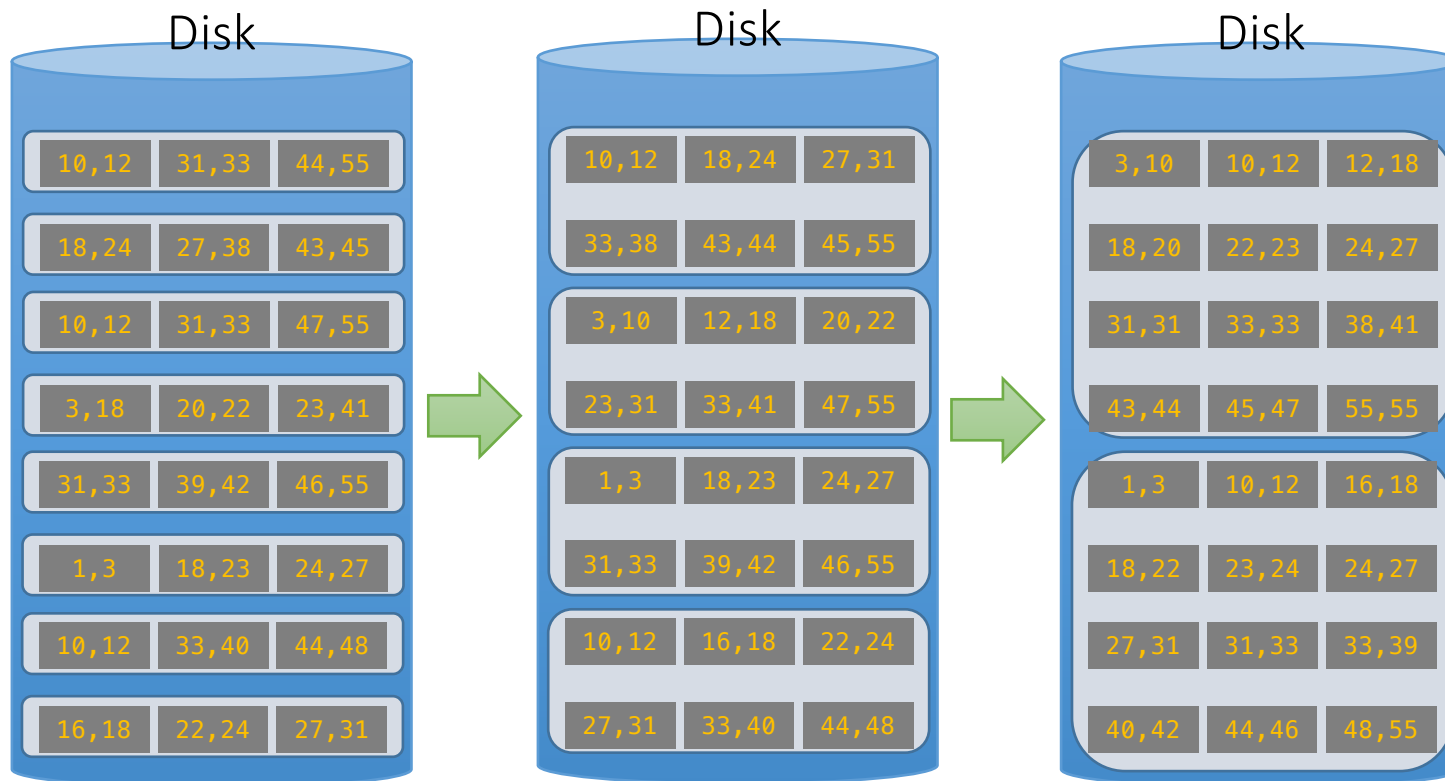
Running External Merge Sort on Larger Files



Assume we still only have 3 buffer pages (*Buffer not pictured*)

