**1**

**Big Scale**

**Roadmap**

Hashing      Sorting      Counting

**Hashing-Sorting-Counting** solves "all" known data scale problems :=)

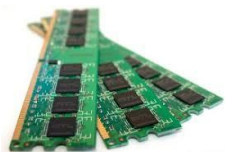+      Boost with a few patterns -- Cache, Parallelize, Pre-fetch

THE BIG IDEA

Note
Works for Relational, noSQL
(e.g.  mySQL, postgres, BigQuery, BigTable, MapReduce, Spark)

**2**

# Big Scale Lego Blocks

# Roadmap

Primary data structures/algorithms

Hashing

Sorting

Counting

HashTables
(hash$_i$(key) --> location)

BucketSort, QuickSort
MergeSort

HashTable + Counter
(hash$_i$(key) --> <count>)

MergeSortedFiles

HashFunctions
(hash$_i$(key) --> location)

MergeSort

MergeSort

HashFunctions
(hash$_i$(key) --> location)

# 3 Why are Sort Algorithms Important?

Why not just use quicksort in main memory??
- How to Sort 10TB - 100 TB of data?
- E.g., with 1GB of RAM, i.e., 0.01-0.001% of data size…

A classic problem in computer science!

Example use cases

1. Query results in sorted order is extremely common
   - e.g., find students in increasing GPA order

2. Core building block for data compression, indexing, joins

# External Merge Sort

**5**
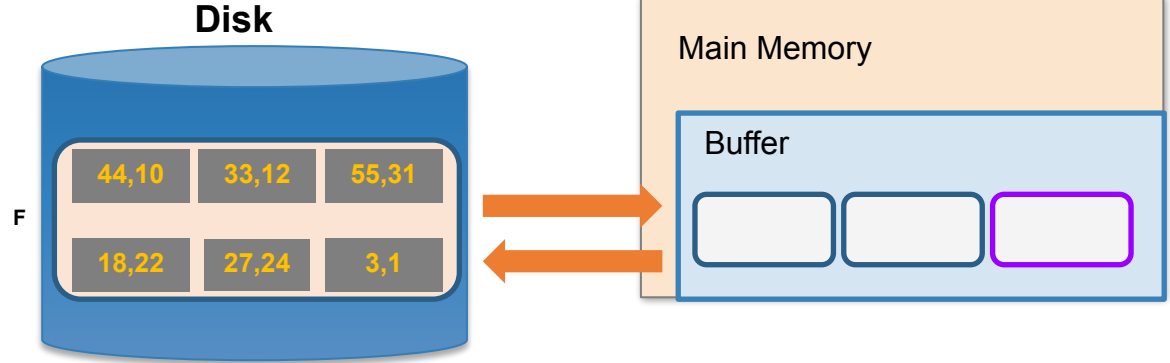
# So how do we sort big files?

1. Split into chunks small enough to sort in memory *("runs")*

2. Merge groups of runs with *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

**Disk**

| 44,10 | 33,12 | 55,31 |
| 18,22 | 27,24 | 3,1 |

F

**Main Memory**

Buffer

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

# External Merge Sort Algorithm

**Disk**

| | | |
|---|---|---|
| $F_1$ | 44,10 | 33,12 | 55,31 |
| $F_2$ | 18,22 | 27,24 | 3,1 |

**Main Memory**

**Buffer**

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

# External Merge Sort Algorithm

**8**



**Disk**

Main Memory

Buffer

F₁

F₂ | 18,22 | 27,24 | 3,1

44,10 | 33,12 | 55,31

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

# External Merge Sort Algorithm

**9**

**Disk**

**Main Memory**

**Buffer**

F₁ 

F₂ | 18,22 | 27,24 | 3,1

10,12 | 31,33 | 44,55

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

# External Merge Sort Algorithm

**10**

**Disk**

**Main Memory**

**Buffer**

F₁: 10,12 | 31,33 | 44,55

Each sorted file is a *run*

F₂: 18,22 | 27,24 | 3,1

Buffer: 1,3 | 18,22 | 24,27

And similarly for $F_2$

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

# External Merge Sort Algorithm
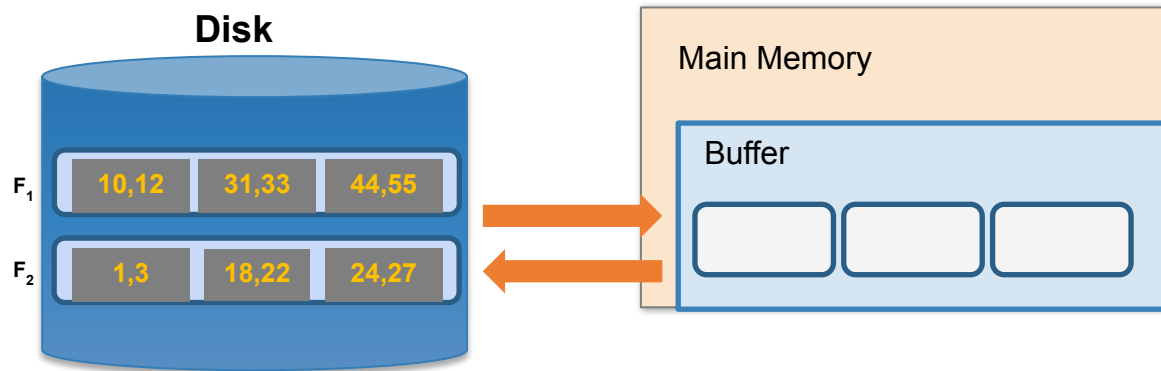


2. Now just run the **external merge** algorithm & we're done!

# Calculating IO Cost

For 3 buffer pages, 6 page file:

