External Sorting

Database System Technology





TA office hours after class



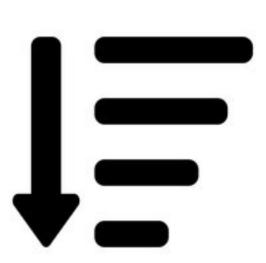
Extra office hours (TBD)

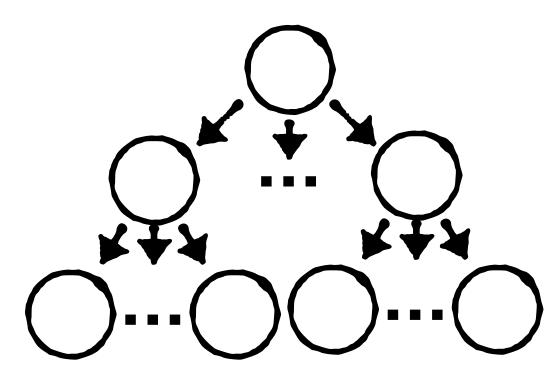
Why do databases need to Sort?

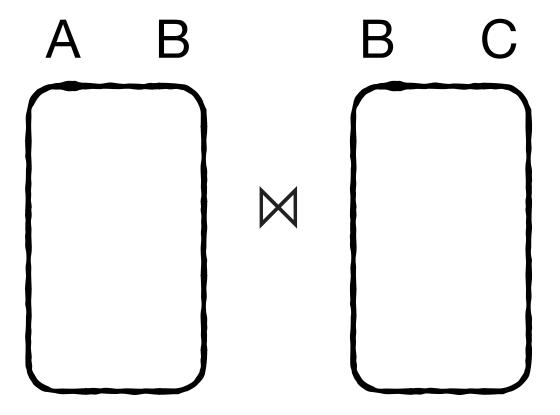
User asking for sorted output (Select...order by...)

Creating an Index

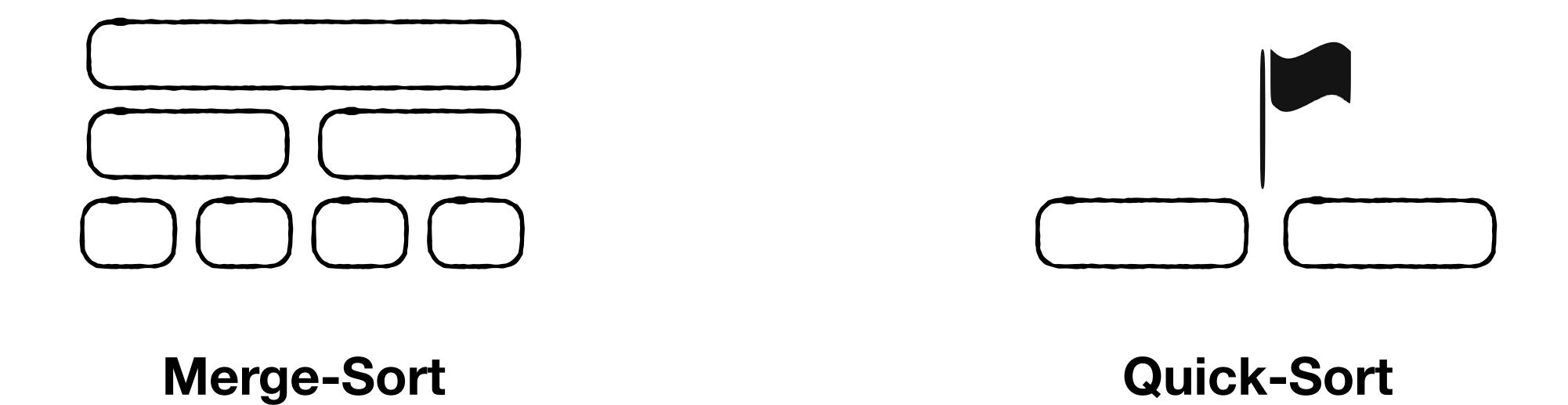
Joining Tables (More later)



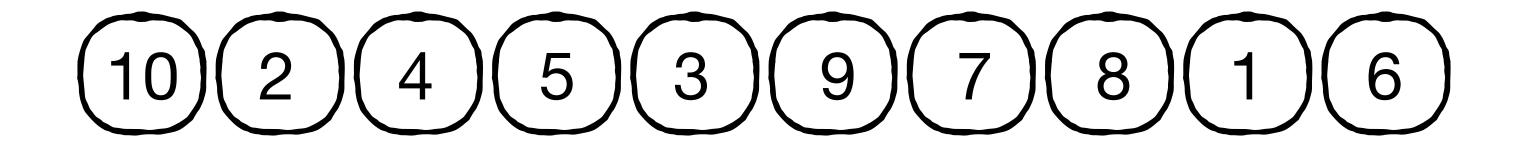




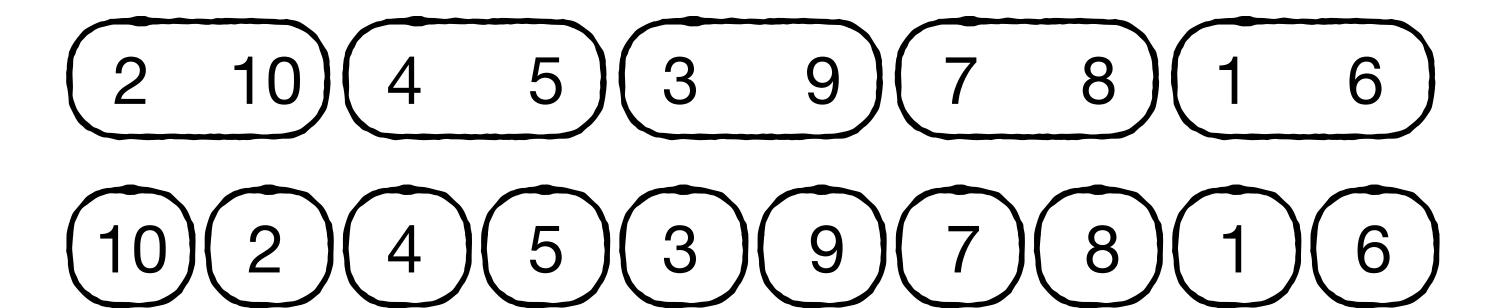
If data fits in memory, traditional sorting algorithms work



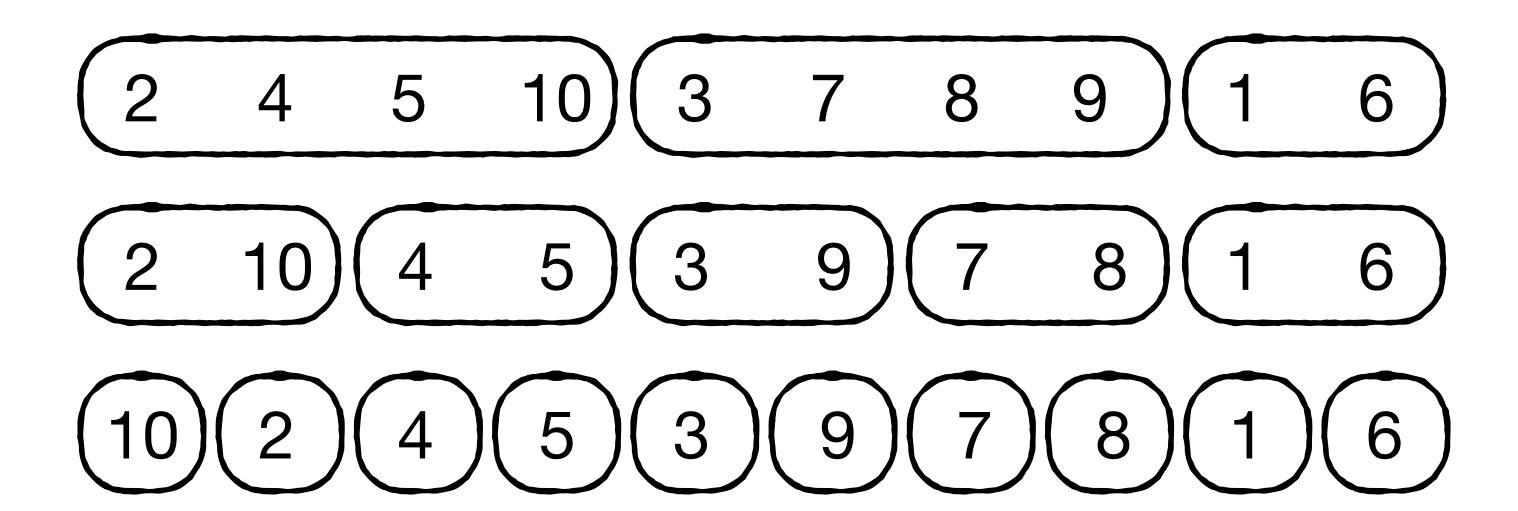
1. Given n entries, divide them into n sublists, each containing one element.



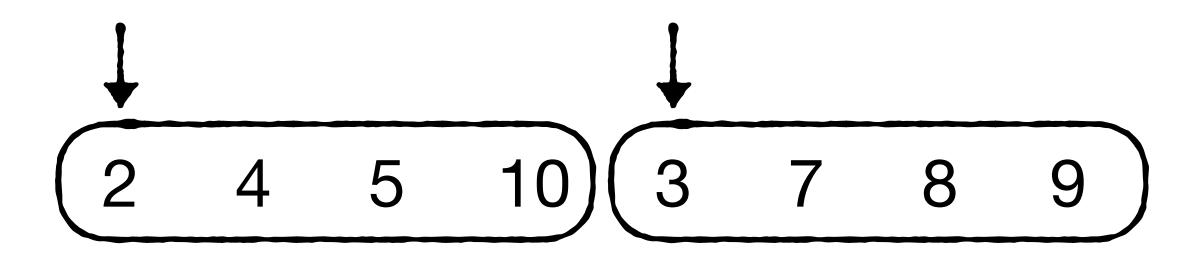
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left



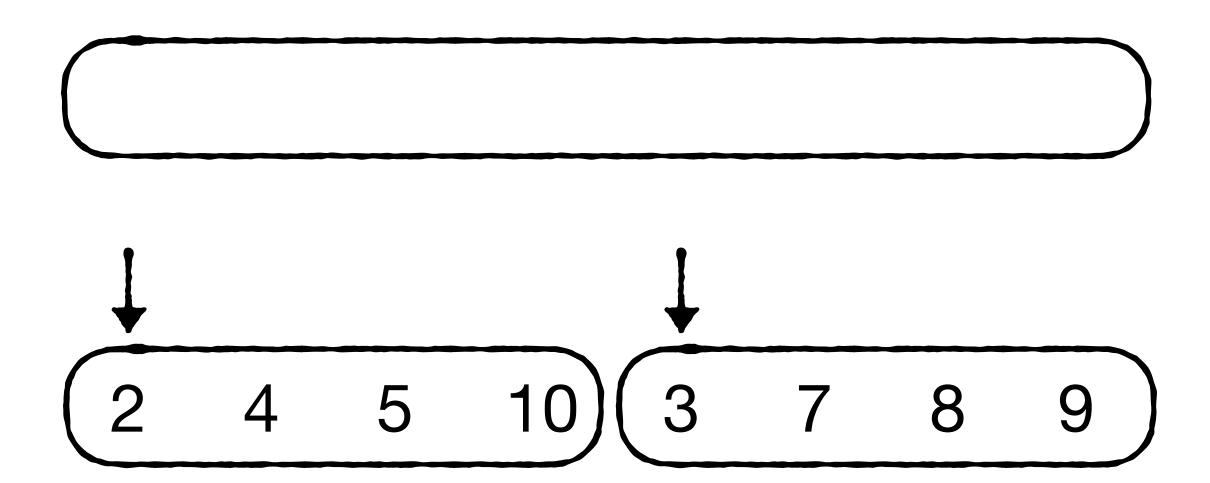
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left



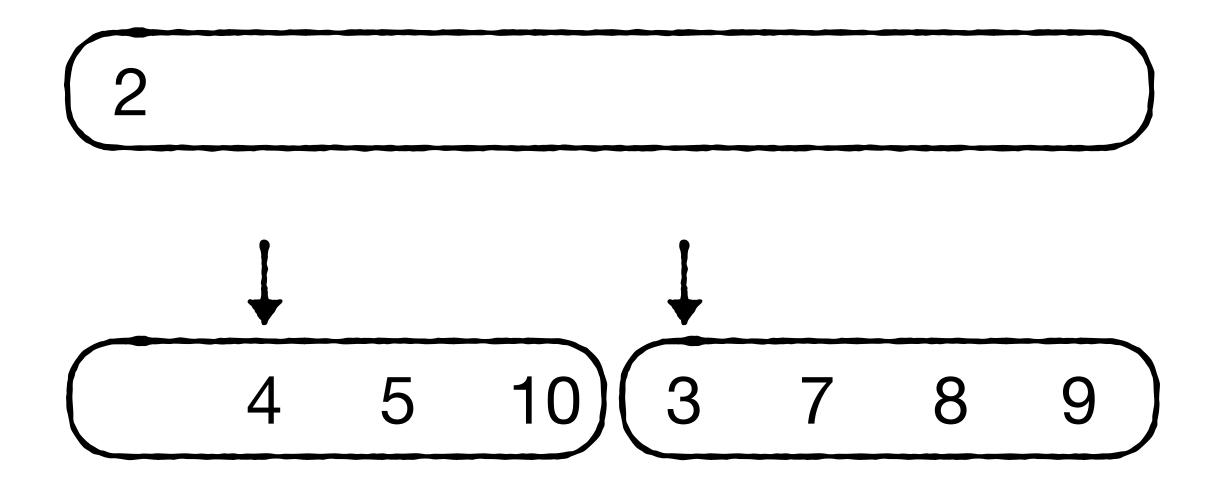
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



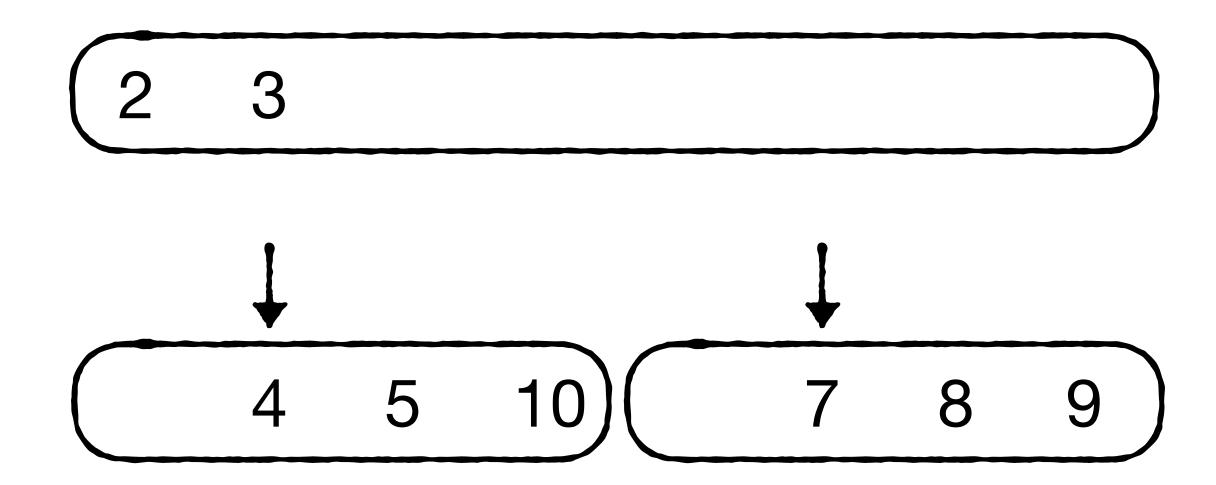
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



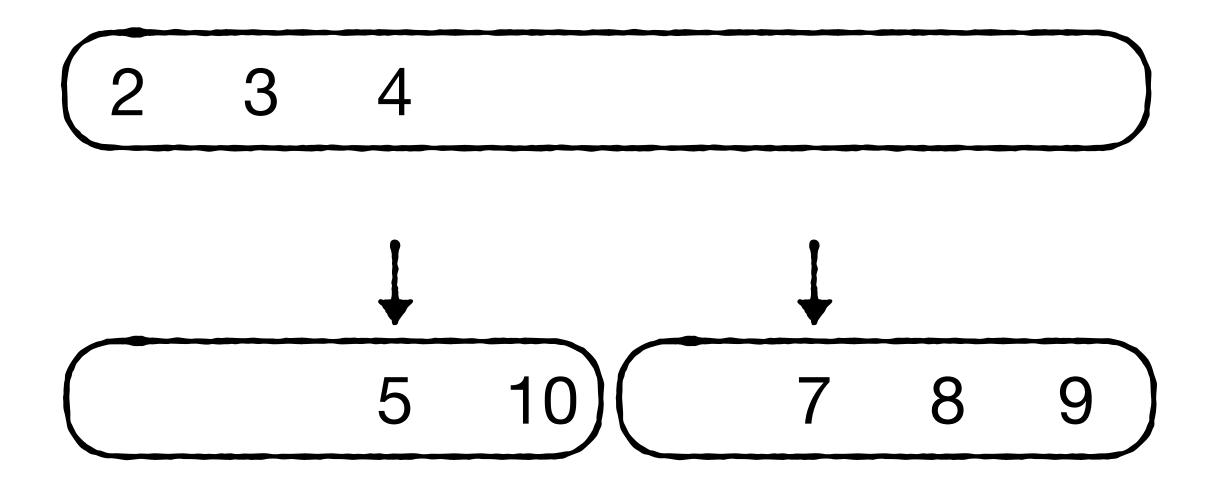
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



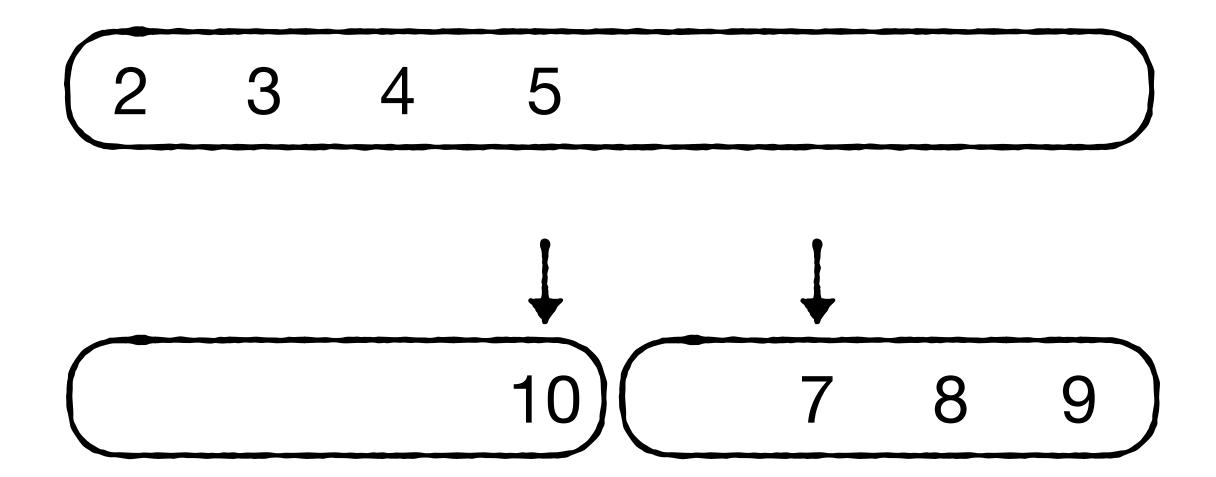
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



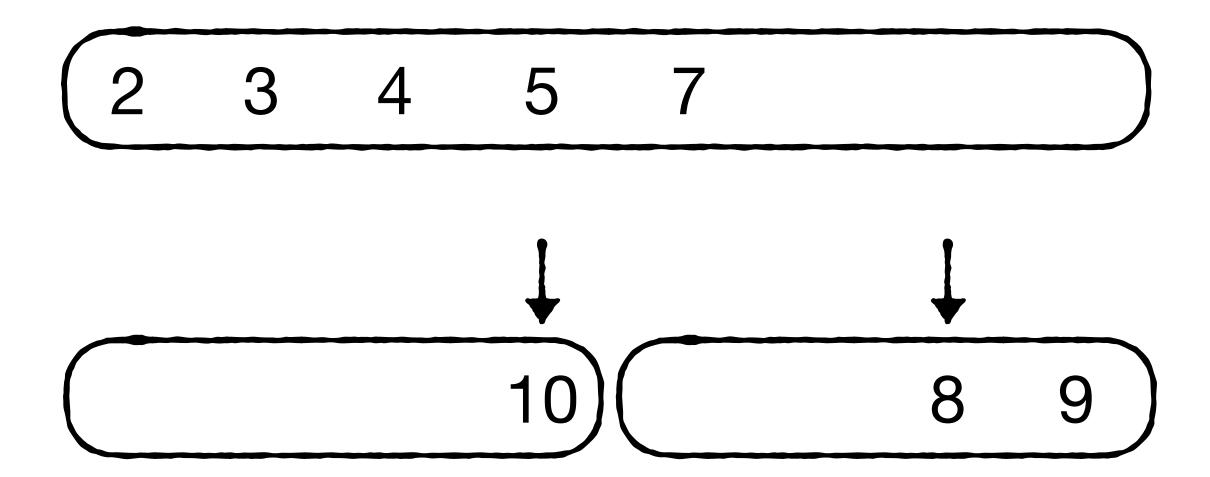
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



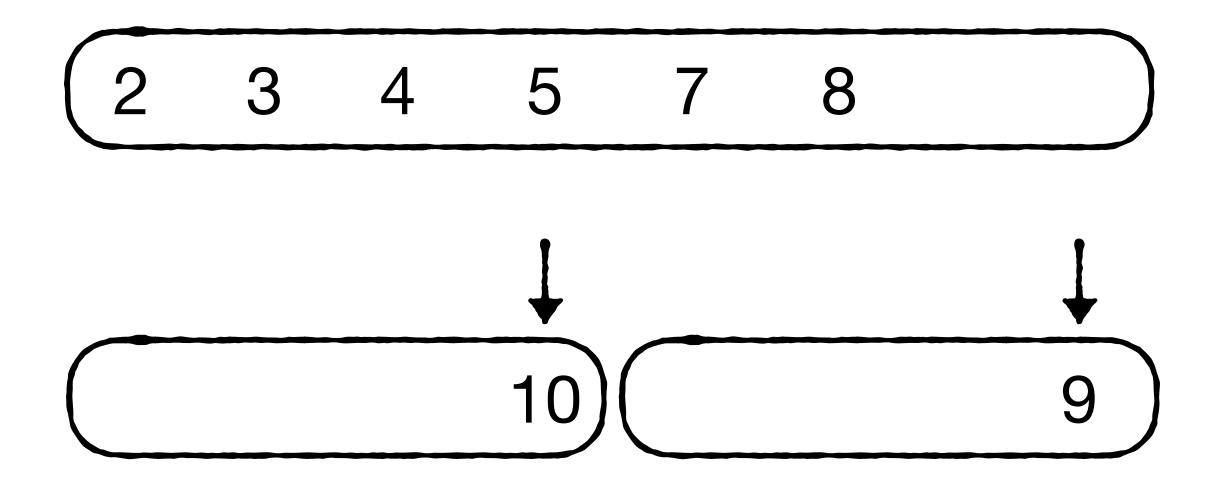
- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step

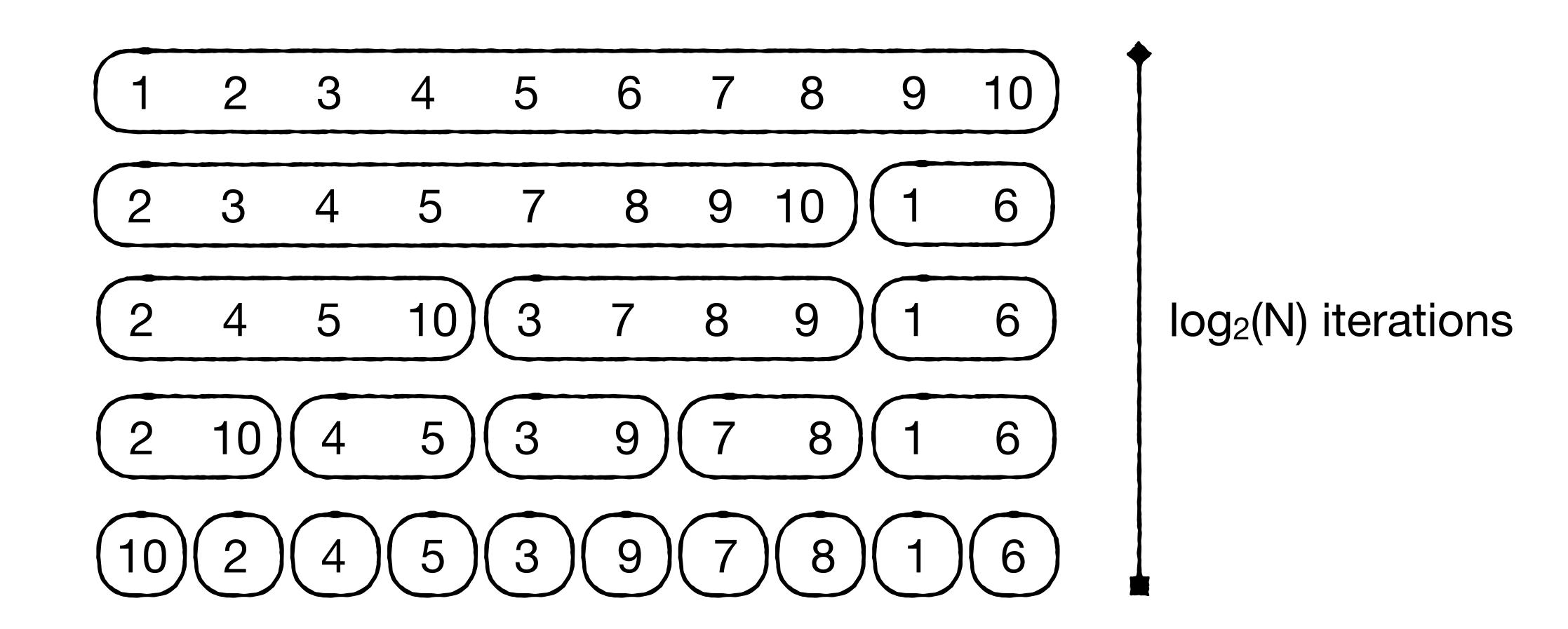


- 1. Given n entries, divide them into n sublists, each containing one element.
- 2. Merge two adjacent sublists at a time until there is one left
- 3. Using cursors, draw the smallest element of the two lists each step



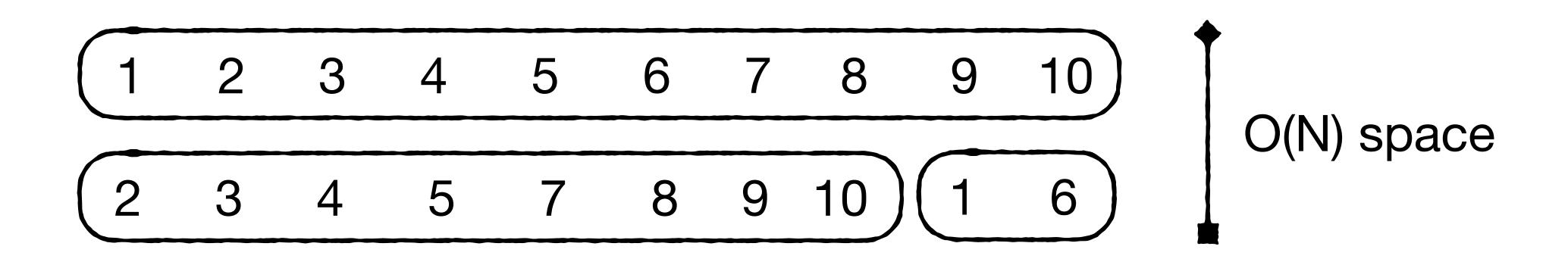
Merge-Sort Analysis

Log(N) iterations, each of which traverses all N elements: O(N log₂(N)) CPU work

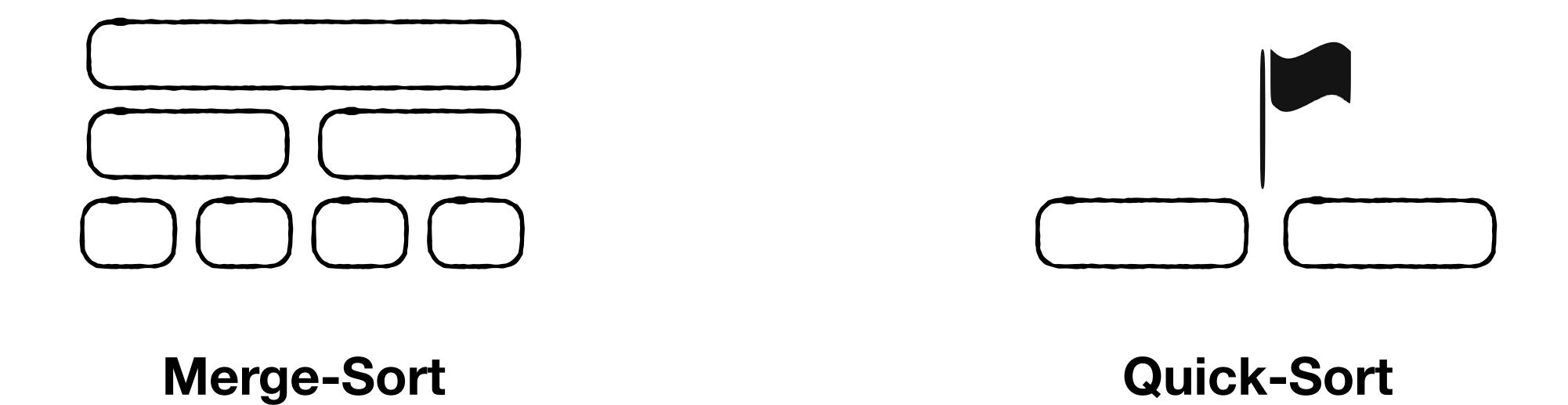


Merge-Sort Analysis

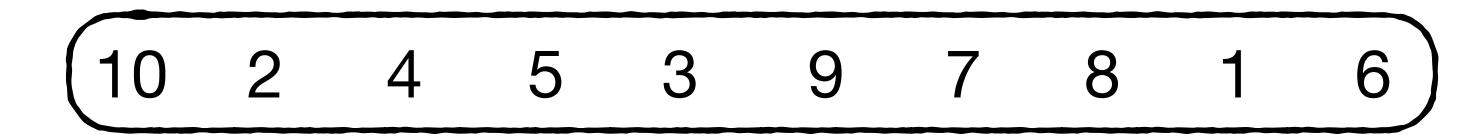
Log(N) iterations, each of which traverses all N elements: O(N log₂(N)) CPU work We need O(N) space by maintaining at most two lists at a time



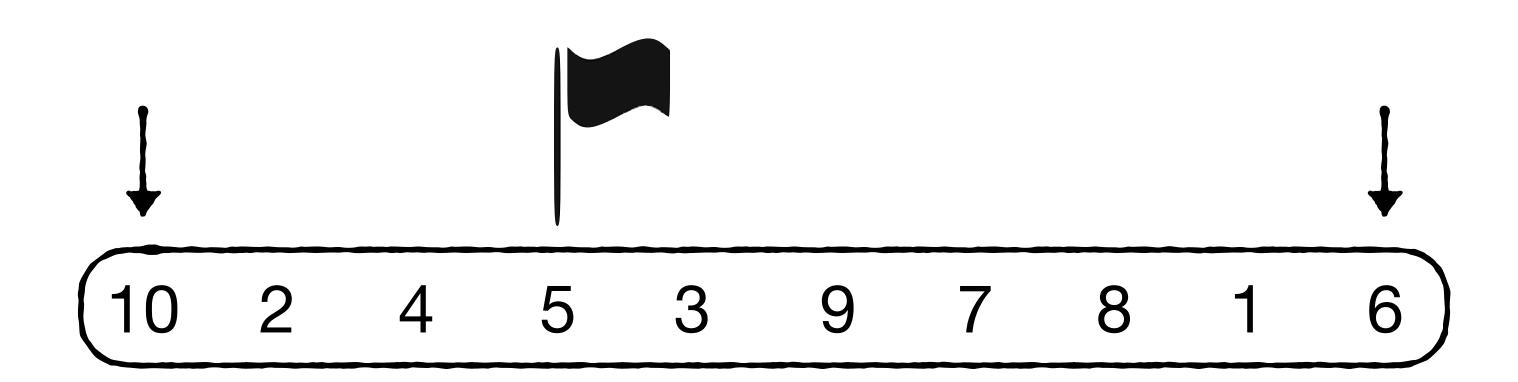
If data fits in memory, traditional sorting algorithms work



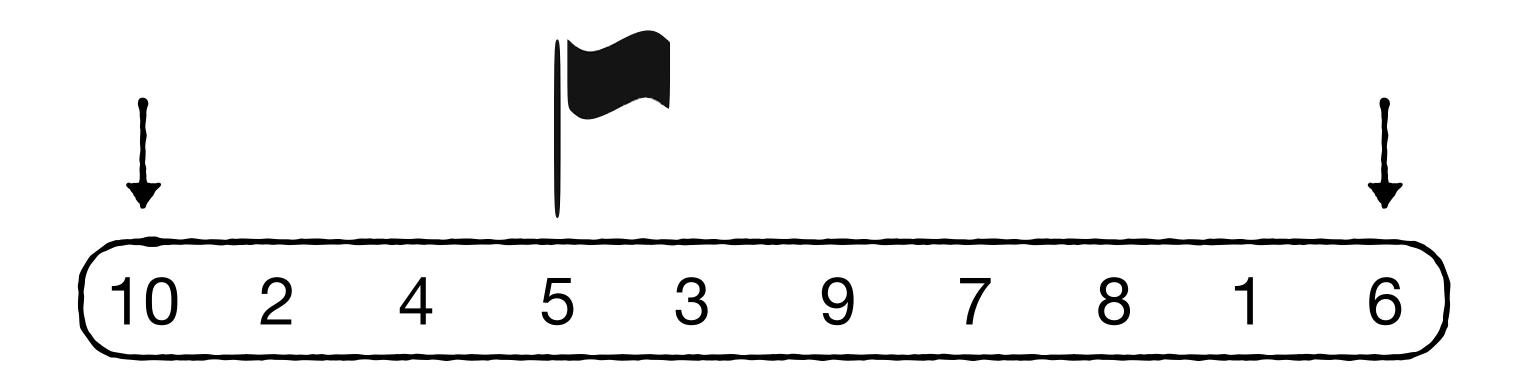
1. Pick random pivot point and initialize two pointers at ends



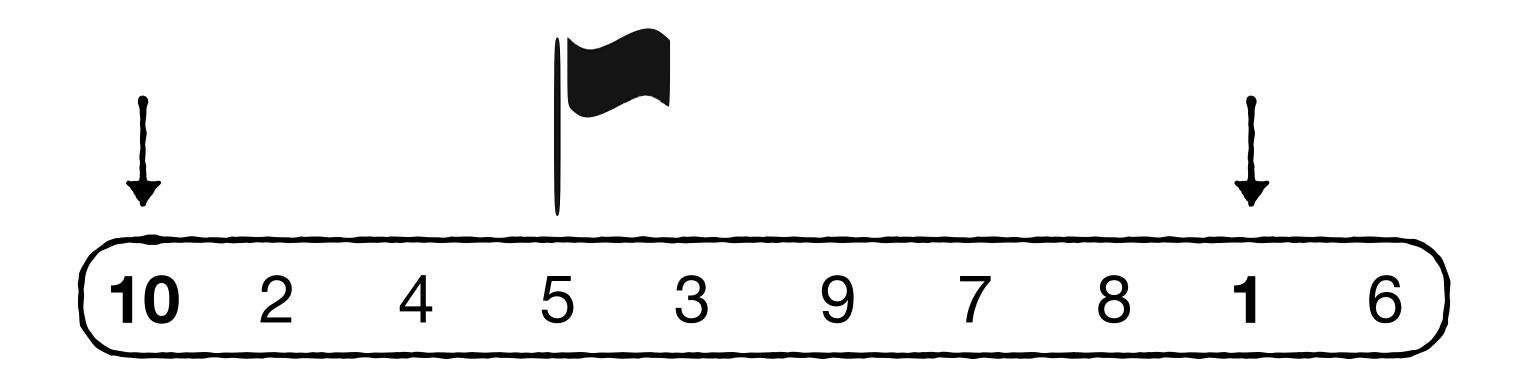
1. Pick random pivot point and initialize two pointers at ends



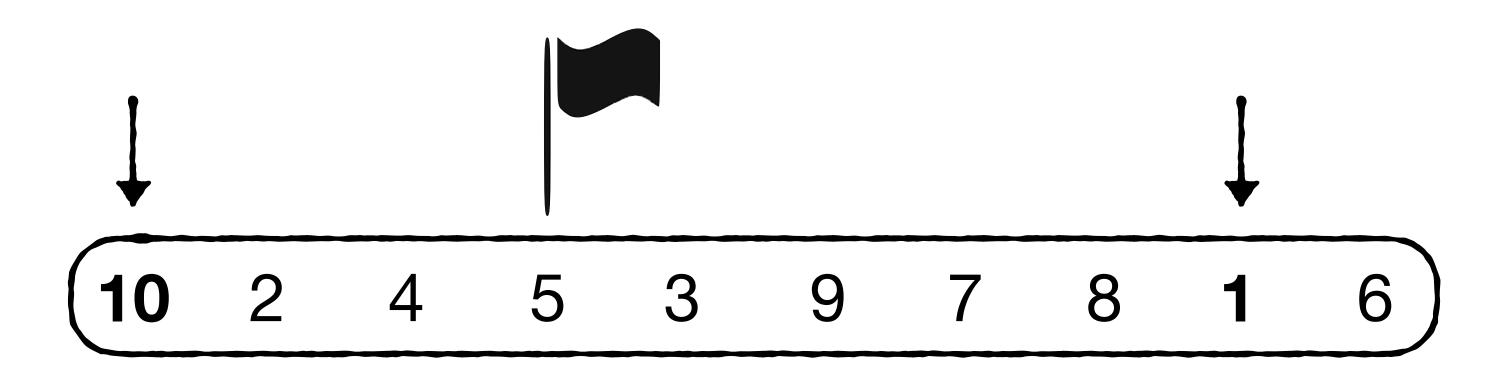
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot



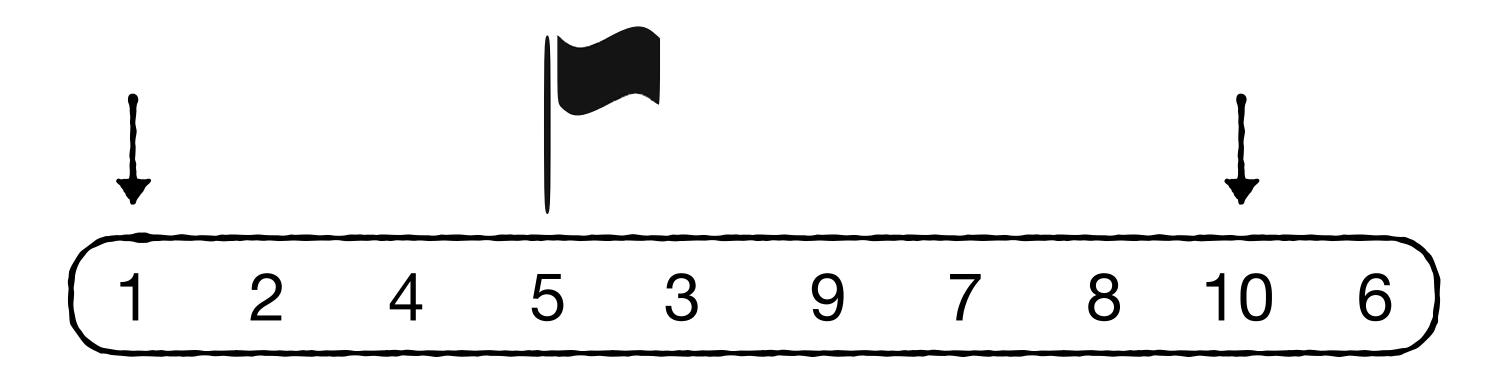
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot



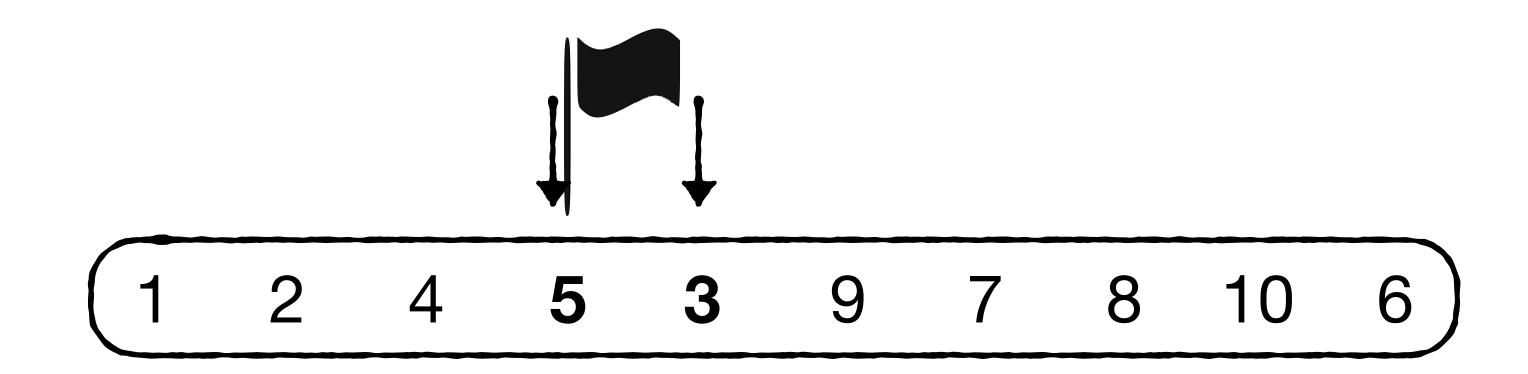
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue



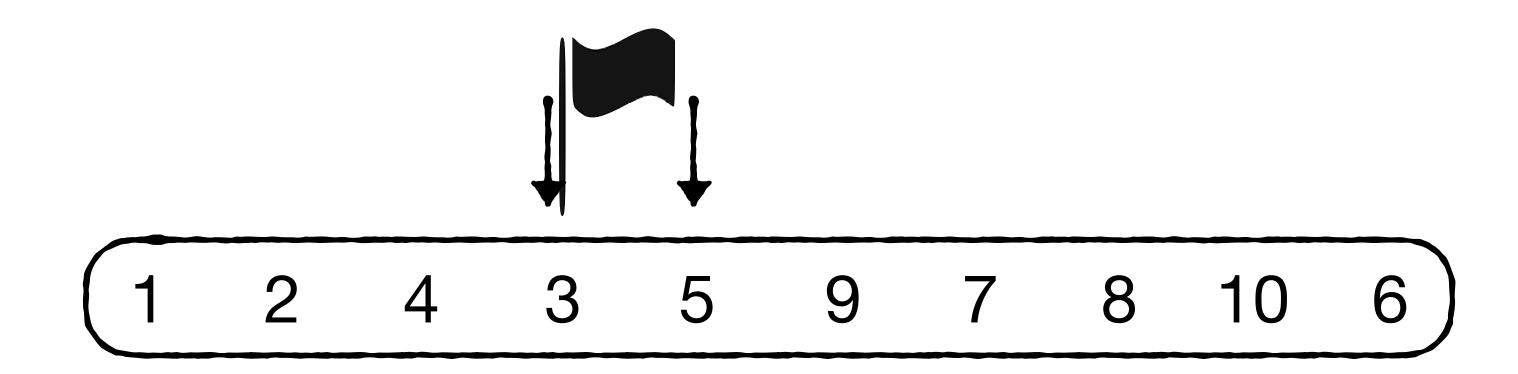
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue



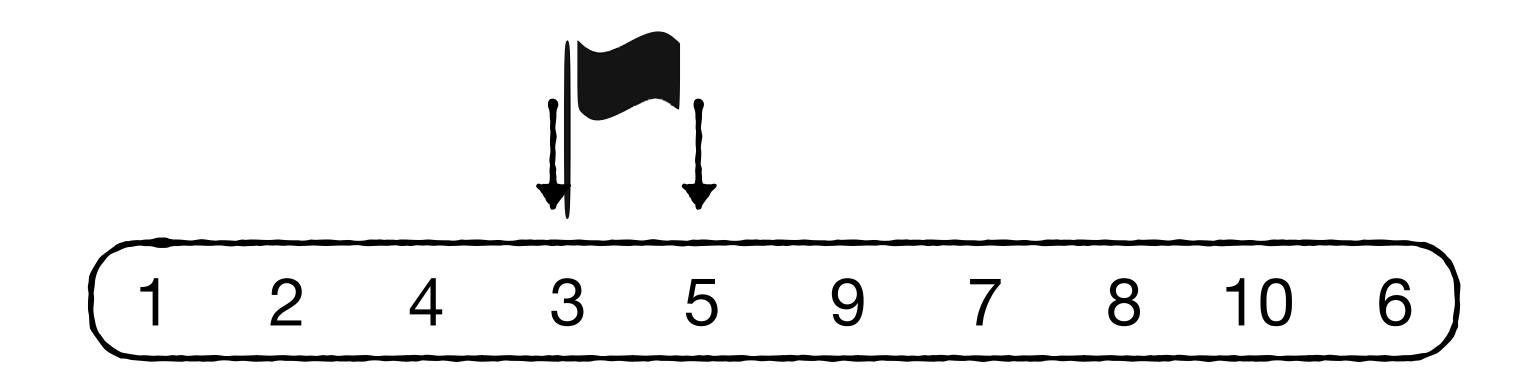
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue



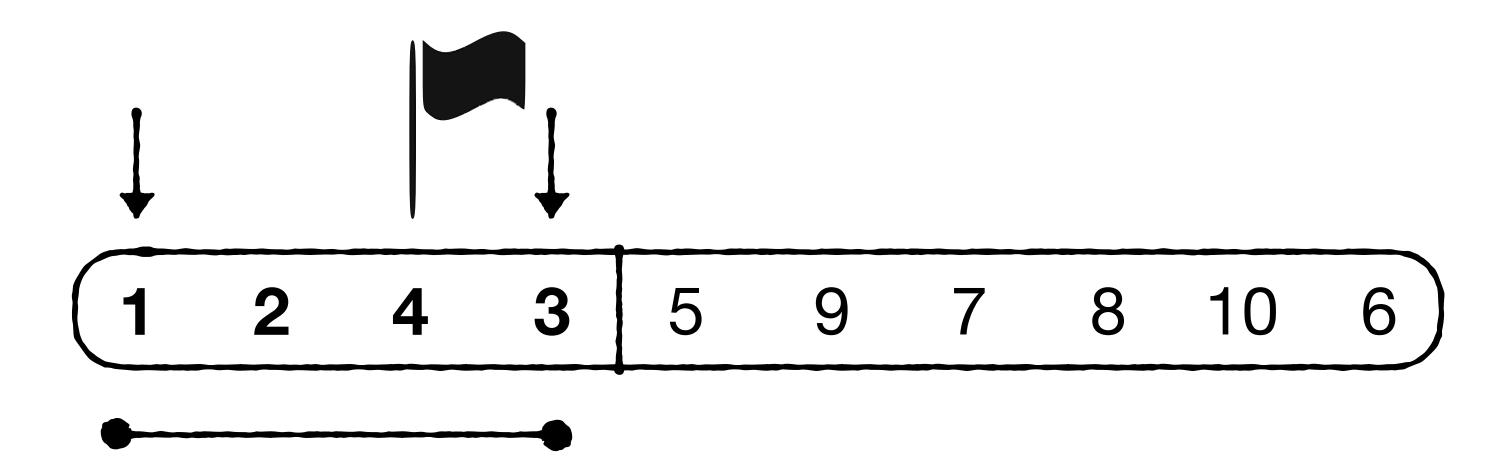
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue



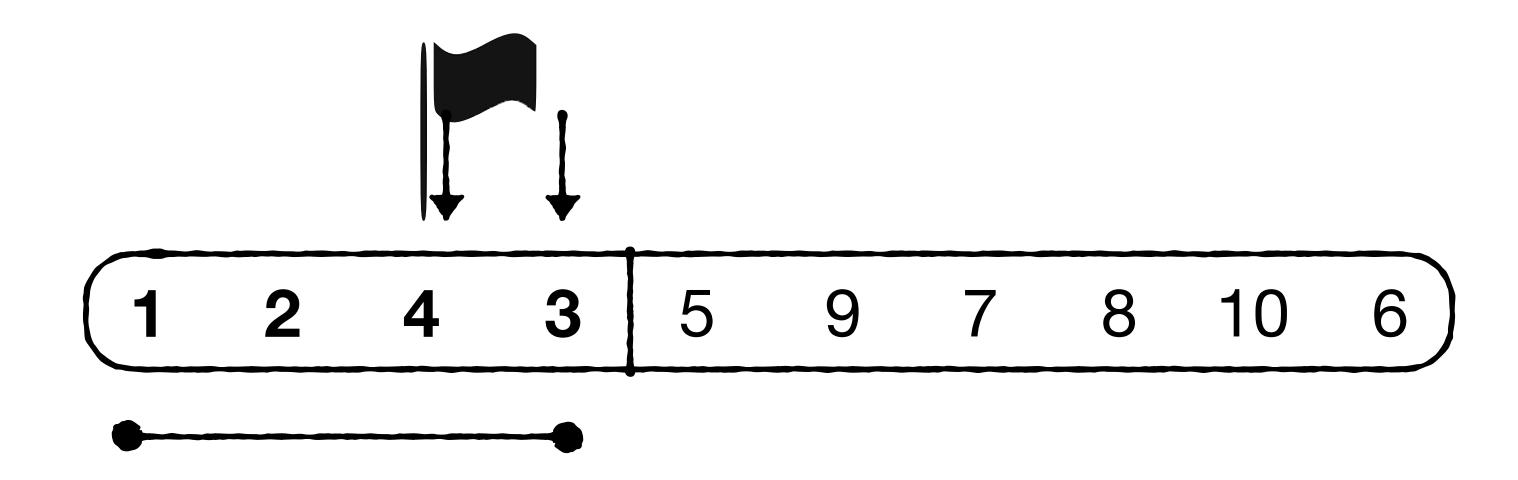
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



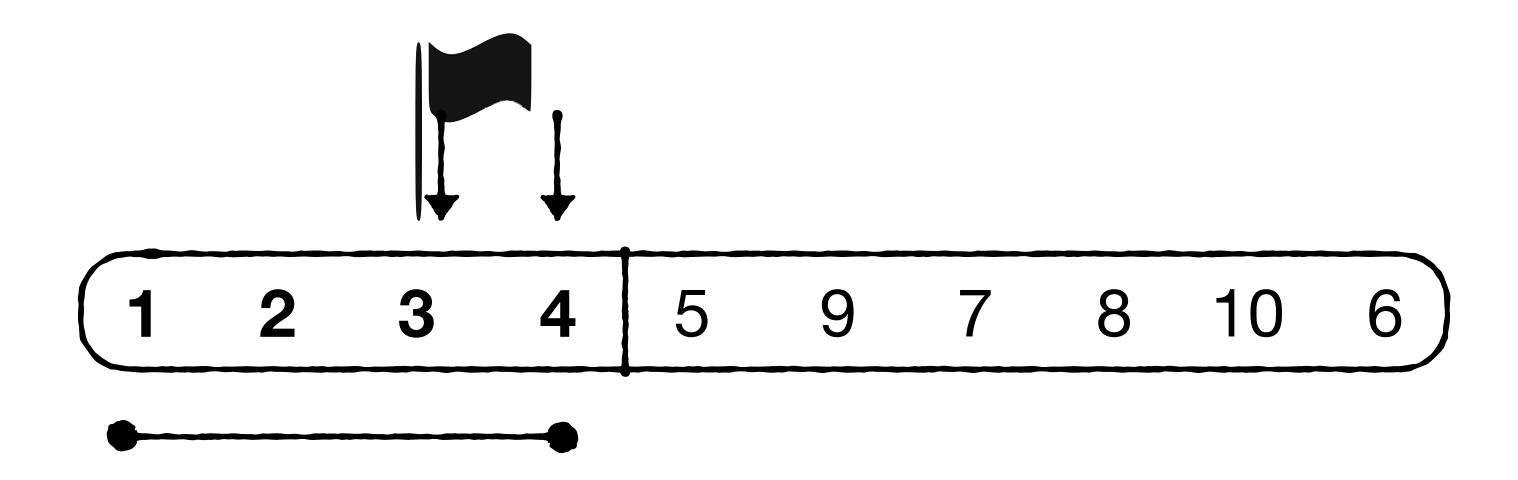
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



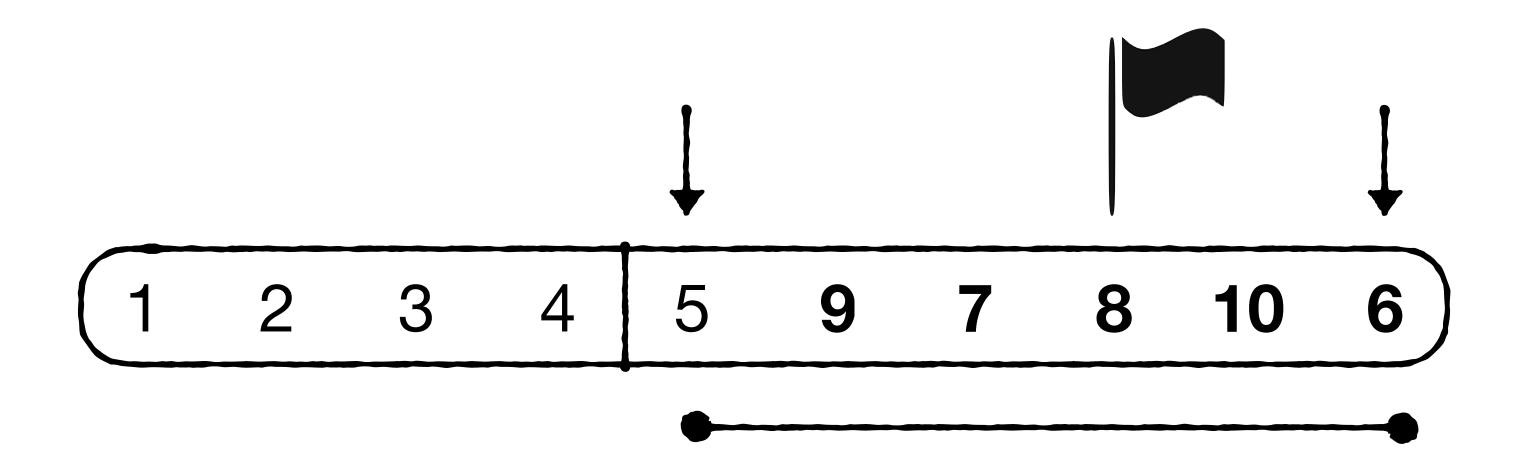
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



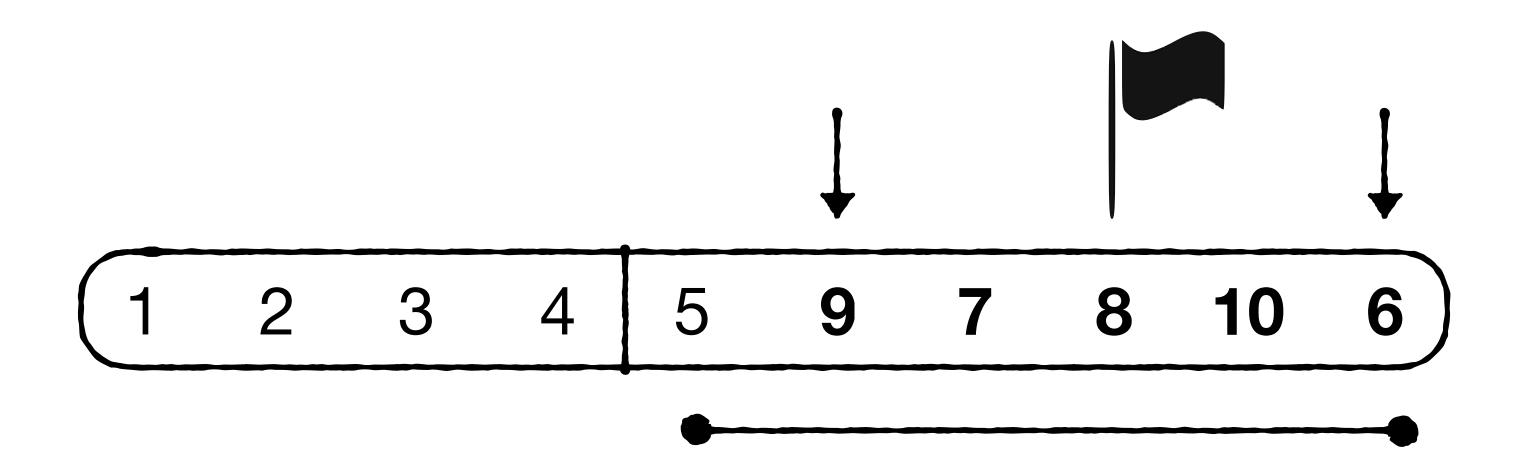
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



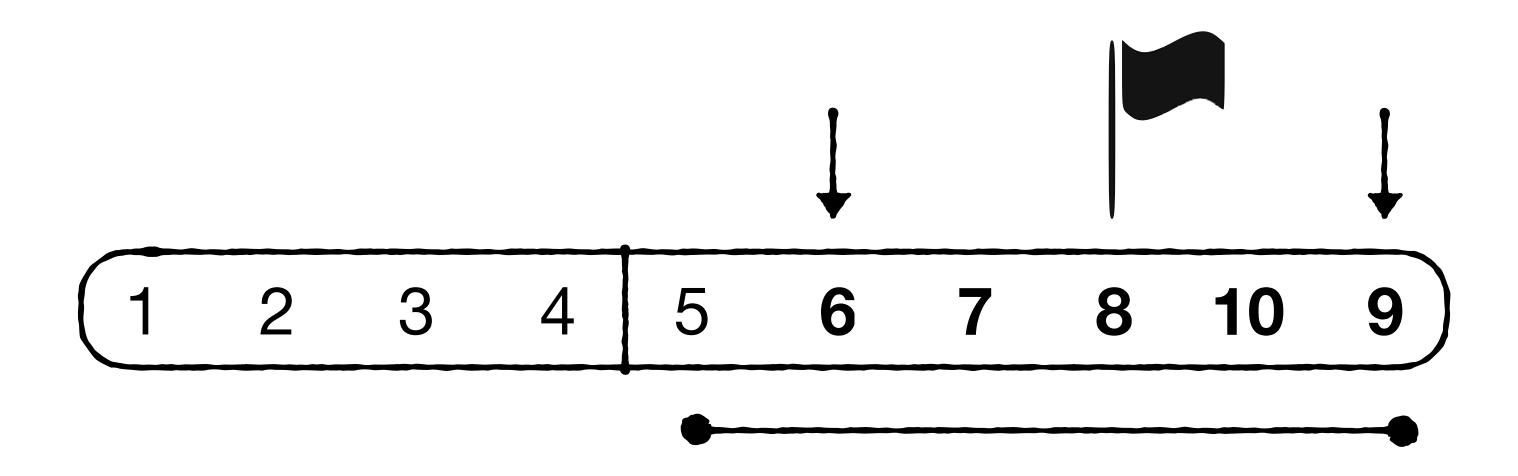
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



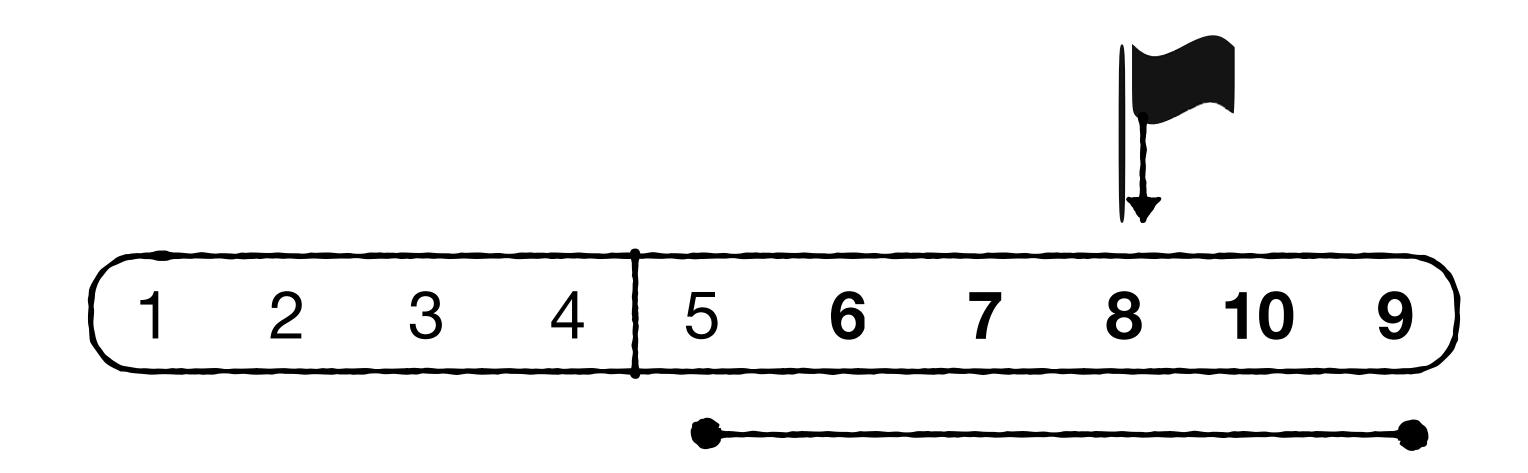
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



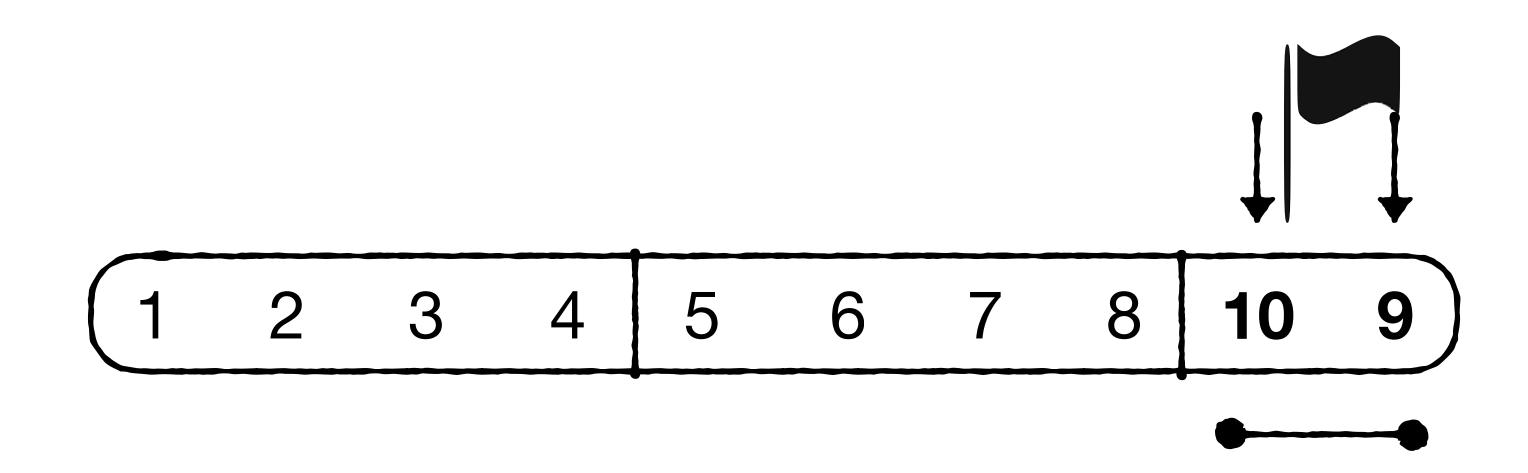
- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



- 1. Pick random pivot point and initialize two pointers at ends
- 2. Move each pointer towards pivot, stopping at first out-of-order value respect to pivot
- 3. Swap pair of out-of-order entries and continue
- 4. Continue recursively on both partitions around pivot.



Properties: Expected O(N log₂ N) worst-case

Properties: Expected O(N log₂ N) worst-case

Can it be worse?

Properties: Expected O(N log₂ N) worst-case

But can be O(N²) with low probability

Properties: Expected O(N log₂ N) worst-case

But can be O(N²) with low probability

(e.g., always pick minimum key as pivot)

Properties: Expected O(N log₂ N) worst-case

But can be O(N²) with low probability

(e.g., always pick minimum key as pivot)

Also uses sequential access, which is fast

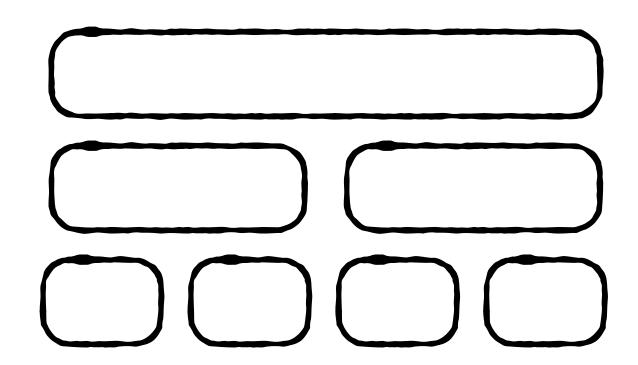
Properties: Expected O(N log₂ N) worst-case

But can be O(N²) with low probability

(e.g., always pick minimum key as pivot)

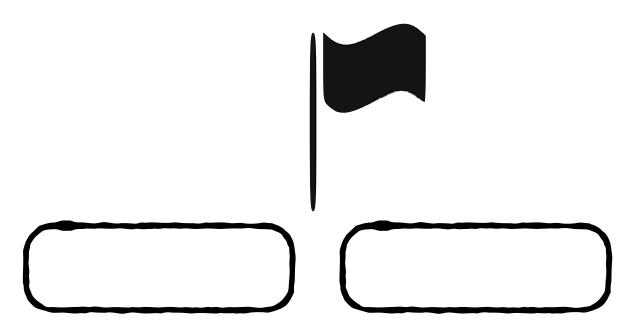
Also uses sequential access, which is fast

In-place algorithm: no need for x2 space like merge-sort



Merge-Sort

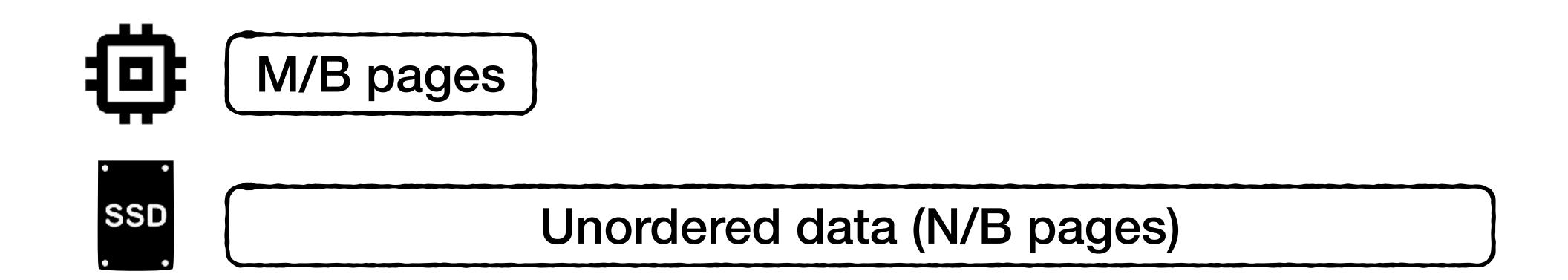
(More robust performance)



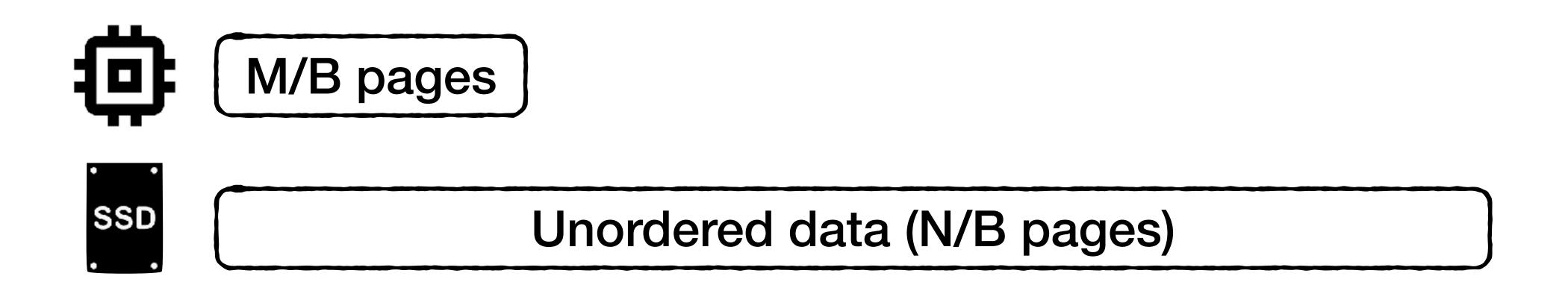
Quick-Sort

(Less space & faster on avg.)

But what if data does not fit in memory?

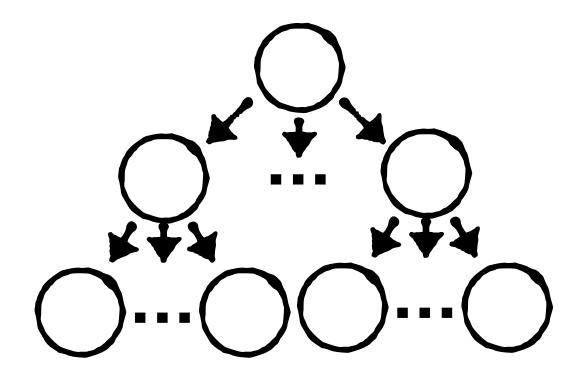


But what if data does not fit in memory?

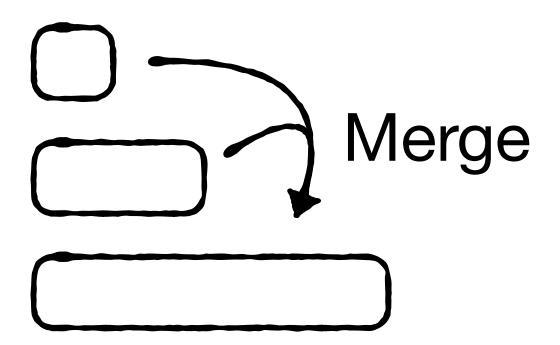


how to sort based on things we've seen in class?

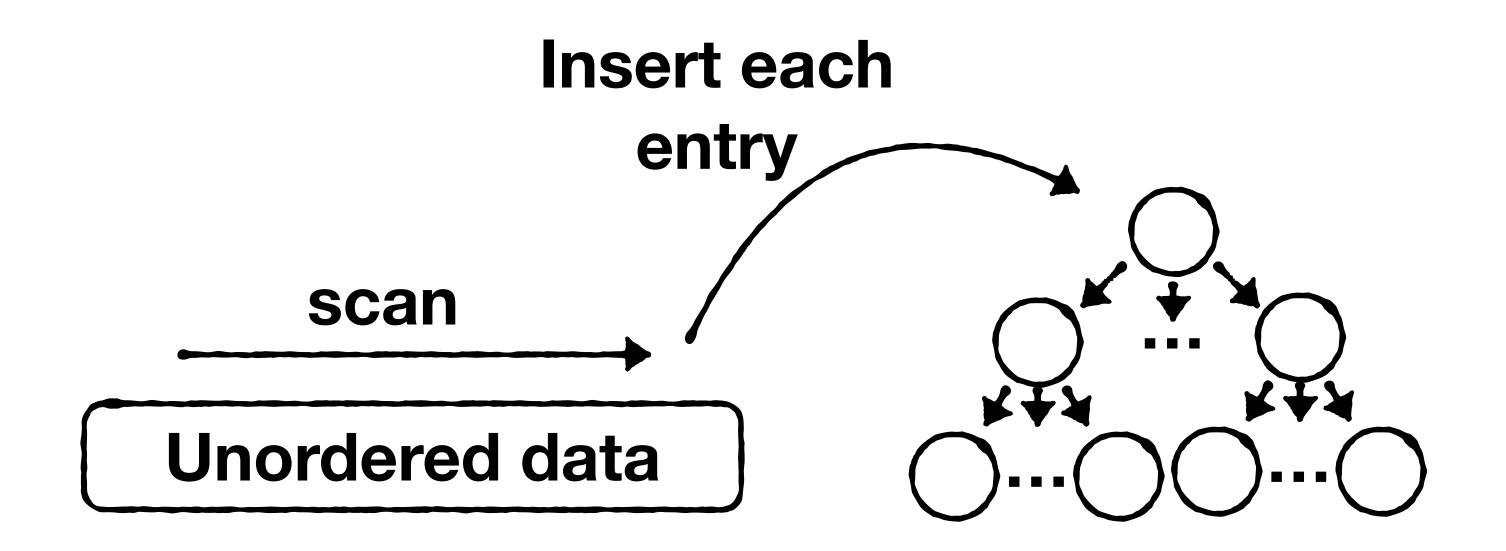
Baselines

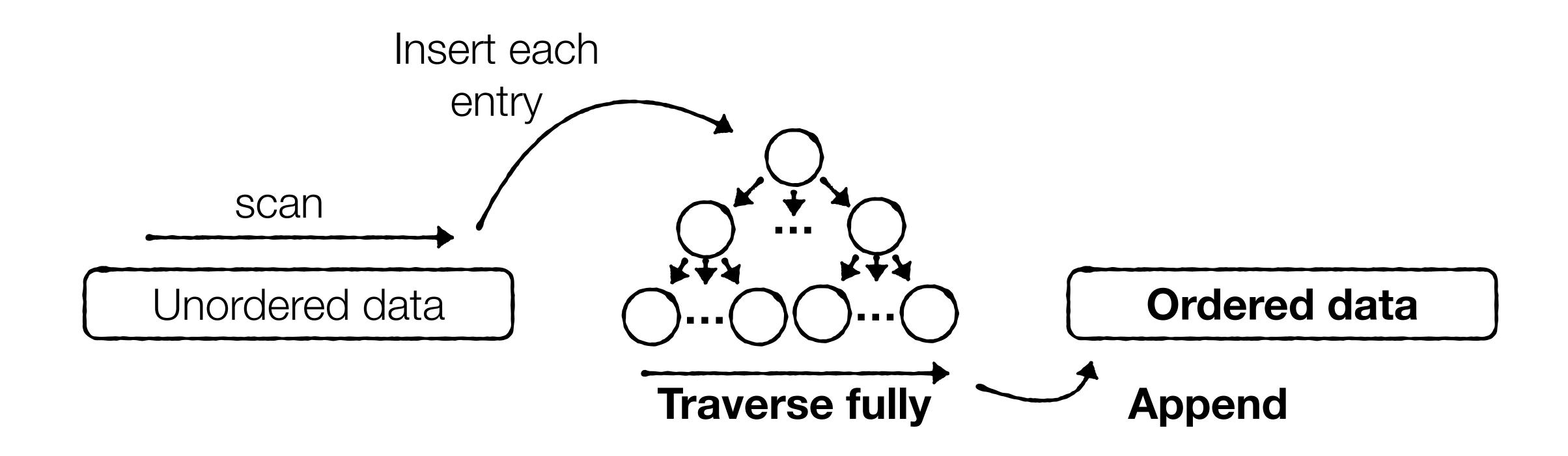


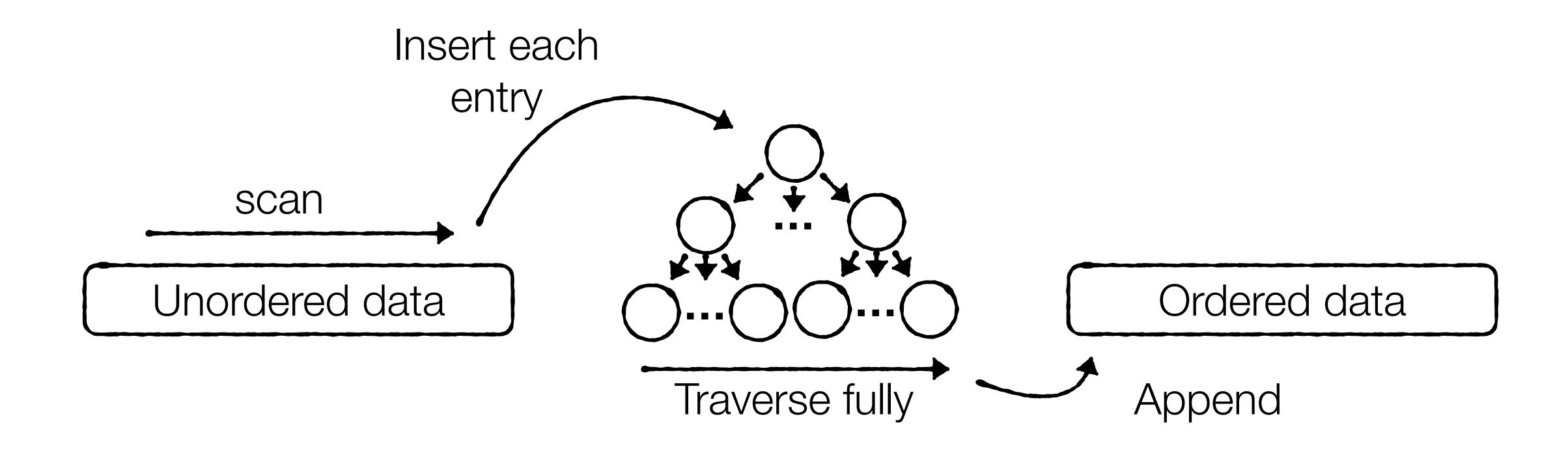
B-Tree



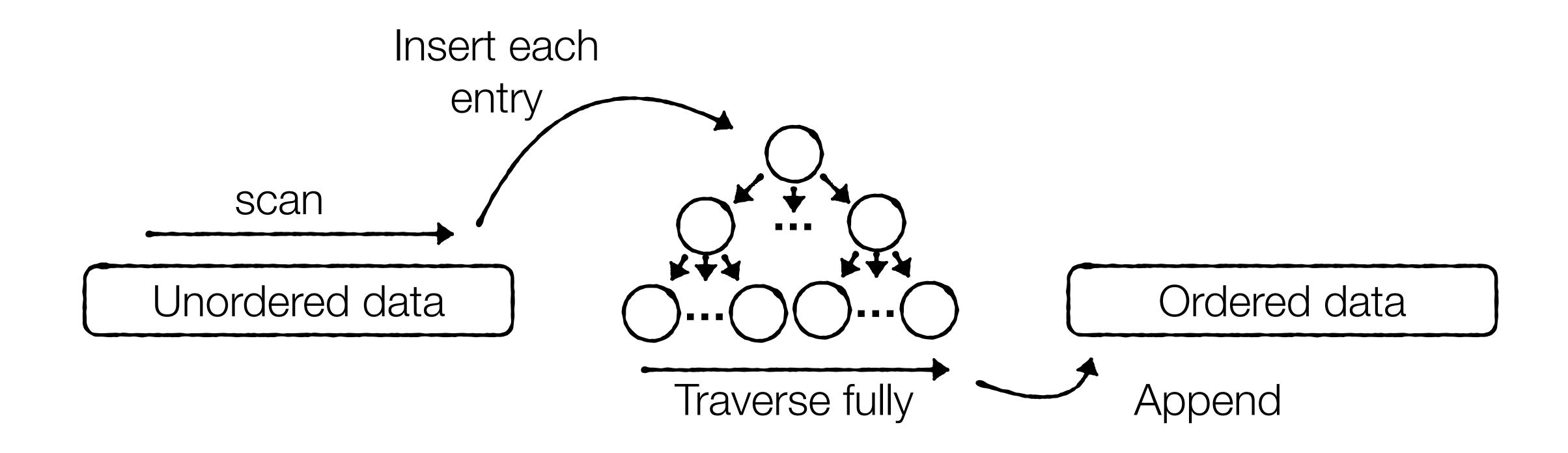
Log-Structured Merge-Tree





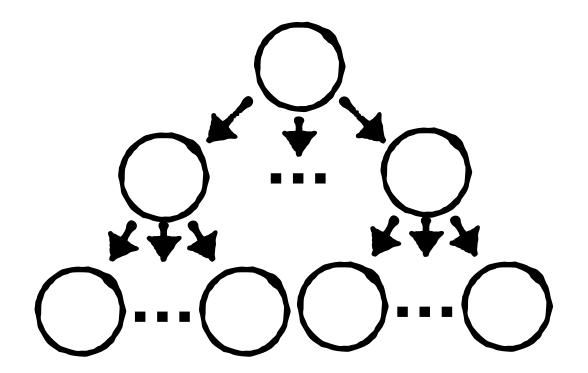


If internal nodes are in storage: O(N log_B N) reads & O(N) writes

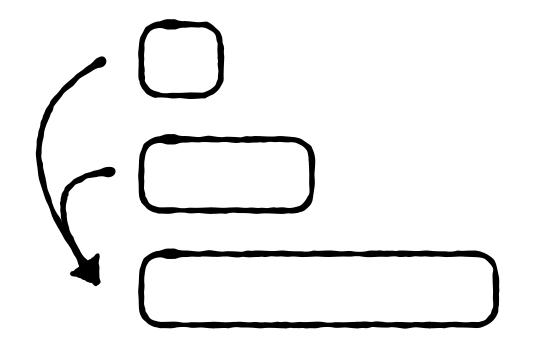


If internal nodes are in memory: O(N) reads & writes

Baselines

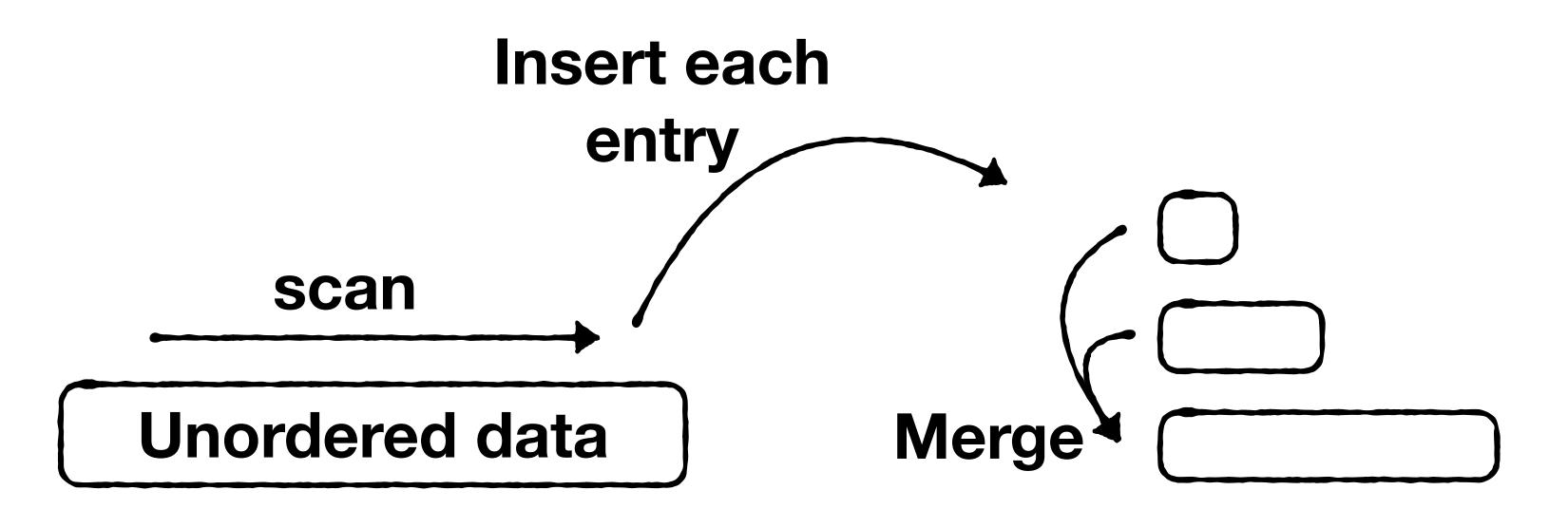


B-Tree

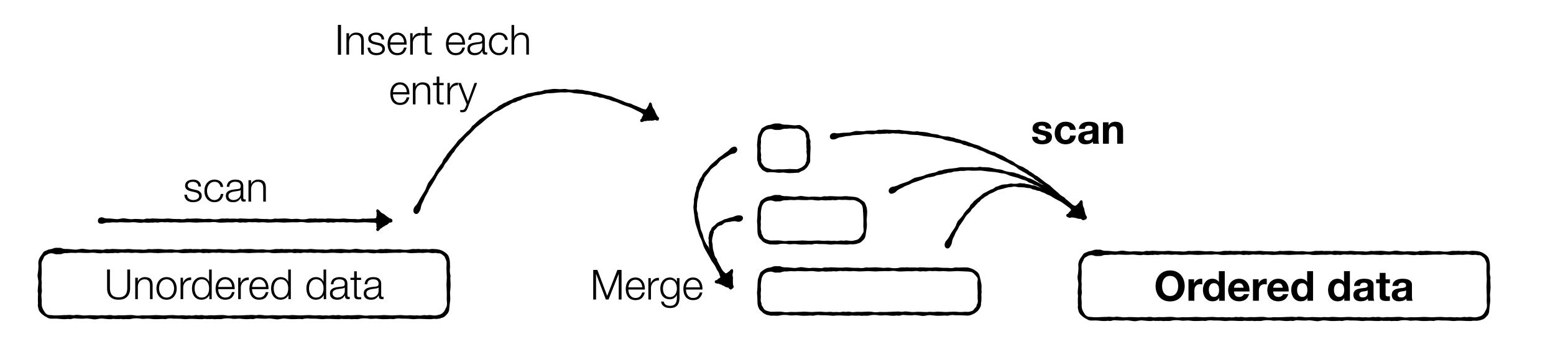


Log-Structured Merge-Tree

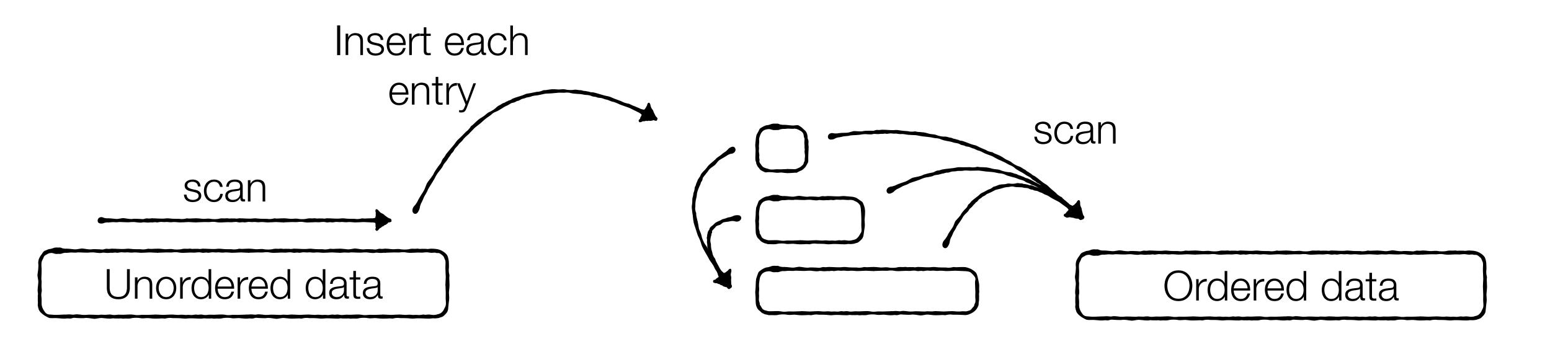
Insert entries into LSM-tree



Insert entries into LSM-tree



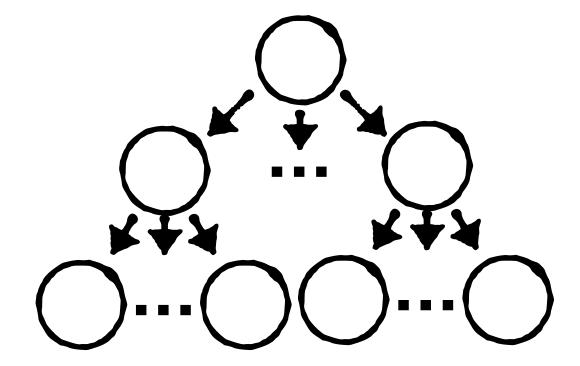
Insert entries into LSM-tree



With basic LSM-tree: O(N/B * log₂(N/P)) reads & writes

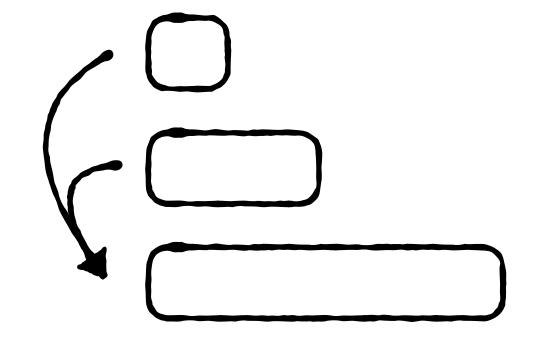
(Regardless of whether internal nodes fit in memory)

B-Tree



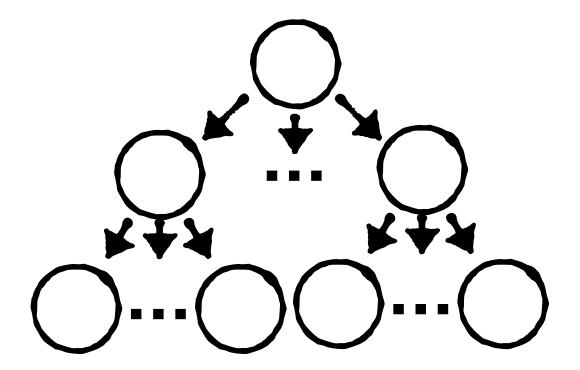
O(N)

LSM-Tree



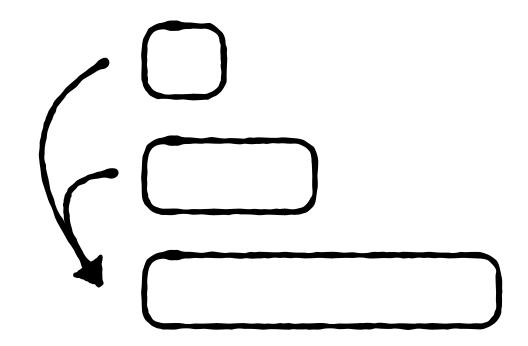
 $O(N/B * log_2(N/P))$

B-Tree



O(N)

LSM-Tree



 $O(N/B * log_2(N/P))$

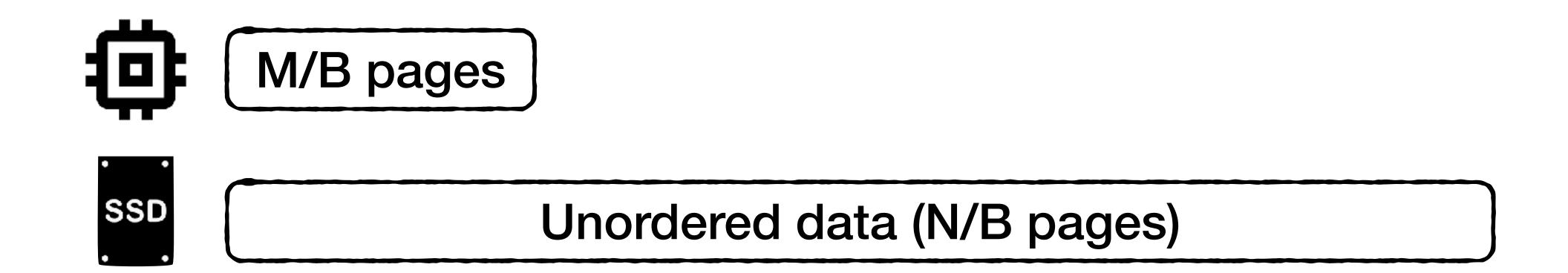
Can we do better?