2. Now merge pairs of (sorted) files... **the resulting files will be sorted!**

Running External Merge Sort on Larger Files

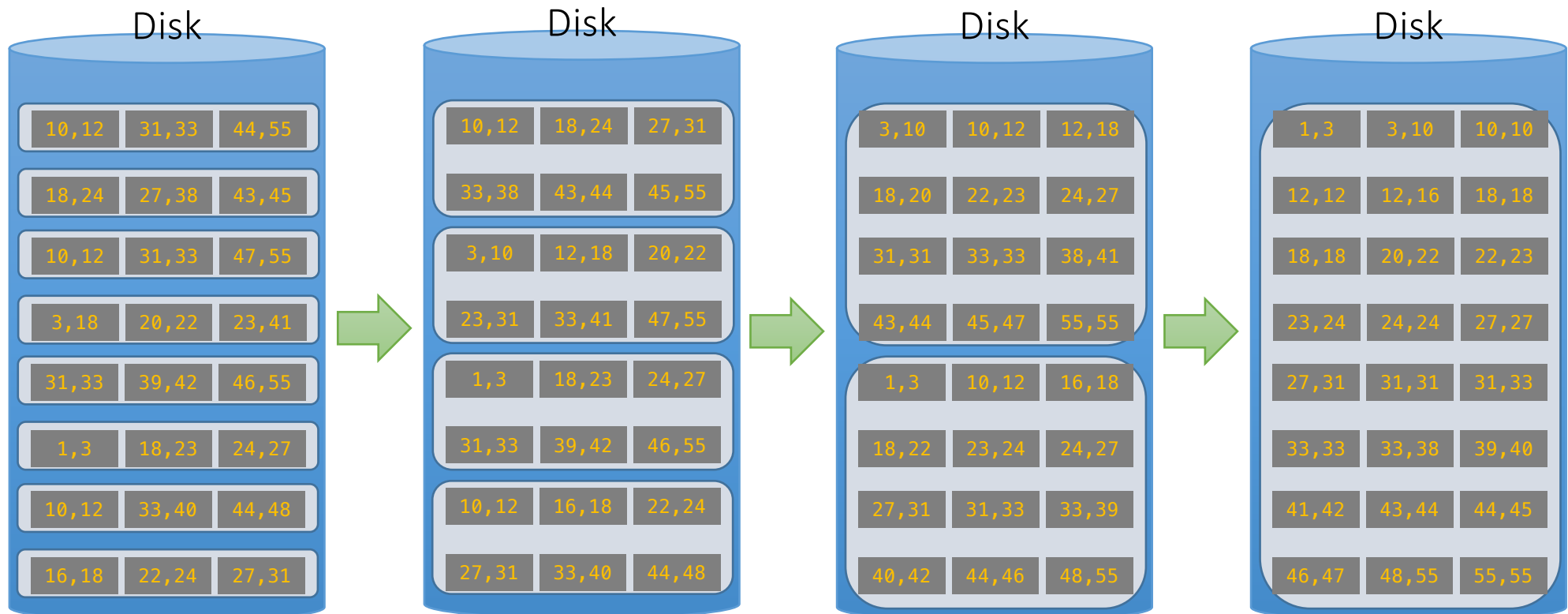


Assume we still only have 3 buffer pages (*Buffer not pictured*)

3. And repeat...

Call each of these steps a *pass*

Running External Merge Sort on Larger Files

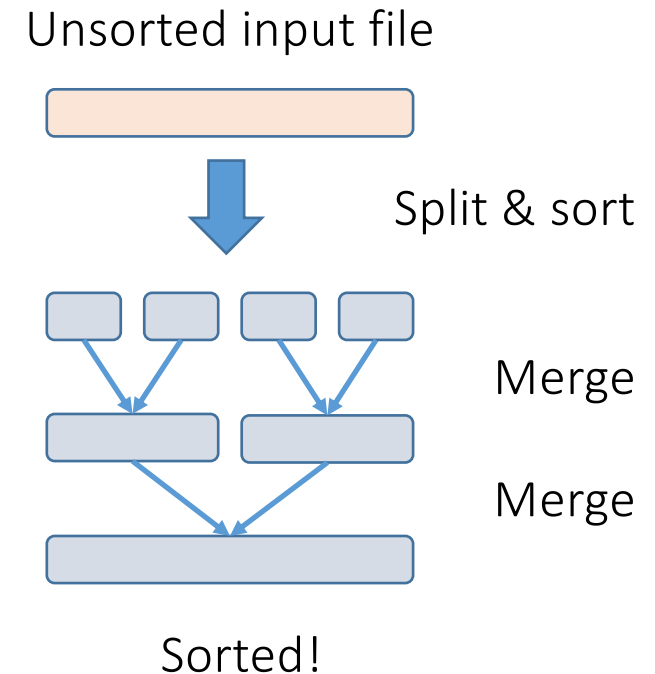


4. And repeat!

Simplified 3-page Buffer Version

Assume for simplicity that we split an N -page file into N single-page **runs** and sort these; then:

- First pass: Merge **$N/2$ pairs of runs** each of length **1 page**
- Second pass: Merge **$N/4$ pairs of runs** each of length **2 pages**
- In general, for N pages, we do $\lceil \log_2 N \rceil$ passes
 - +1 for the initial split & sort
- Each pass involves reading in & writing out all the pages = **$2N$ IO**



→ $2N * (\lceil \log_2 N \rceil + 1)$ total IO cost!

3. Optimizations for sorting

Using $B+1$ buffer pages to reduce # of passes

Suppose we have $B+1$ buffer pages now; we can:

1. Increase length of initial runs. Sort $B+1$ at a time!

At the beginning, we can split the N pages into runs of length $B+1$ and sort these in memory

IO Cost:

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



$$2N(\lceil \log_2 \frac{N}{B+1} \rceil + 1)$$

Starting with runs
of length 1

Starting with runs of
length $B+1$

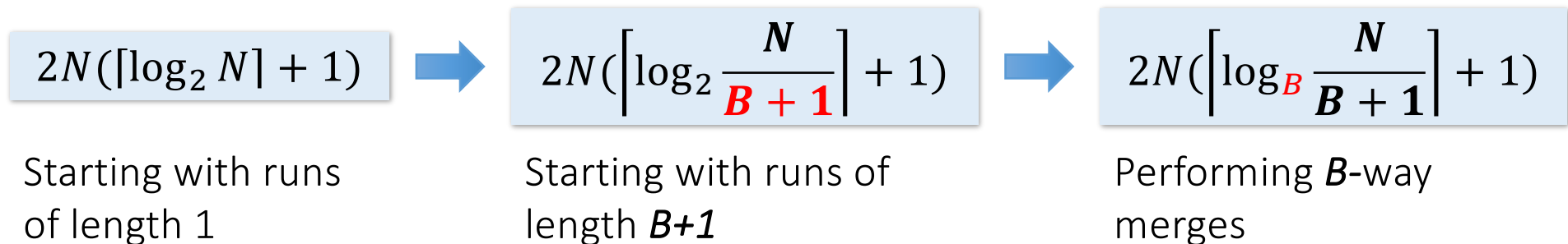
Using $B+1$ buffer pages to reduce # of passes

Suppose we have $B+1$ buffer pages now; we can:

2. Perform a B -way merge.

On each pass, we can merge groups of B runs at a time (vs. merging pairs of runs)!