1.  Split into <u>two 3-page files</u> and sort in memory
    = 1 R + 1 W per page = 2*(3 + 3) = 12 IOs

2.  Merge each pair of sorted chunks with *external merge algorithm*
    = 2*(3 + 3) = 12 IOs *[Recall: ExtMergeSort in 2*(M+N)]*

3.  Total cost = 24 IO

<u>Note</u>: What's "IO" and how does it map to seek/scans?
●  24 IOs = 24 disk block read/writes
●  Are disk blocks contiguous?
    ○  Cost = 1 seek + time to scan 24 blocks
    ○  Else, cost = 24 seeks + scan 24 blocks

⇒ For such problems, we'll use IO units for simplicity

# 13 Running External Merge Sort on Larger Files

**Disk**

| | | |
|---|---|---|
| 10,12 | 31,33 | 44,55 |
| 45,38 | 18,43 | 24,27 |
| 10,12 | 31,33 | 47,55 |
| 41,3 | 18,22 | 23,20 |
| 42,46 | 31,33 | 39,55 |
| 1,3 | 18,23 | 24,27 |
| 10,12 | 48,33 | 44,40 |
| 16,31 | 18,22 | 24,27 |

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

# Running External Merge Sort on Larger Files

**Disk**

| 10,12 | 31,33 | 44,55 |

| 45,38 | 18,43 | 24,27 |

| 10,12 | 31,33 | 47,55 |

| 41,3 | 18,22 | 23,20 |

| 42,46 | 31,33 | 39,55 |

| 1,3 | 18,23 | 24,27 |

| 10,12 | 48,33 | 44,40 |

| 16,31 | 18,22 | 24,27 |

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

# 15

# Running External Merge Sort on Larger Files

**Disk**

| | | |
|---|---|---|
| 10,12 | 31,33 | 44,55 |
| 18,24 | 27,38 | 43,45 |
| 10,12 | 31,33 | 47,55 |
| 3,18 | 20,22 | 23,41 |
| 31,33 | 39,42 | 46,55 |
| 1,3 | 18,23 | 24,27 |
| 10,12 | 33,40 | 44,48 |
| 16,18 | 22,24 | 27,31 |

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

Each sorted file is a ***run***

# Running External Merge Sort on Larger Files



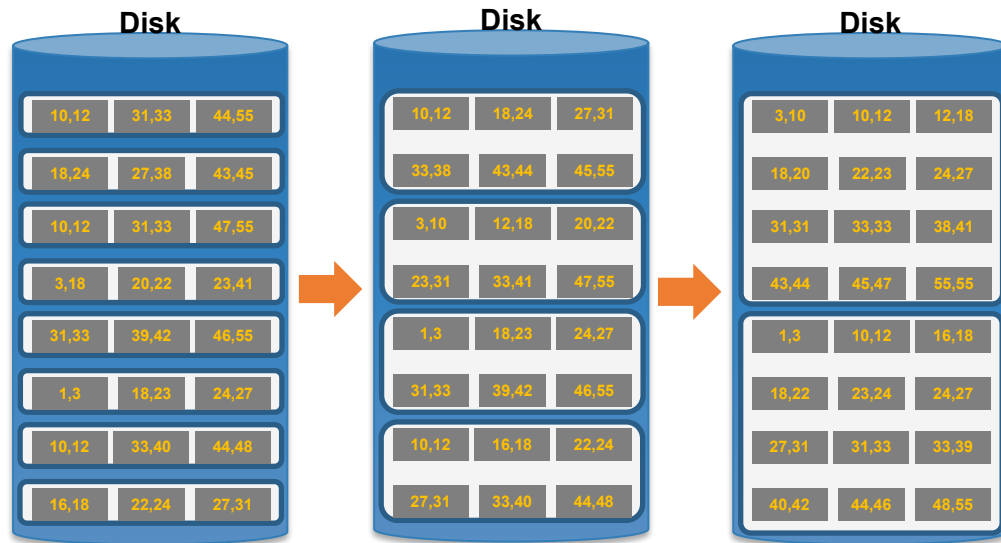**Disk**

| 10,12 | 31,33 | 44,55 |
| 18,24 | 27,38 | 43,45 |
| 10,12 | 31,33 | 47,55 |
| 3,18 | 20,22 | 23,41 |
| 31,33 | 39,42 | 46,55 |
| 1,3 | 18,23 | 24,27 |
| 10,12 | 33,40 | 44,48 |
| 16,18 | 22,24 | 27,31 |

**Disk**

| 10,12 | 18,24 | 27,31 |
| 33,38 | 43,44 | 45,55 |
| 3,10 | 12,18 | 20,22 |
| 23,31 | 33,41 | 47,55 |
| 1,3 | 18,23 | 24,27 |
| 31,33 | 39,42 | 46,55 |
| 10,12 | 16,18 | 22,24 |
| 27,31 | 33,40 | 44,48 |