Ideally we want O(N/B)

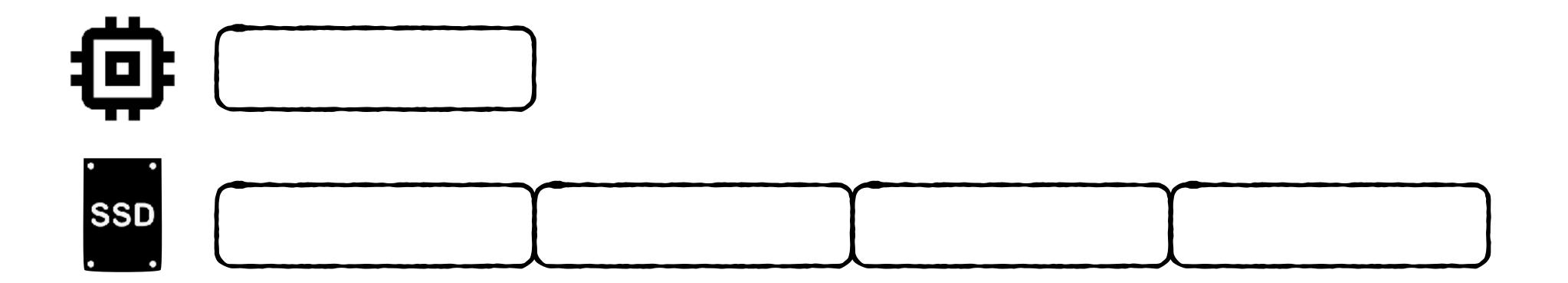
How can we efficiently sort data that doesn't fit in memory?



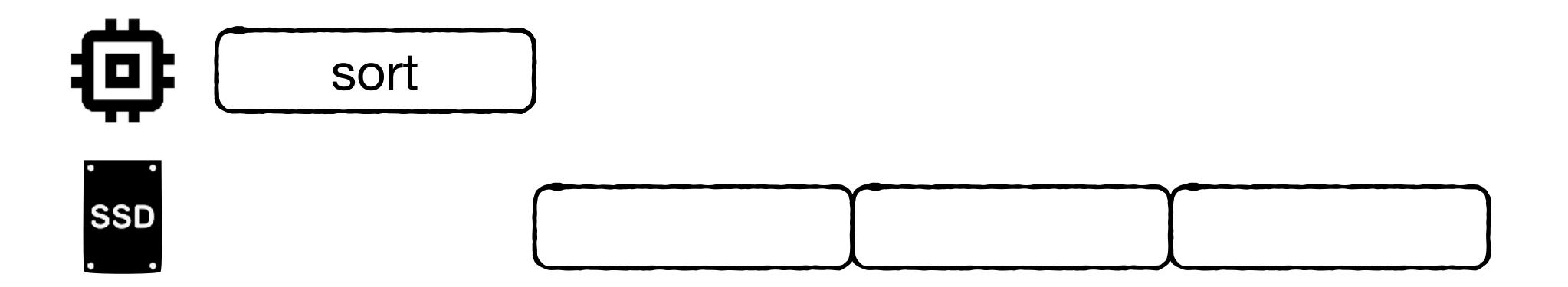
Think & then discuss with your neighbors



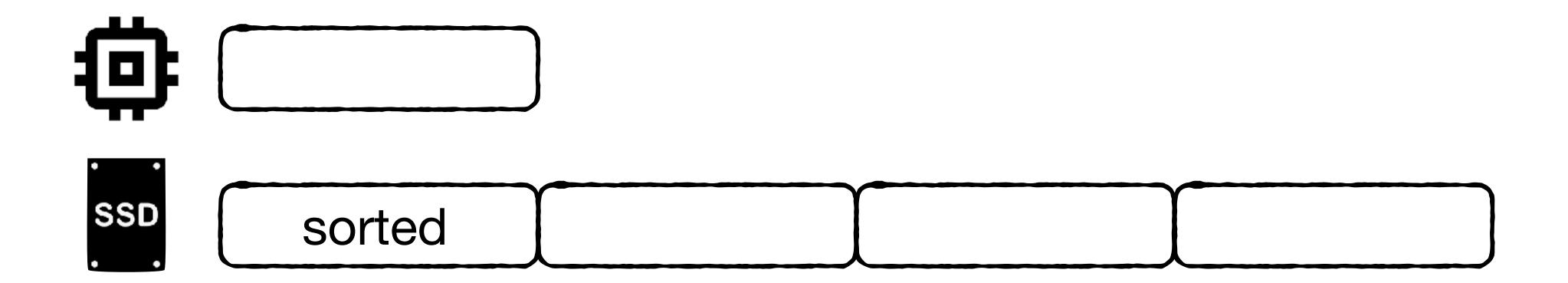
Multi-Way Merge-Sort

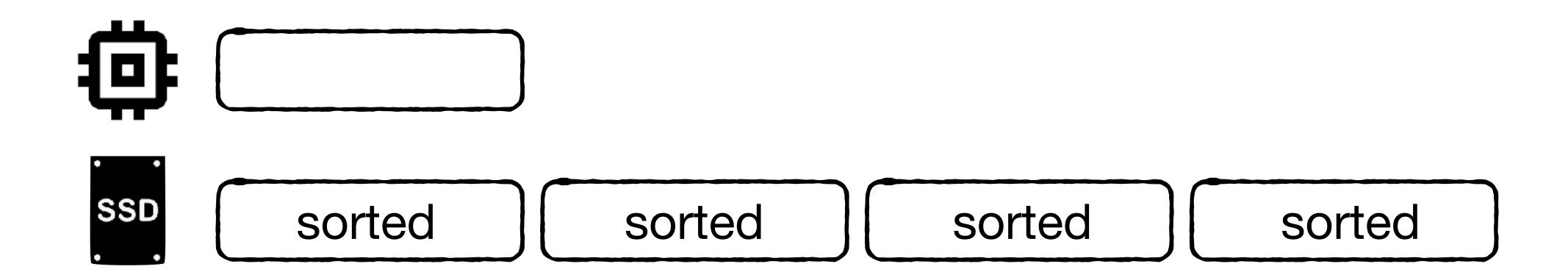


Multi-Way Merge-Sort

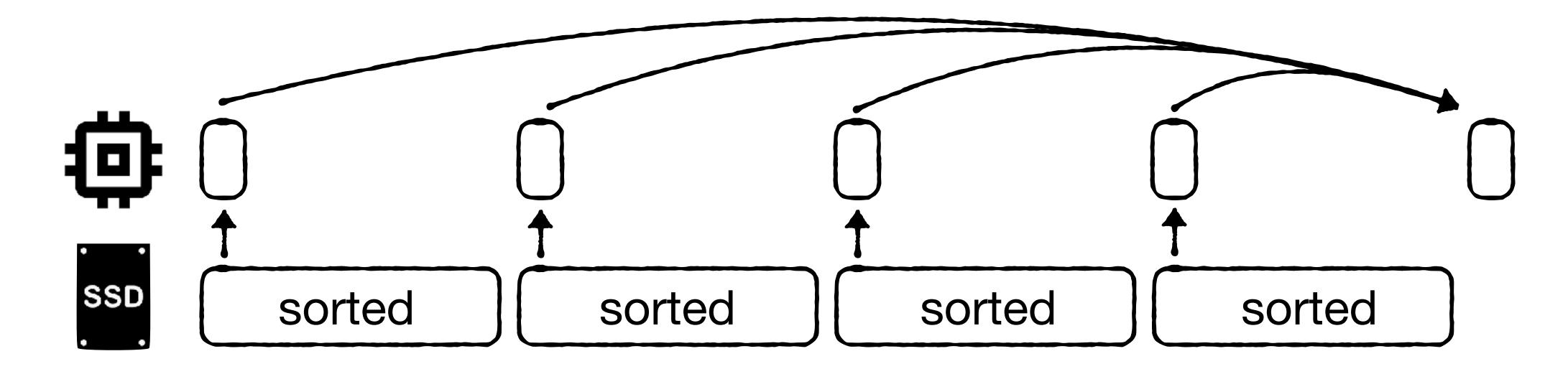


Multi-Way Merge-Sort

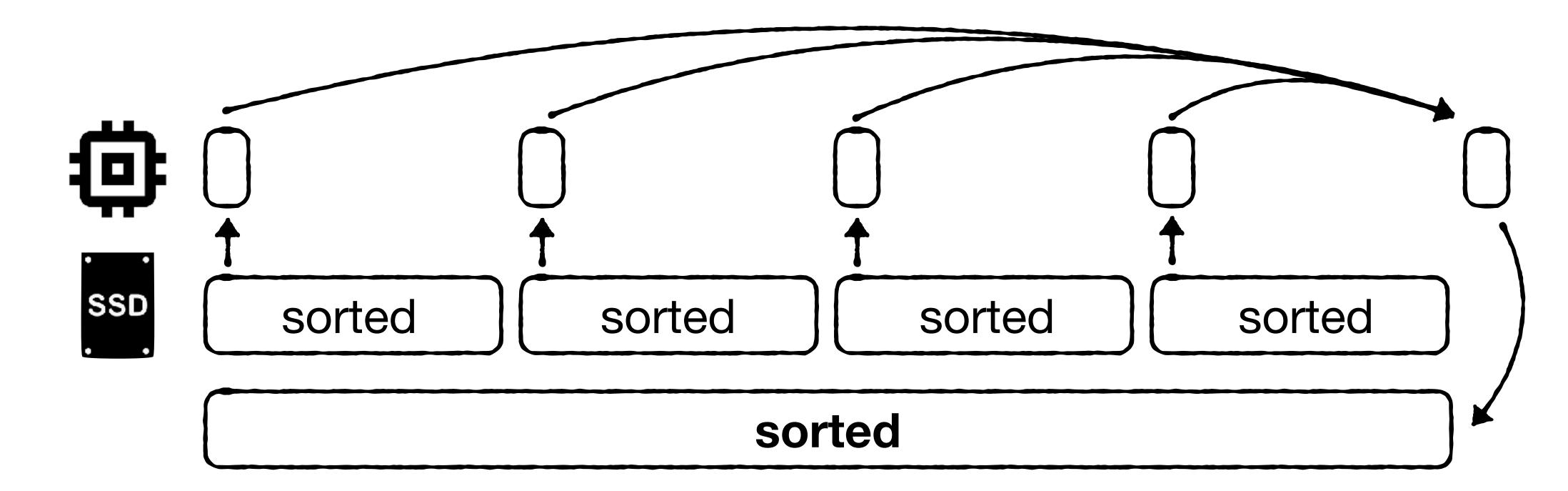




- 1. Sequentially read chunks that fit in memory, sort, and store back as temporary files
- 2. Allocate a buffer for each input file and merge into output stream.

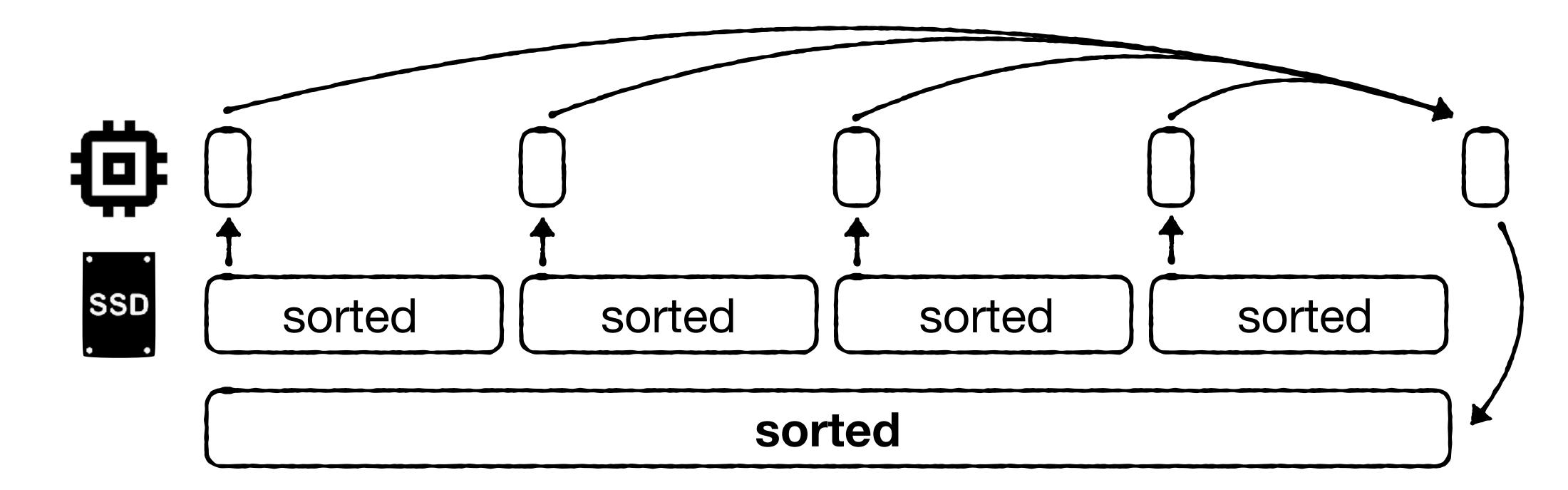


- 1. Sequentially read chunks that fit in memory, sort, and store back as temporary files
- 2. Allocate a buffer for each input file and merge into output stream.

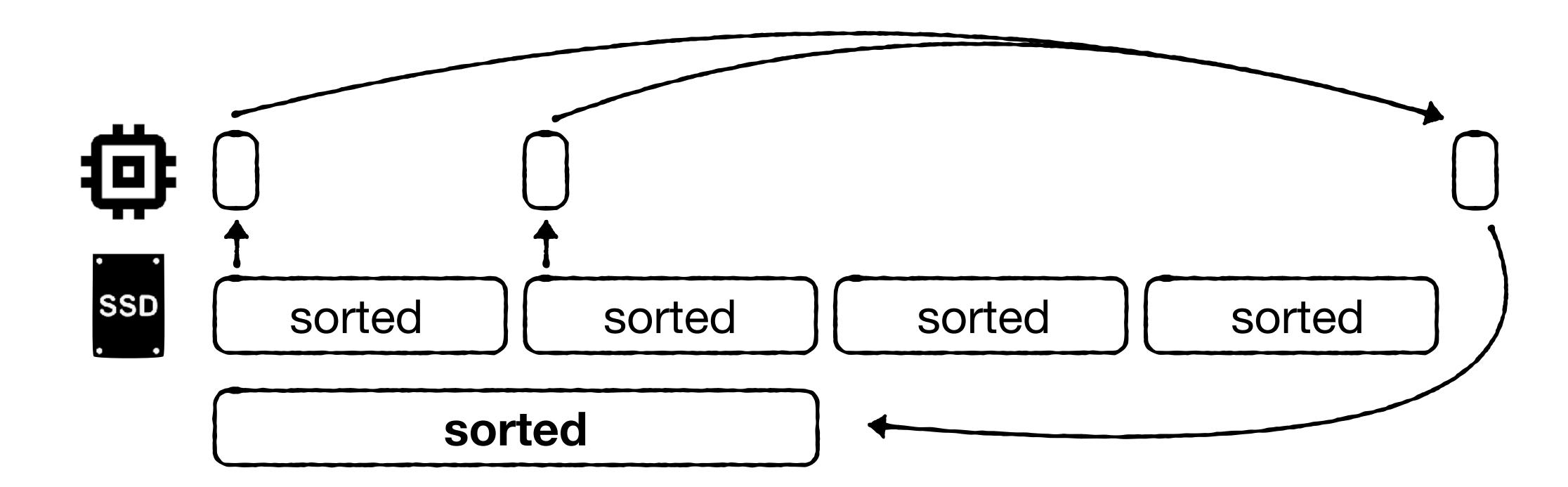


Problem?

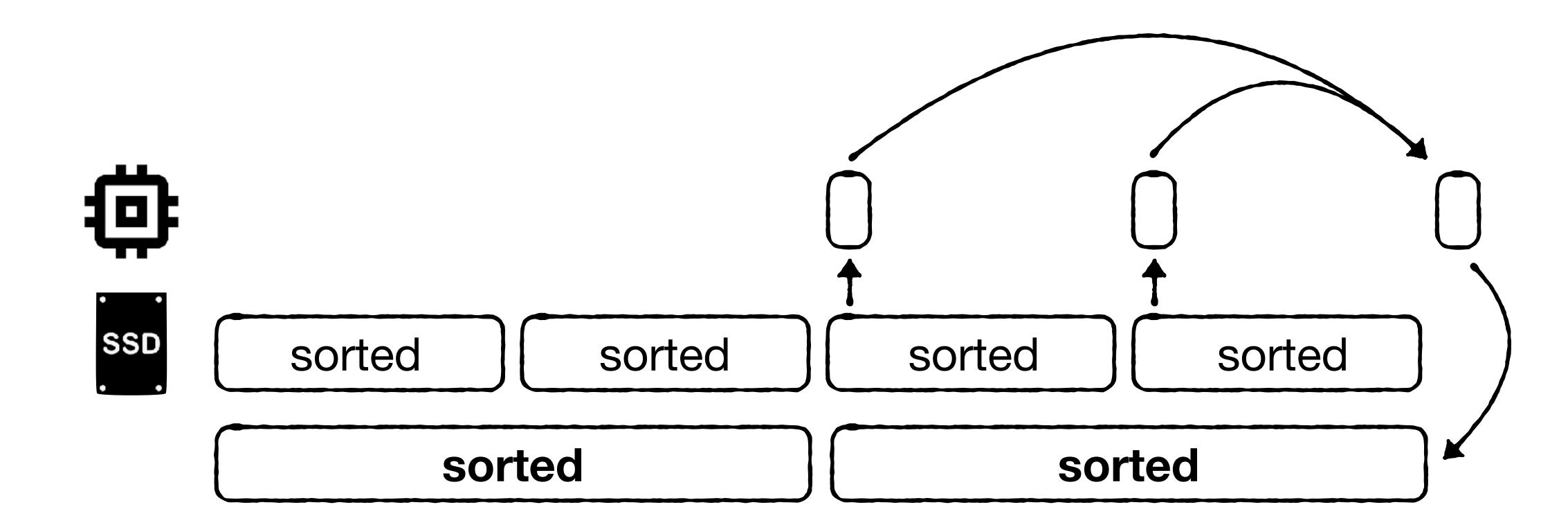
- 1. Sequentially read chunks that fit in memory, sort, and store back as temporary files
- 2. Allocate a buffer for each input file and merge into output stream.



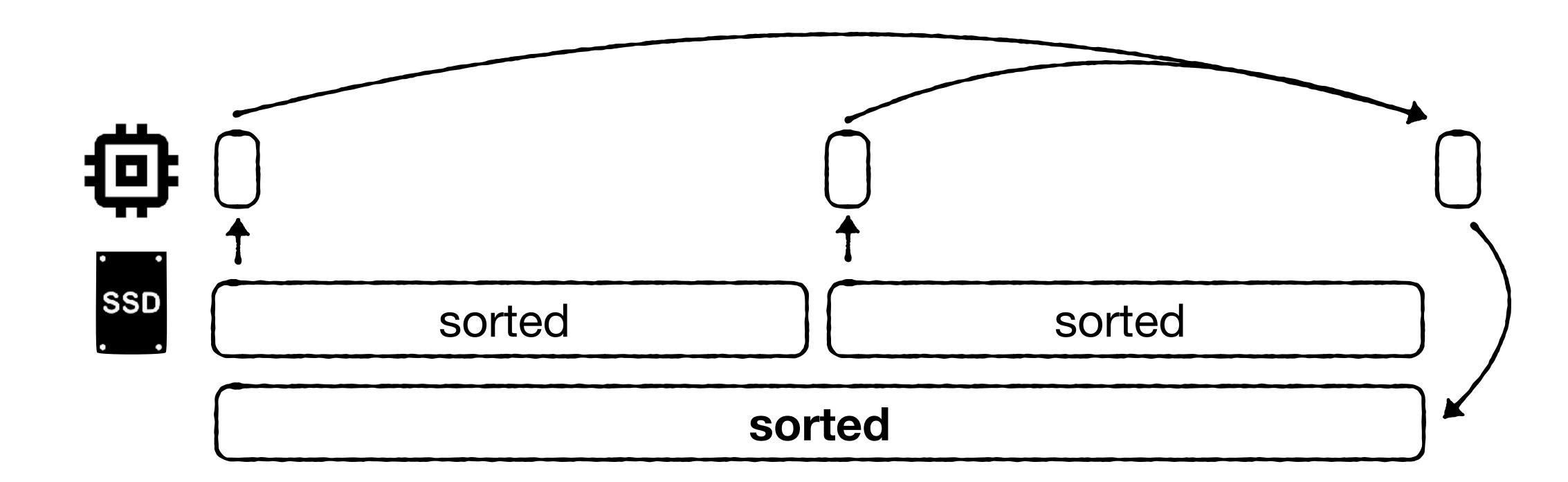
3. If we have little memory, may have to do multiple passes



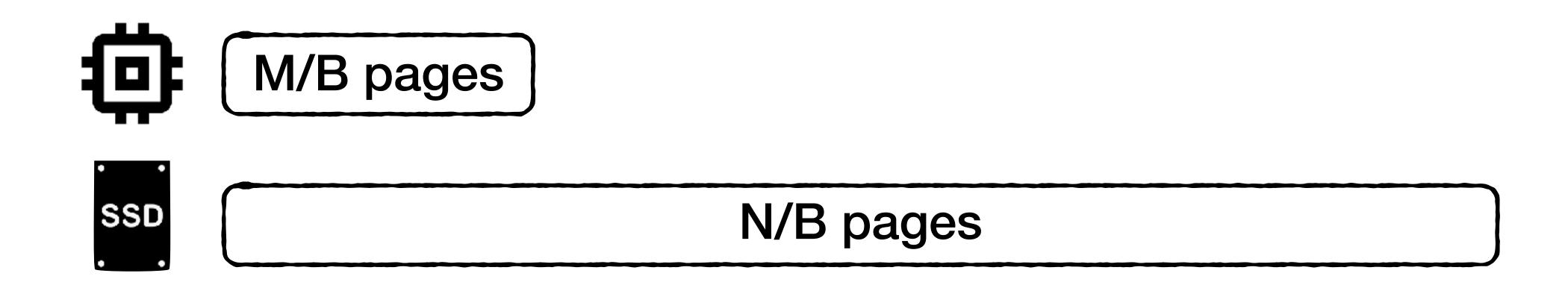
3. If we have little memory, may have to do multiple passes



3. If we have little memory, may have to do multiple passes

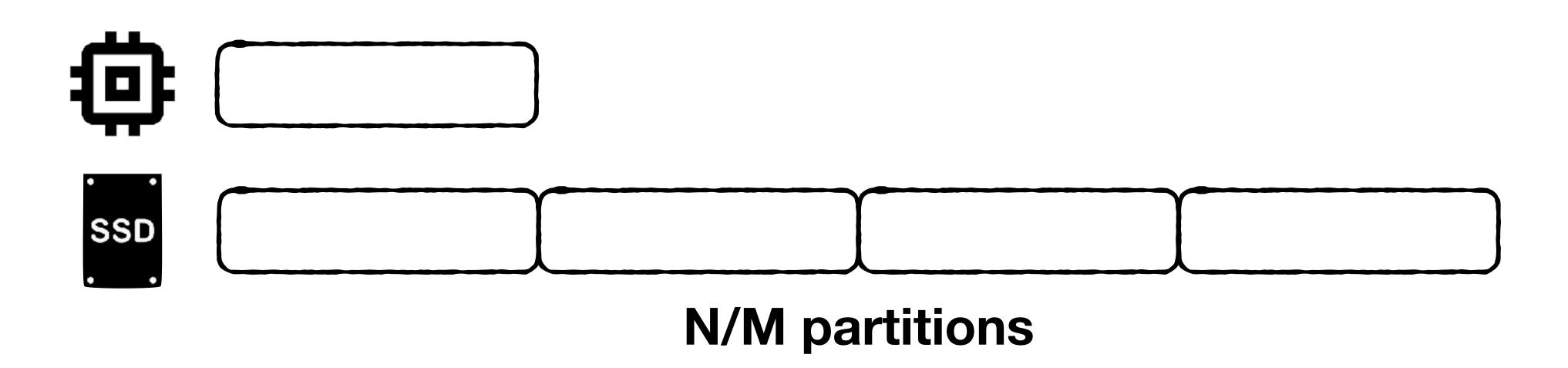


How many merging iterations (passes) must we do?



How many merging iterations (passes) must we do?

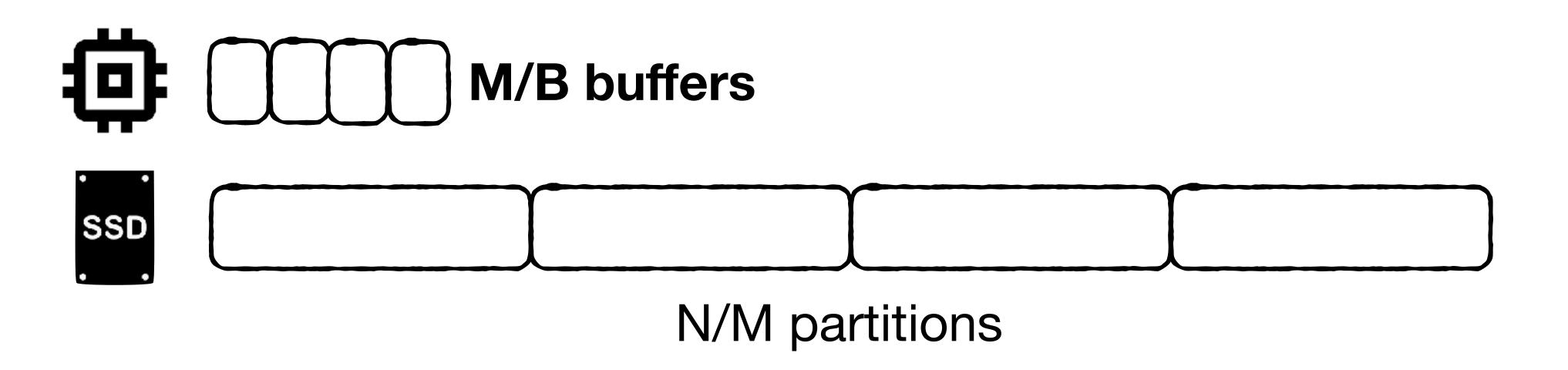
Data is divided into N/M partitions



How many merging iterations (passes) must we do?

Data is divided into N/M partitions

We have M/B buffers, so each iteration merges M/B partitions

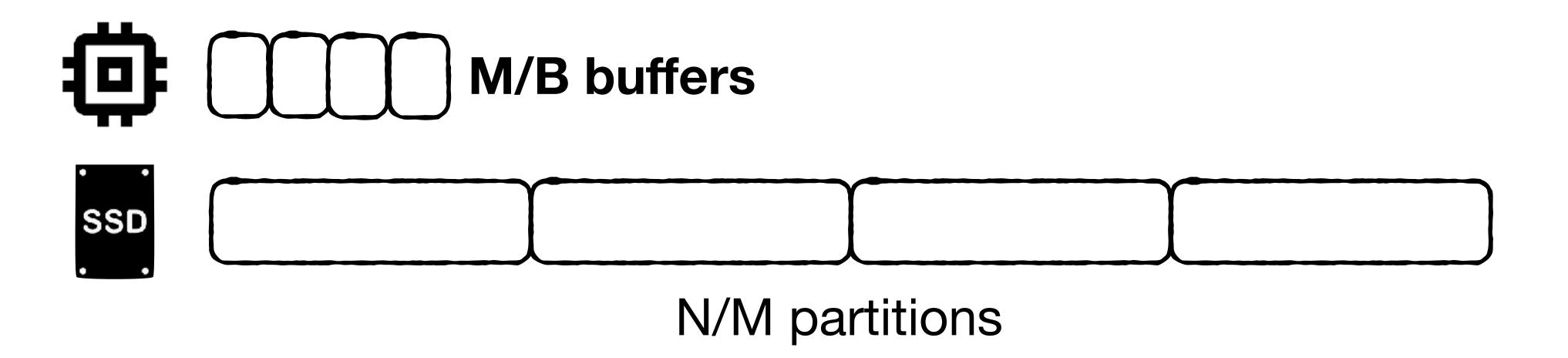


How many merging iterations (passes) must we do?

Data is divided into N/M partitions

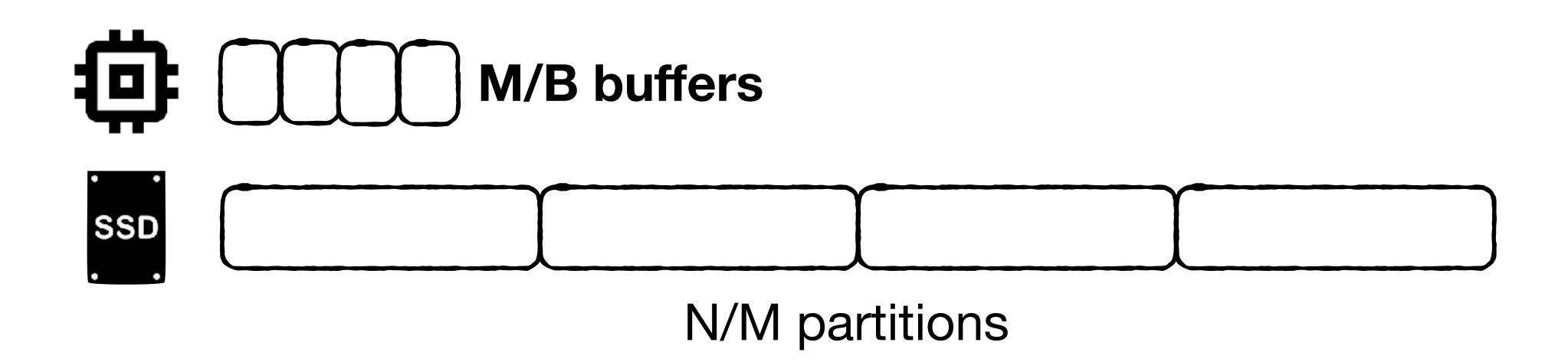
We have M/B buffers, so each iteration merges M/B partitions

Iterations = $log_{M/B}(N/M)$



Iterations = $log_{M/B}(N/M)$

Each iteration does a full pass over the data costing O(N/B)

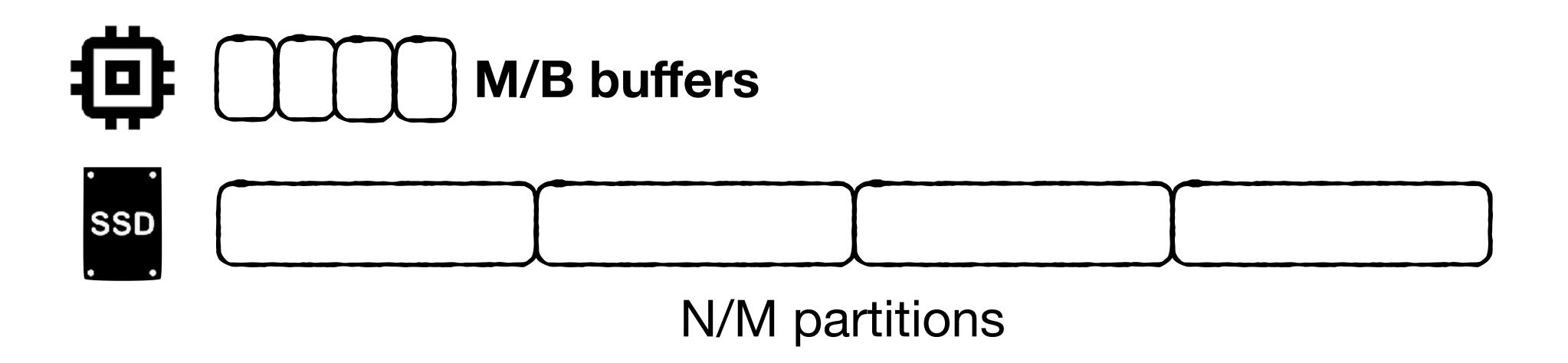


Analysis

Iterations = $log_{M/B}(N/M)$

Each iteration does a full pass over the data costing O(N/B)

Cost of Merging Phase: N/B · [log_{M/B}(N/M)]

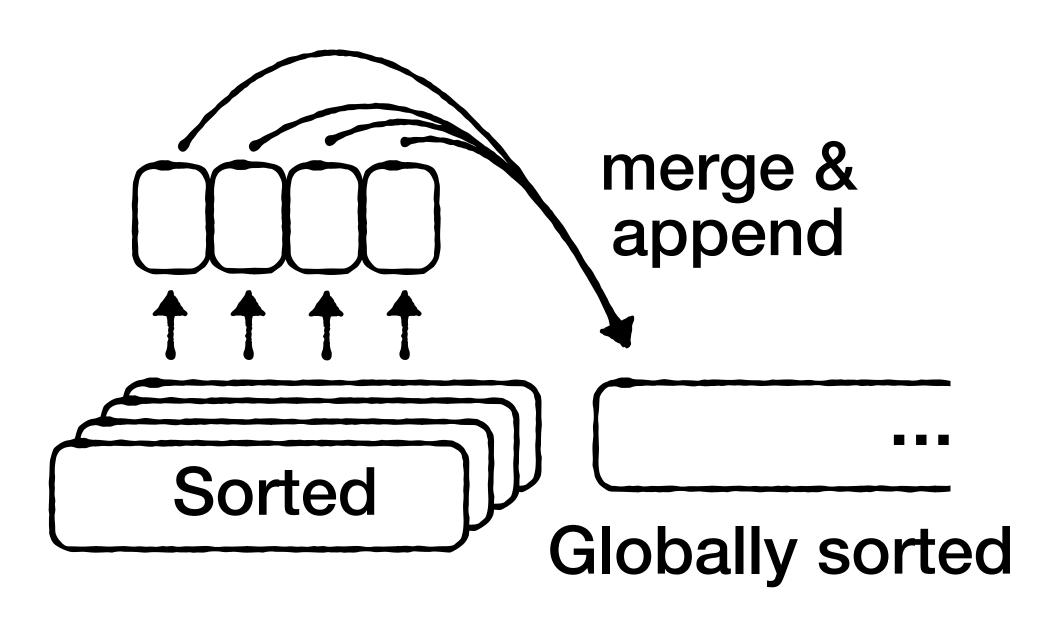


Analysis

Partitioning Phase N/B I/Os

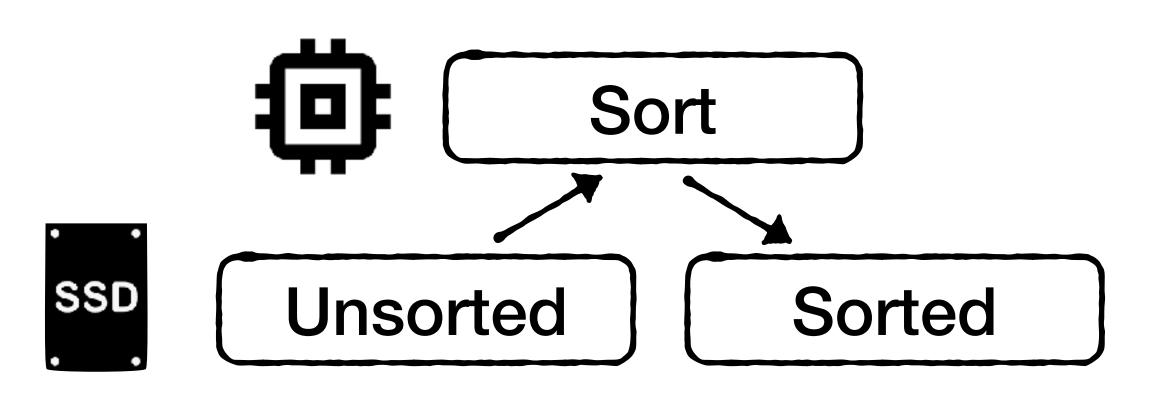
Sort
Unsorted Sorted

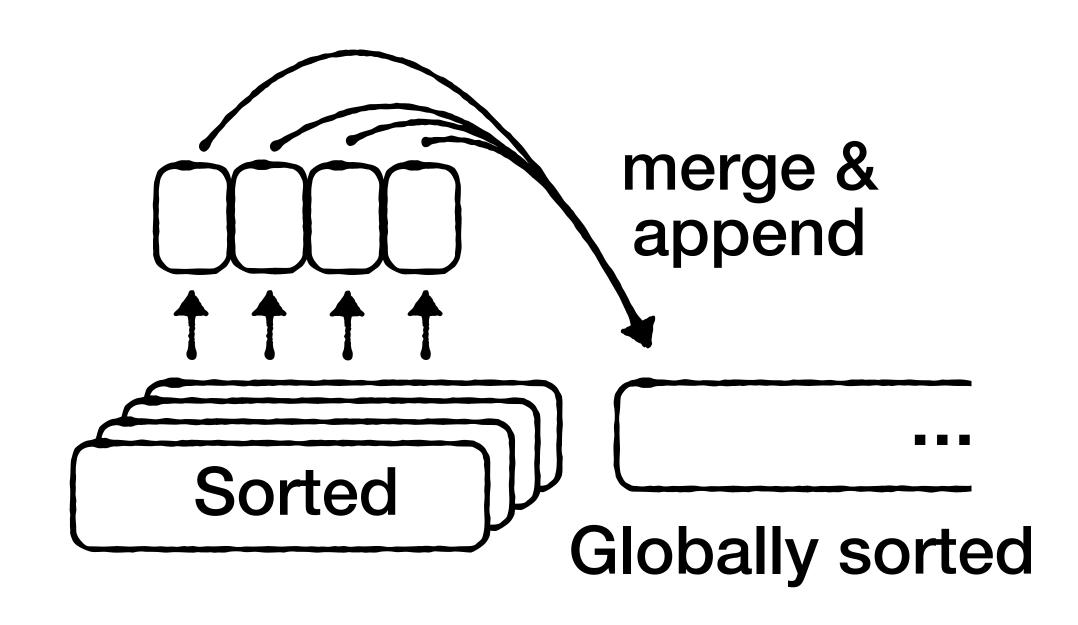
Merging Phase N/B · [log_{M/B}(N/M)]



Analyzing I/O for External Sort

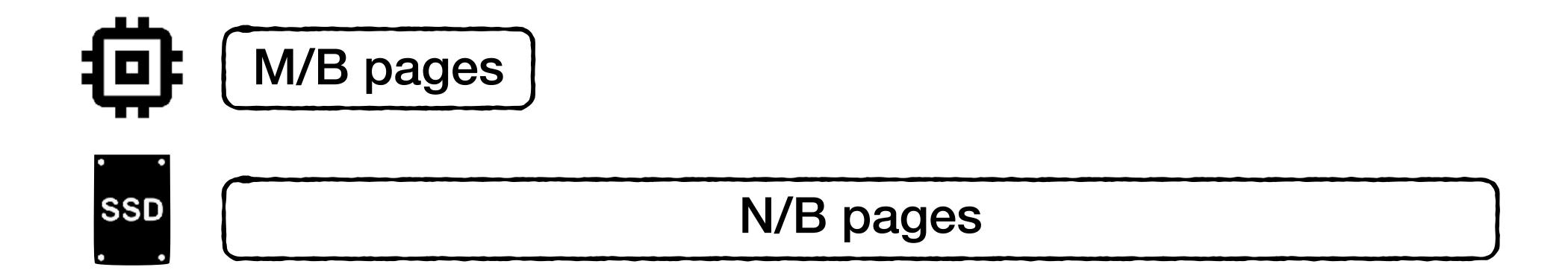
Partitioning Phase Merging Phase Total N/B I/Os + $N/B \cdot \lceil \log_{M/B}(N/M) \rceil = O(N/B \cdot \log_{M/B}(N/B))$