IO Cost:



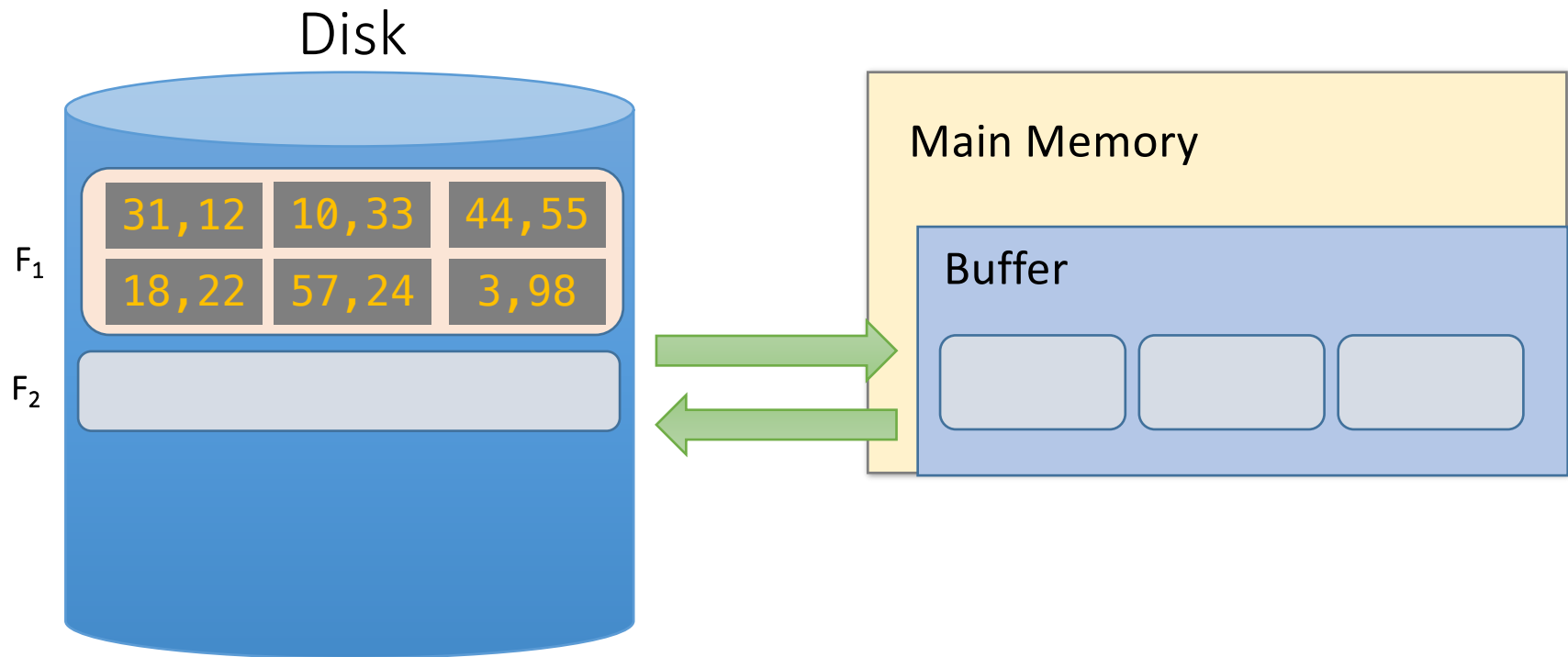
Repacking

Repacking for even longer initial runs

- With $B+1$ buffer pages, we can now start with ***$B+1$ -length initial runs*** (and use ***B -way merges***) to get $2N(\lceil \log_B \frac{N}{B+1} \rceil + 1)$ IO cost...
- Can we reduce this cost more by getting even longer initial runs?
- Use **repacking**- produce longer initial runs by “merging” in buffer as we sort at initial stage

Repacking Example: 3 page buffer

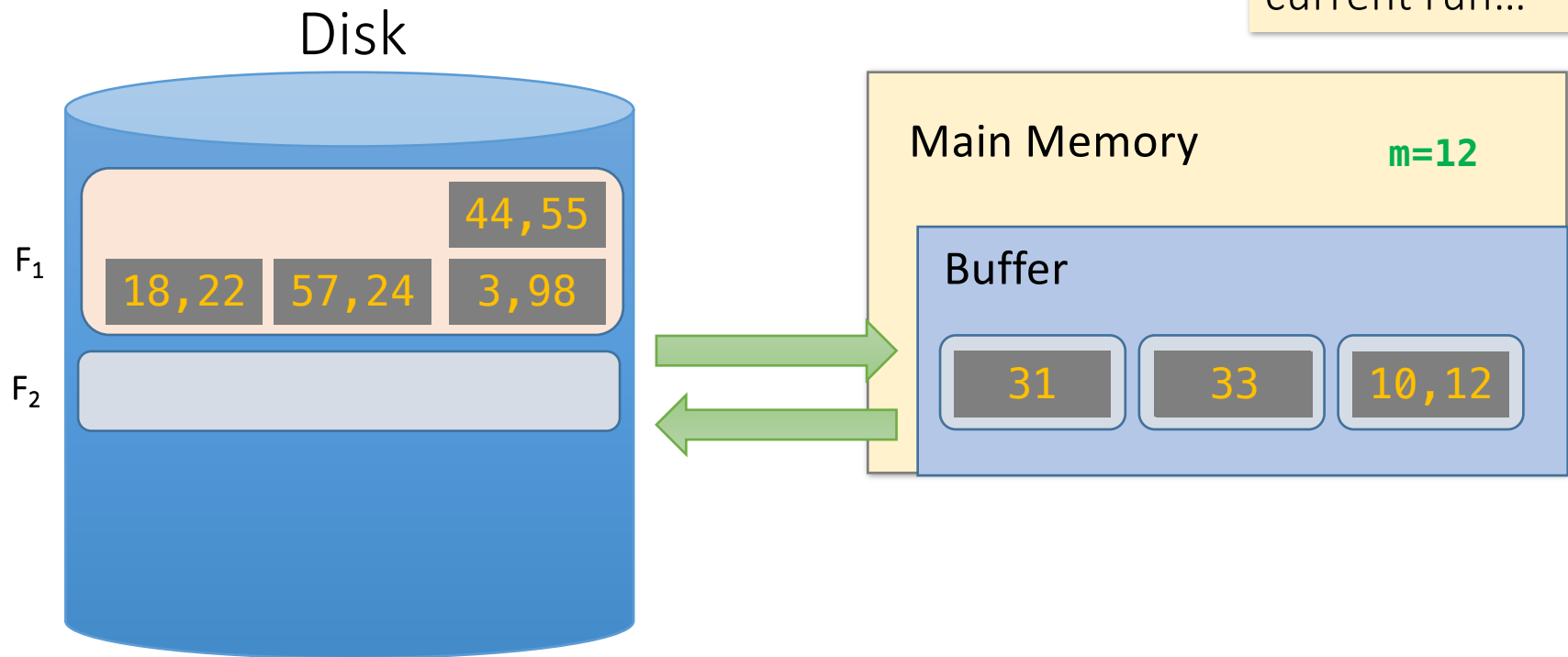
- Start with unsorted single input file, and load 2 pages