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

# Running External Merge Sort on Larger Files

**Disk**

| 10,12 | 31,33 | 44,55 |
| 18,24 | 27,38 | 43,45 |
| 10,12 | 31,33 | 47,55 |
| 3,18 | 20,22 | 23,41 |
| 31,33 | 39,42 | 46,55 |
| 1,3 | 18,23 | 24,27 |
| 10,12 | 33,40 | 44,48 |
| 16,18 | 22,24 | 27,31 |

**Disk**

| 10,12 | 18,24 | 27,31 |
| 33,38 | 43,44 | 45,55 |
| 3,10 | 12,18 | 20,22 |
| 23,31 | 33,41 | 47,55 |
| 1,3 | 18,23 | 24,27 |
| 31,33 | 39,42 | 46,55 |
| 10,12 | 16,18 | 22,24 |
| 27,31 | 33,40 | 44,48 |

**Disk**

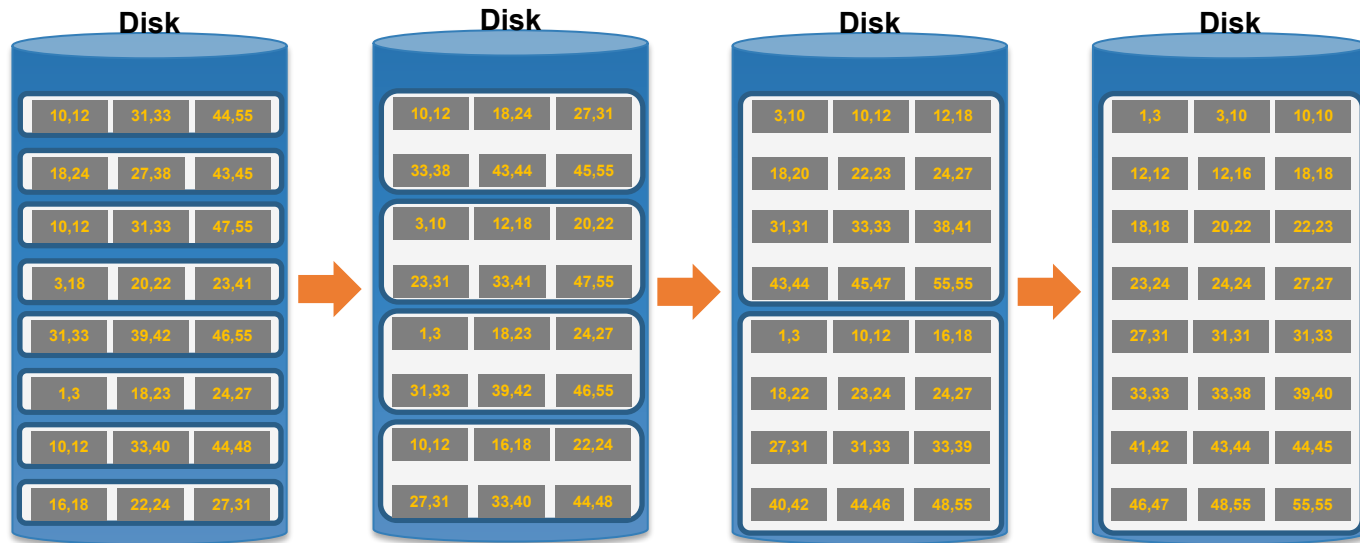| 3,10 | 10,12 | 12,18 |
| 18,20 | 22,23 | 24,27 |
| 31,31 | 33,33 | 38,41 |
| 43,44 | 45,47 | 55,55 |
| 1,3 | 10,12 | 16,18 |
| 18,22 | 23,24 | 24,27 |
| 27,31 | 31,33 | 33,39 |
| 40,42 | 44,46 | 48,55 |

3. And repeat…

Call each of these steps a ***pass***
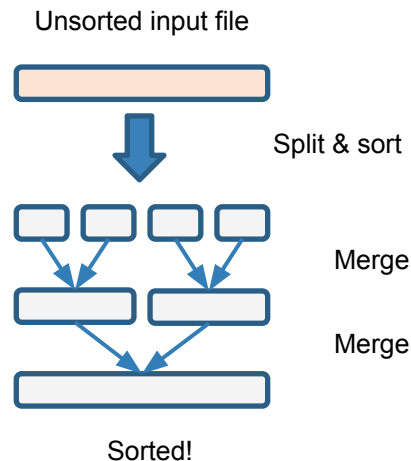
# Running External Merge Sort on Larger Files



4. And repeat!

# Simplified 3-page Buffer Version

For N page file, we do
- Sort step: Sort in 2N IOs
- Merge steps:
  - $\lceil \log_2 N \rceil$ passes
  - 2N IOs/pass
    (each page is read+write once)

Unsorted input file

Split & sort

Merge

Merge

Sorted!