Impact of More Memory

 $O(N/B \cdot \log_{M/B}(N/B))$

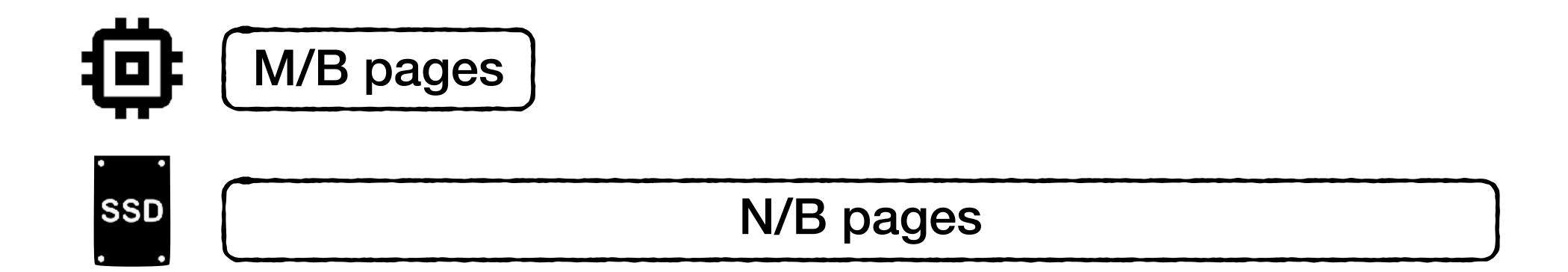


Impact of More Memory

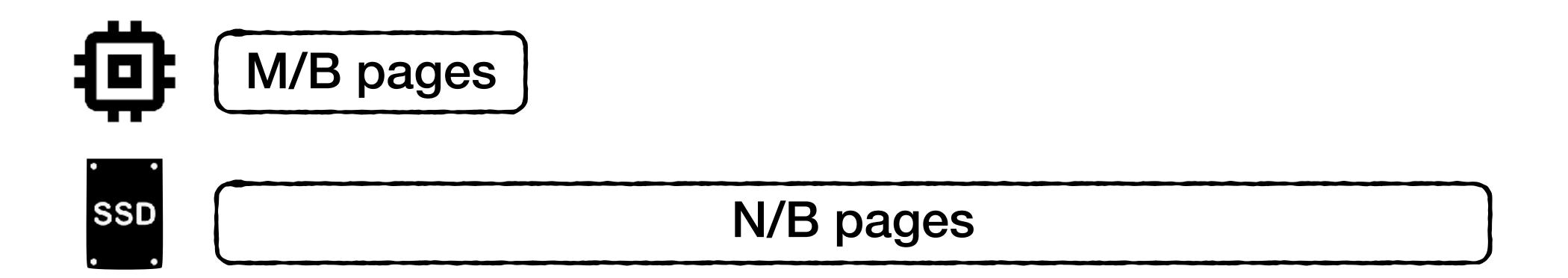
 $O(N/B \cdot \log_{M/B}(N/B))$



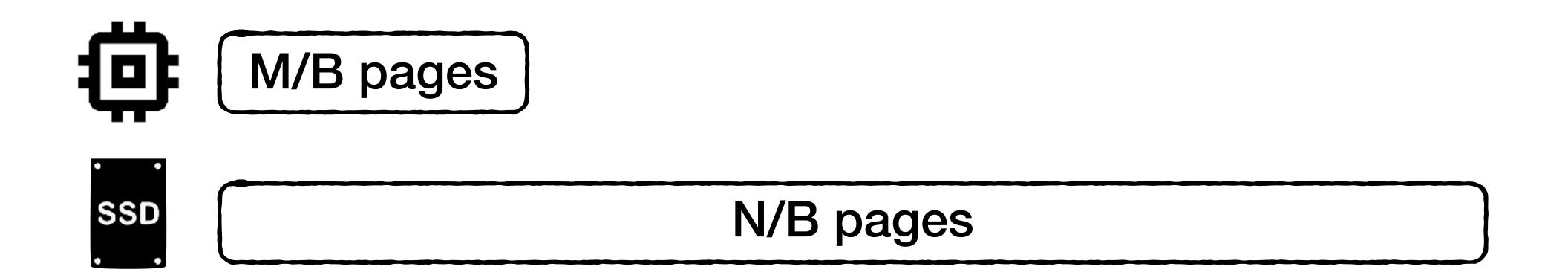
Merging more partitions per pass



 $O(N/B \cdot \log_{M/B}(N/B))$

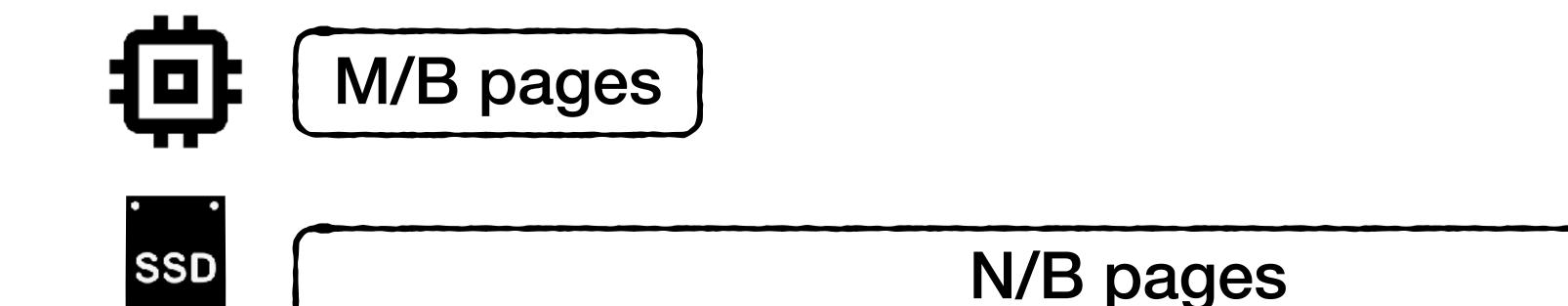


Let $log_{M/B}(N/B) = 2$ and solve for M



Let $log_{M/B}(N/B) = 2$ and solve for M

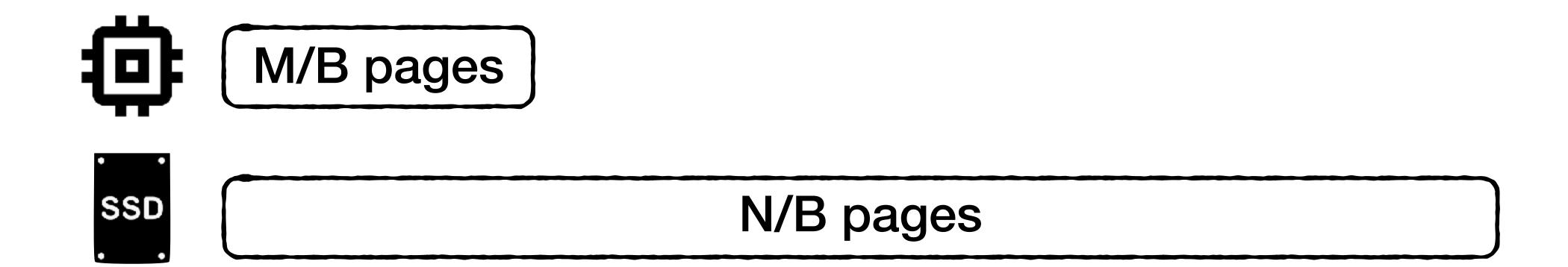
We get: $M = \sqrt{N \cdot B}$ (Measured in entries)



Let $log_{M/B}(N/B) = 2$ and solve for M

We get:
$$M = \sqrt{N \cdot B}$$
 (Measured in entries)

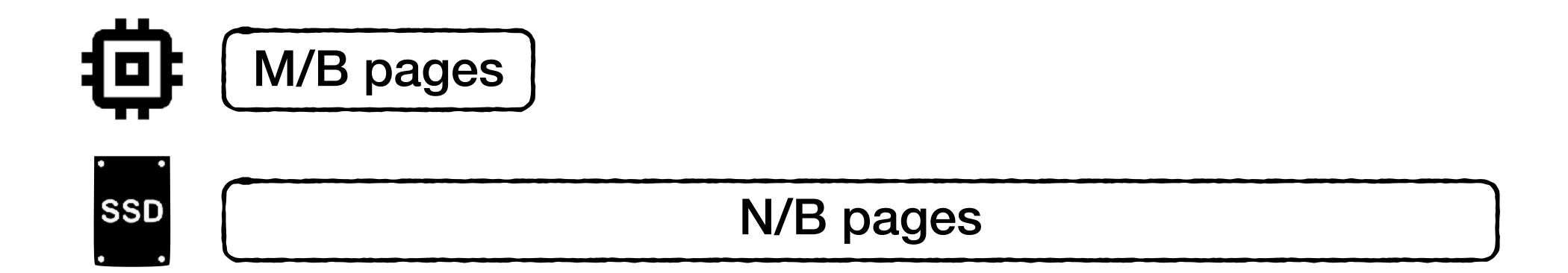
Hence, memory can accommodate $\sqrt{N/B}$ buffers



Two-Pass Merge Sort Algorithm

Use at least $M = \sqrt{N \cdot B}$ memory to partition the data.

This creates at most $N/M = \sqrt{N/B}$ sorted partitions

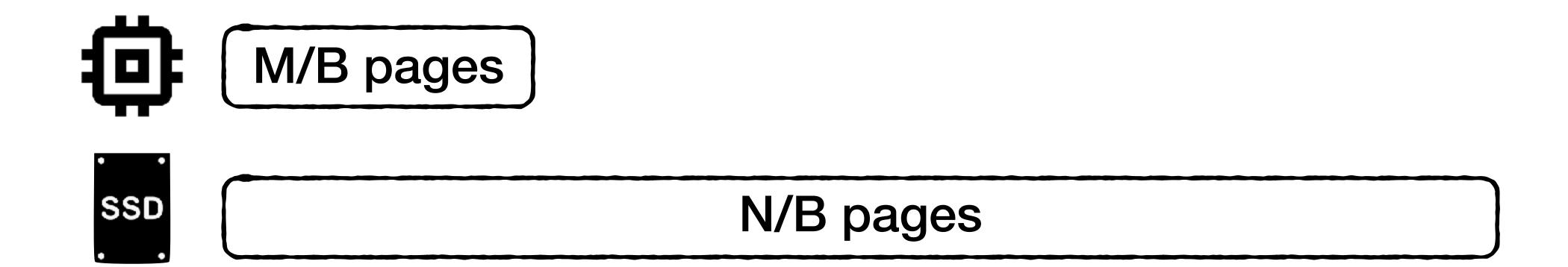


Two-Pass Merge Sort Algorithm

Use at least $M = \sqrt{N \cdot B}$ memory to partition the data.

This creates at most $N/M = \sqrt{N/B}$ sorted partitions

Then merge in one pass using at most $\sqrt{N/B}$ input buffers



Two-Pass Merge Sort Algorithm

Use at least $M = \sqrt{N \cdot B}$ memory to partition the data.

This creates at most $N/M = \sqrt{N/B}$ sorted partitions

Then merge in one pass using at most $\sqrt{N/B}$ input buffers

Cost: O(N/B)





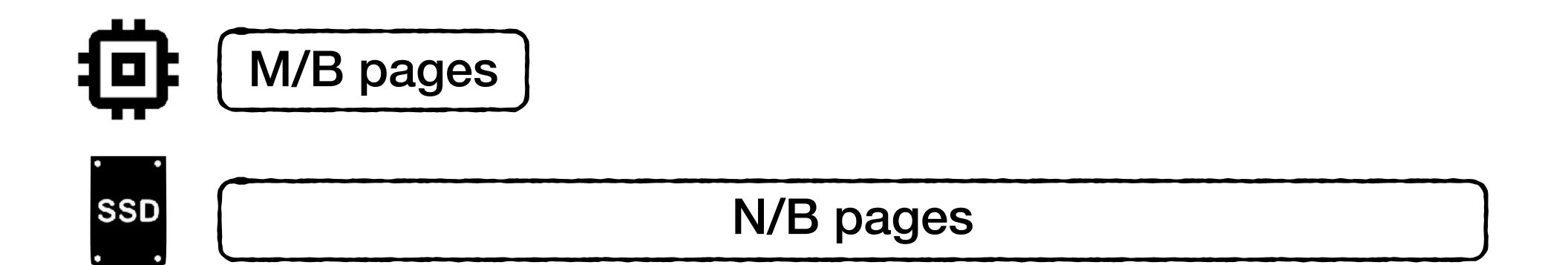
N/B pages

How much memory do we need in practice?

Assume 1 TB, 16 byte entries, and 4KB pages

N=2³⁶ entries. And with 4KB pages, B=2⁸ entries.

We need
$$M = \sqrt{N \cdot B} = \sqrt{2^{44}} = 2^{22}$$
 entries in memory, or 64 MB



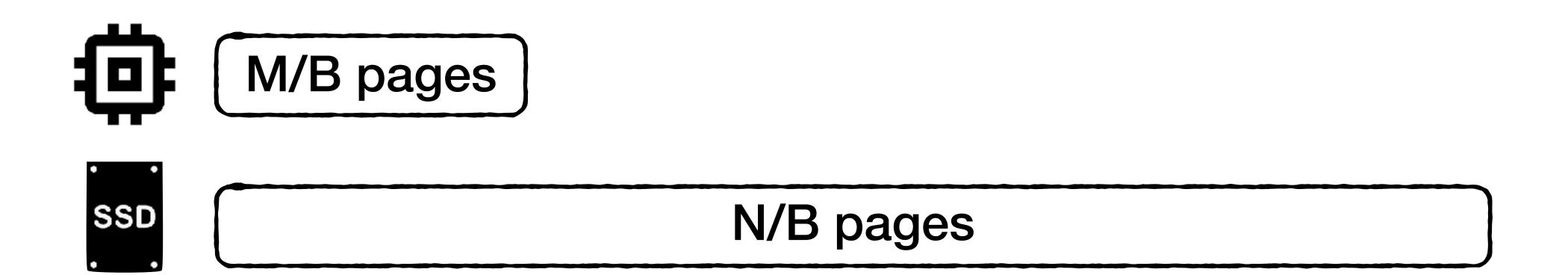
How much memory do we need in practice?

Assume 1 TB, 16 byte entries, and 4KB pages

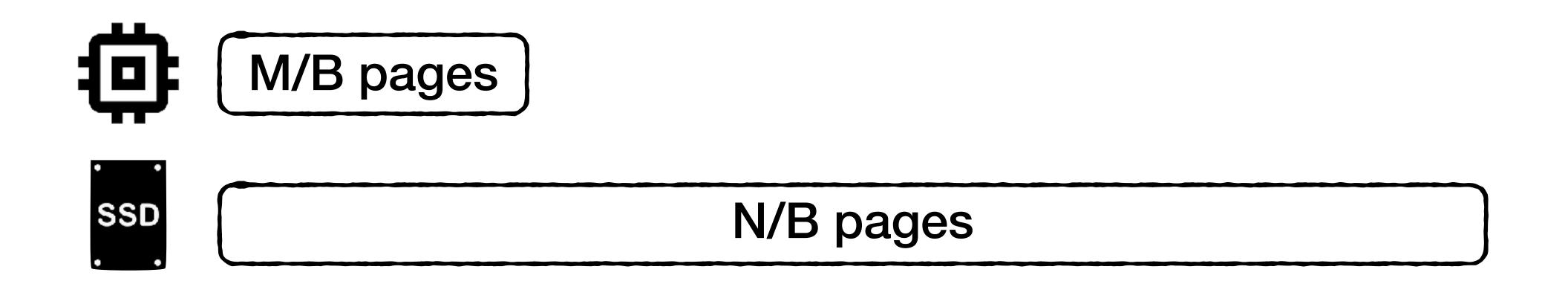
N=2³⁶ entries. And with 4KB pages, B=2⁸ entries.

We need
$$M = \sqrt{N \cdot B} = \sqrt{2^{44}} = 2^{22}$$
 entries in memory, or 64 MB

Hence, for all practical purposes, a 2 pass sorting algorithm is practical.



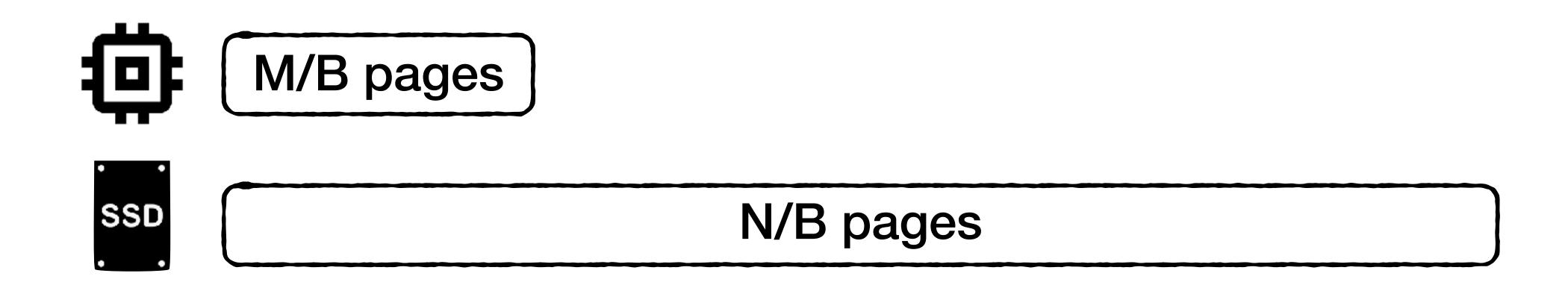
Achieved our goal of sorting using O(N/B) I/Os using little memory:)



Achieved our goal of sorting using O(N/B) I/Os using little memory:)

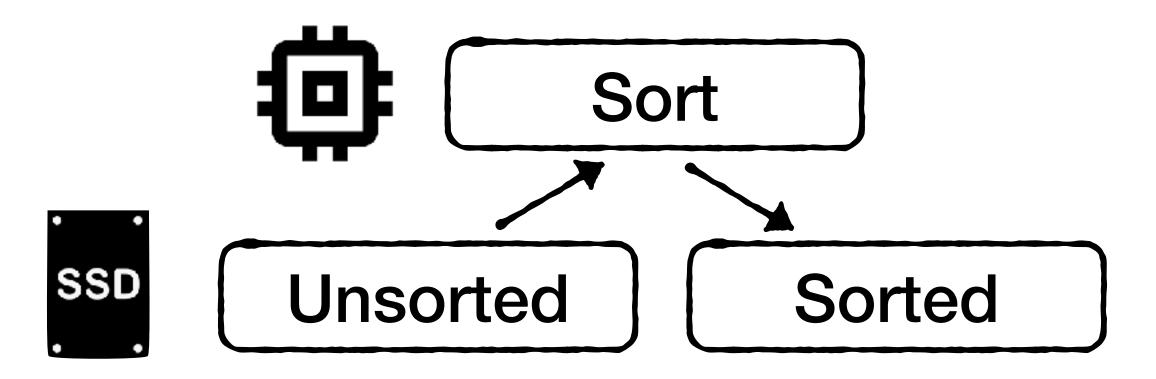
But how about CPU overheads?

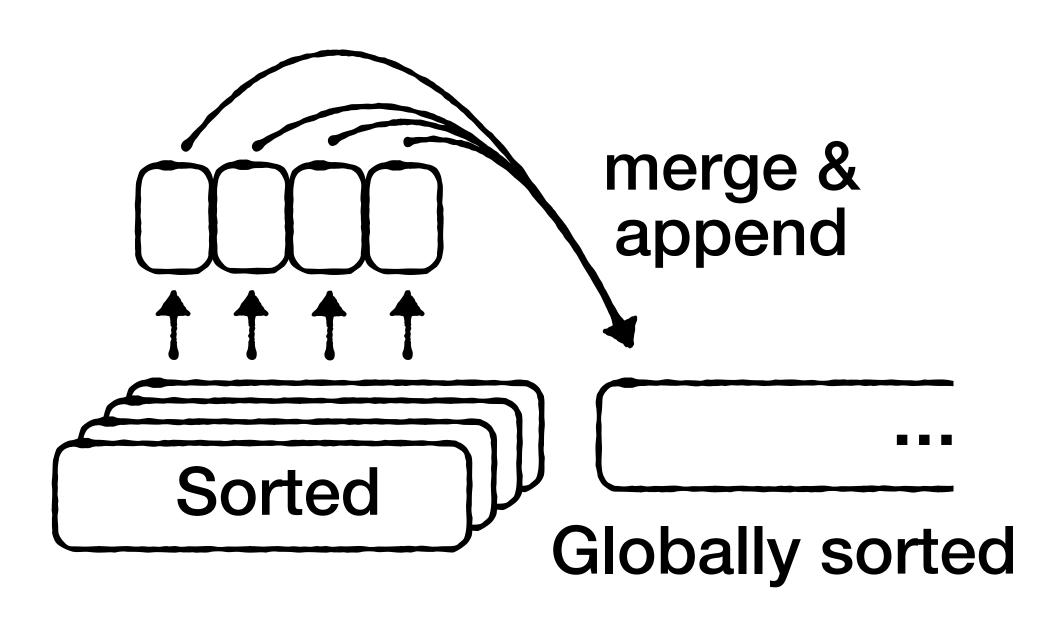
We still expect O(N · log₂ N). Do we achieve it?



Analyzing CPU

Partitioning Phase



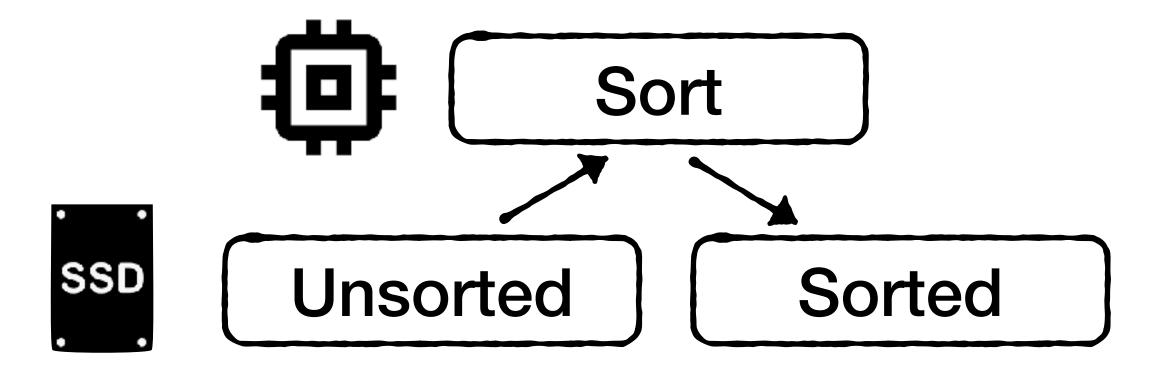


Partitioning Phase

Each chunk contains M entries

Need O(M · log₂ M) CPU cycles to merge-sort it in-memory

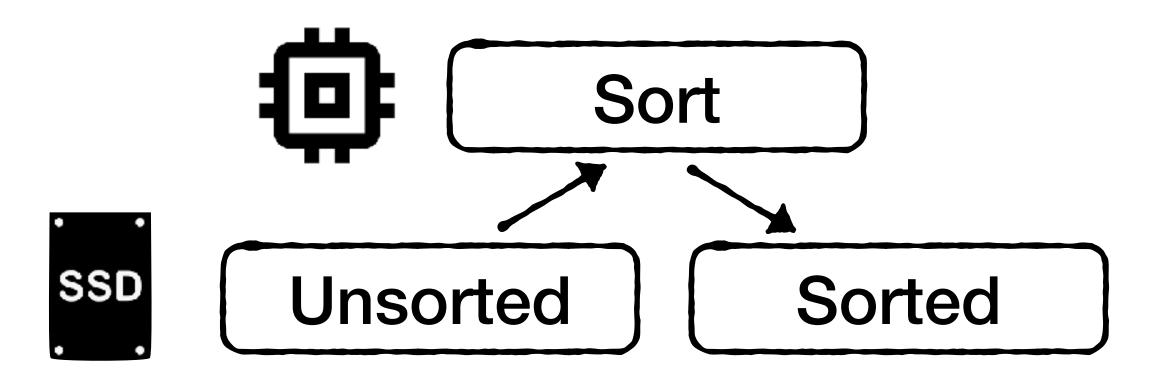
Doing this for all N/M chunk takes O(N · log₂ M) CPU

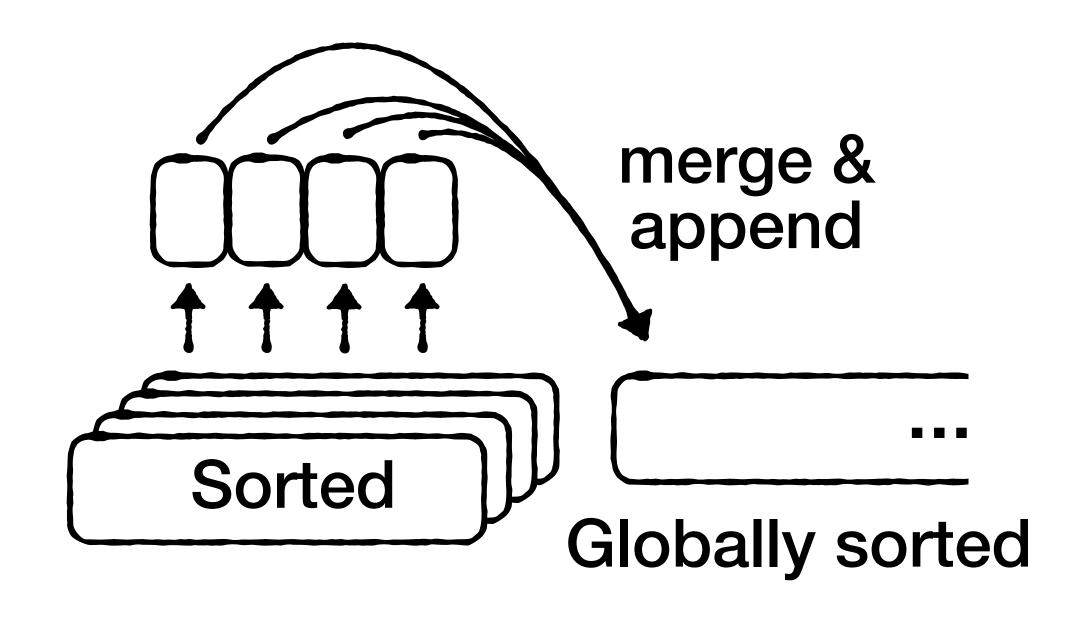


Analyzing CPU

Partitioning Phase

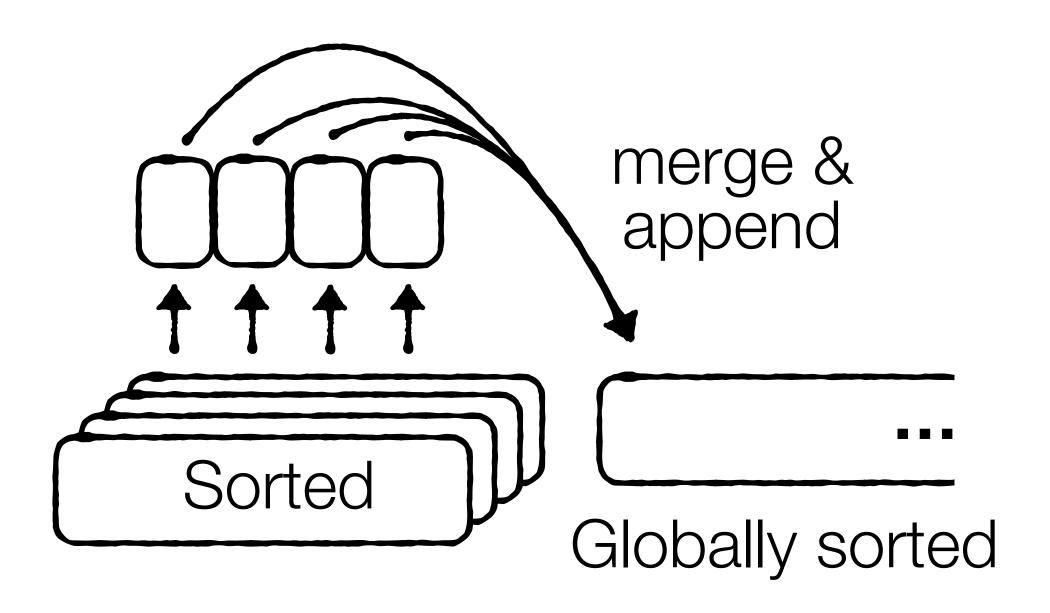
 $O(N \cdot log_2 M)$

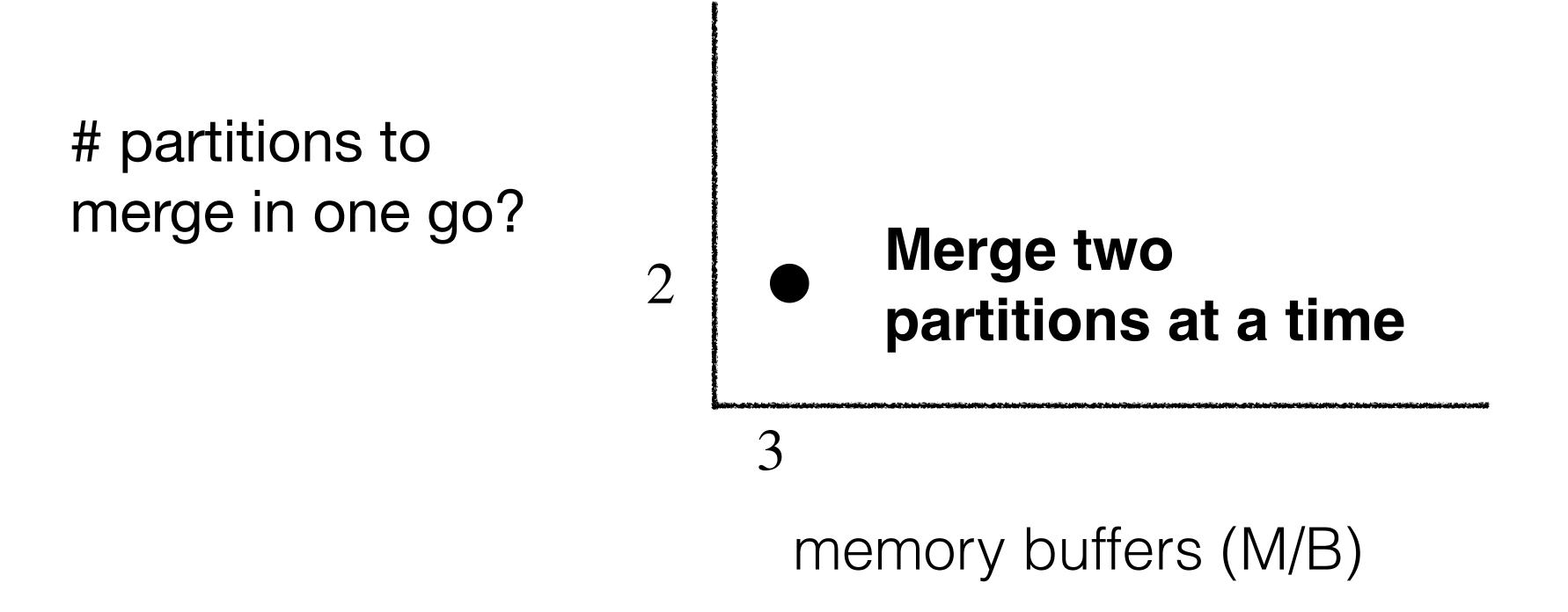


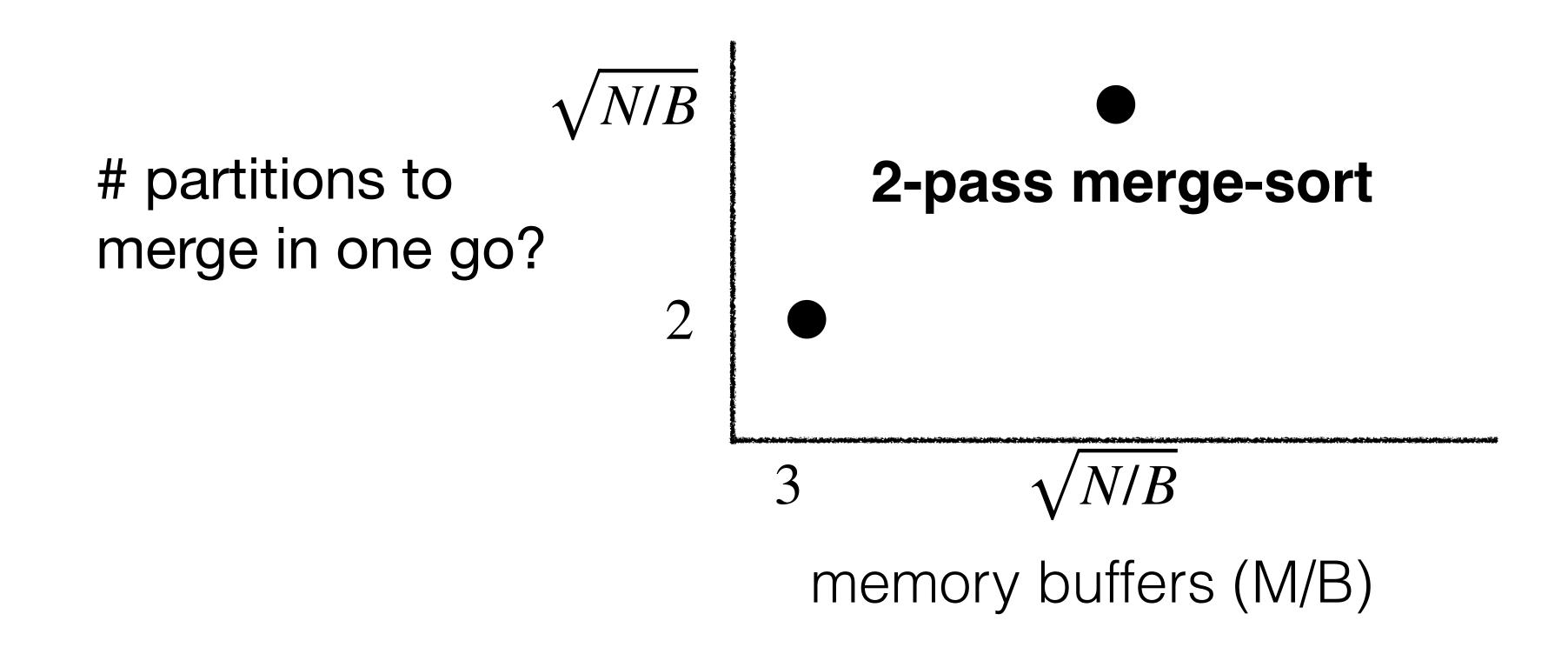


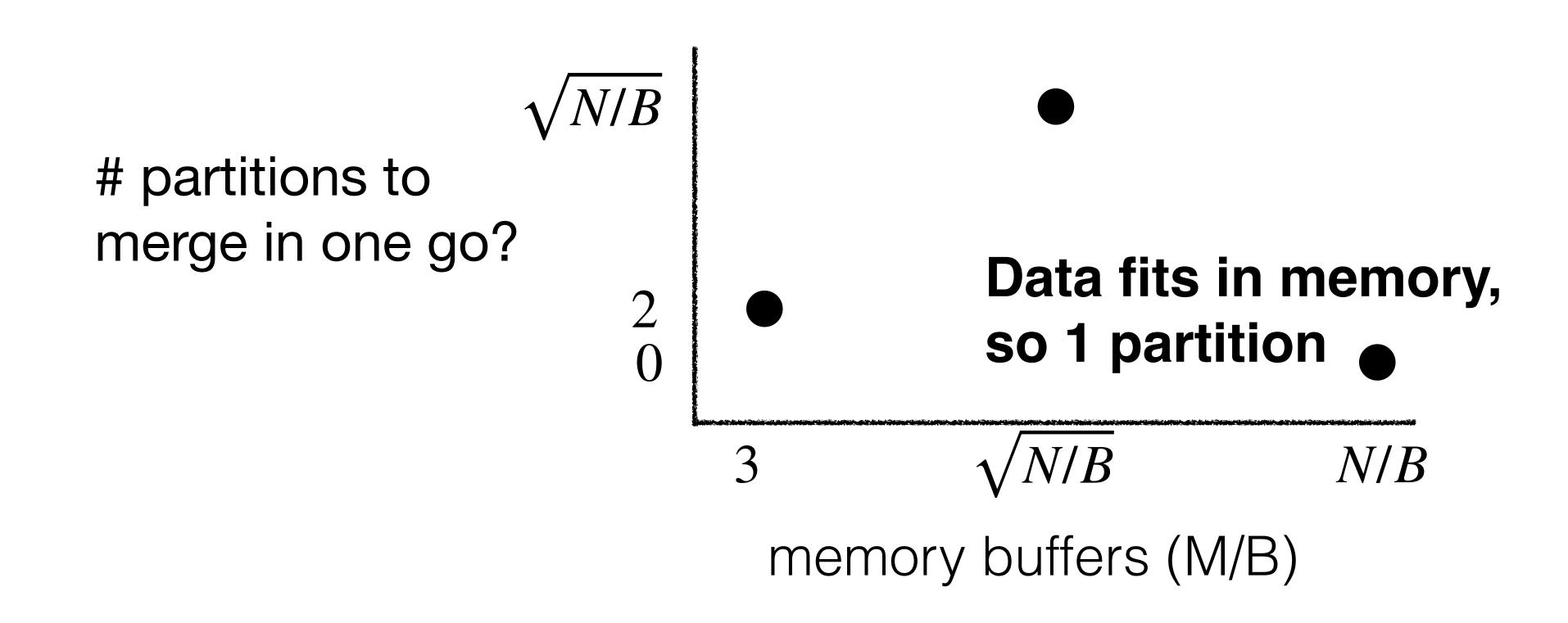
maximum # partitions to merge in one go?

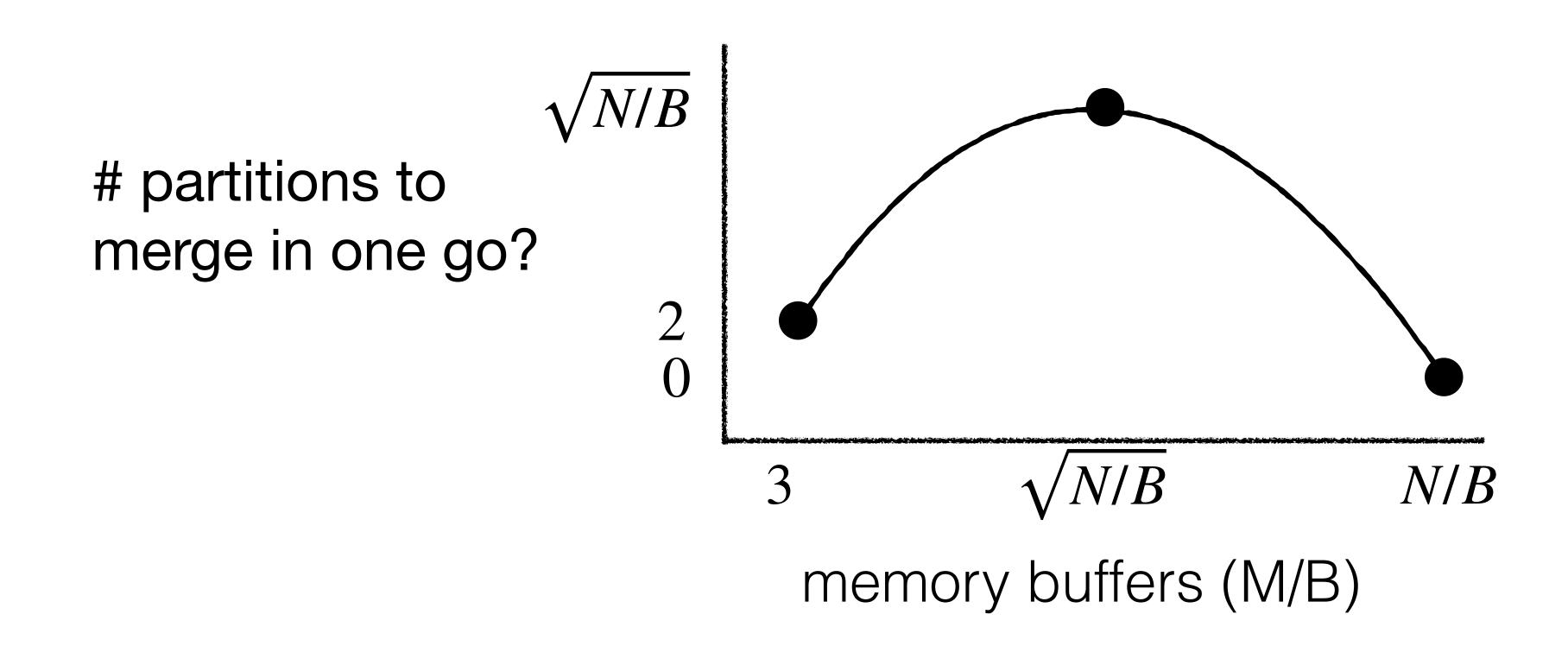
How do we merge them?