Repacking Example: 3 page buffer

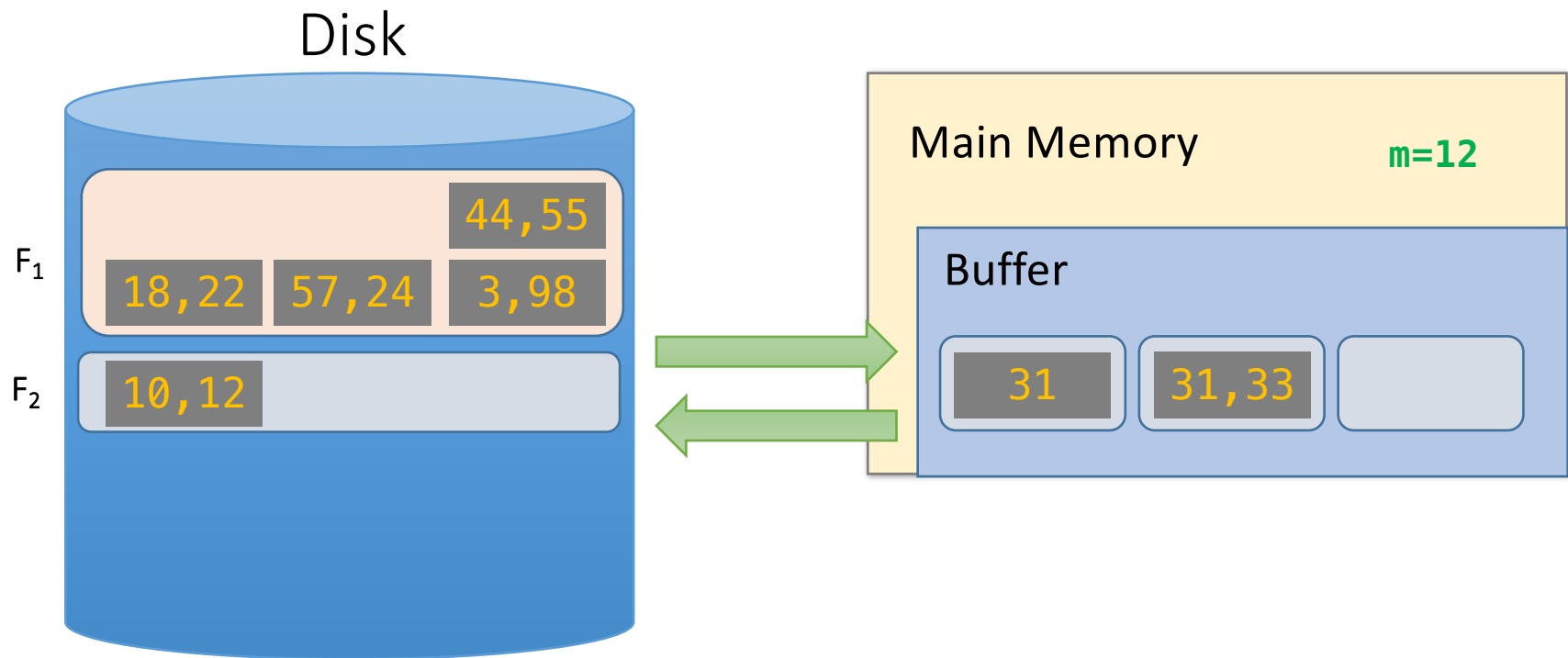
- Take the minimum two values, and put in output page

Also keep track of max (last) value in current run...