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

# 20 External Merge Sort: Optimizations

Now assume we have **B+1** buffer pages (vs 3 pages in examples so far)

Three optimizations:

1. Increase the length of initial runs
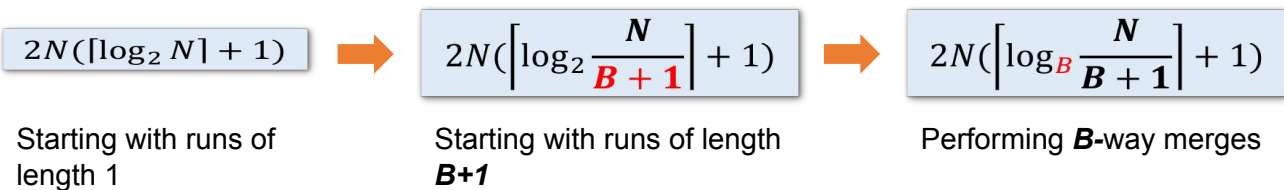
2. B-way merges

3. Repacking

# 21 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!
Split the N pages into runs of length B+1 and sort these in memory

2. **Perform a B-way merge.**
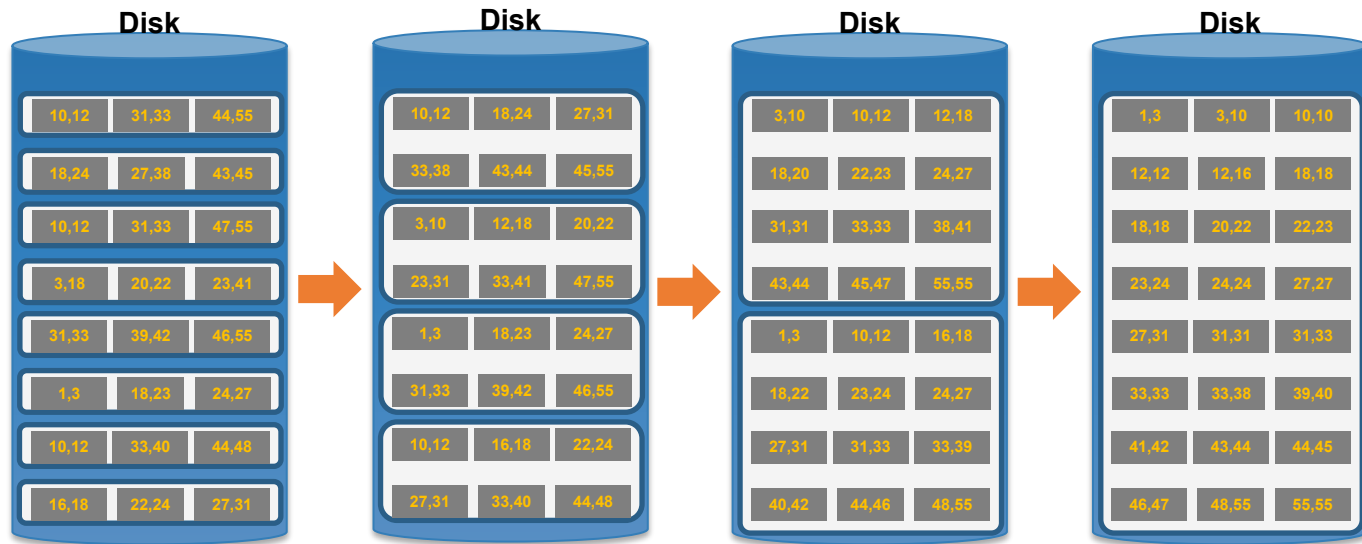On each pass, merge groups of *B* runs at a time (vs. merging pairs of runs)!

IO Cost:

$$2N(\lceil \log_2 N \rceil + 1)$$  ➡️  $$2N\left(\left\lceil \log_2 \frac{N}{B+1} \right\rceil + 1\right)$$  ➡️  $$2N\left(\left\lceil \log_B \frac{N}{B+1} \right\rceil + 1\right)$$

Starting with runs of length 1

Starting with runs of length **B+1**

Performing **B**-way merges

Pretty fast IO aware sort !!

# Repacking for longer runs (Optimization)

**Disk**

| | | |
|---|---|---|
| 10,12 | 31,33 | 44,55 |
| 18,24 | 27,38 | 43,45 |
| 10,12 | 31,33 | 47,55 |
| 3,18 | 20,22 | 23,41 |
| 31,33 | 39,42 | 46,55 |
| 1,3 | 18,23 | 24,27 |
| 10,12 | 33,40 | 44,48 |
| 16,18 | 22,24 | 27,31 |

**Disk**

| | | |
|---|---|---|
| 10,12 | 18,24 | 27,31 |
| 33,38 | 43,44 | 45,55 |
| 3,10 | 12,18 | 20,22 |
| 23,31 | 33,41 | 47,55 |
| 1,3 | 18,23 | 24,27 |
| 31,33 | 39,42 | 46,55 |
| 10,12 | 16,18 | 22,24 |
| 27,31 | 33,40 | 44,48 |

**Disk**

| | | |
|---|---|---|
| 3,10 | 10,12 | 12,18 |
| 18,20 | 22,23 | 24,27 |
| 31,31 | 33,33 | 38,41 |
| 43,44 | 45,47 | 55,55 |
| 1,3 | 10,12 | 16,18 |
| 18,22 | 23,24 | 24,27 |
| 27,31 | 31,33 | 33,39 |
| 40,42 | 44,46 | 48,55 |

**Disk**

| | | |
|---|---|---|
| 1,3 | 3,10 | 10,10 |
| 12,12 | 12,16 | 18,18 |
| 18,18 | 20,22 | 22,23 |
| 23,24 | 24,24 | 27,27 |
| 27,31 | 31,31 | 31,33 |
| 33,33 | 33,38 | 39,40 |
| 41,42 | 43,44 | 44,45 |
| 46,47 | 48,55 | 55,55 |

Idea: What if it's already 'nearly' or 'partly' sorted?