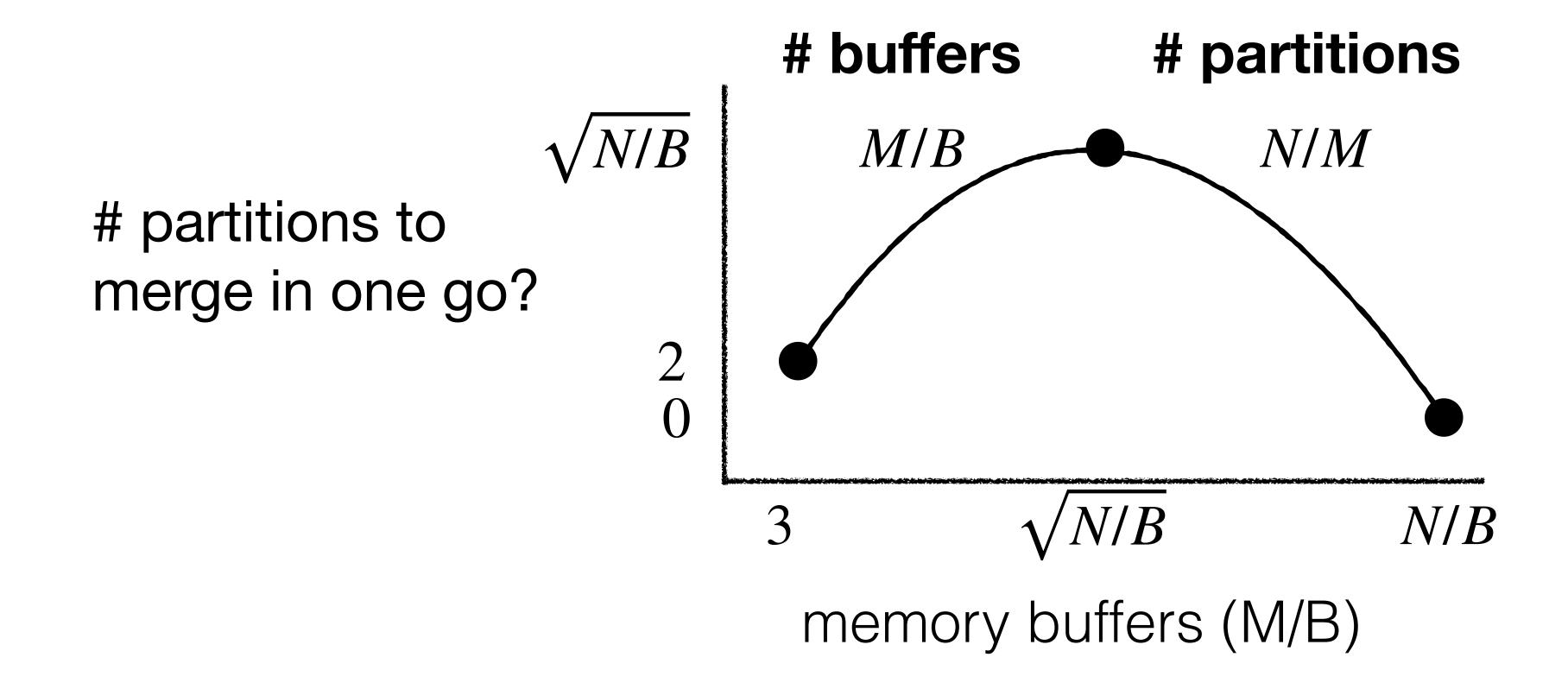










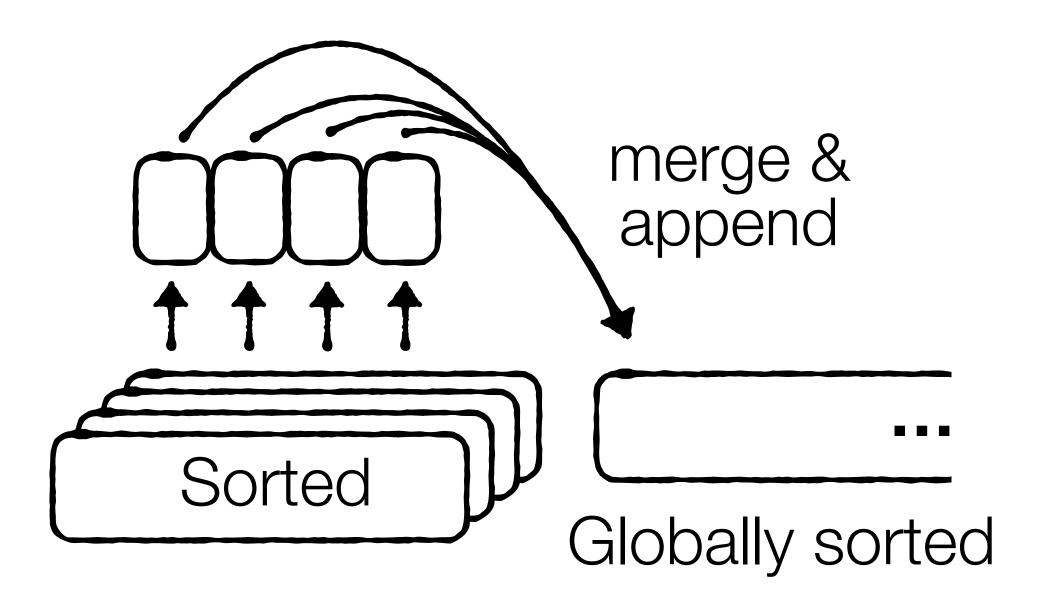


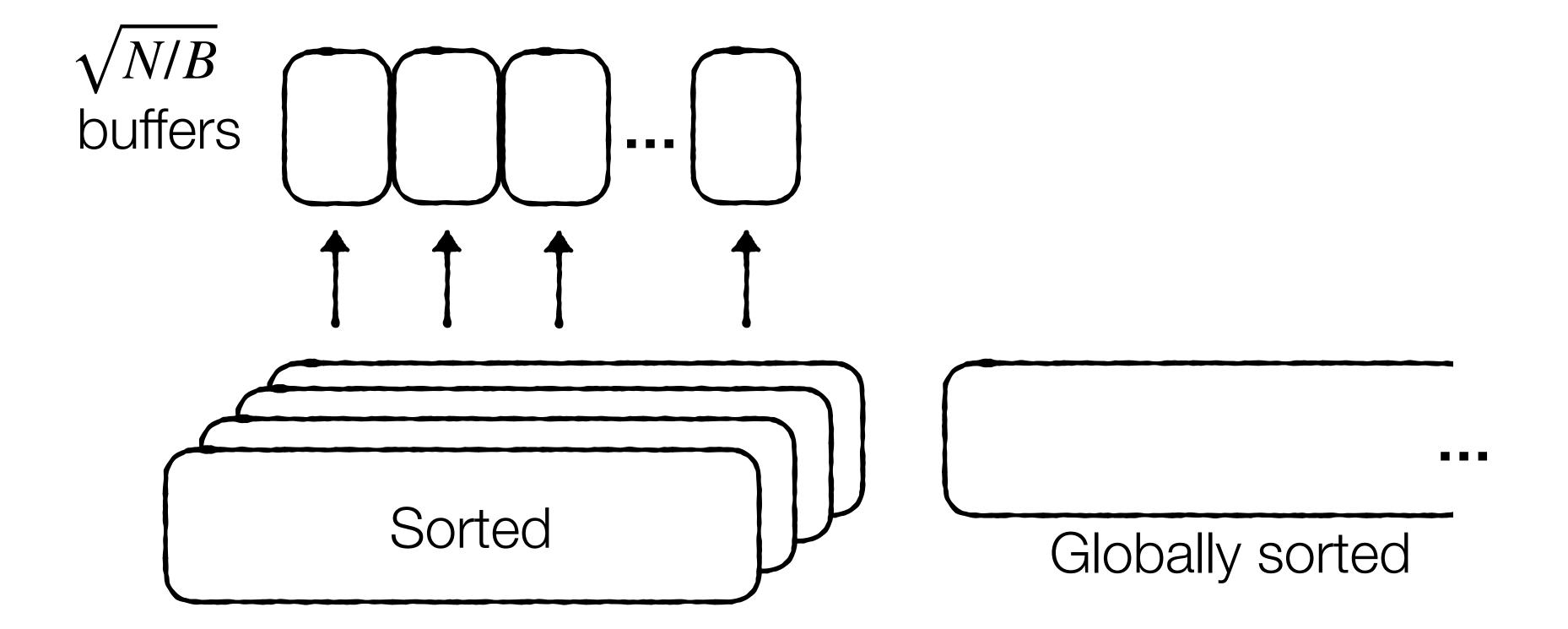
memory buffers (M/B)

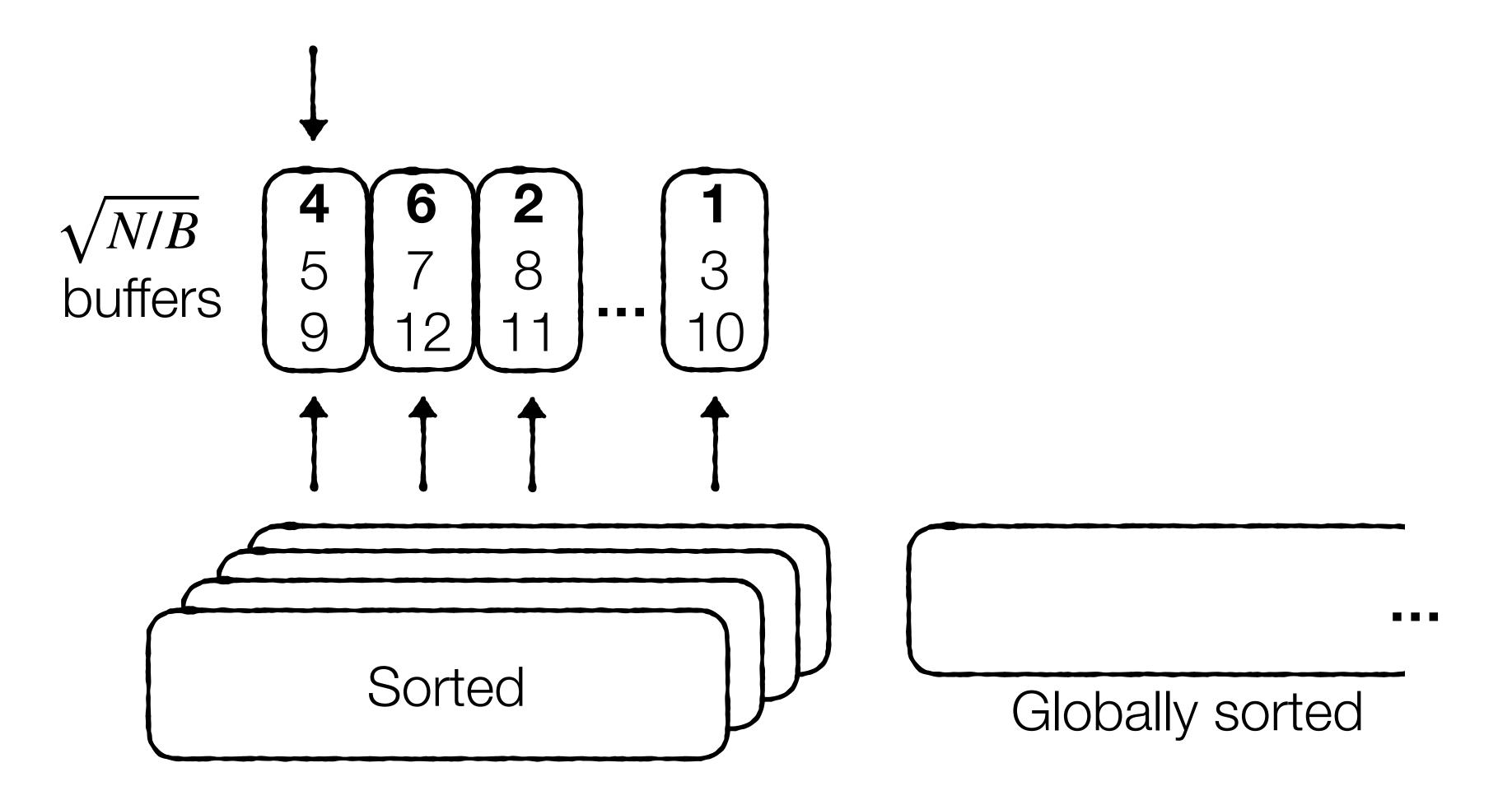
maximum # partitions to merge in one go?

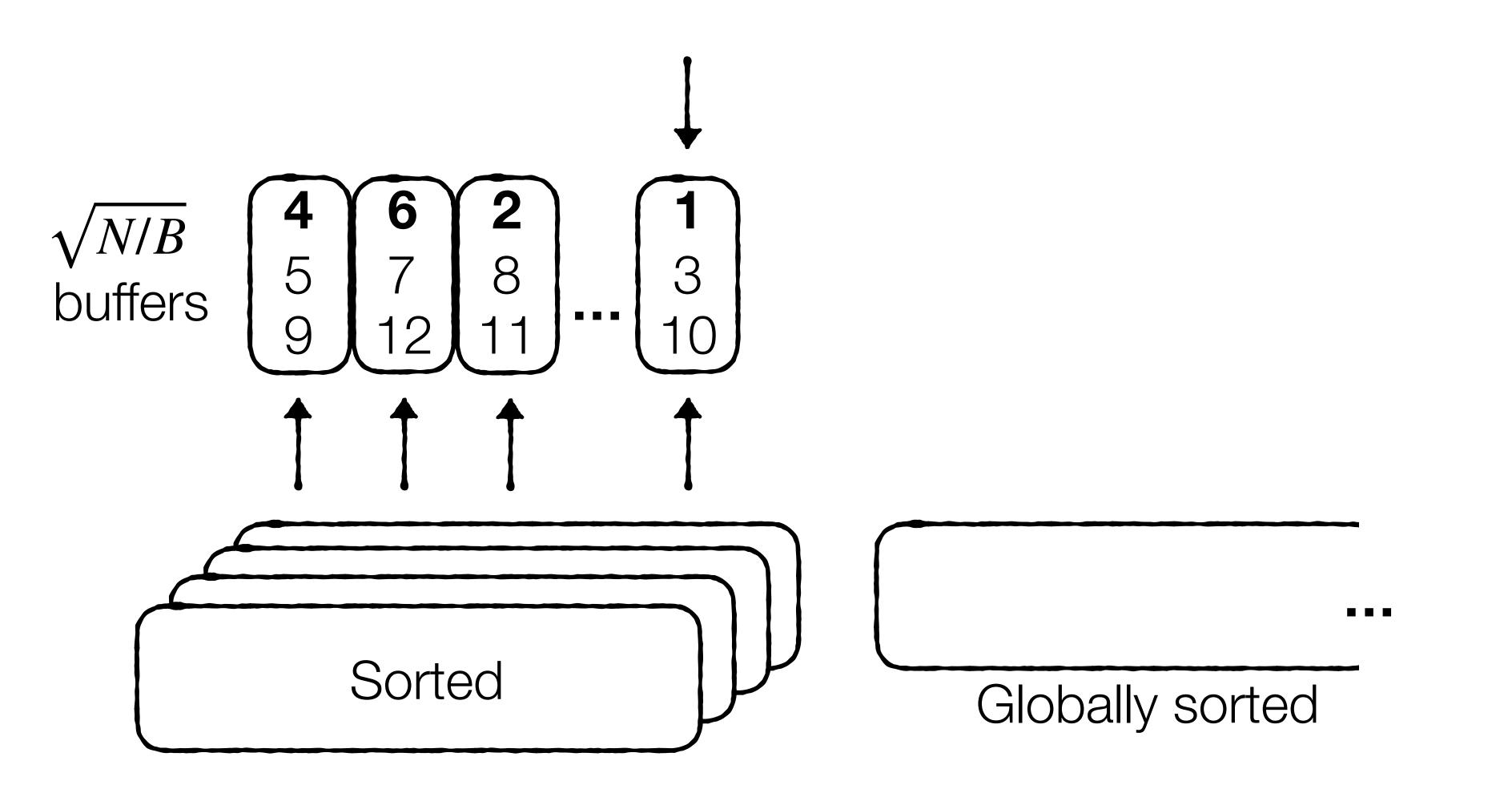
$$\sqrt{N/B}$$

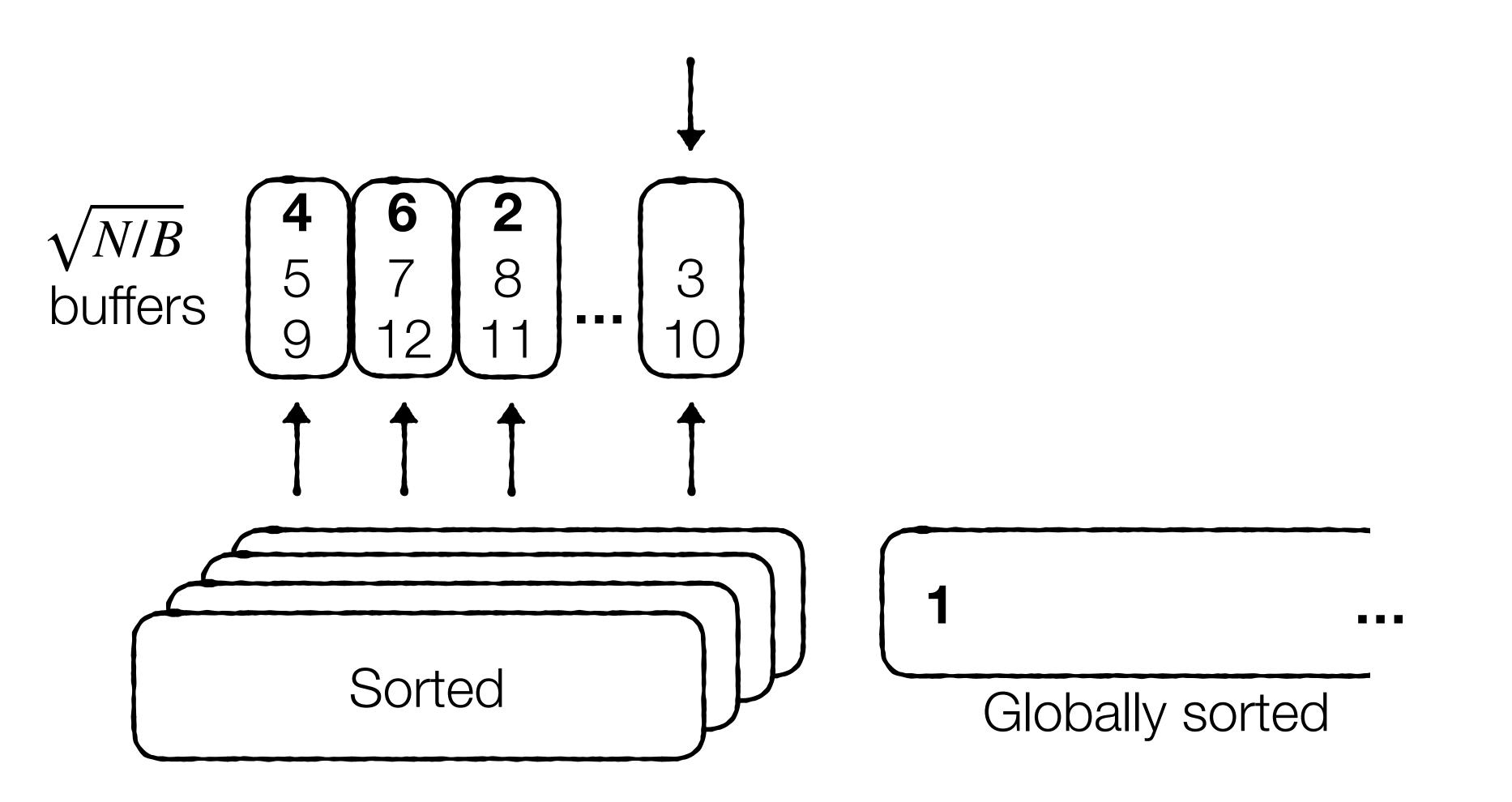
How to merge partitions?

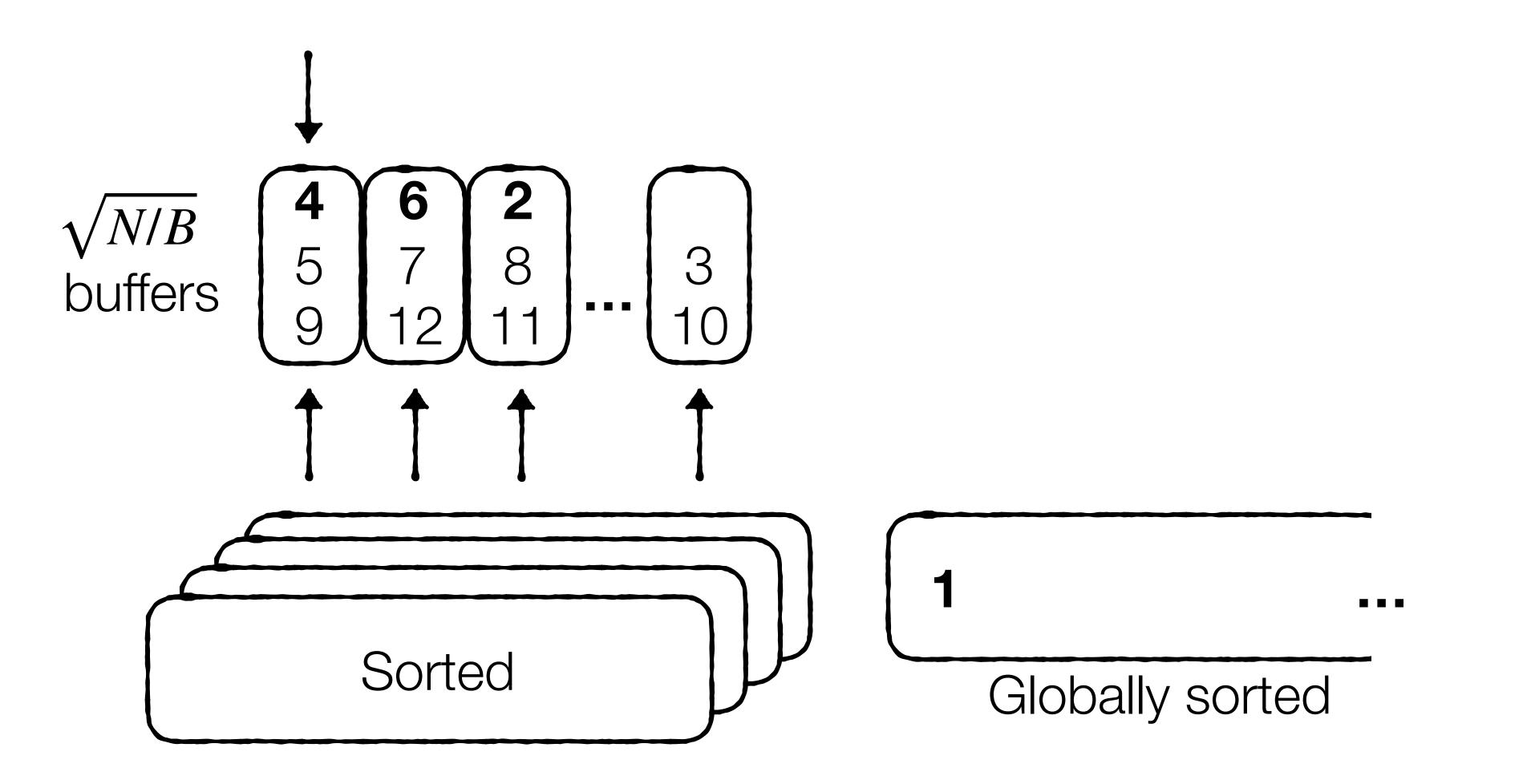


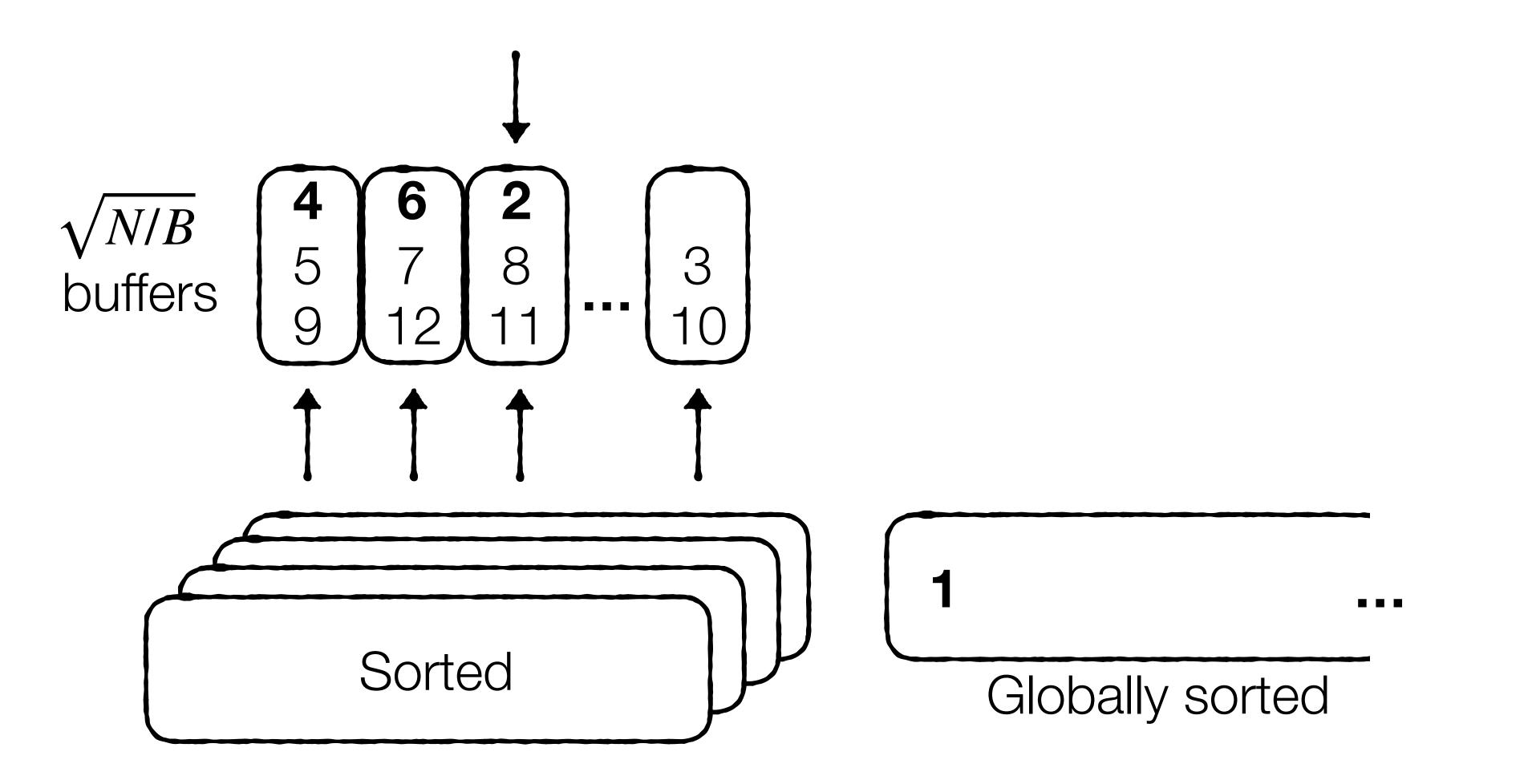


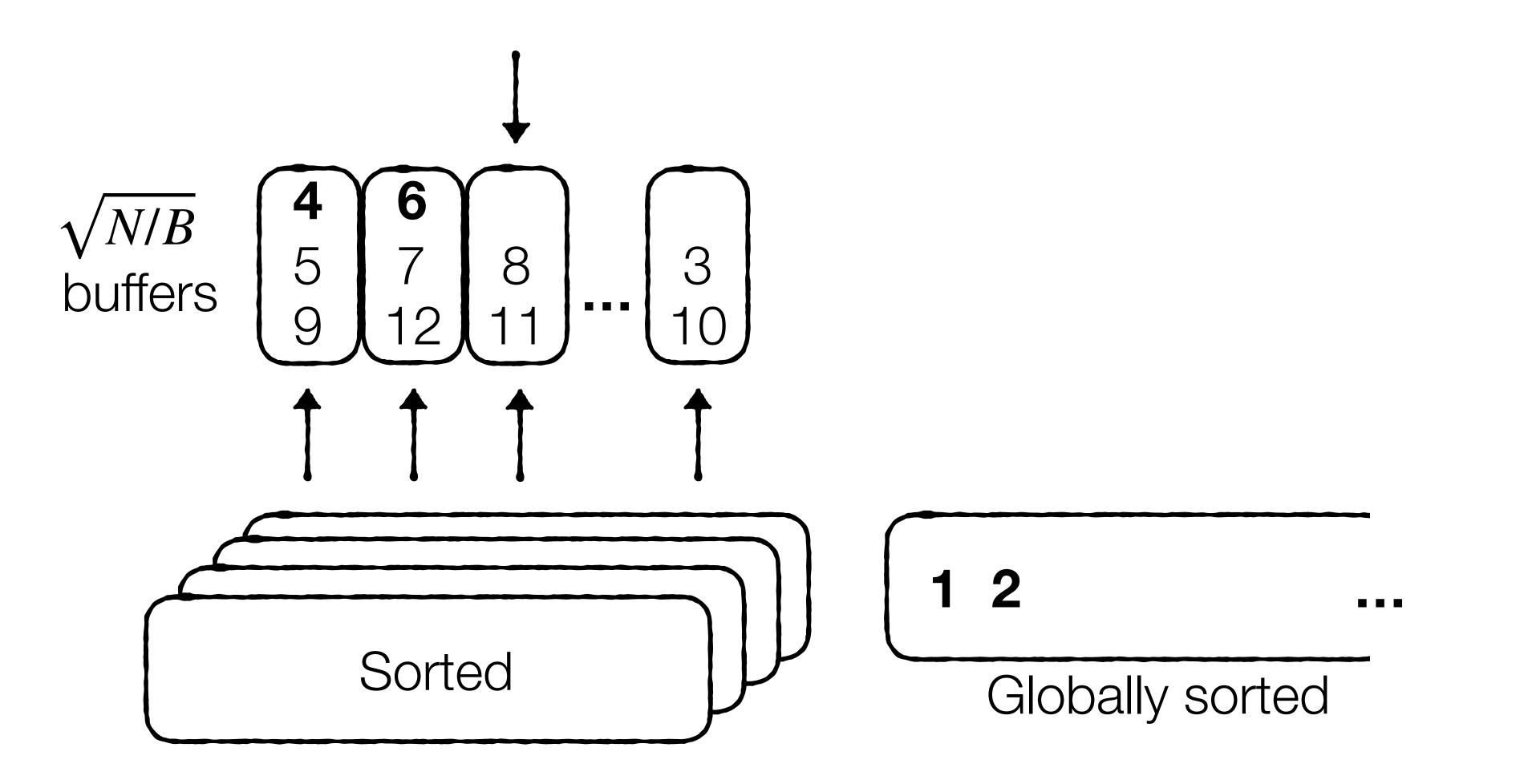


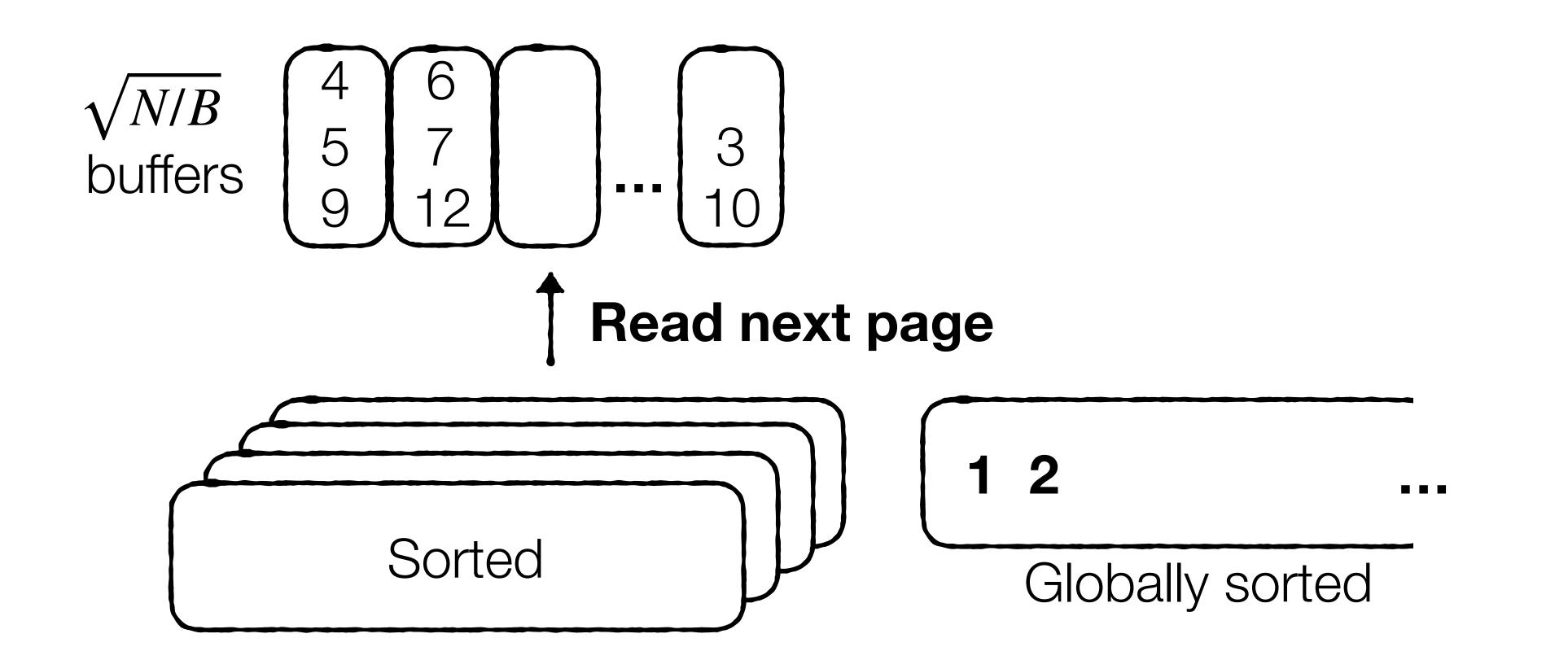


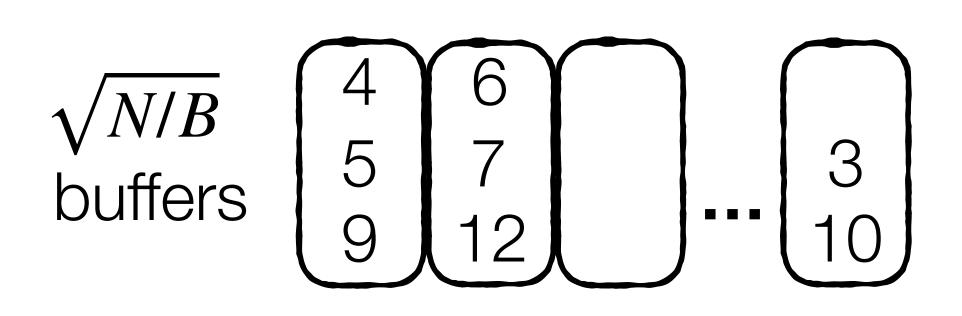




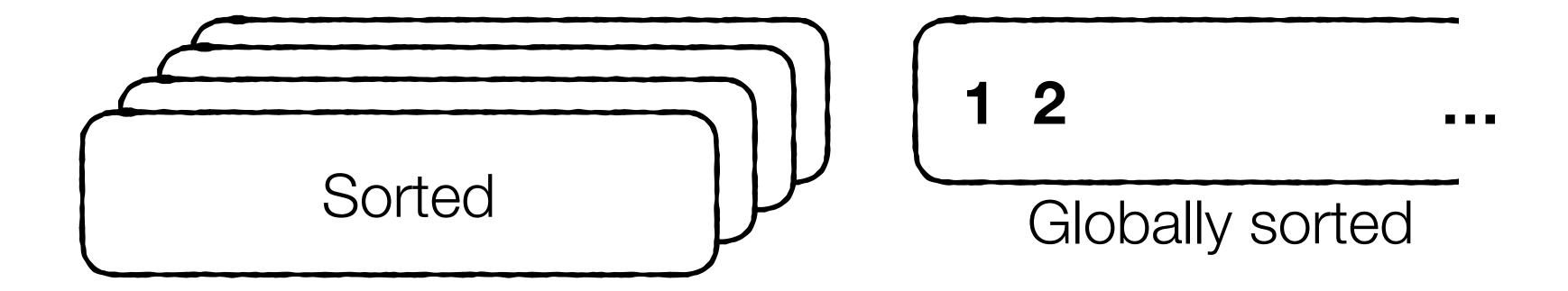




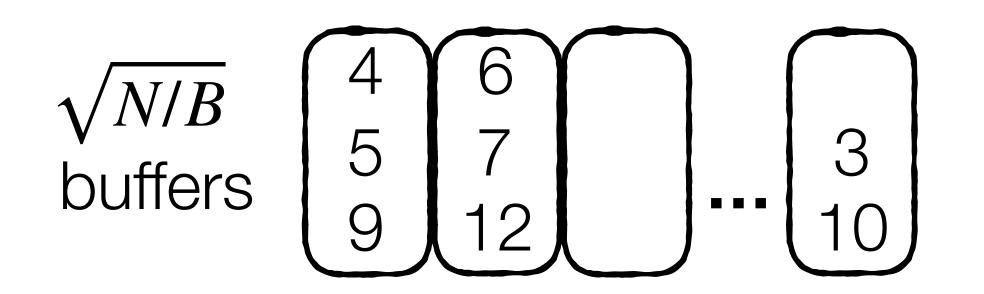


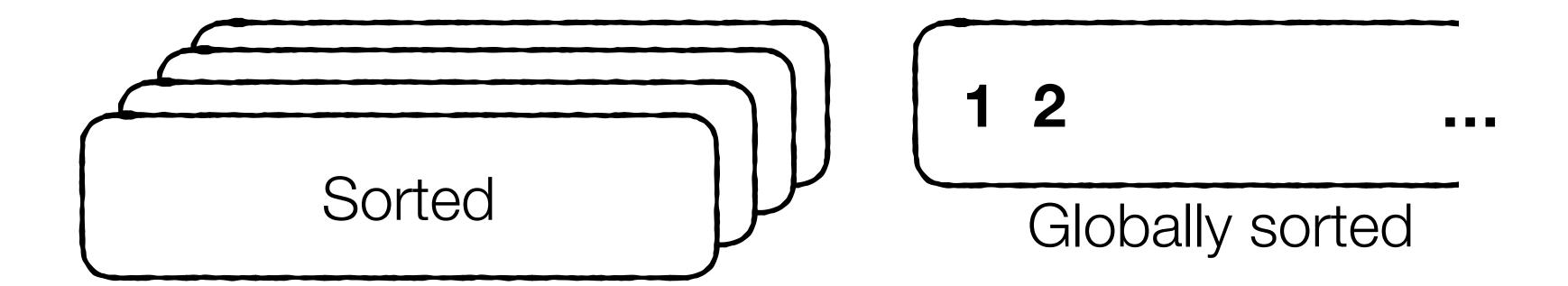


$$O(\sqrt{N/B} \cdot N)$$
 comparisons



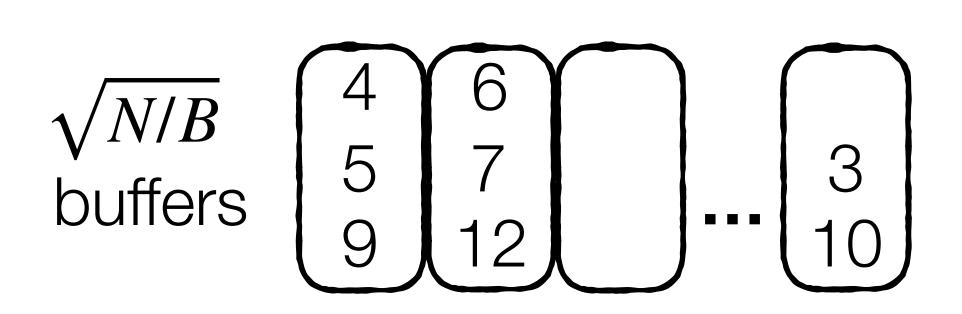
How to merge partitions more efficiently?



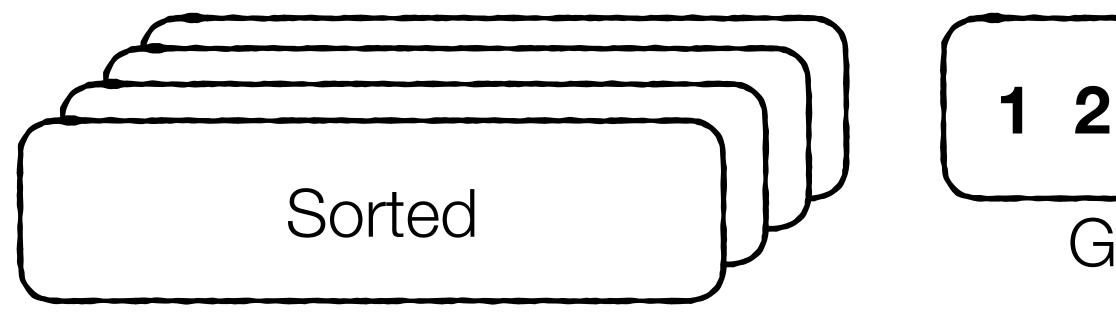


How to merge partitions more efficiently?

Discuss with your neighbors (2 min)



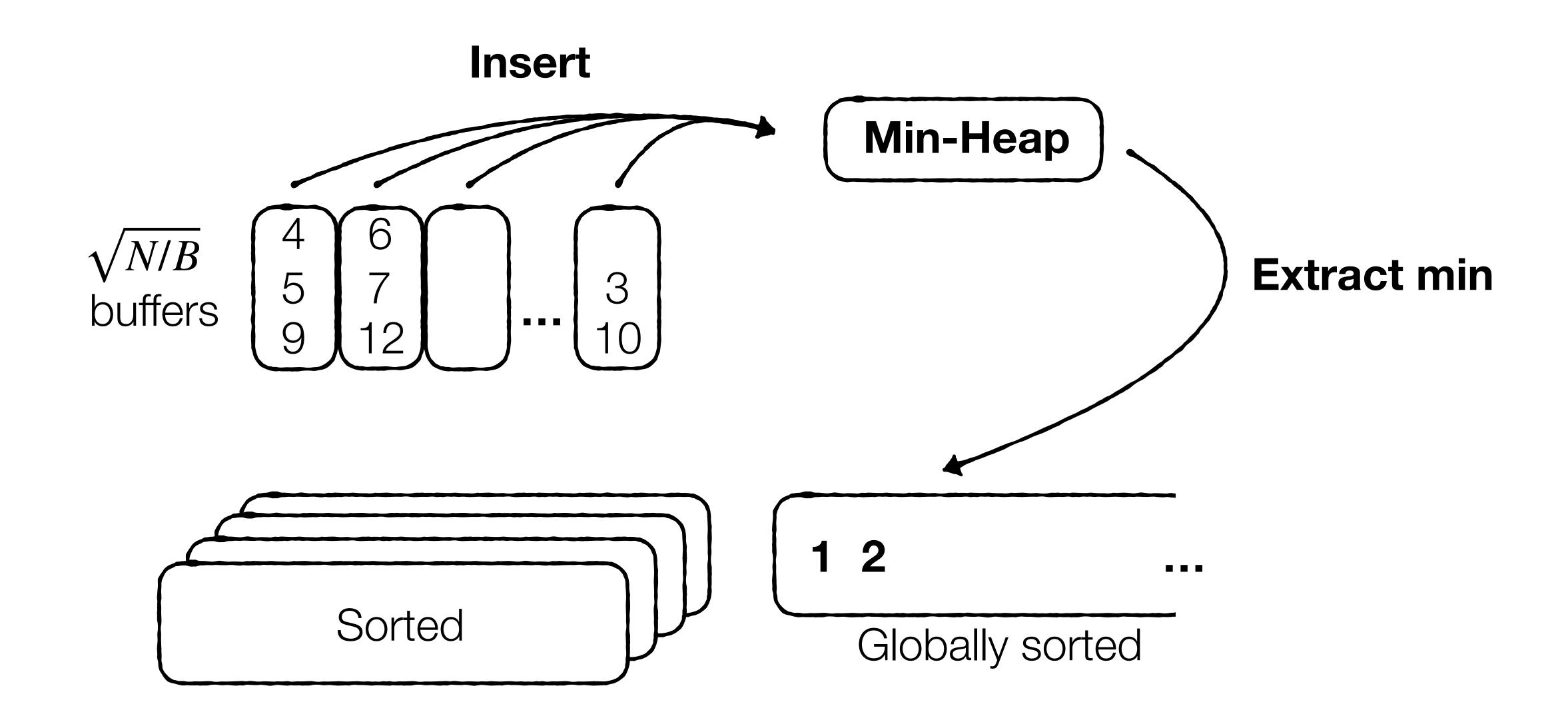




1 2

Globally sorted

How to merge partitions more efficiently?