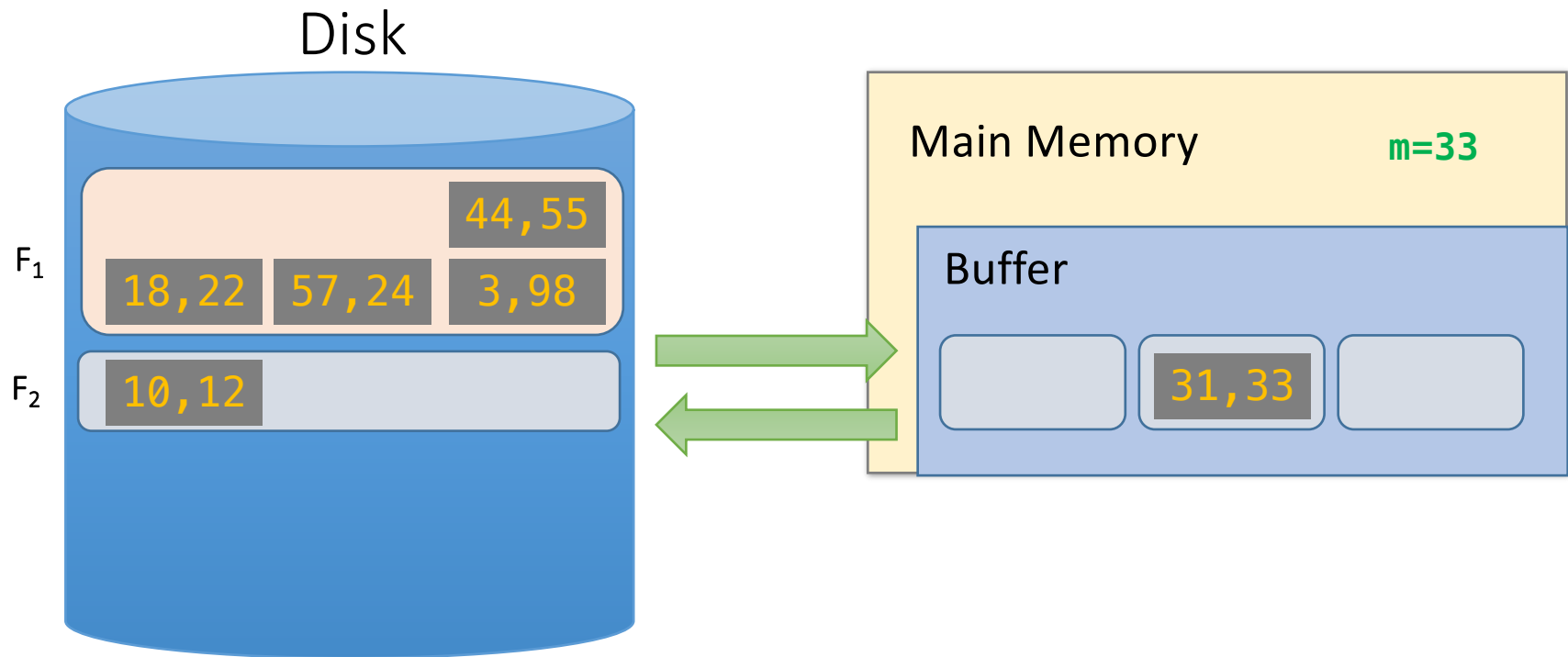
Repacking Example: 3 page buffer

- Next, **repack**



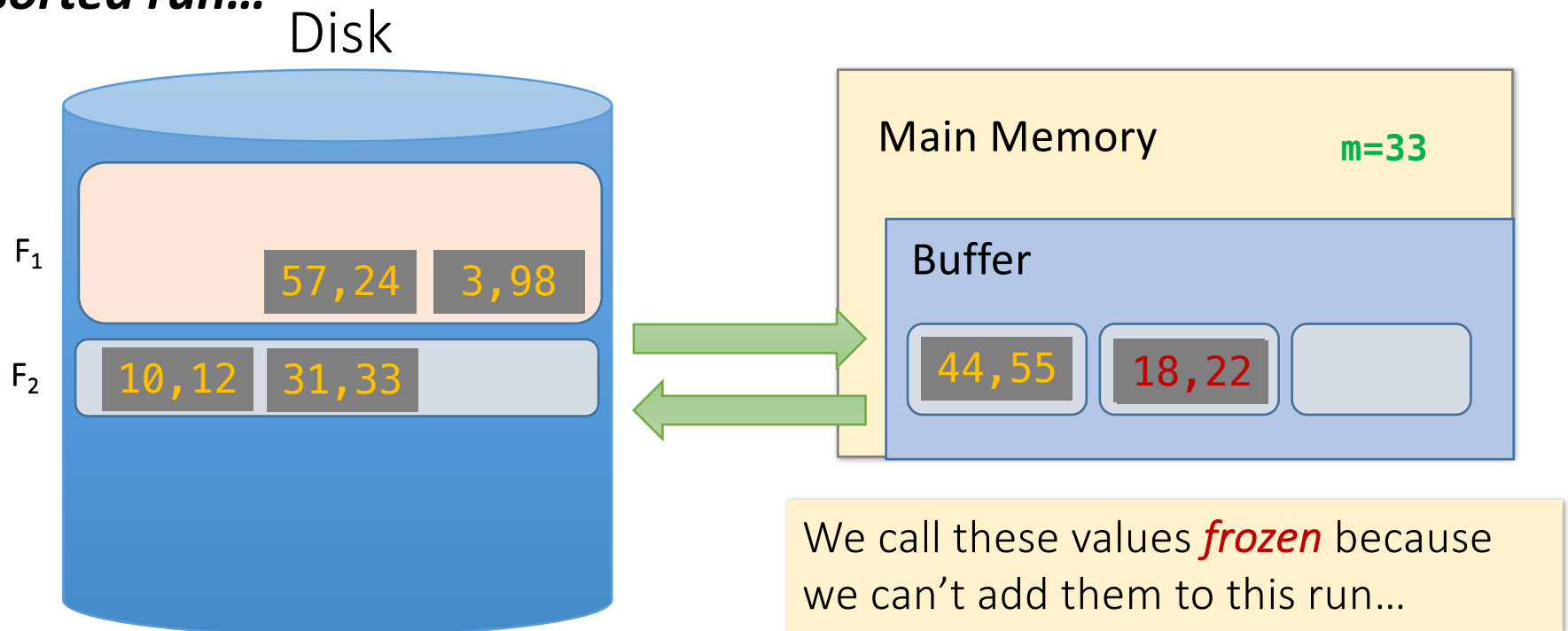
Repacking Example: 3 page buffer

- Next, **repack**, then load another page and continue!