Can we be smarter with buffer? **Optimistic sorting**

# Repacking Example: 3 page buffer

**23**

Start with unsorted single input file, and load 2 pages

**Disk**

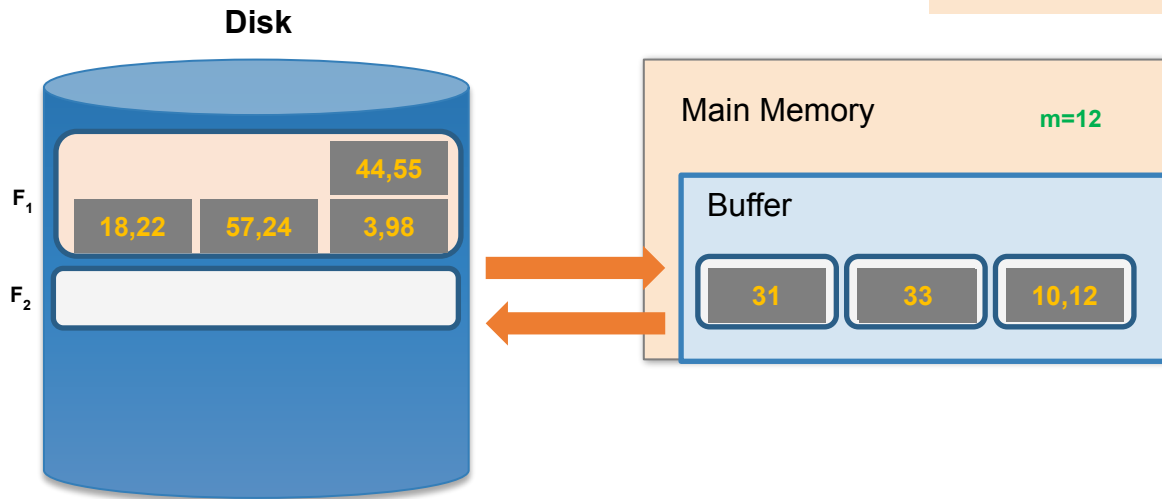| | | |
|---|---|---|
| 31,12 | 10,33 | 44,55 |
| 18,22 | 57,24 | 3,98 |

F₁

F₂

Main Memory

Buffer

# Repacking Example: 3 page buffer

Take the minimum two values, and put in output page
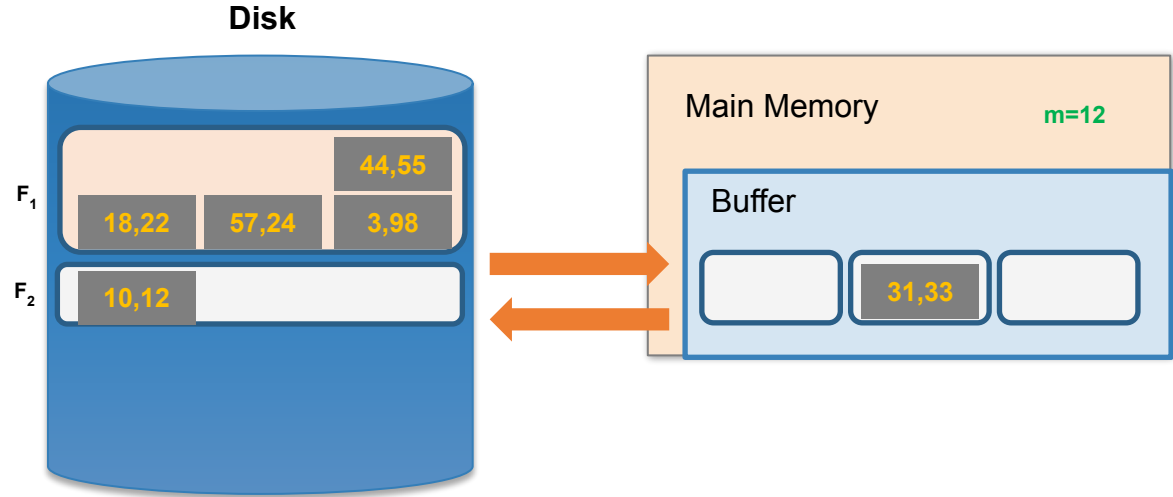
Also keep track of max (last) value in current run…

**Disk**

F₁: 44,55 | 18,22 | 57,24 | 3,98

F₂:

**Main Memory** m=12

**Buffer:** 31 | 33 | 10,12

# Repacking Example: 3 page buffer

- Next, *repack*

**Disk**

**Main Memory**

**m=12**

**Buffer**

$F_1$

44,55

18,22 | 57,24 | 3,98

$F_2$

10,12

31,33

# Repacking Example: 3 page buffer

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

**Disk**

**F₁**

| 44,55 |
|---|
| 18,22 | 57,24 | 3,98 |

**F₂**

| 10,12 |
|---|

**Main Memory**  m=33

**Buffer**

| | 31,33 | |
|---|---|---|

# Repacking Example: 3 page buffer

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

**Disk**

F₁ | 57,24 | 3,98

F₂ | 10,12 | 31,33 | 44,55

Main Memory          **m=55**

Buffer

| 18,22 |

We call these values ***frozen*** because we can't add them to this run…

# 29

## Repacking Example: 3 page buffer

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

**Disk**

Main Memory      **m=55**

$F_1$

3,98

$F_2$    10,12    31,33    44,55

Buffer

57,24    18,22

# Repacking Example: 3 page buffer

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

**Disk**

$F_1$

$F_2$ | 10,12 | 31,33 | 44,55 |

**Main Memory**  m=55

Buffer

| 57,24 | 18,22 | 3,98 |

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



**Disk**

$F_1$

$F_2$ | 10,12 | 31,33 | 44,55 |
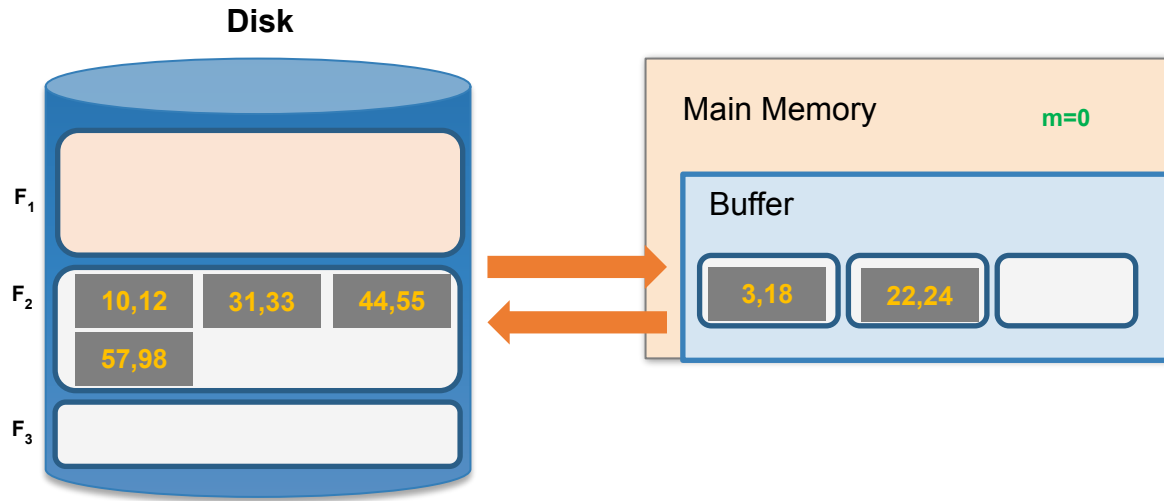
Main Memory    **m=55**

Buffer

| 3,24 | 18,22 | 57,98 |

# 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

**Disk**

Main Memory     $m=0$

Buffer

F$_1$

F$_2$   10,12   31,33   44,55

57,98

F$_3$

3,24   18,22

# 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

**Disk**

# Repacking

- Note that, for buffer with B+1 pages:
  - Best case: If input file is sorted → nothing is frozen → we get a single run!
  - Worst case: If input file is reverse sorted → everything is frozen → we get runs of length B+1

- In general, with repacking we do <u>no worse</u> than without it!

- Engineer's approximation: runs will have ~2(B+1) length

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

# 35 10 TB Sorting Example

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

Sort 10 TB file with 1 GB of RAM
- I.e., File has 156.25 Million Disk Blocks, RAM: 15625 Pages
- I.e., N = 156.25 Million, B = 15624

⇒ Log $_{15624}$ (N/[2(B+1)]) ~= Log $_{15625}$ (5000) = 0.88

⇒ Sort cost = 2N (ceil[0.88]+1) = 4*N IOs

That's AMAZING!!!

Algorithm sorts BIG files (10,000x bigger than RAM) with a small constant factor (4x) on data size

Sorting

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

Sort N pages with B+1 buffer size

*(vs n log n, for n tuples in RAM. Negligible for large data, vs IO -- much, much slower)*

~ 2 N

Sort N pages when N ~= B

     *(because (Log$_B$ 0.5) < 0)*

~ 4 N
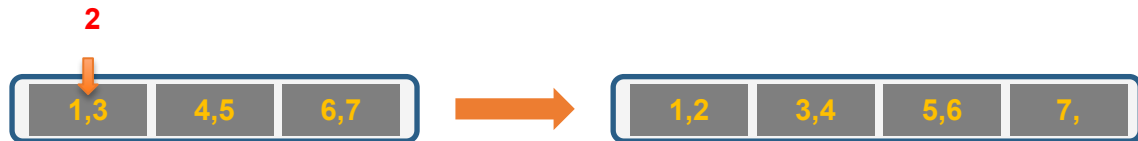
Sort N pages when N ~= 2*B^2

     *(because (Log$_B$ B) = 1)*

We assume cost = 1 IO for read and 1 IO for write.
Alternative IO model (e.g, SSDs in HW#2): 1 IO for read and 8 IOs for write?