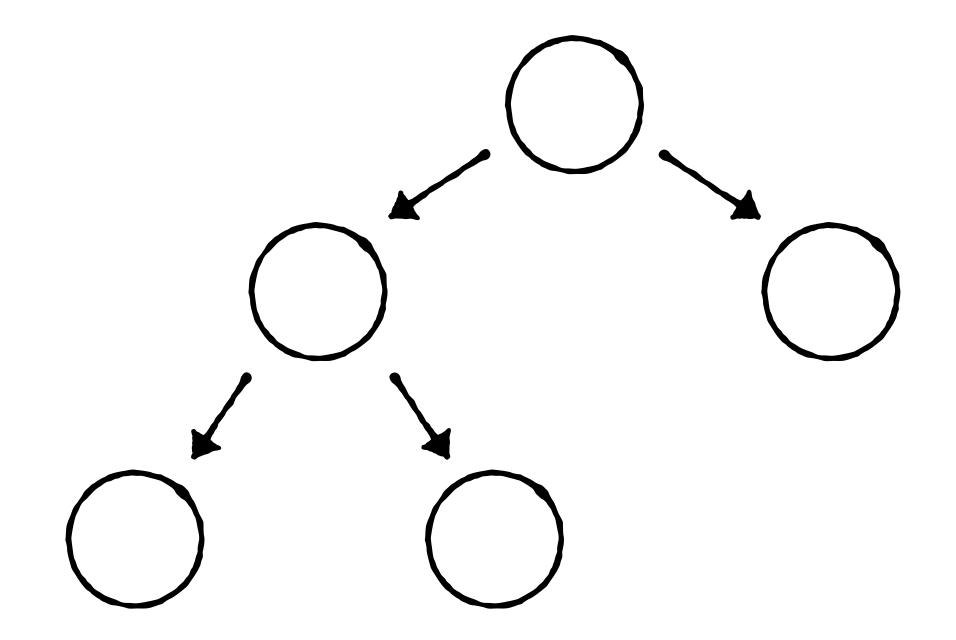
Well-known data structure that efficiently extracts the minimum value in a collection of data items

APIRuntimeInsert(key) $O(log_2 N)$ Key = extract_min() $O(log_2 N)$

Well-known data structure that efficiently extracts the minimum value in a collection of data items

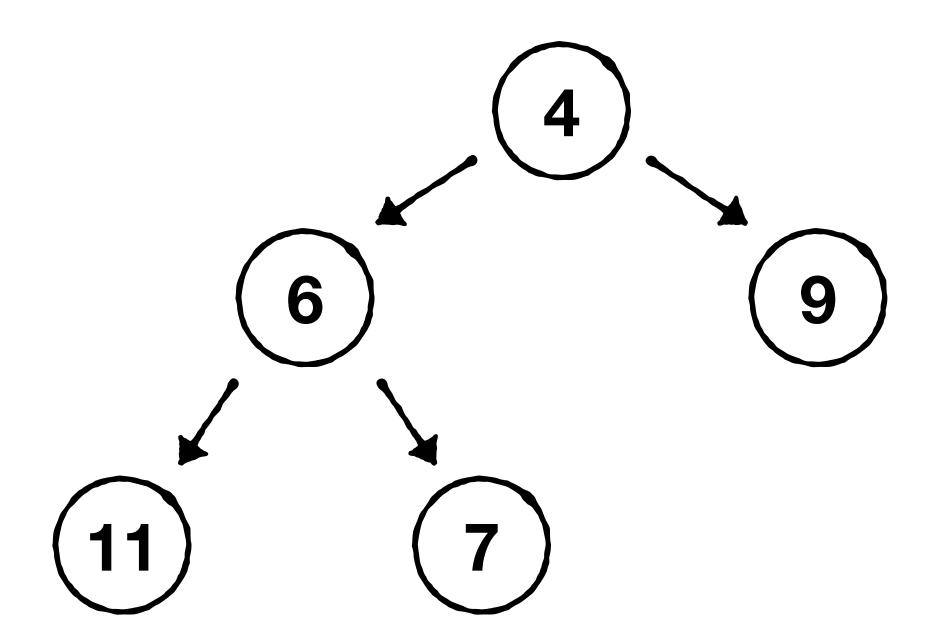
API	Runtime
Insert(key)	O(log ₂ N)
Key = extract_min()	O(log ₂ N)
<pre>min_key = insert_and_extract(new_key) (efficiently combines both operations)</pre>	O(log ₂ N)

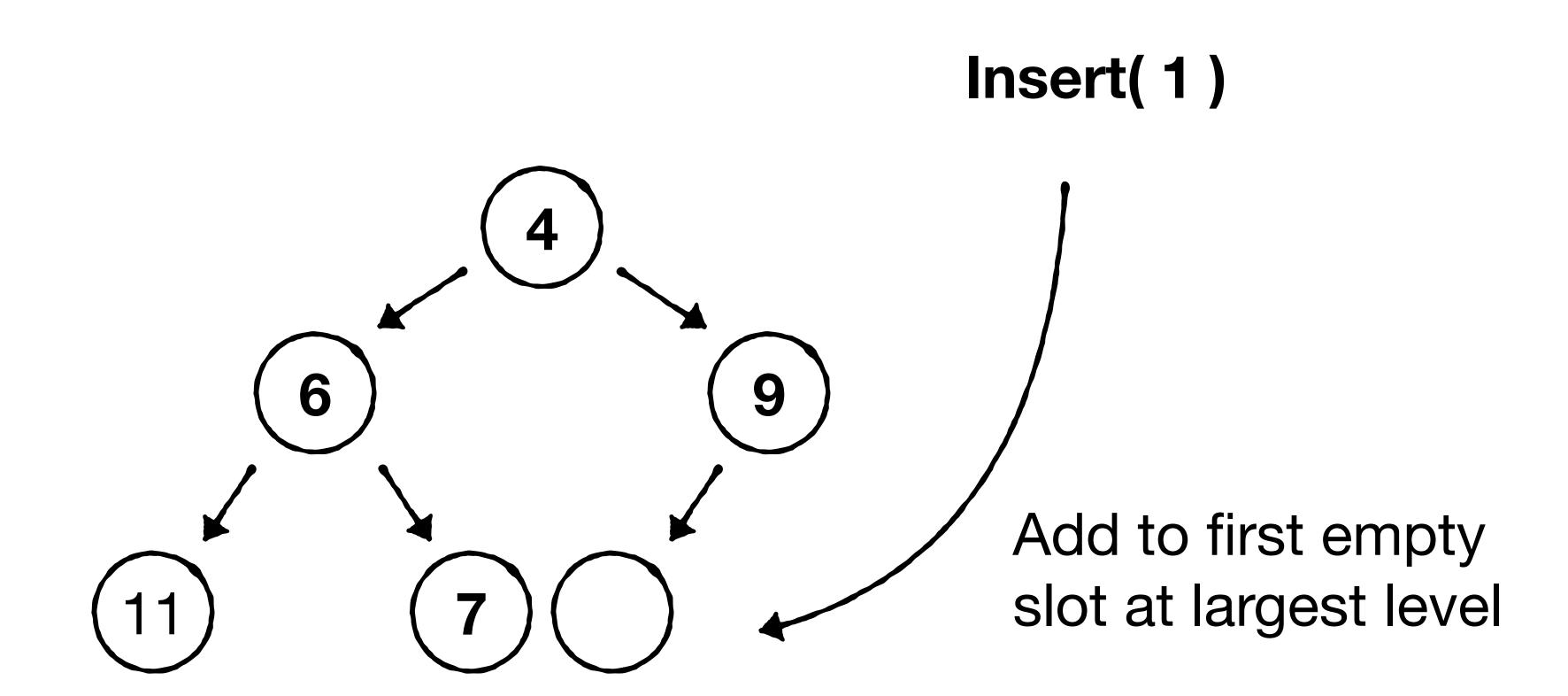
(1) Complete binary tree (All levels are full & largest level is full from left to right)

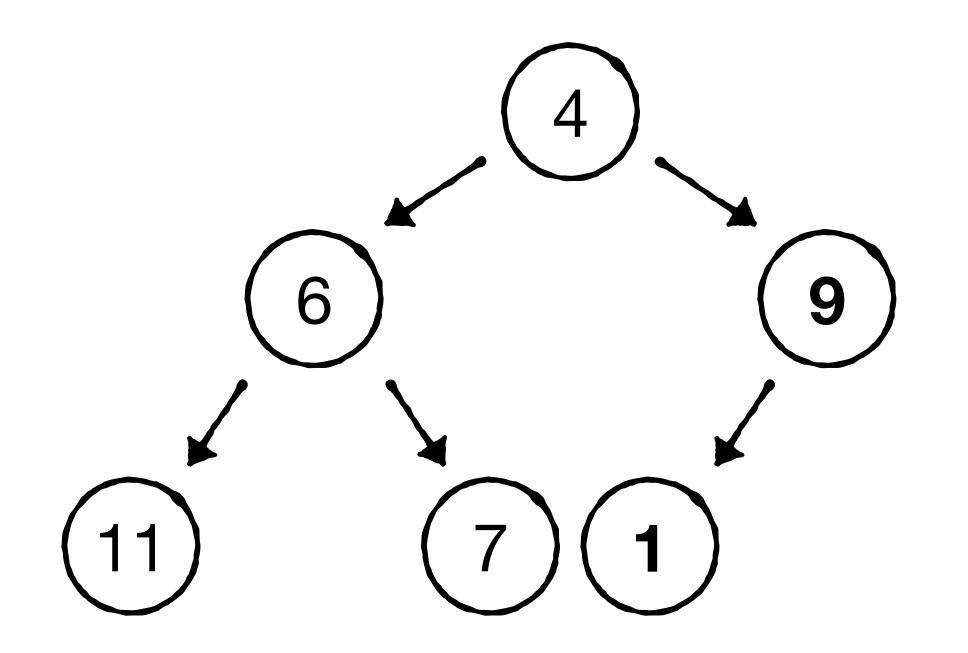


(1) Complete binary tree(All levels are full & largest level is full from left to right)

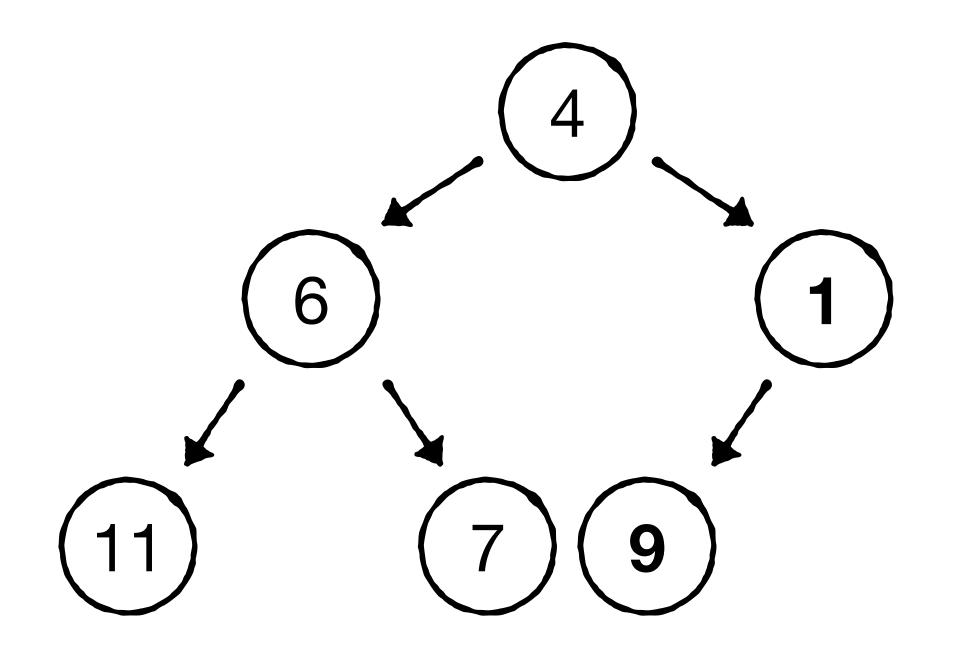
(2) Parent key always smaller than children's.





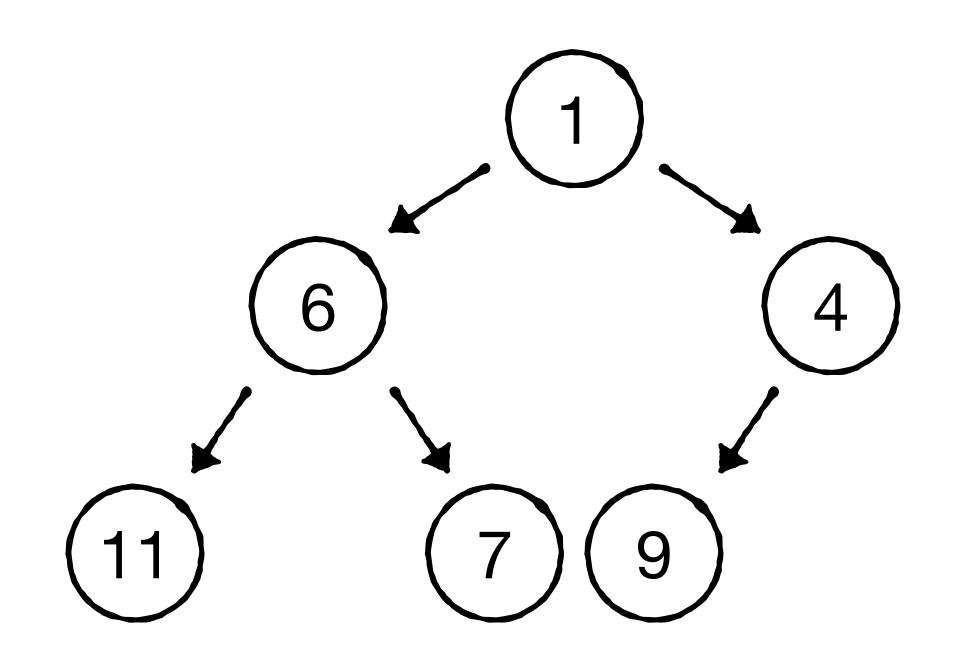


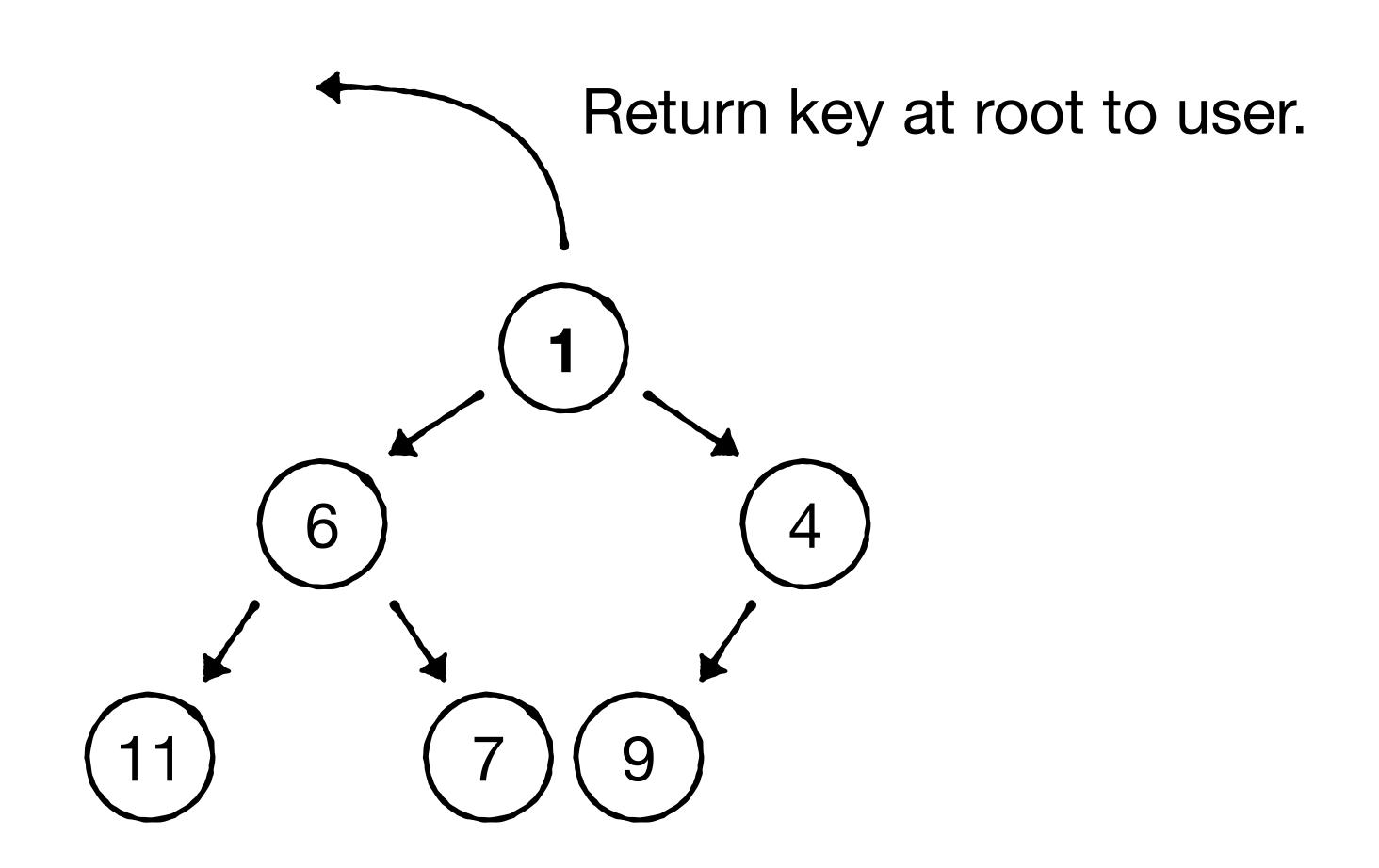
Swap until parent is always smaller

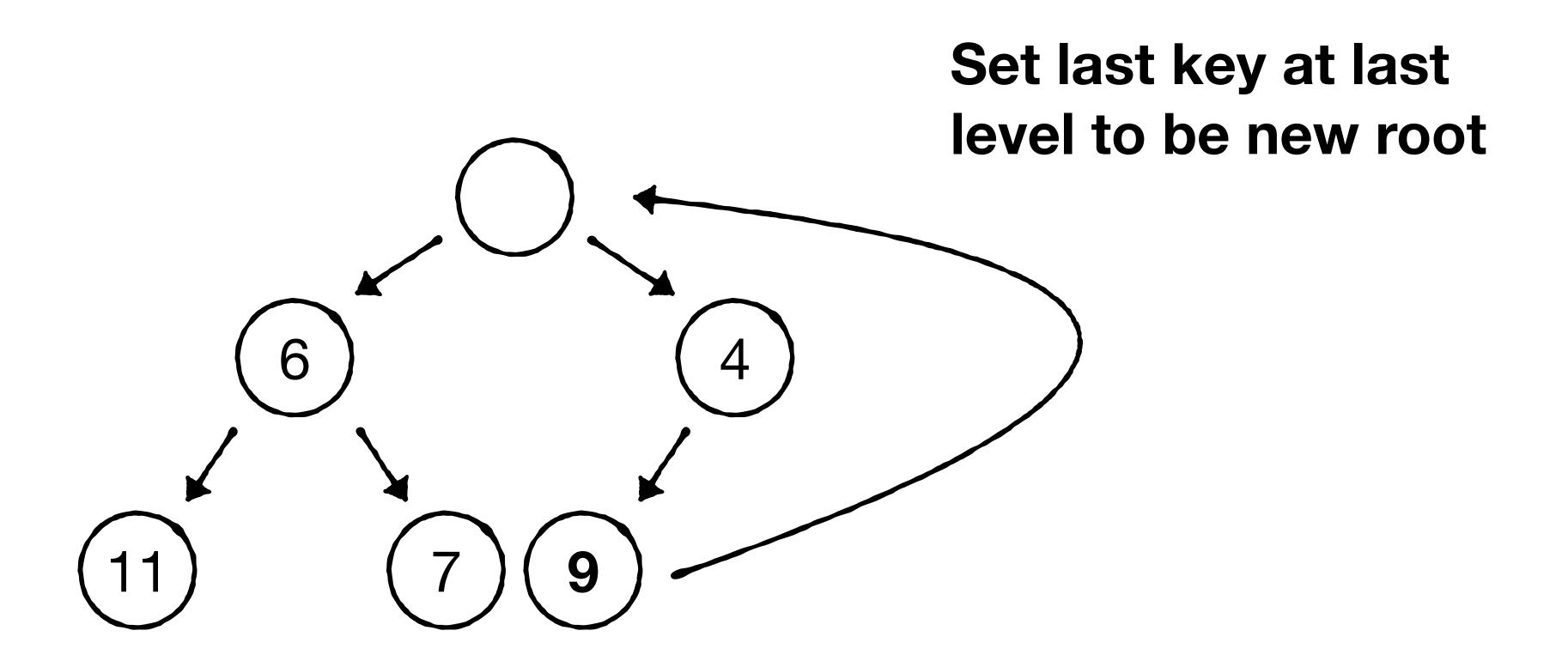


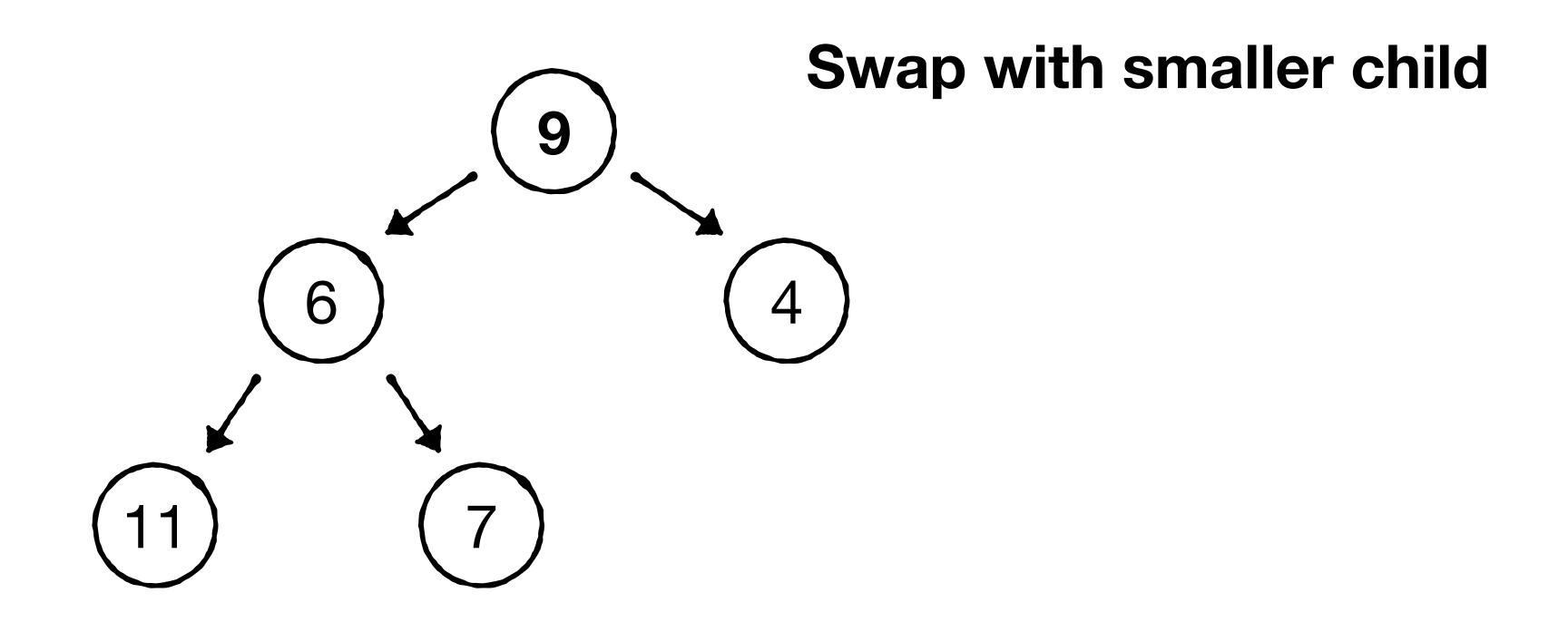
Swap until parent is always smaller

O(log₂ N) insertion cost

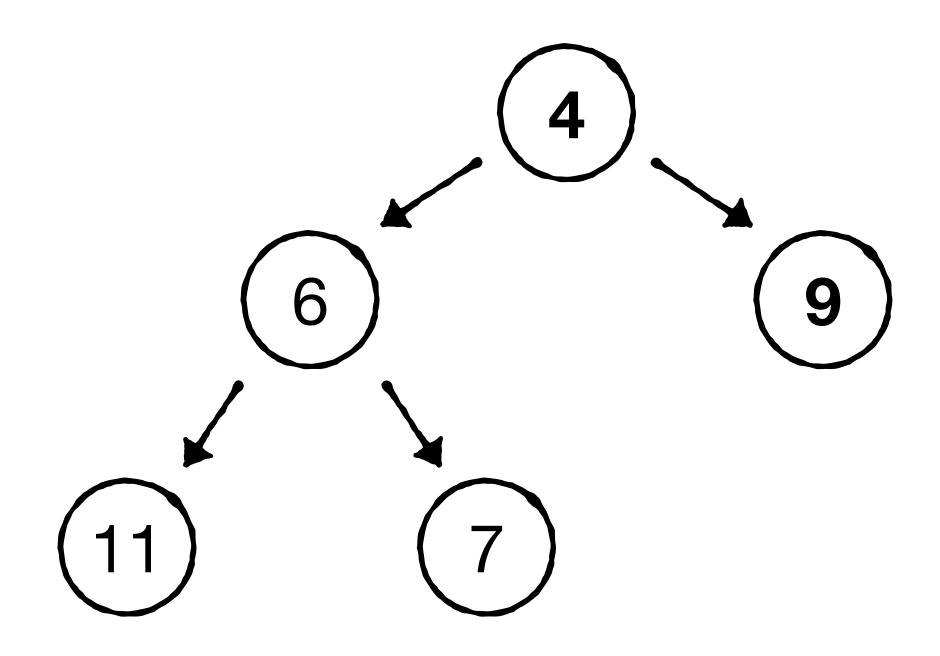




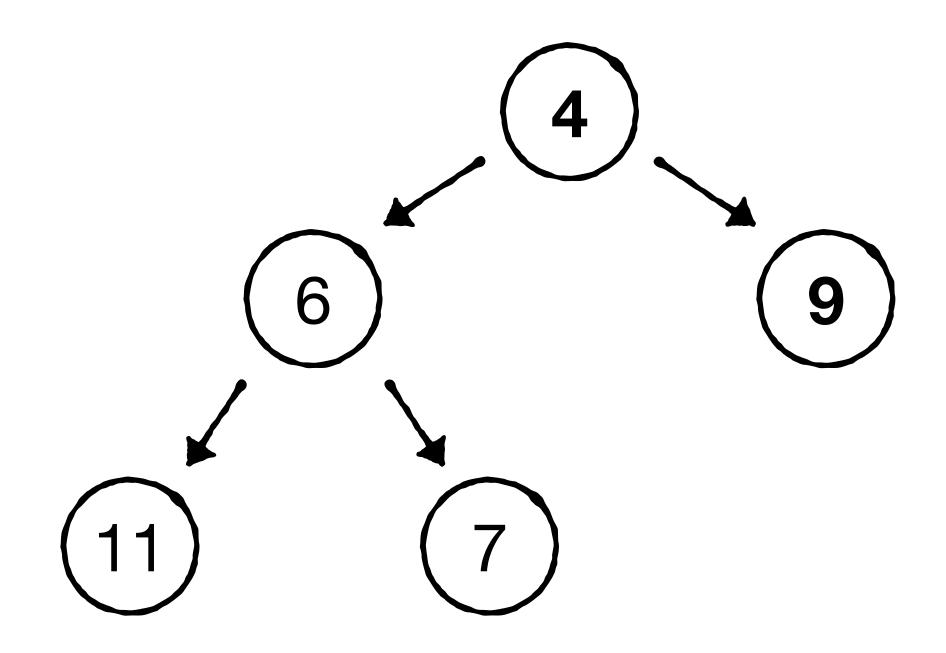


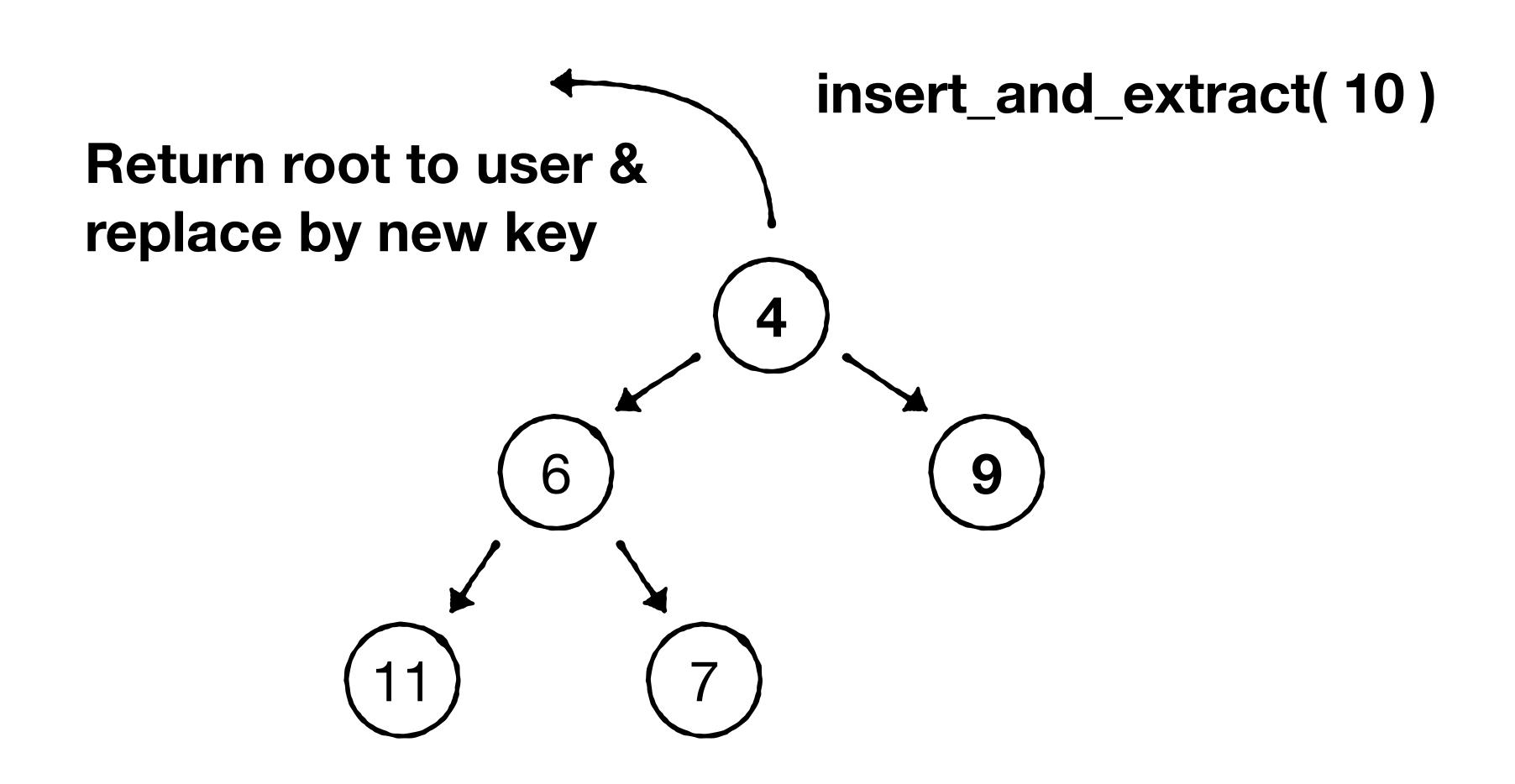


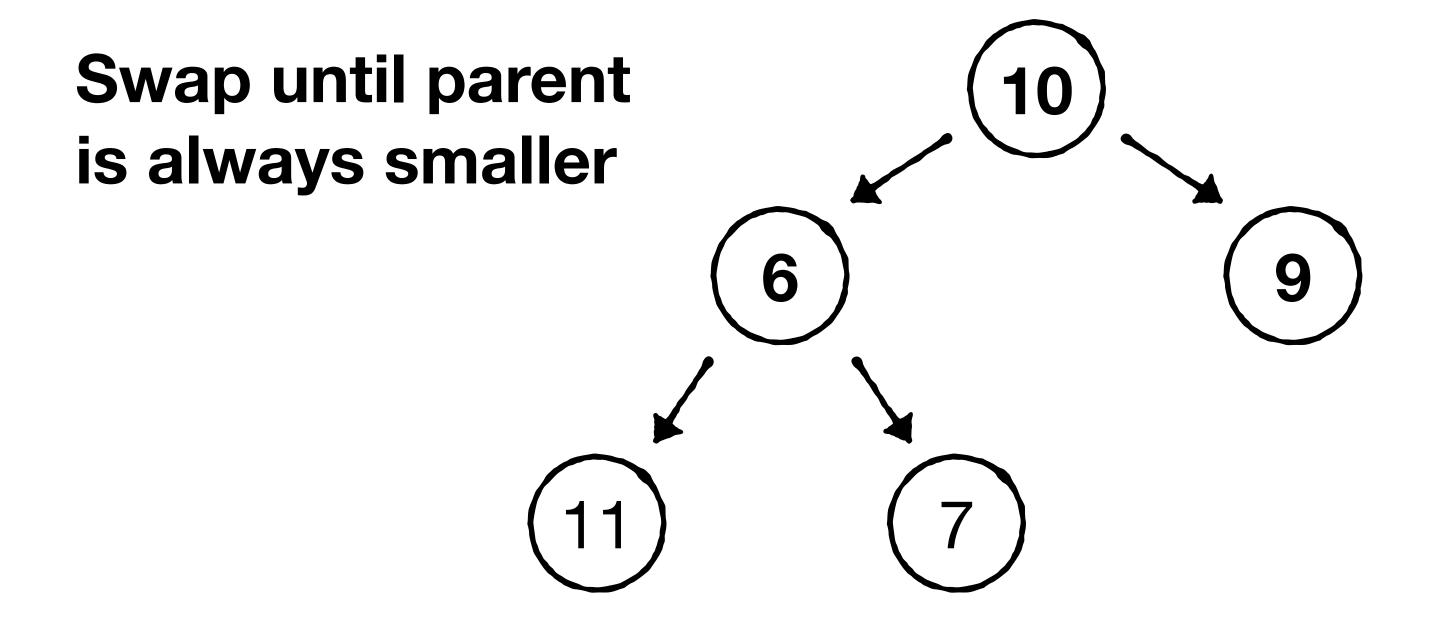
O(log₂ N) min extraction cost

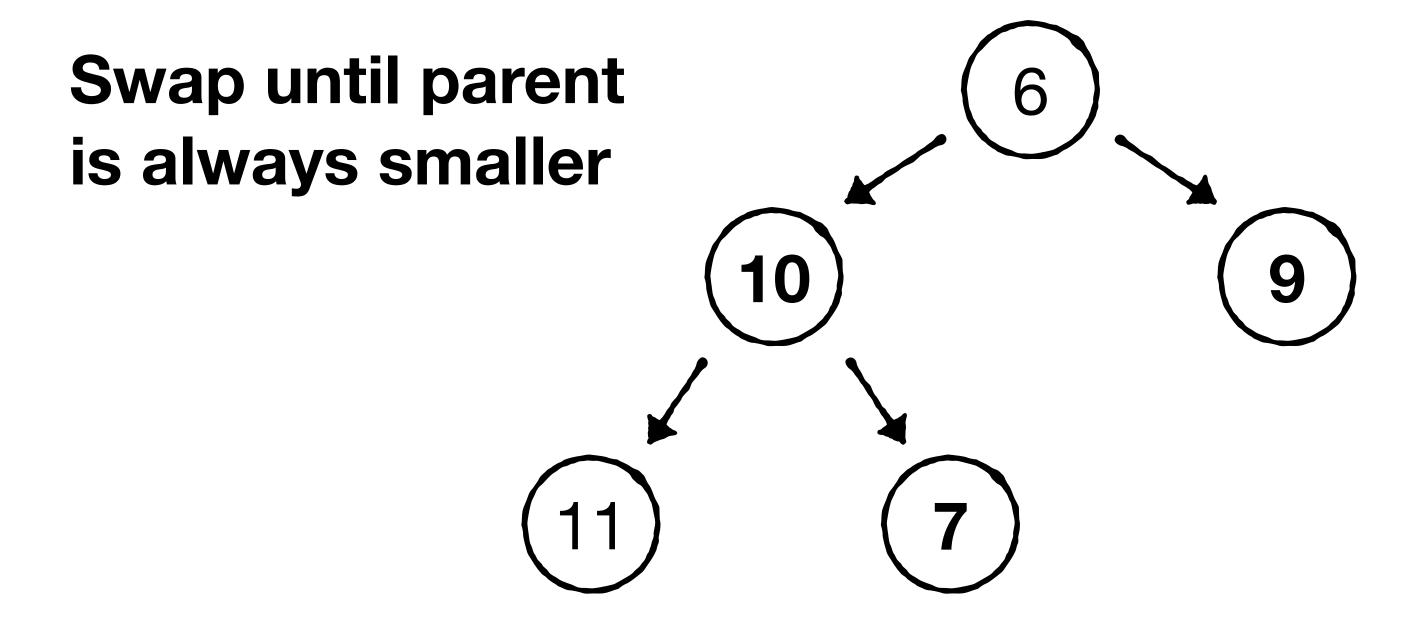


min_key = insert_and_extract(new_key)

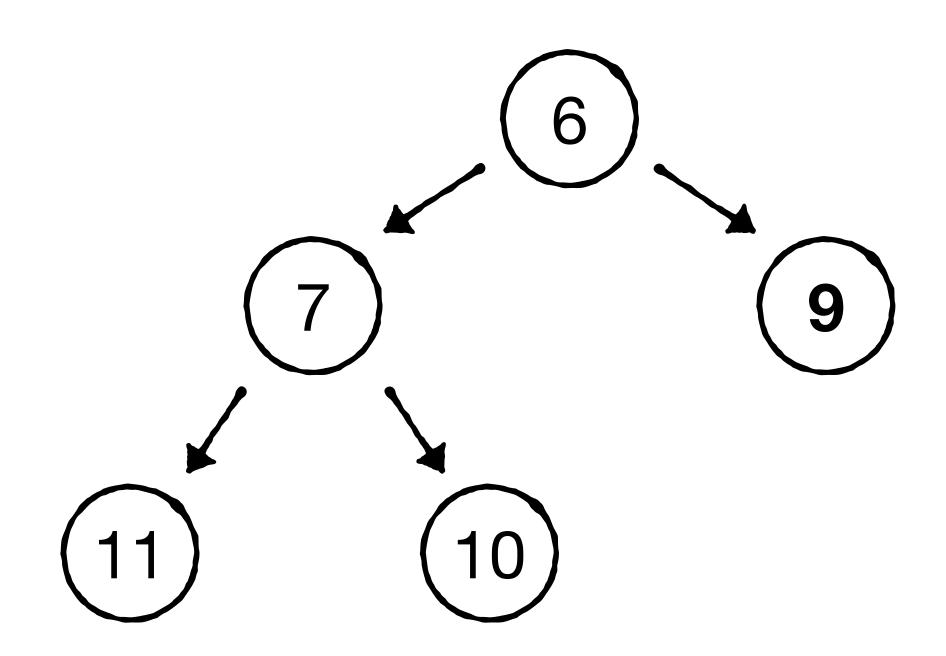






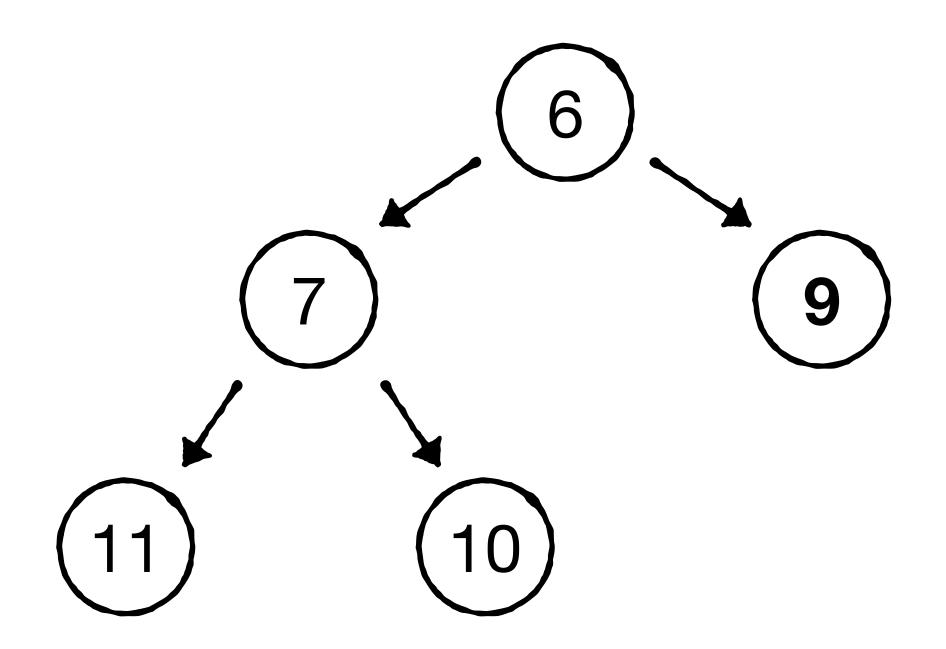


O(log₂ N) for insert_and_extract

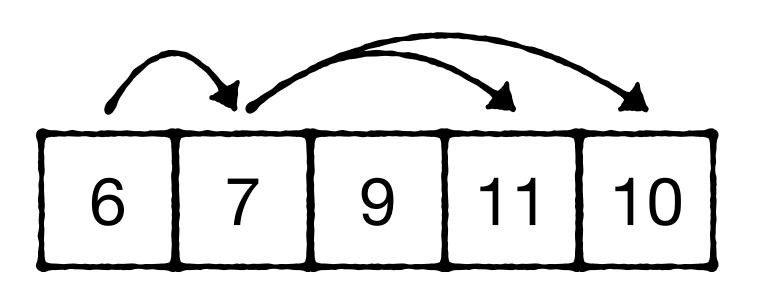


Binary Min-Heap Implementation

Tree with Pointers



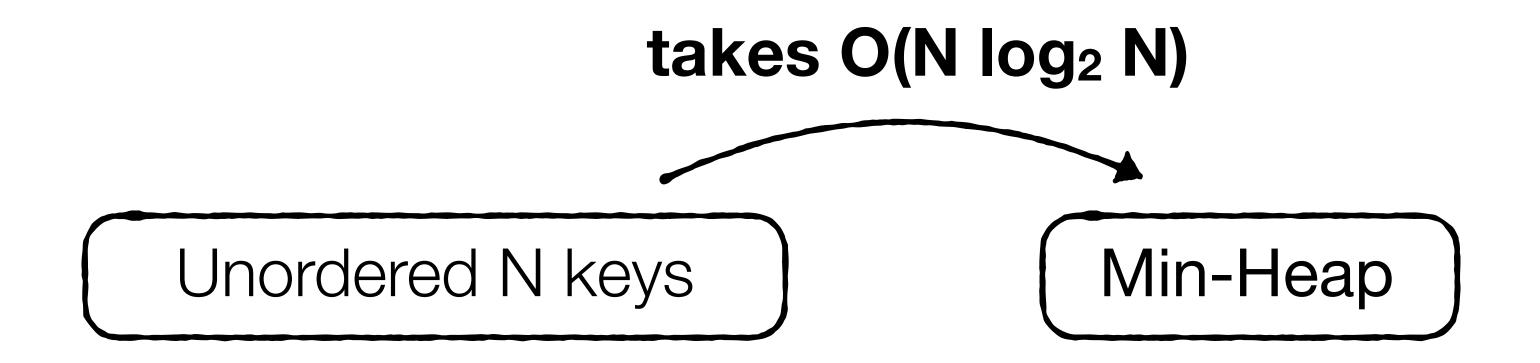
Array



Compact since the binary tree is complete. Avoids overhead of pointers.