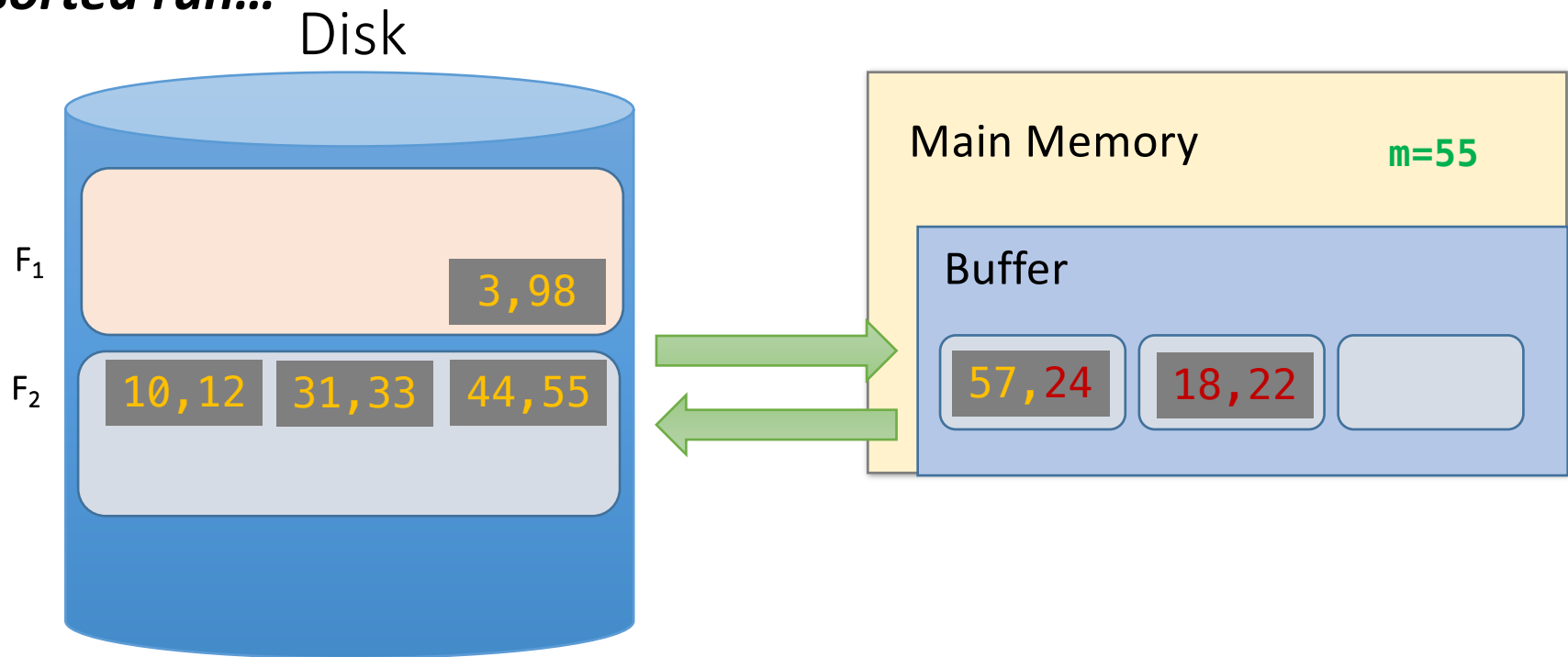
Repacking Example: 3 page buffer

- Now, however, ***the smallest values are less than the largest (last) in the sorted run...***



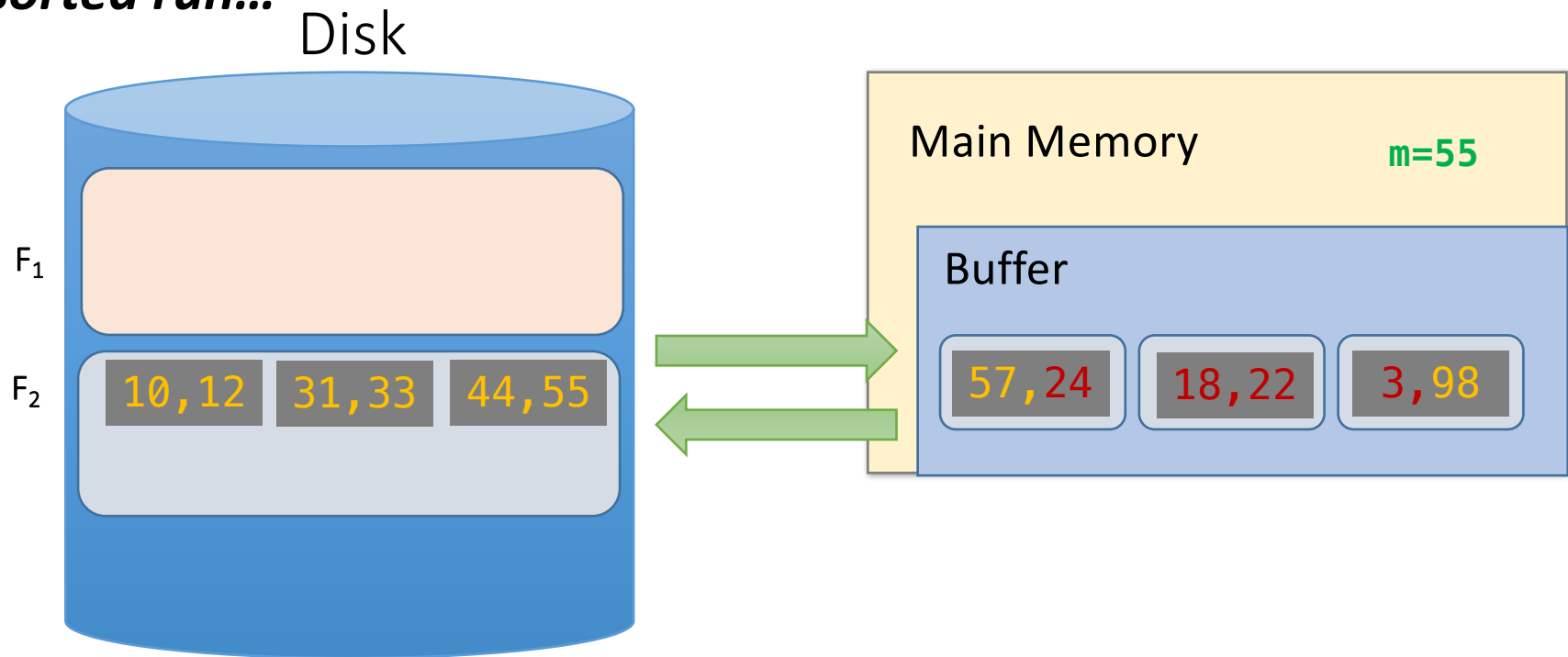
Repacking Example: 3 page buffer

- Now, however, ***the smallest values are less than the largest (last) in the sorted run...***