# Sorting, with insertions?

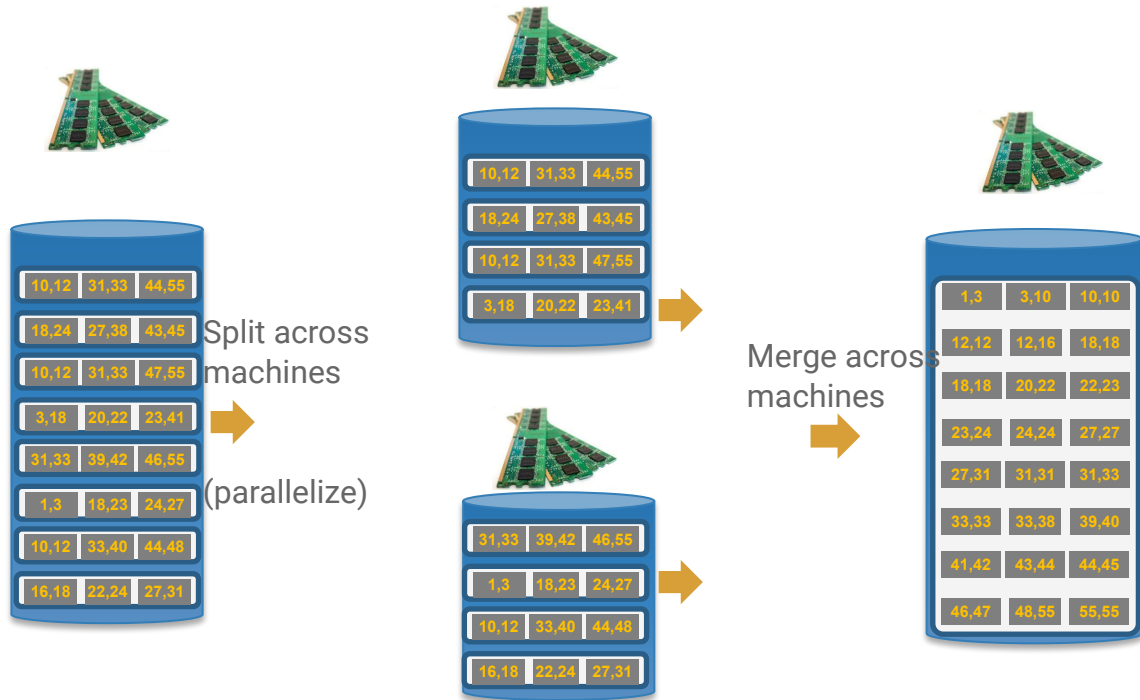- What if we want to **insert** a new person, but keep list sorted?

**2**

| 1,3 | 4,5 | 6,7 |

➡️

| 1,2 | 3,4 | 5,6 | 7, |

- We would have to potentially shift **N** records, requiring **~ 2*N/P** IO (worst case) operations (where P = # of records per page)!
  - We could leave some "slack" in the pages…

Could we get faster insertions?
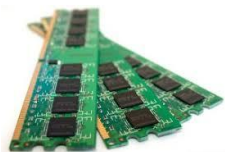(next section)

**38**

Scaling, Speeding Sort (in Cluster)

Split across machines

(parallelize)

Merge across machines

MergeSort locally in each machine (in parallel)

Notes
- Use N machines (N >= 2)
- Could reuse machines
- Speedup at cost of network bandwidth (especially with current data centers)

**39**

Big Scale
Lego Blocks

Roadmap

Primary data structures/algorithms

Hashing

Sorting

Counting

HashTables
(hash$_i$(key) --> location)

BucketSort, QuickSort
MergeSort

HashTable + Counter
(hash$_i$(key) --> <count>)

HashFunctions
(hash$_i$(key) --> location)

MergeSortedFiles

MergeSort

HashFunctions
(hash$_i$(key) --> location)

MergeSort

# Let's build Indexes

# Example [Reminder]

**41**

**CName_Index**

Block #

| CName |
|-------|
| AAPL |
| AAPL |
| AAPL |
| GOOG |
| GOOG |
| GOOG |
| Alibaba |
| Alibaba |

Block #

**Company**

| CName | Date | Price | Country |
|-------|------|-------|---------|
| AAPL | Oct1 | 101.23 | USA |
| AAPL | Oct2 | 102.25 | USA |
| AAPL | Oct3 | 101.6 | USA |
| GOOG | Oct1 | 201.8 | USA |
| GOOG | Oct2 | 201.61 | USA |
| GOOG | Oct3 | 202.13 | USA |
| Alibaba | Oct1 | 407.45 | China |
| Alibaba | Oct2 | 400.23 | China |

How?

**PriceDate_Index**

| Date | Price | Block # |
|------|-------|---------|
| Oct1 | 101.23 | |
| Oct2 | 102.25 | |
| Oct3 | 101.6 | |
| Oct1 | 201.8 | |
| Oct2 | 201.61 | |
| Oct3 | 202.13 | |
| Oct1 | 407.45 | |
| Oct2 | 400.23 | |

1. Index contains search values + Block #: e.g., DB block number.
   - In general, "pointer" to where the record is stored (e.g., RAM page, DB block number or even machine + DB block)
   - Index is conceptually a table. In practice, implemented very efficiently (see how soon)

2. Can have multiple indexes to support multiple search keys

**42**

How to build?