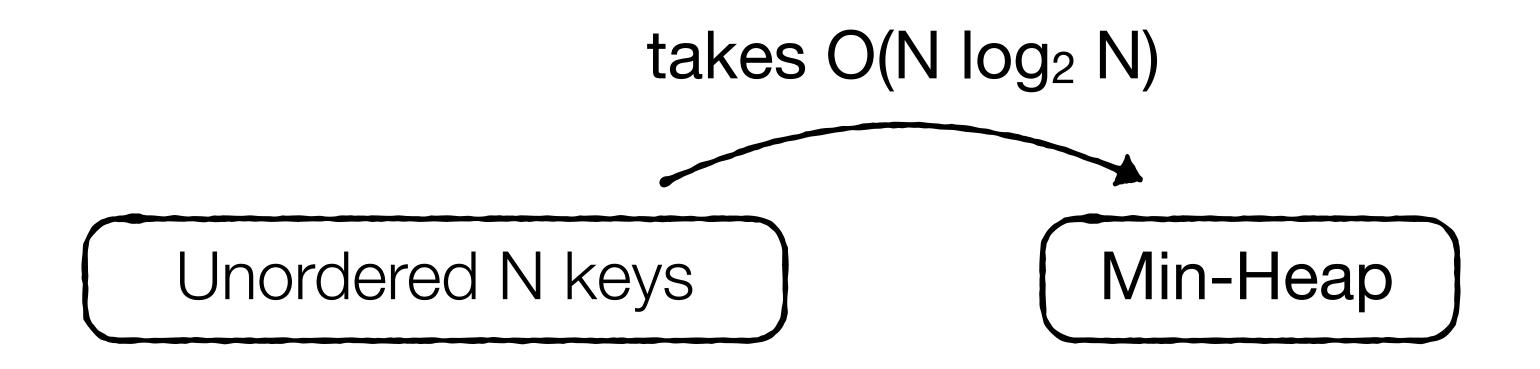
Binary Min-Heap Construction

We can construct a heap for N entries using normal insertions



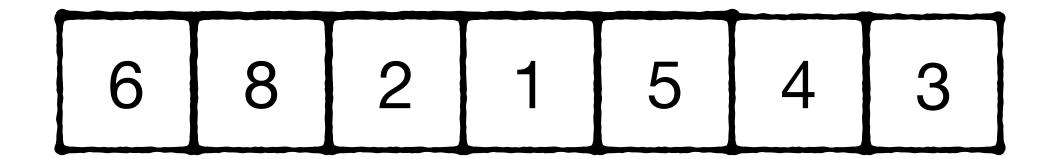
Binary Min-Heap Construction

We can construct a heap for N entries using normal insertions



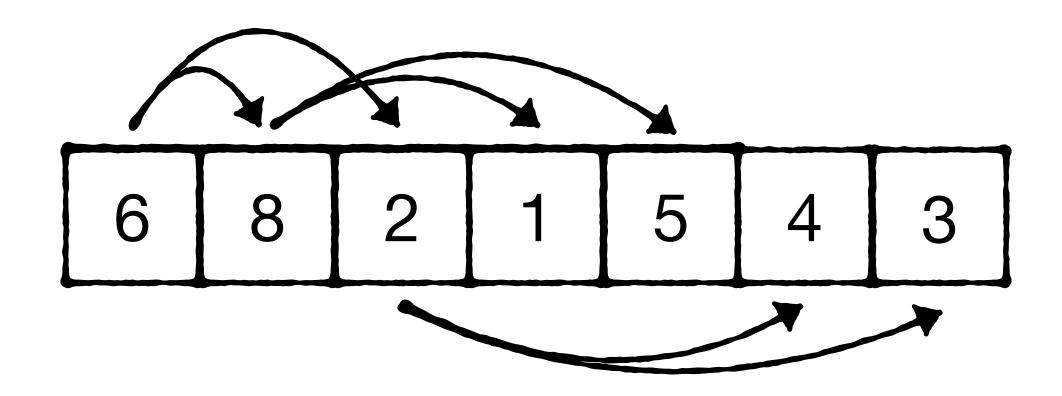
We can do better:)

Binary Min-Heap Efficient Array Construction

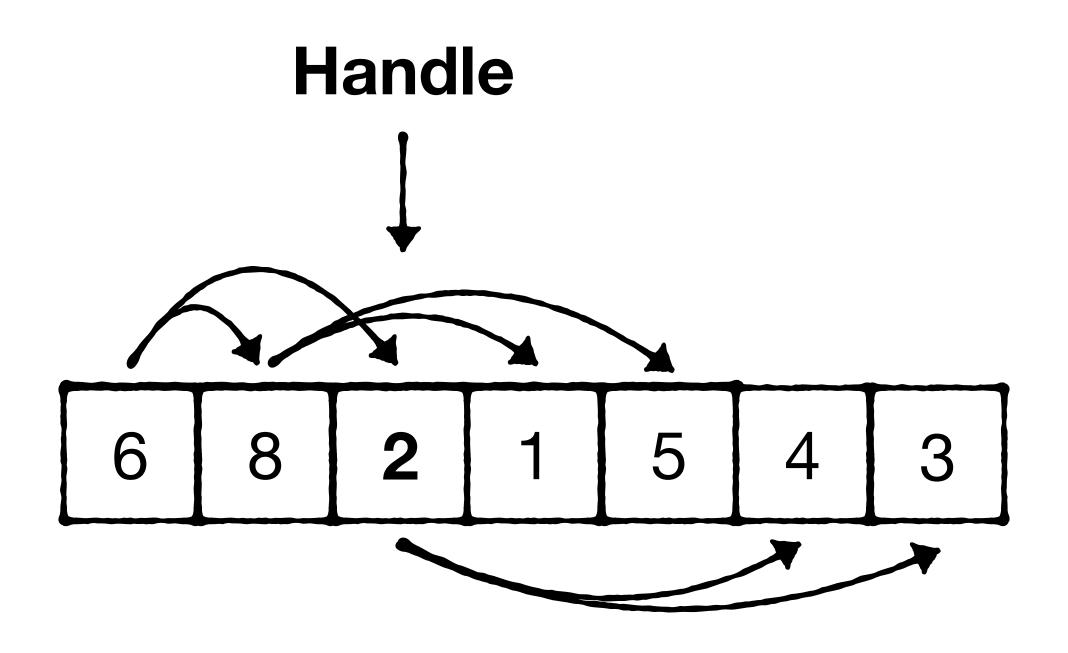


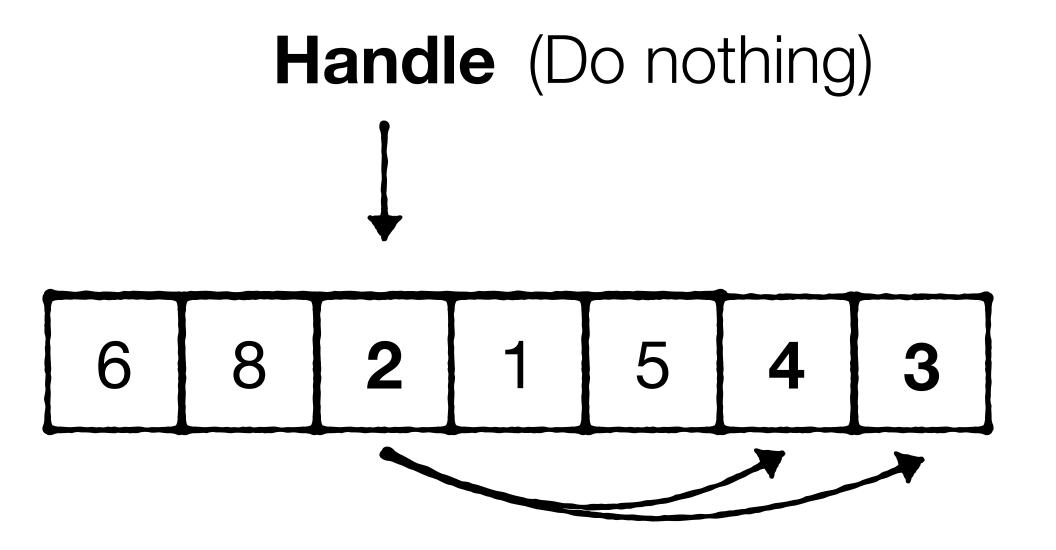
Start with Unordered N keys

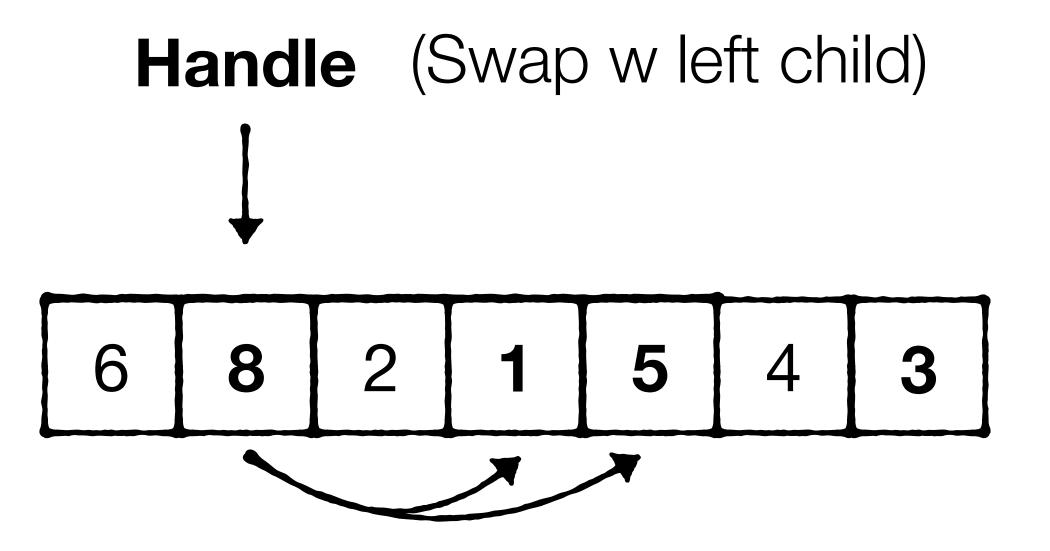
Binary Min-Heap Efficient Array Construction

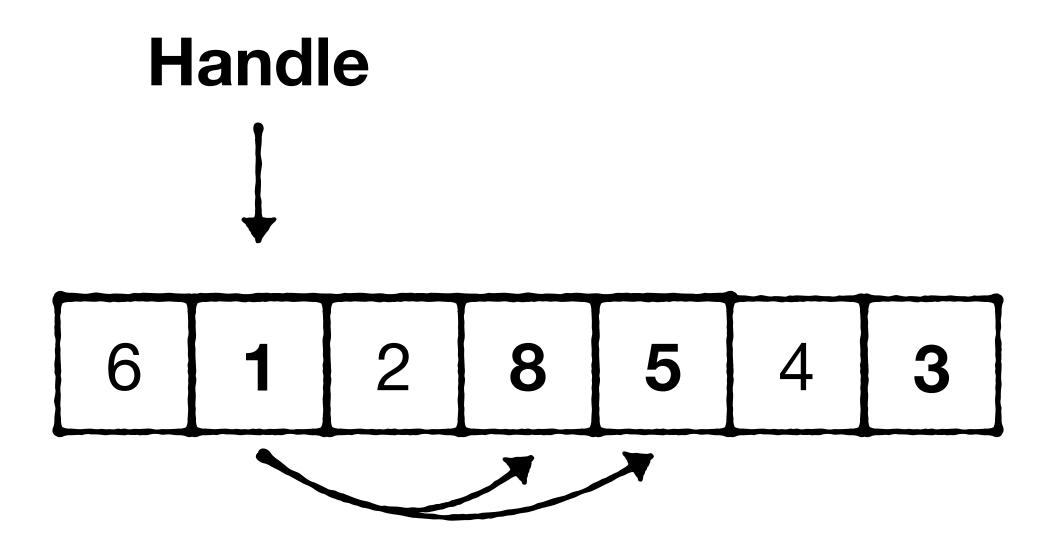


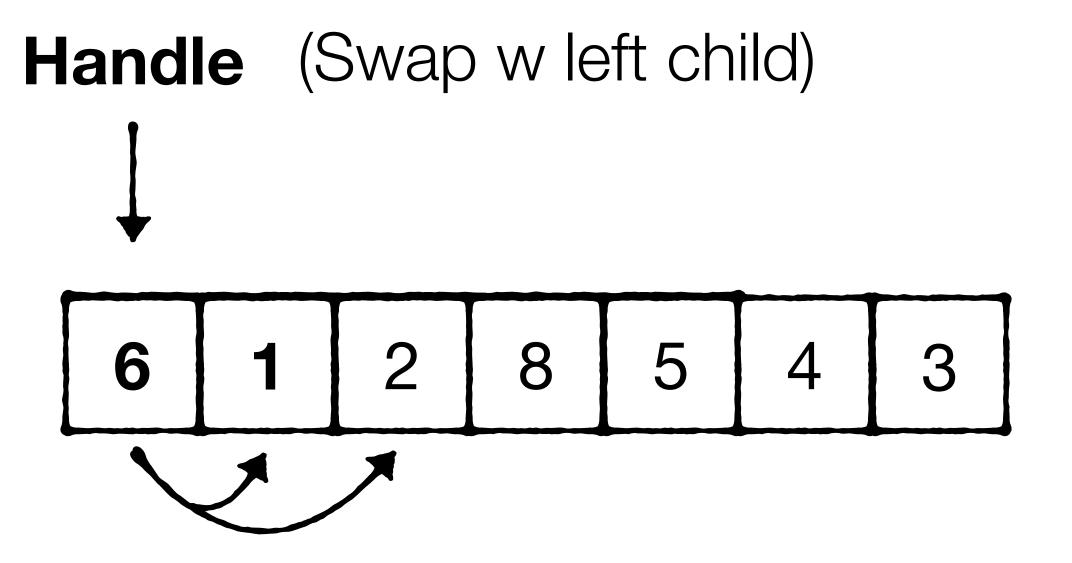
We can infer which slots should be the parent of which other slots.

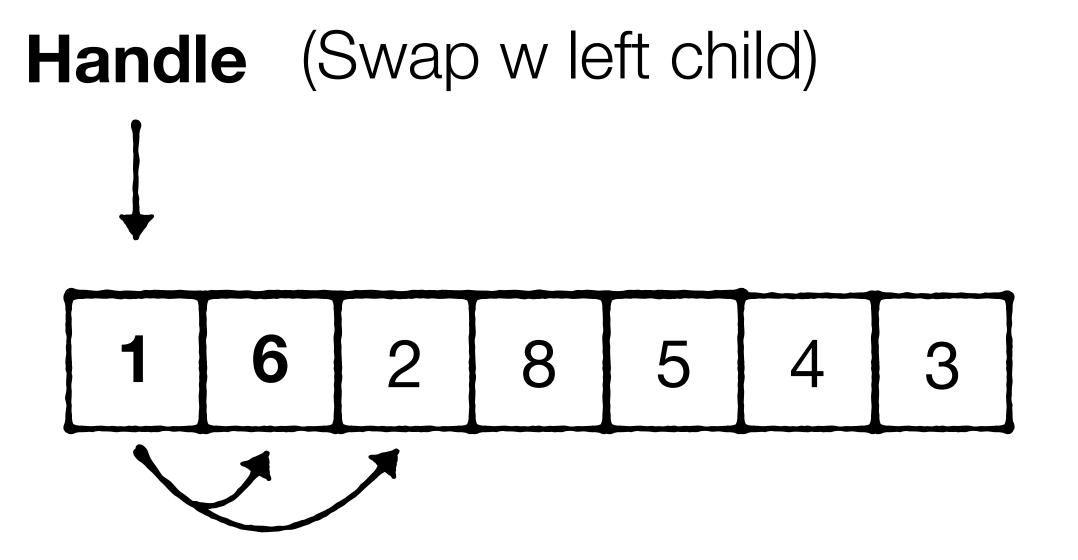


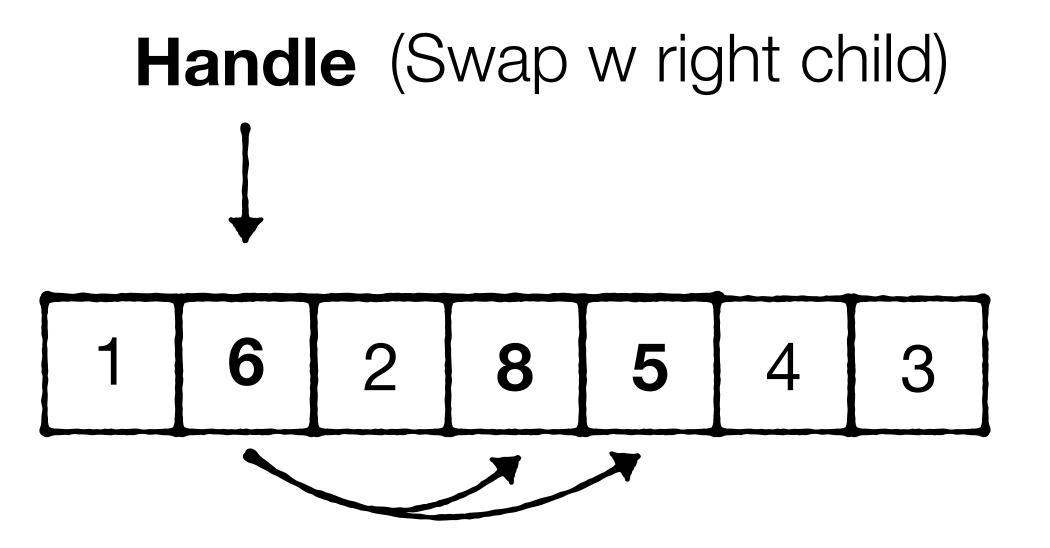


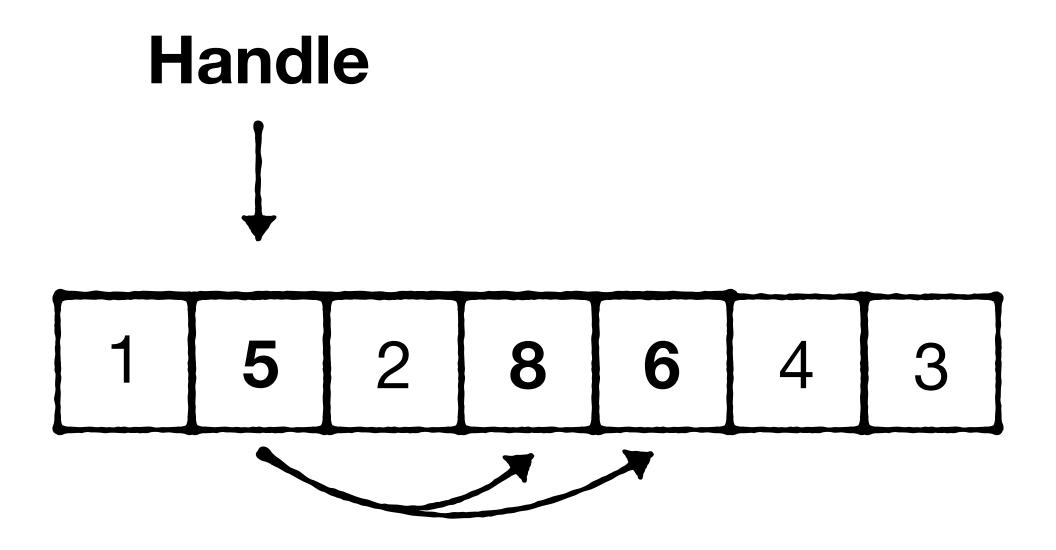






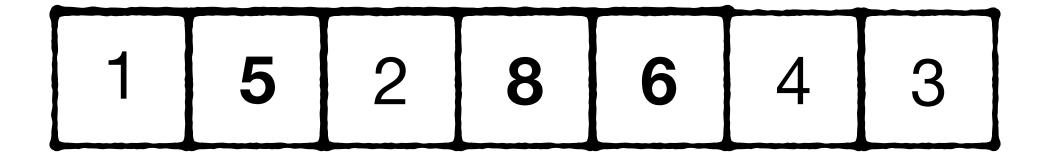






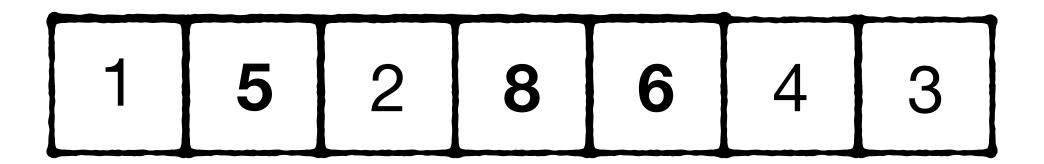
 1
 5
 2
 8
 6
 4
 3

Max swaps: 2 1 1 0 0 0 0



Max swaps: 2 1 1 0 0 0 0

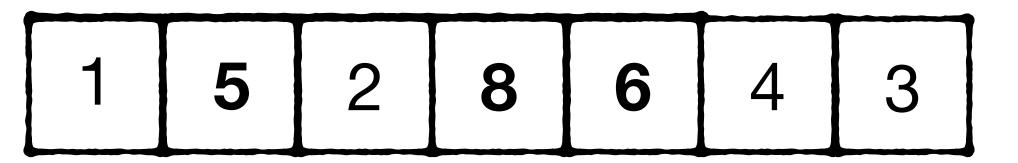
Max swaps for a sub-tree is equal to its depth



Max swaps:

$$2 + 1 + 1 + 0 + 0 + 0 + 0 = O(N)$$

Max swaps for a sub-tree is equal to its depth

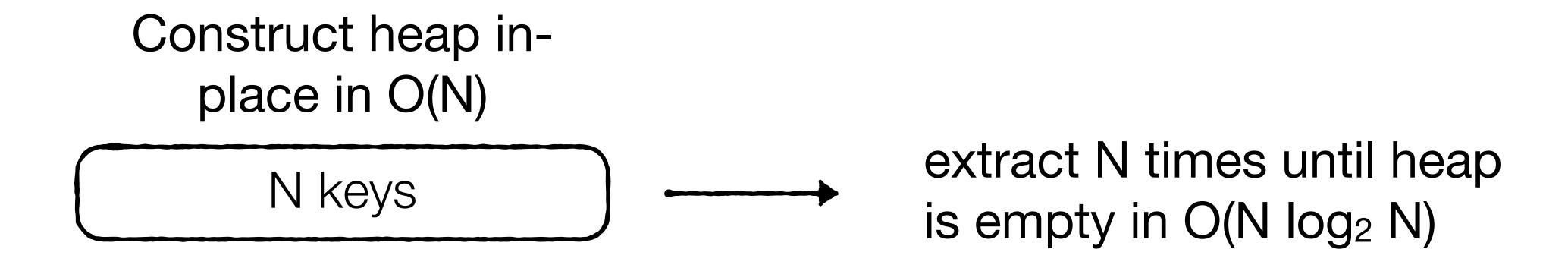


 $O(N) < O(N \log_2 N)$

from before with pure insertions

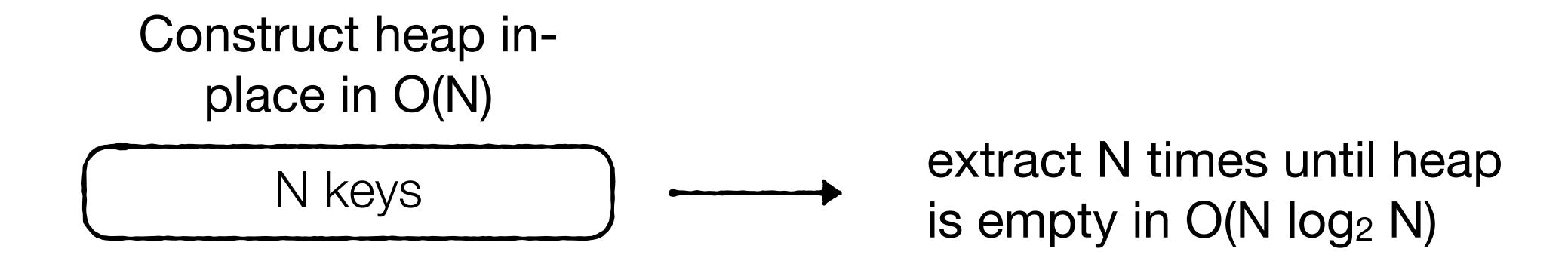
Side-Note on Heap-Sort

If data fits in memory, we can sort it using a heap

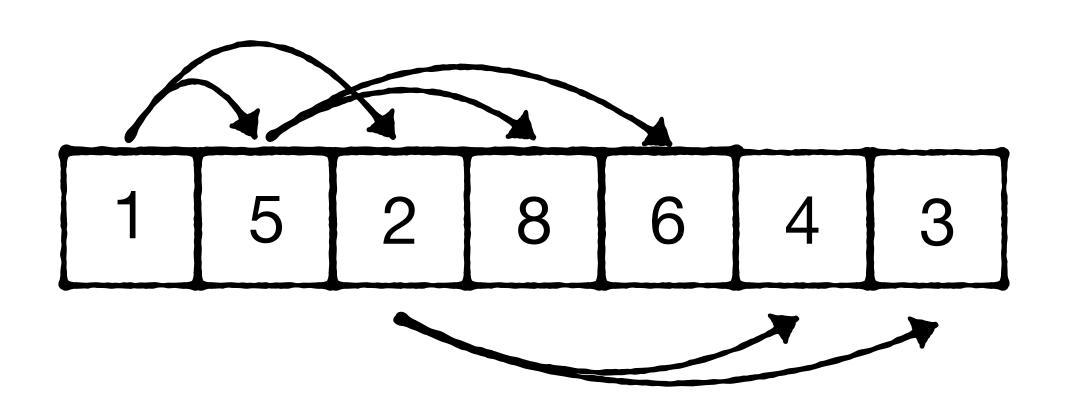


Side-Note on Heap-Sort

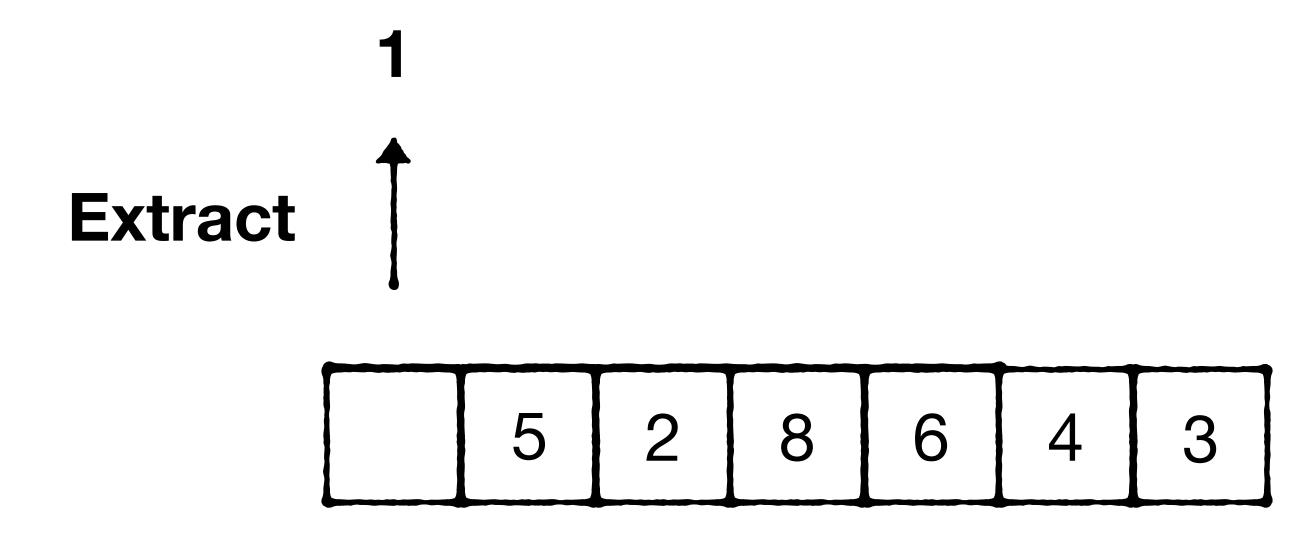
If data fits in memory, we can sort it using a heap

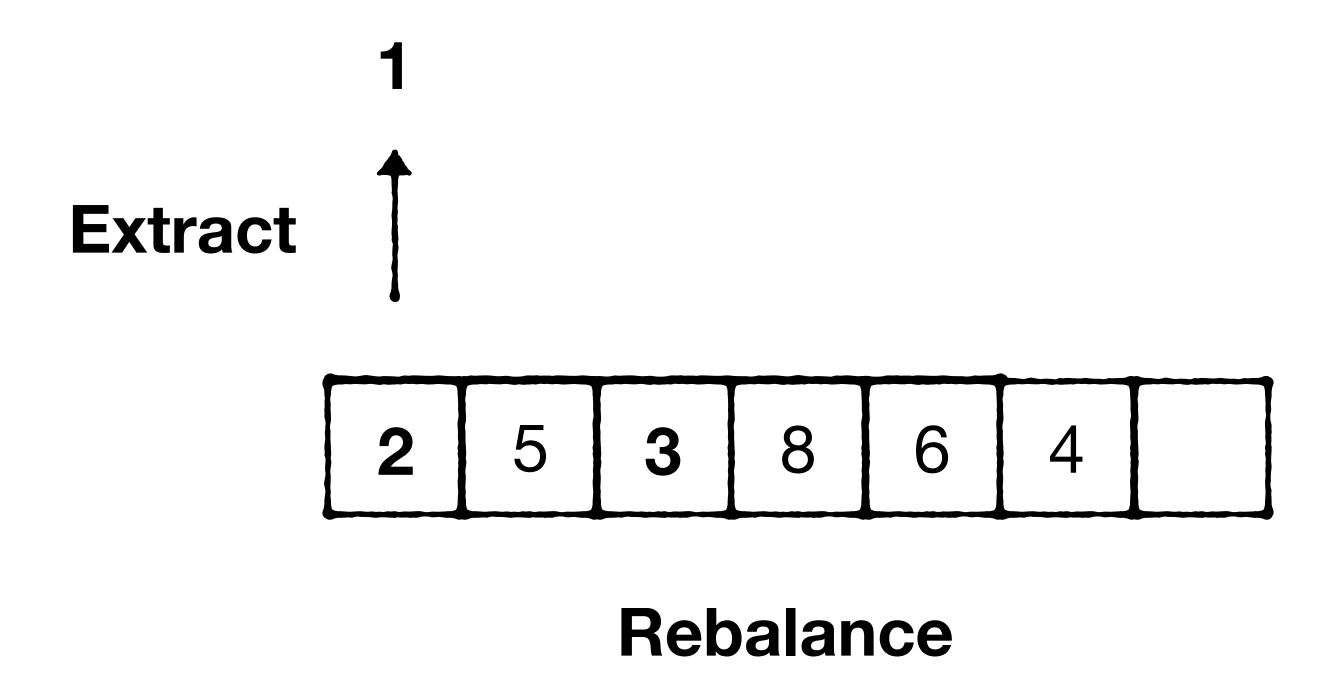


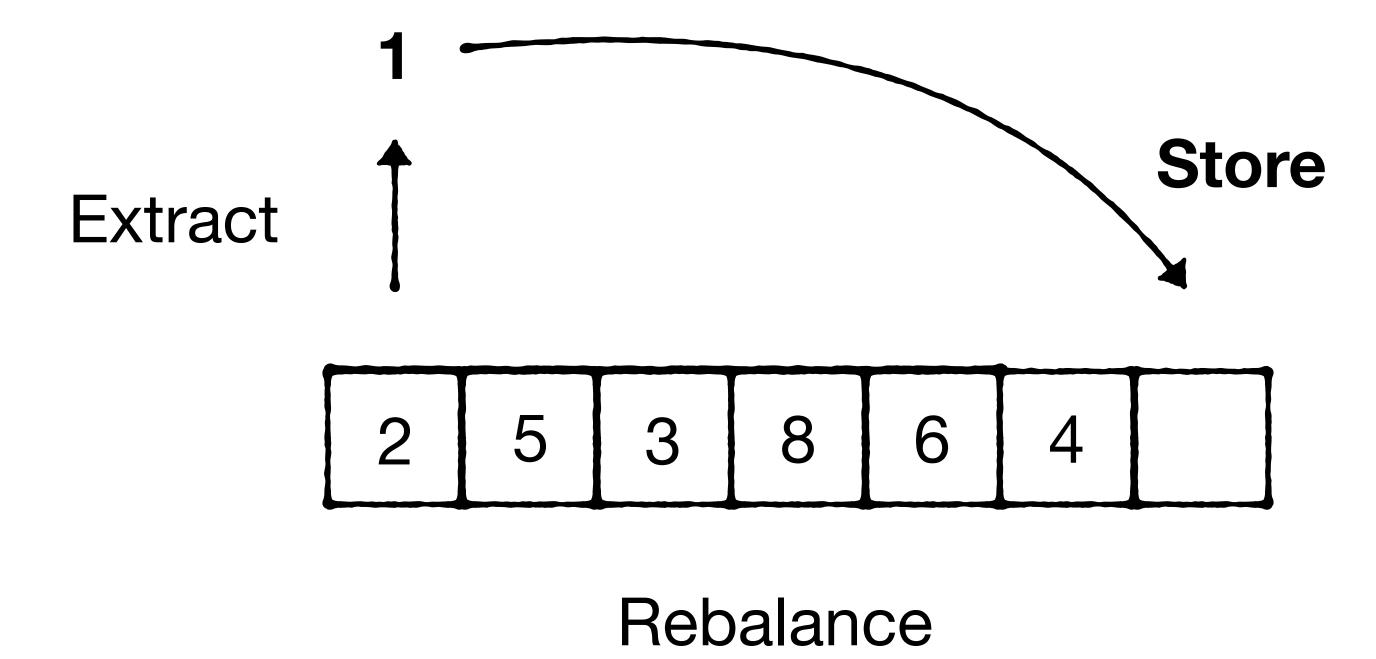
Can you do this in-place?:)

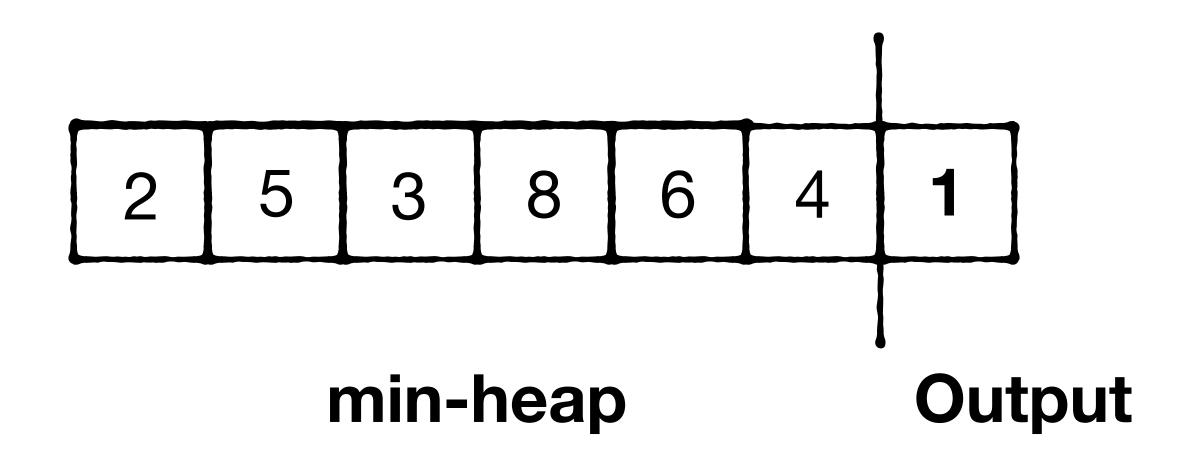


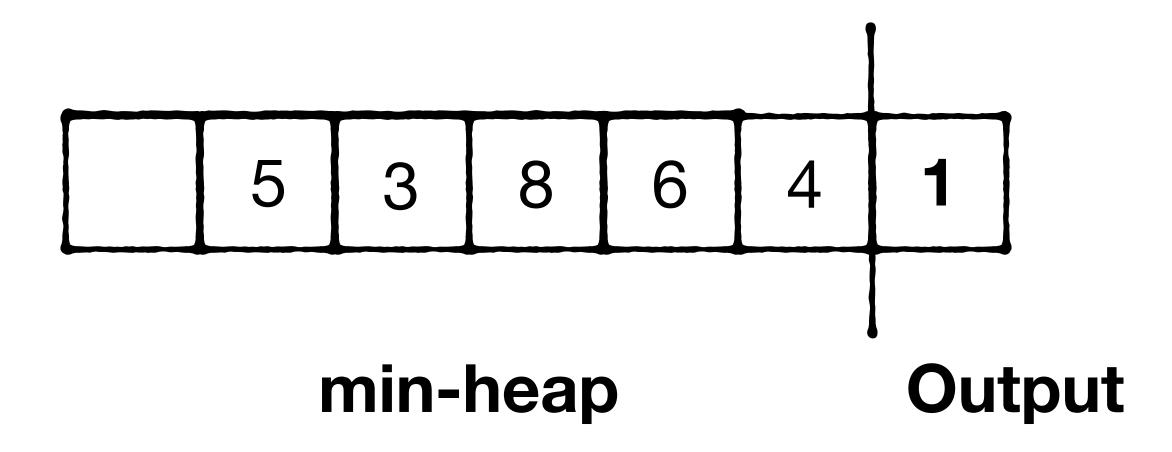
Recall min-heap can be stored as compact array

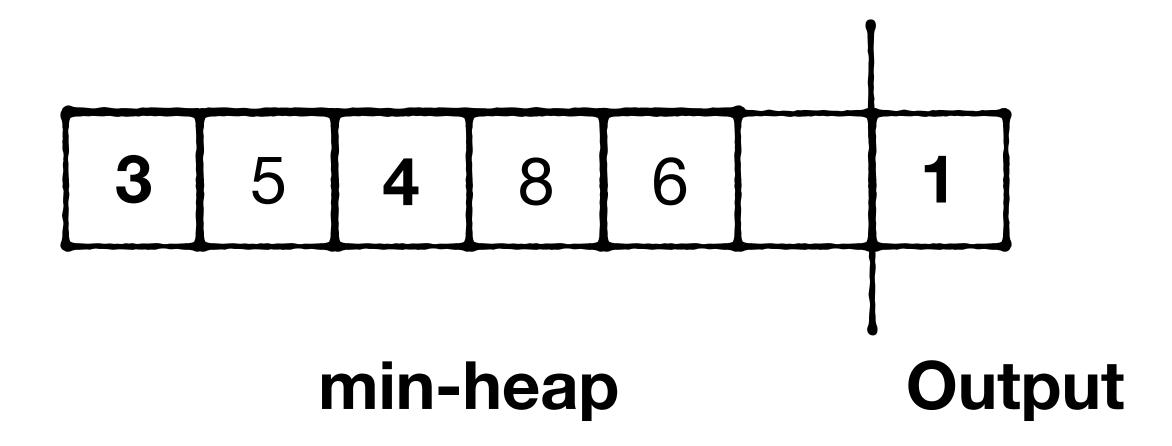


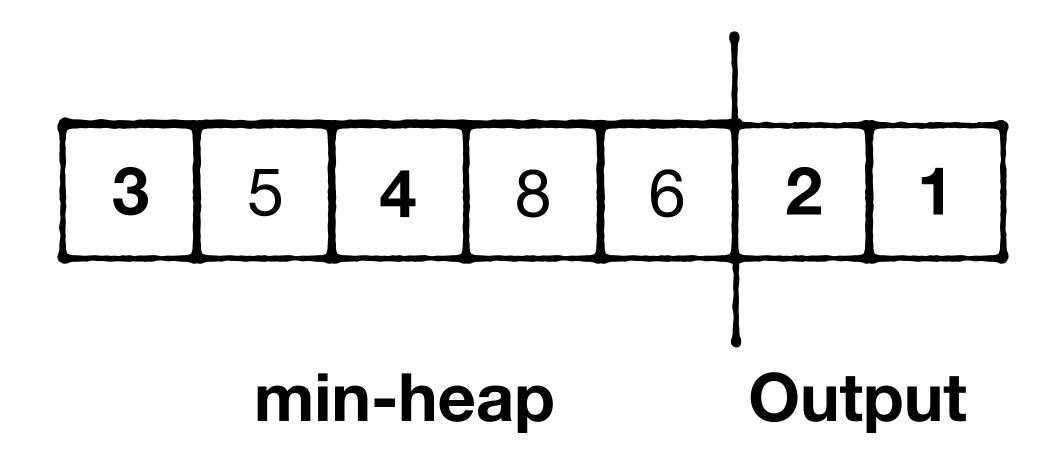




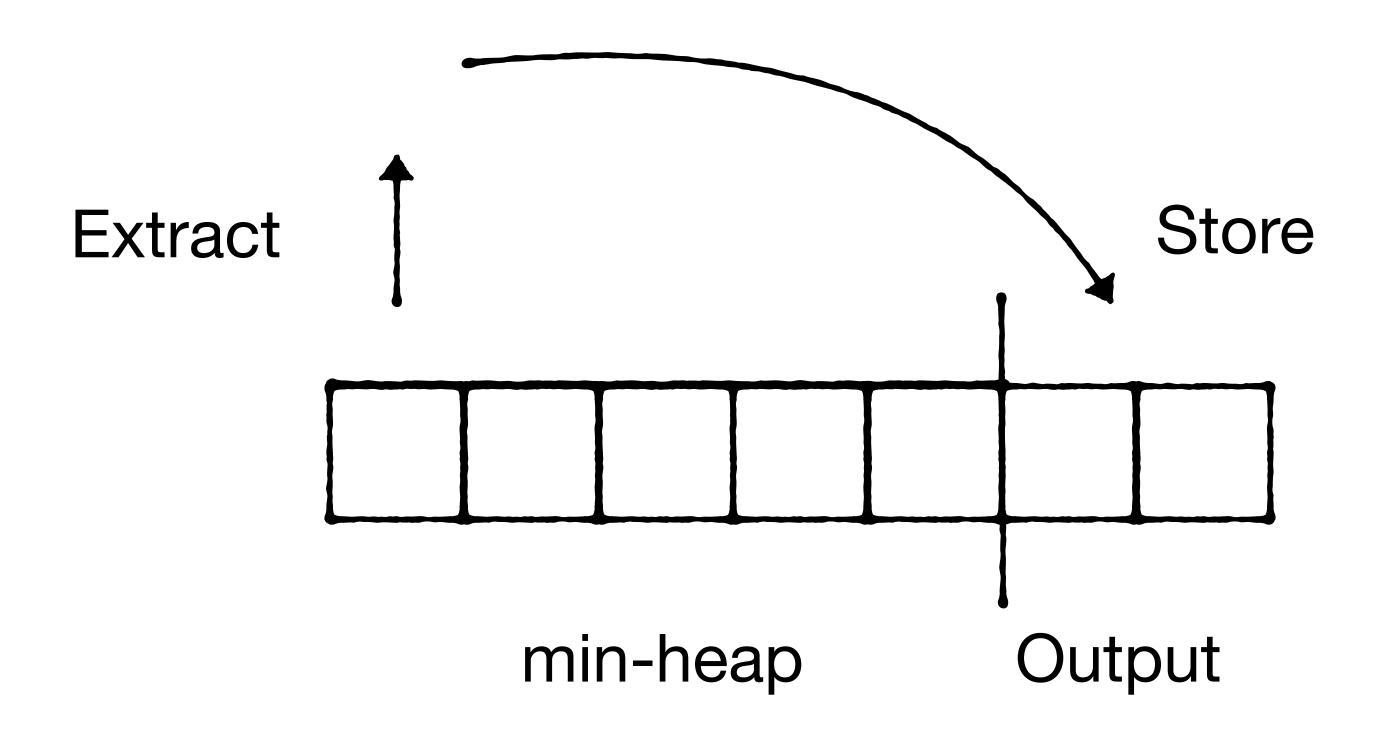




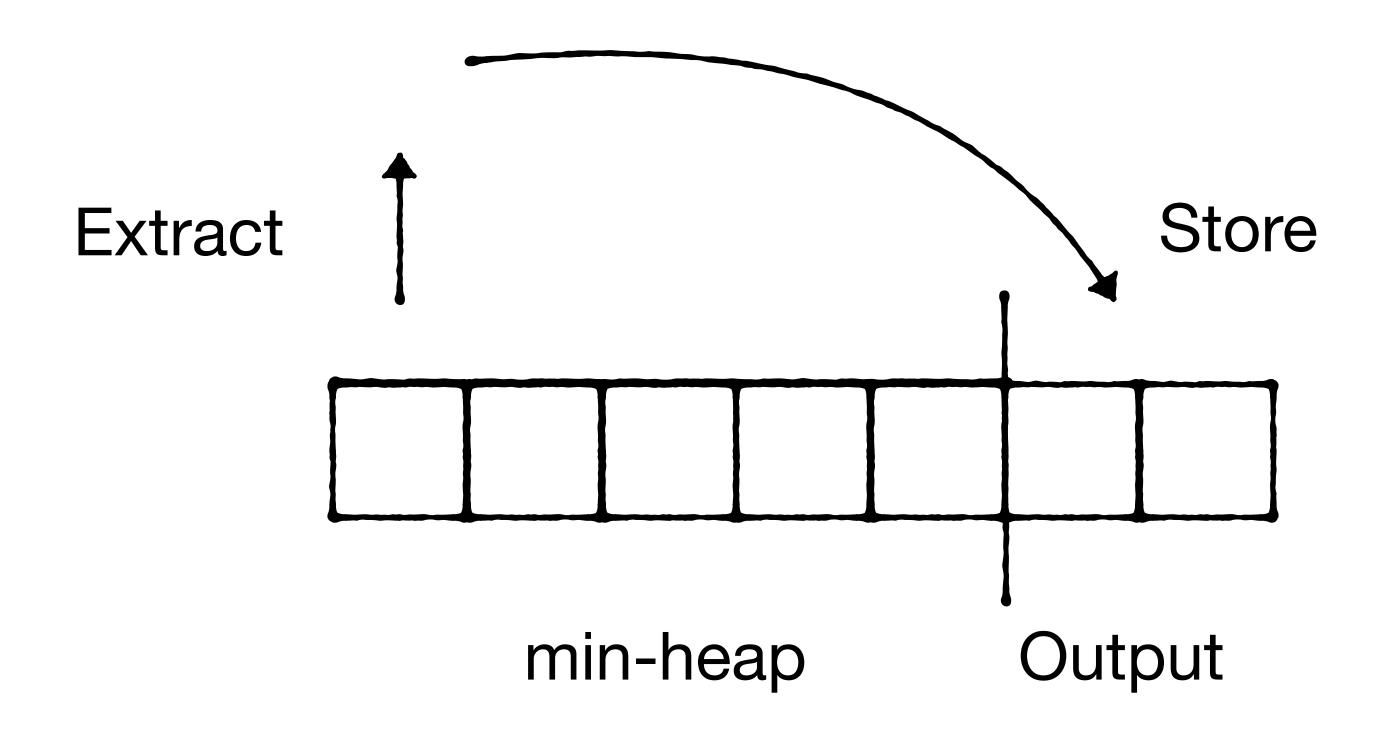




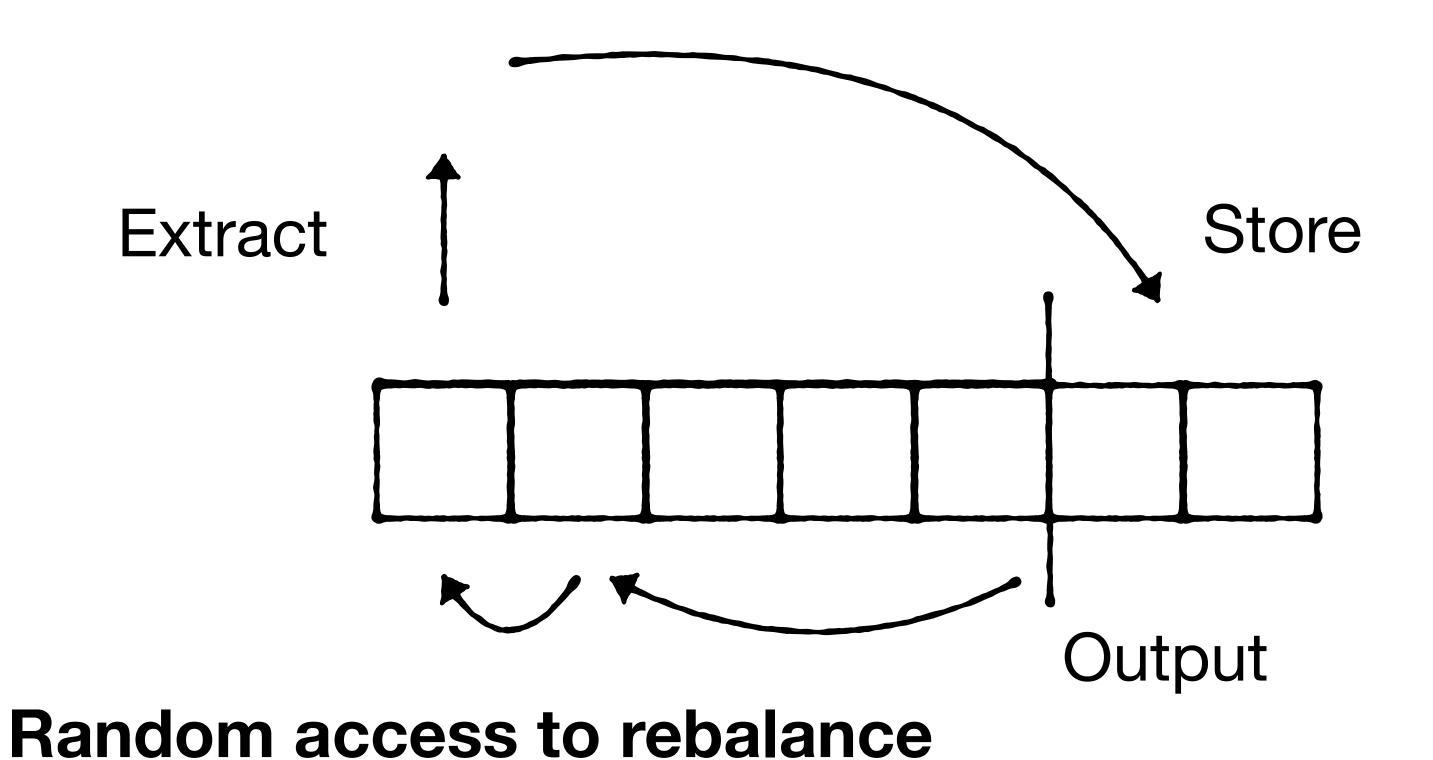
in-place algorithm:)



Downsides compared to quick-sort or merge-sort?