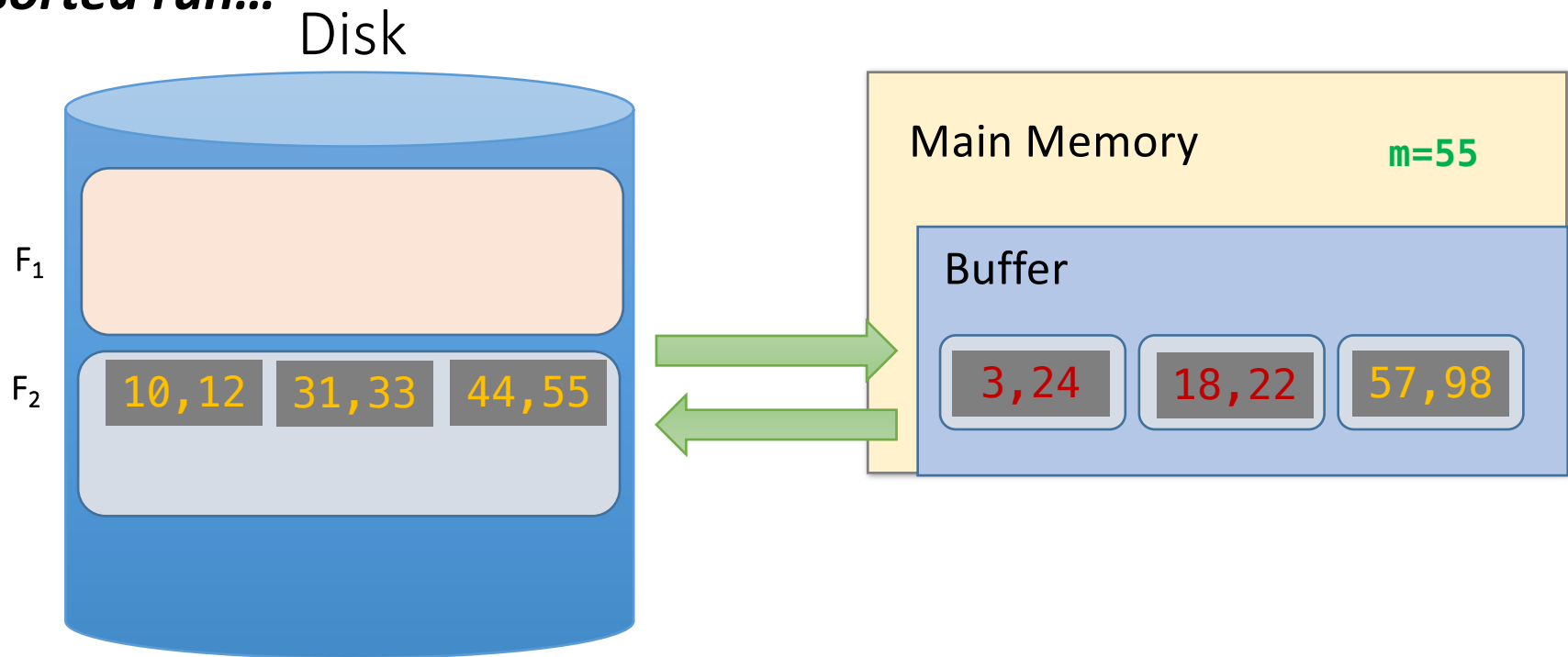
Repacking Example: 3 page buffer

- Now, however, ***the smallest values are less than the largest (last) in the sorted run...***