1. How is data organized?

   ▷ Is data in Row or Column store?

   ▷ Is data sorted or not?

2. How do we organize search values?

# 43

## Recall Data Layout

Company(CName, StockPrice, Date, Country)
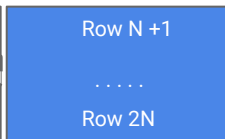
Logical Table

| Company | | Col3 | |
|---------|------|-------|---------|
| CName | Date | Price | Country |
| AAPL | Oct1 | 101.23 | USA |
| AAPL | Oct2 | 102.25 | USA |
| AAPL | Oct3 | 101.6 | USA |
| GOOG | Oct1 | 201.8 | USA |
| GOOG | Oct2 | 201.61 | USA |
| GOOG | Oct3 | 202.13 | USA |
| Alibaba | Oct1 | 407.45 | China |
| Alibaba | Oct2 | 400.23 | China |

Row1, Row3, Row5, Row8

### Row based storage (aka Row Store)

Page — Row1 ..... Row N
Page — Row N +1 ..... Row 2N
Page … n (RAM/DIsk)

### Column based storage (aka Column Store)

Page — CName
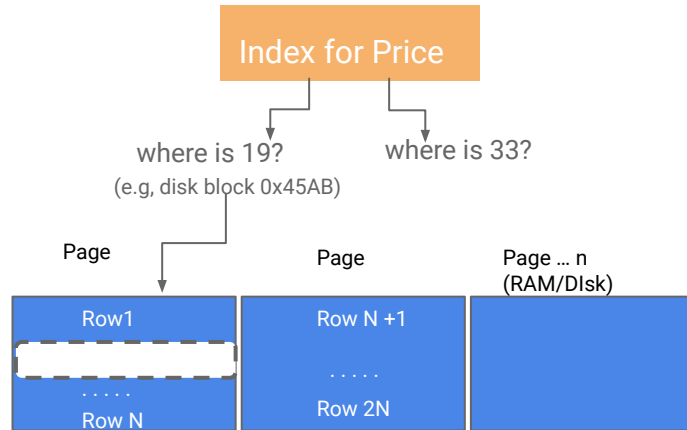Page — Date
Page — Price
Page — Company

# Index on row store

**4**

Query: Search for cname with specific price?
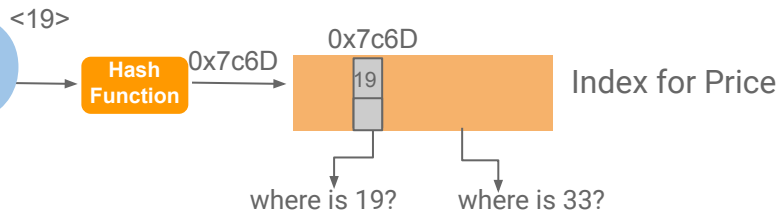⇒ If 'price' is an <u>indexed column</u>, query will be fast.
⇒ 'Price' is <u>search key</u>. Values in price column are <u>search values</u>.

**4**

Index for Price

where is 19?
(e.g, disk block 0x45AB)

where is 33?

Page

Page

Page ... n
(RAM/DIsk)

| Row1 | Row N +1 | |
| | . . . . . | |
| . . . . . | Row 2N | |
| Row N | | |

"Real" data layout, with full records
(including cname, prices, etc.)

# Our 1st Hashing index

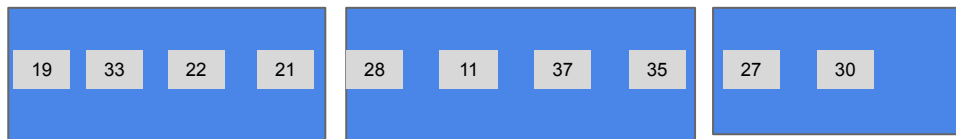**45**

<19>

Hash Function → 0x7c6D

0x7c6D

| 19 |

Index for Price

**Goal of Index: where is location of each value?**

For simplicity, we'll only show the underline{search values} for Price index

Example row: <'goog', price=19, date=Oct 1, ...> as <price=19> or <19>

where is 19?     where is 33?

| 19 | 33 | 22 | 21 |   | 28 | 11 | 37 | 35 |   | 27 | 30 |

Maintain locations of N values

Sorted
| 11 | 19 | 21 | 22 |   | 27 | 28 | 30 |   | 33 | 35 | 37 |

If sorted, will need to maintain locations only of smallest value in each block.

How it works in practice?

1. Schema designer picks a column to keep data sorted by (e.g., price). Index for that column is cheap.

2. For other columns, index will be bigger (e.g., CName)

# Index Types

- Hash Tables
  - IO-aware hashing (e.g., *linear* or *extendible hashing*)

These data structures are "IO aware"

- B-Trees *(covered next)*
  - Very good for range queries, sorted data
  - Some old databases only implemented B-Trees
  - *We will look at a variant called B+ Trees*

**Real difference between structures**:
costs of ops *determines which index you pick and why*