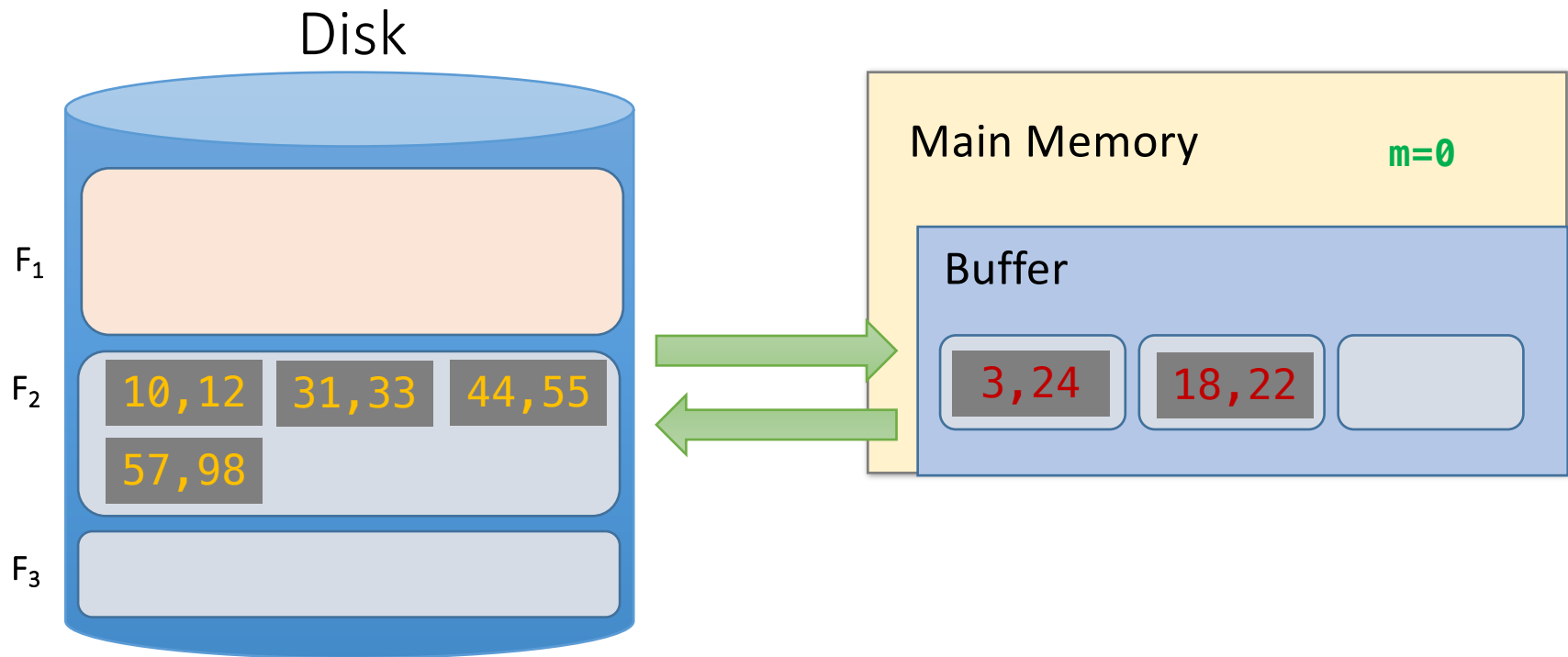
Repacking Example: 3 page buffer

- Now, however, ***the smallest values are less than the largest (last) in the sorted run...***



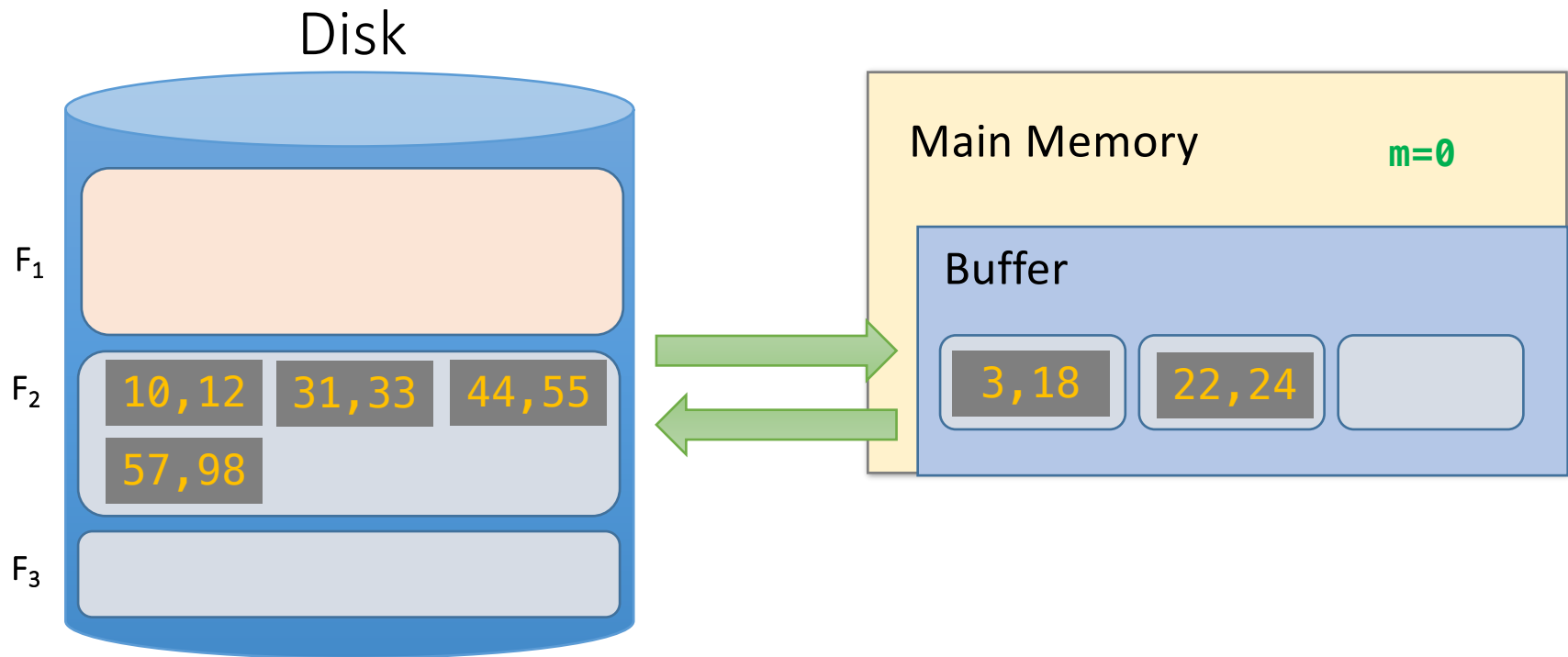
Repacking Example: 3 page buffer

- Once **all buffer pages have a frozen value**, or input file is empty, start new run with the frozen values



Repacking Example: 3 page buffer

- Once ***all buffer pages have a frozen value***, or input file is empty, start new run with the frozen values



Repacking

- Note that, for buffer with $B+1$ pages:
 - If input file is sorted \rightarrow nothing is frozen \rightarrow we get **a single** run!
 - If input file is reverse sorted (worst case) \rightarrow everything is frozen \rightarrow we get runs of length **$B+1$**
- In general, with repacking we do **no worse** than without it!
- What if the file is already sorted?
- Engineer's approximation: runs will have **$\sim 2(B+1)$** length

$$\sim 2N \left(\left\lceil \log_B \frac{N}{2(B+1)} \right\rceil + 1 \right)$$

Summary

- Basics of IO and buffer management.
 - See notebook for more fun! (Learn about *sequential flooding*)
- We introduced the IO cost model using **sorting**.
 - Saw how to do merges with few IOs,
 - Works better than main-memory sort algorithms.
- Described a few optimizations for sorting