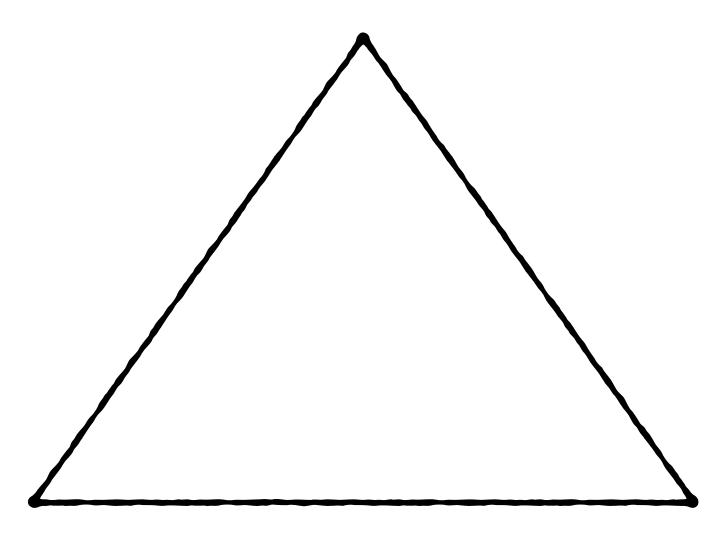
Downsides compared to quick-sort or merge-sort?



Heap-Sort

Worst case O(N log N)
In-place

More random access



Merge-Sort

Worst case O(N log N)

2x more space

Sequential access

Quick-Sort

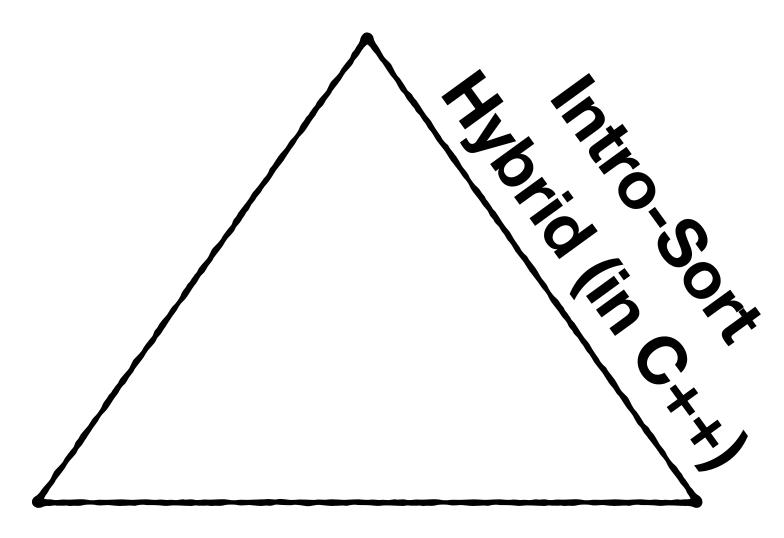
Avg. worst case O(N log N)
In-place

Sequential access

Heap-Sort

Worst case O(N log N)
In-place

More random access



Merge-Sort

Worst case O(N log N)

2x more space

Sequential access

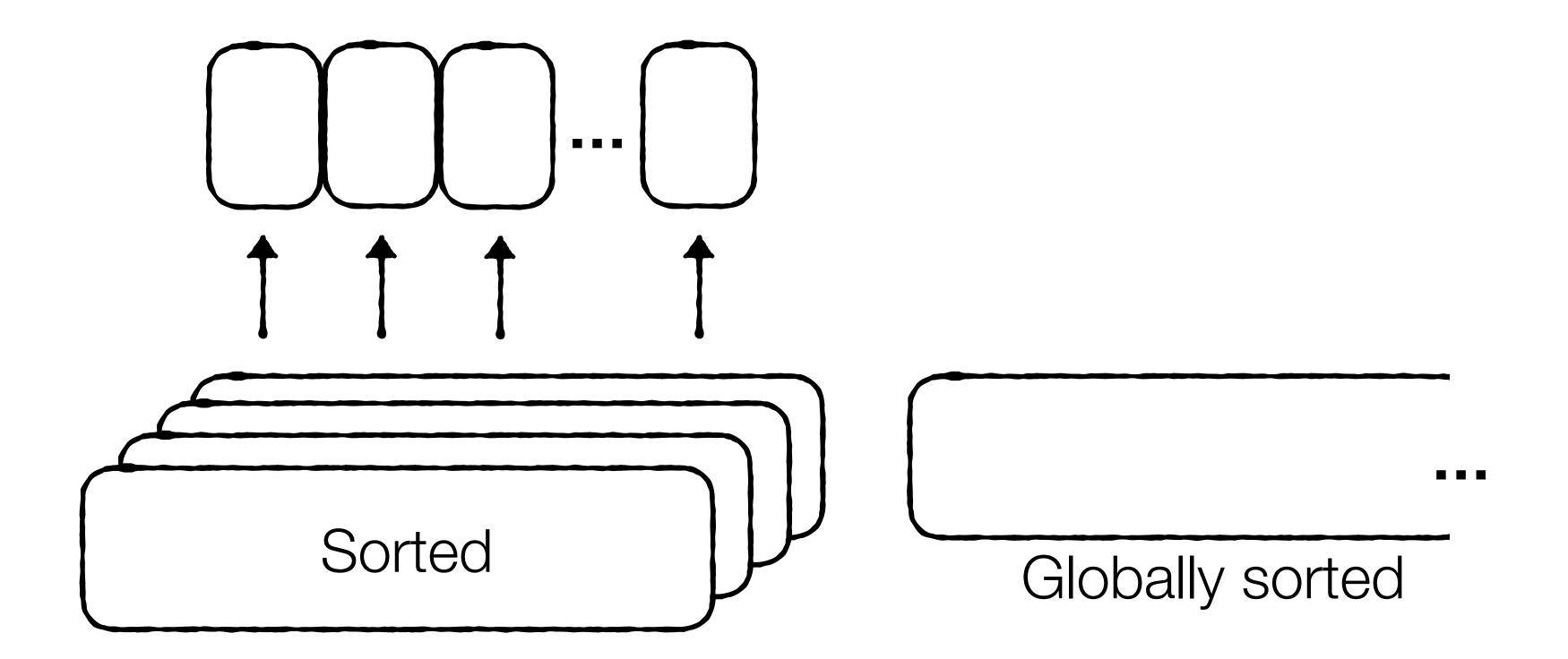
Quick-Sort

Avg. worst case O(N log N)
In-place

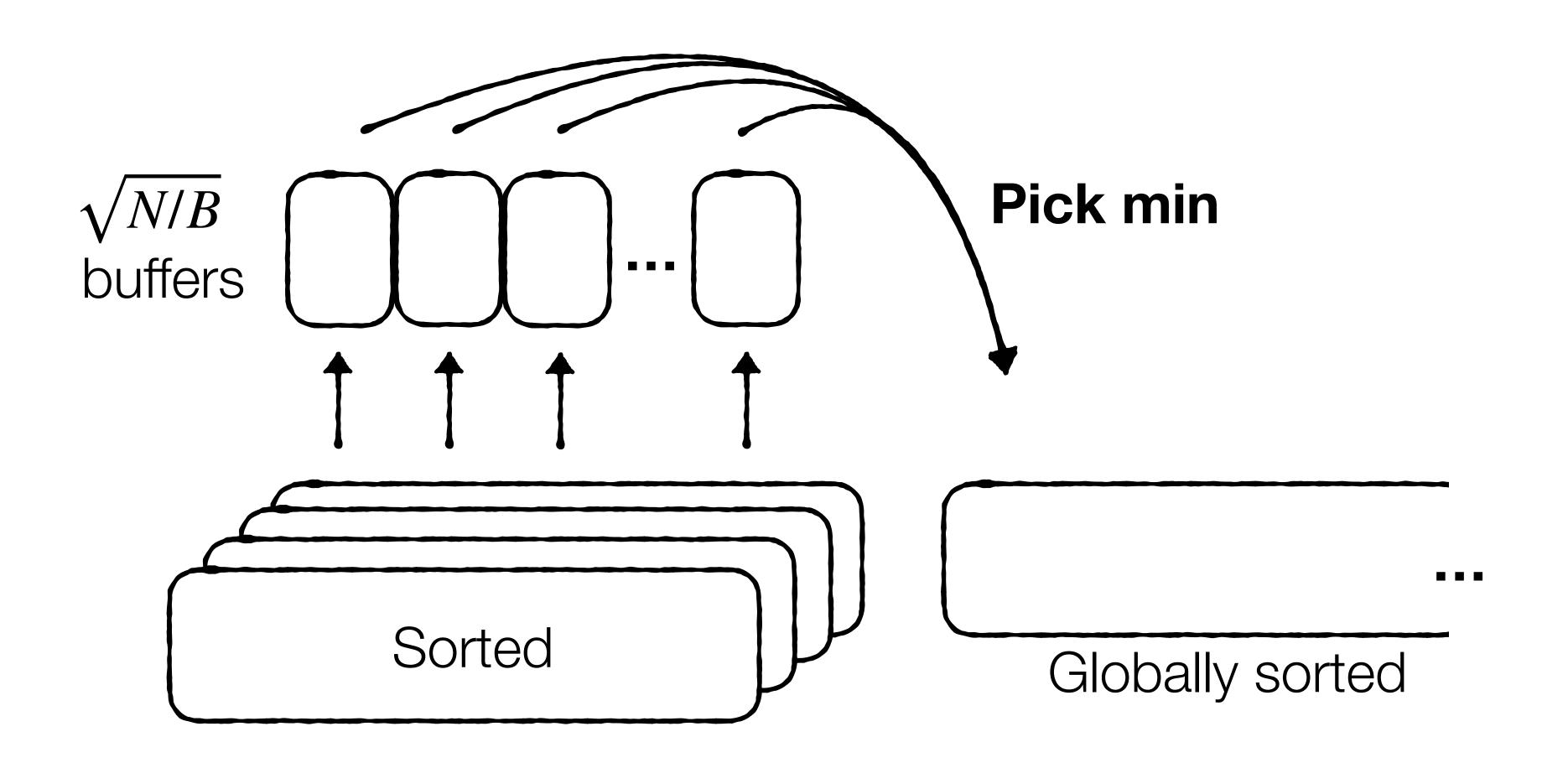
Sequential access



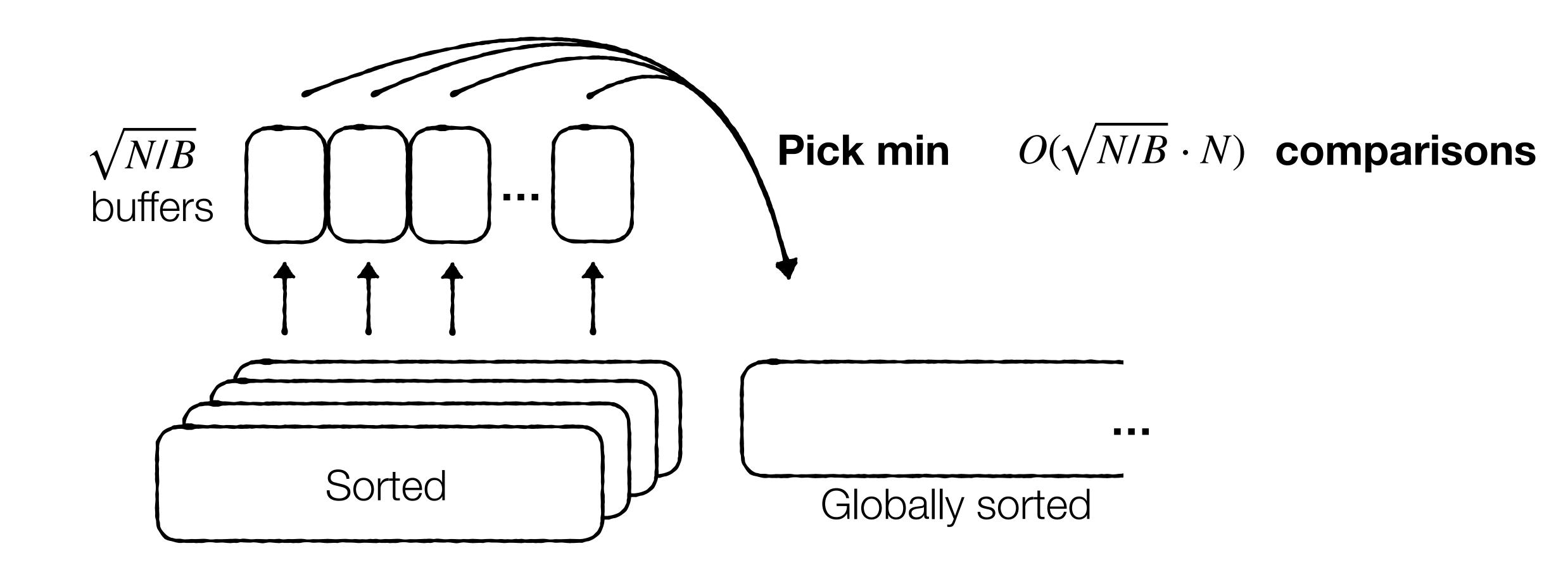
How to merge partitions?



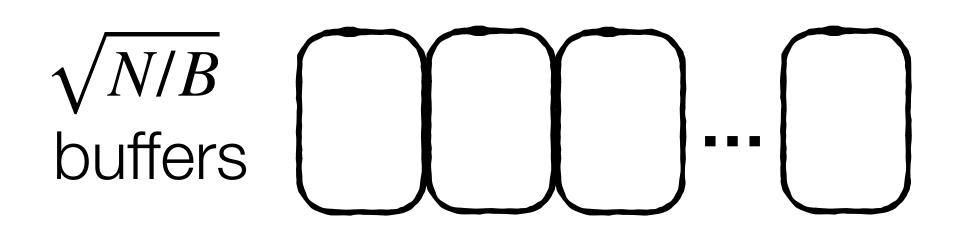
How to merge partitions?

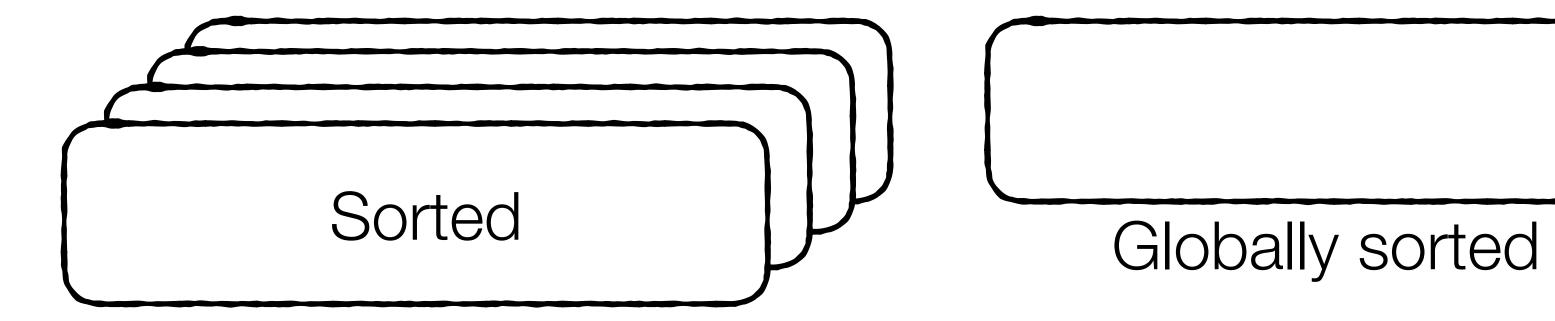


How to merge partitions?



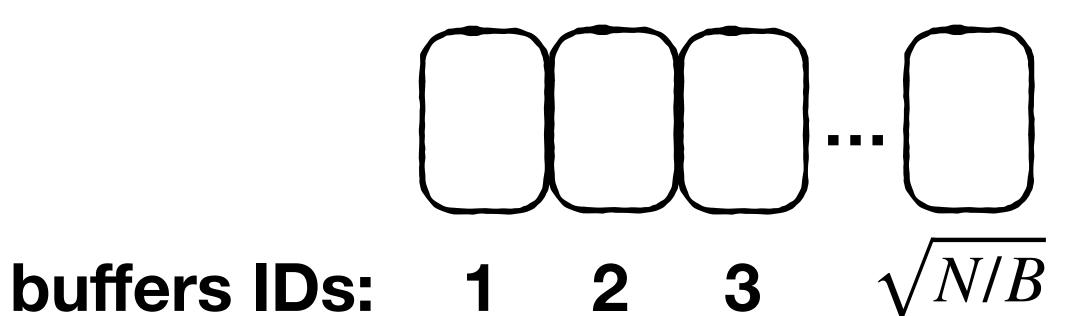
Min-Heap

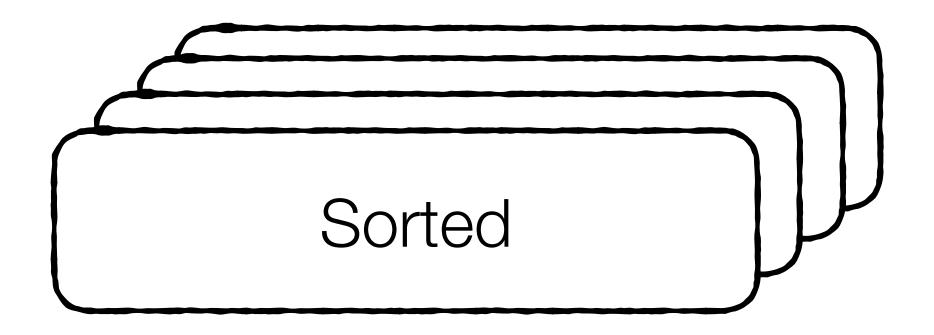




pair<key, buffer ID>

Min-Heap

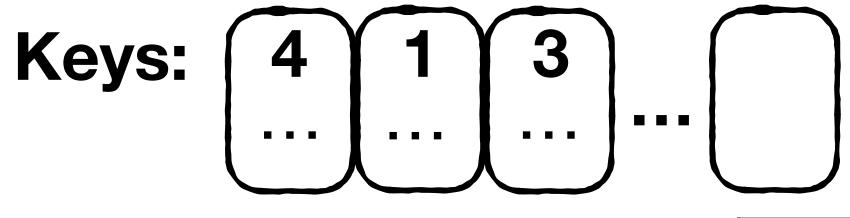




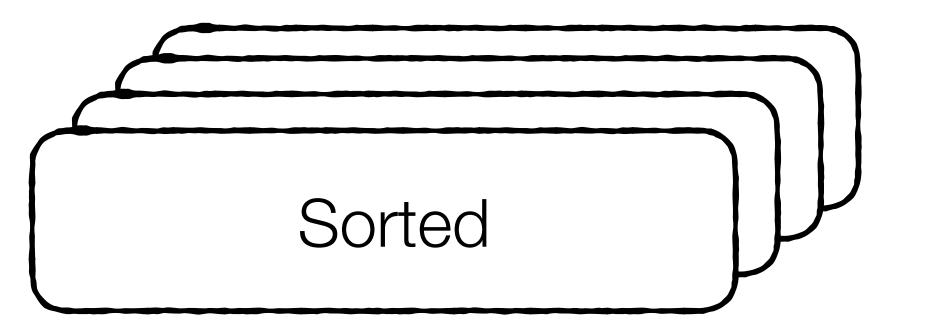


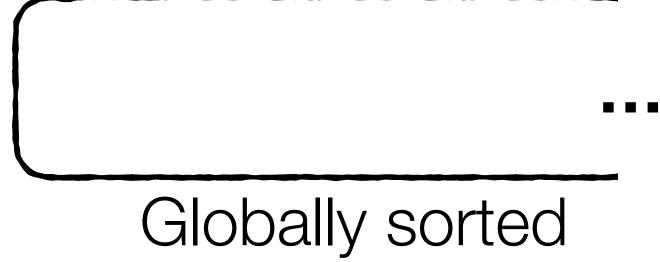
pair<key, buffer ID>

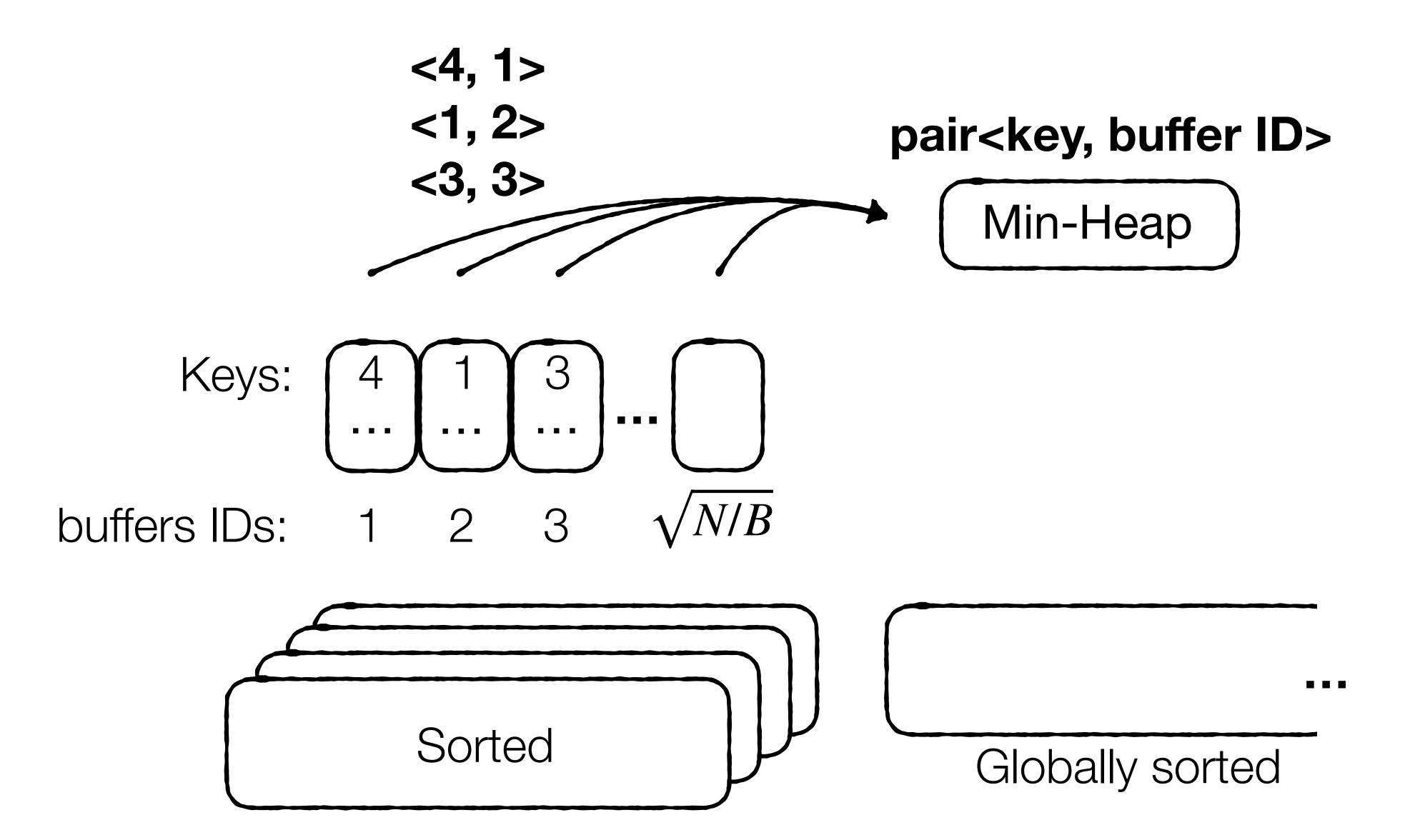
Min-Heap



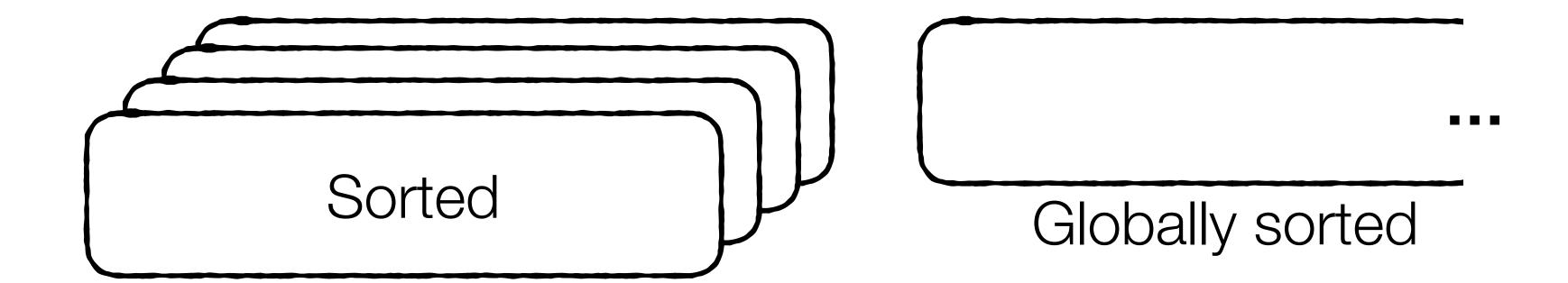
buffers IDs: 1 2 3 $\sqrt{N/B}$

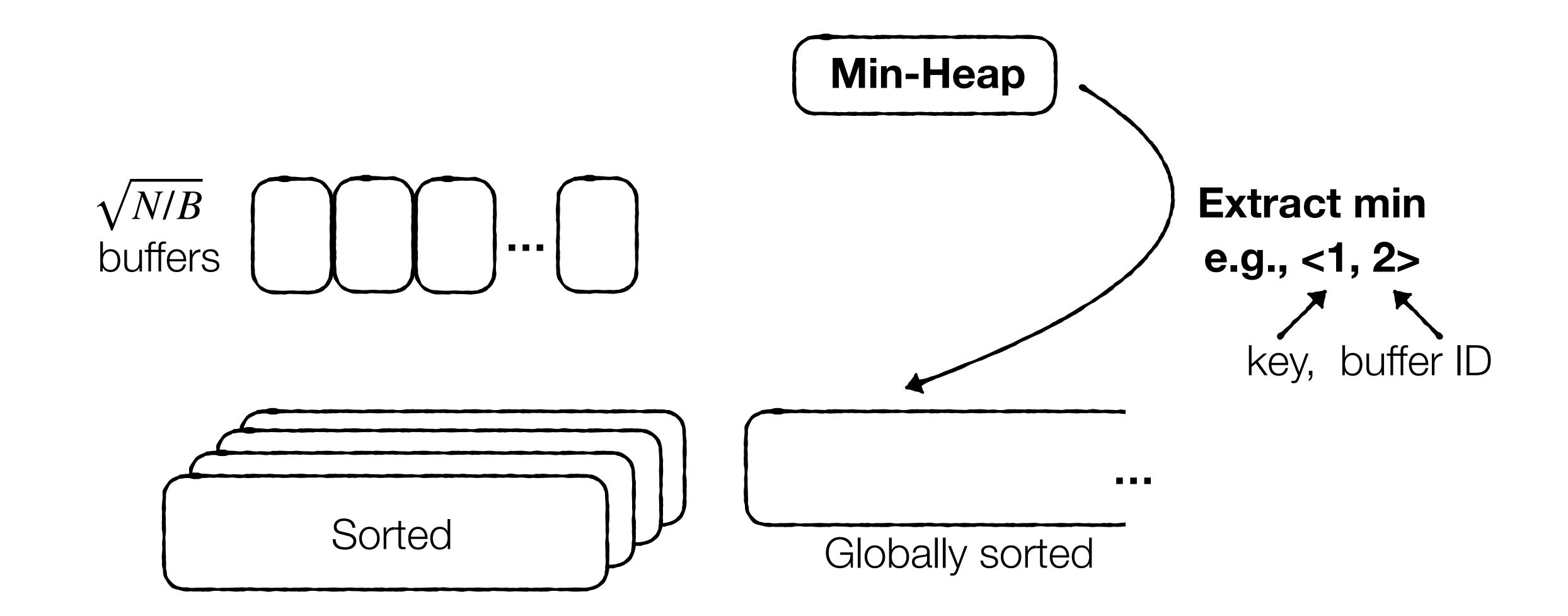


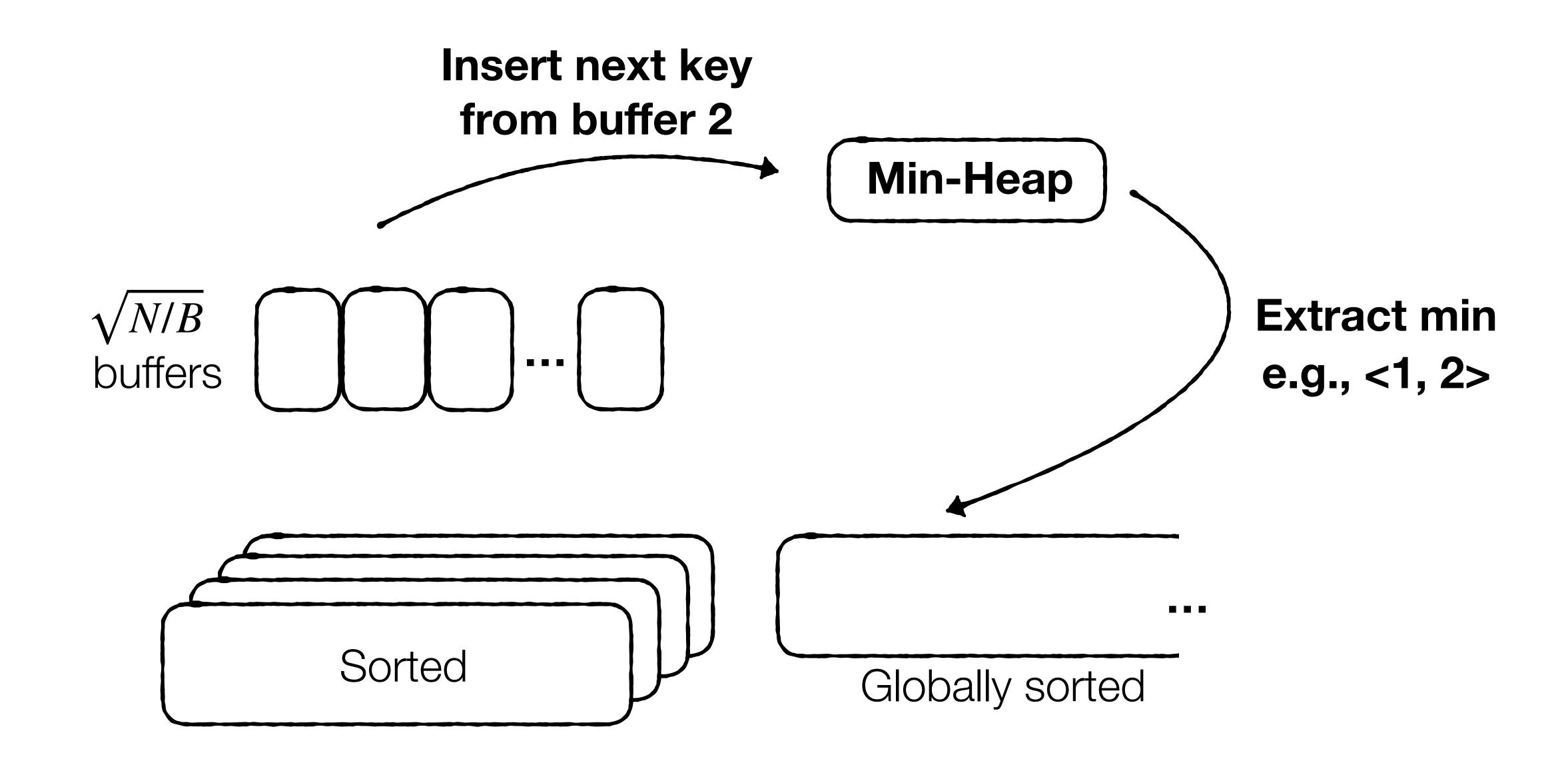


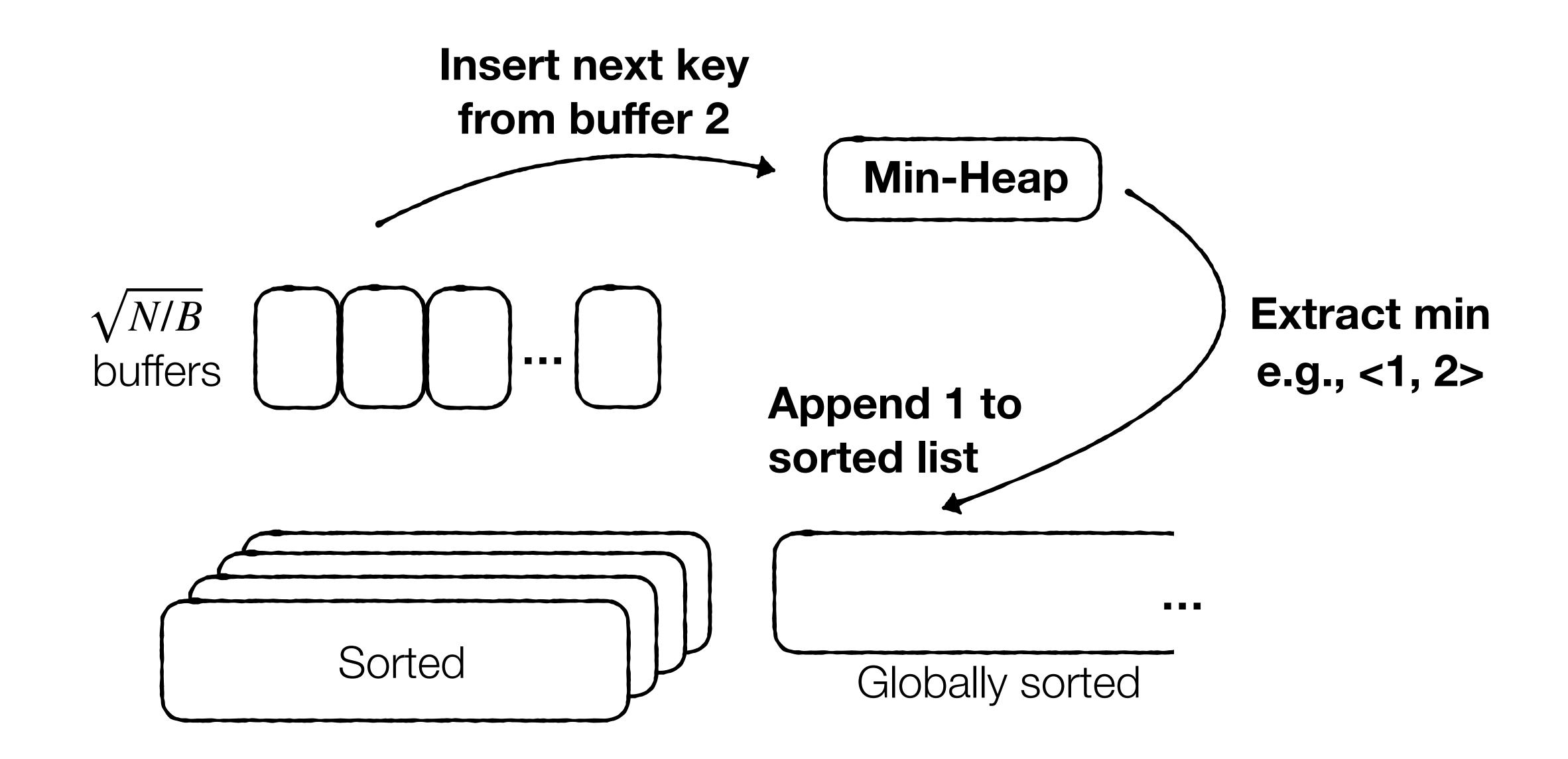


Construct in $O(\sqrt{N/B})$ $\sqrt{N/B}$ buffers ...

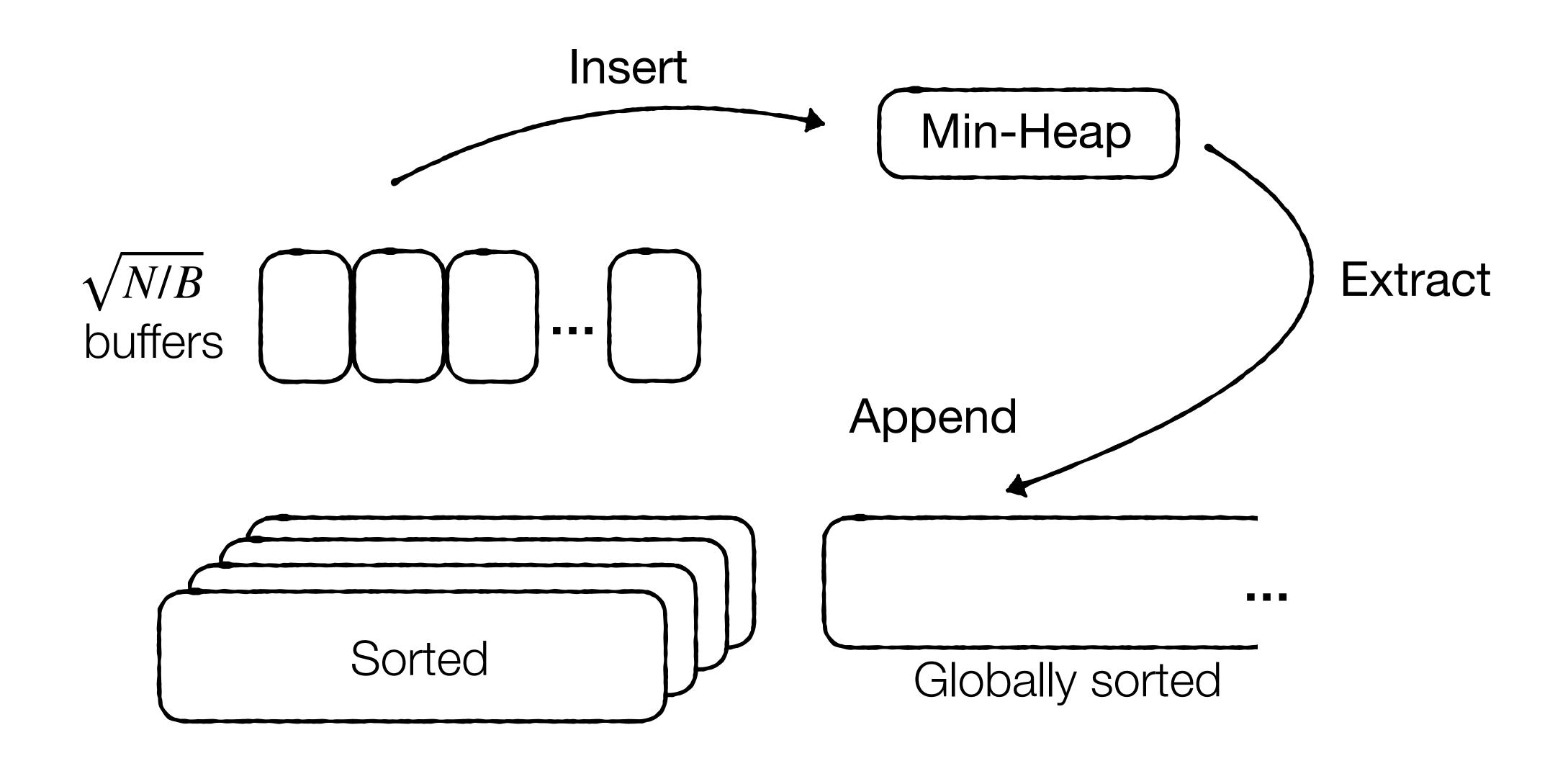




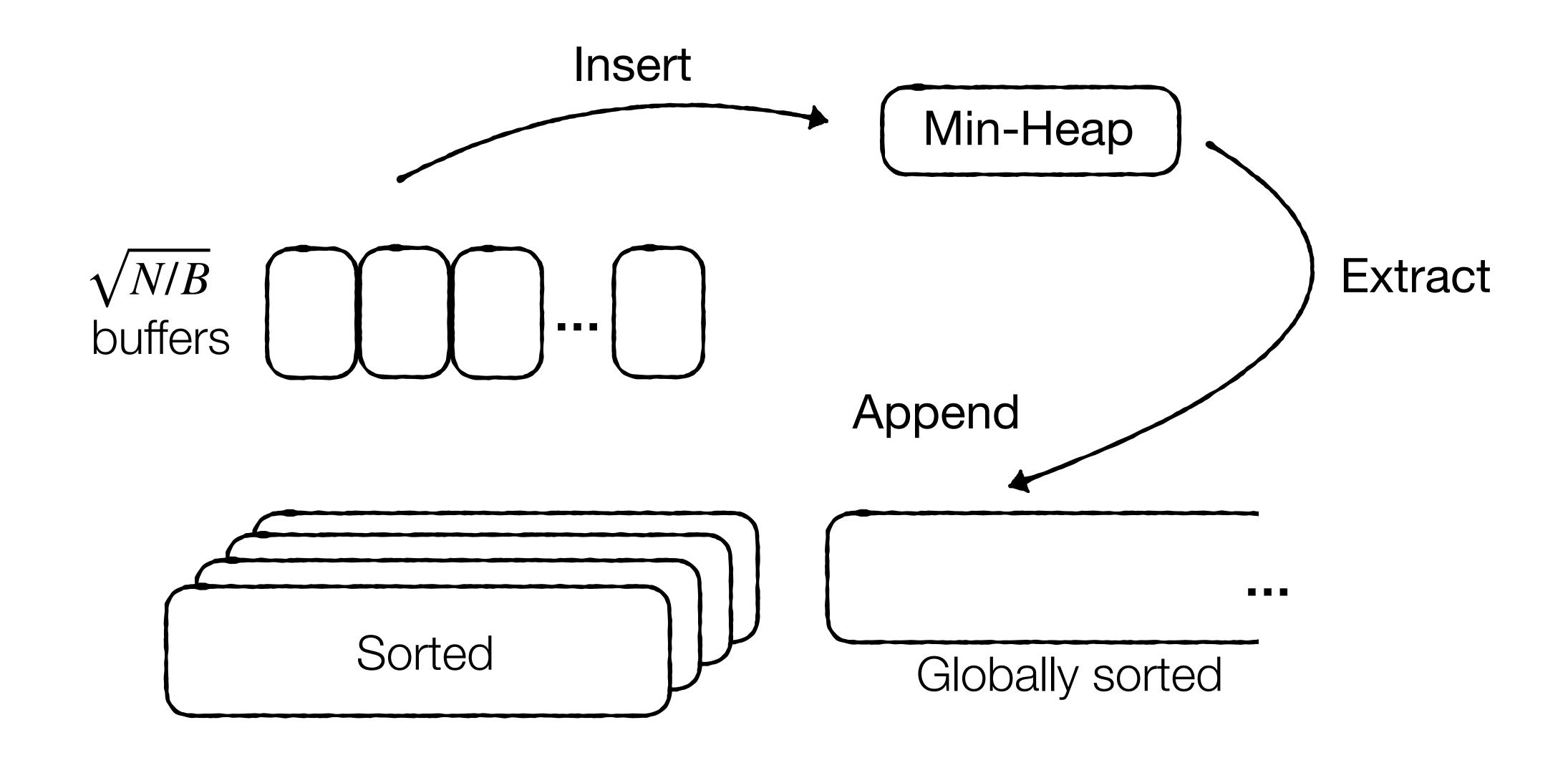




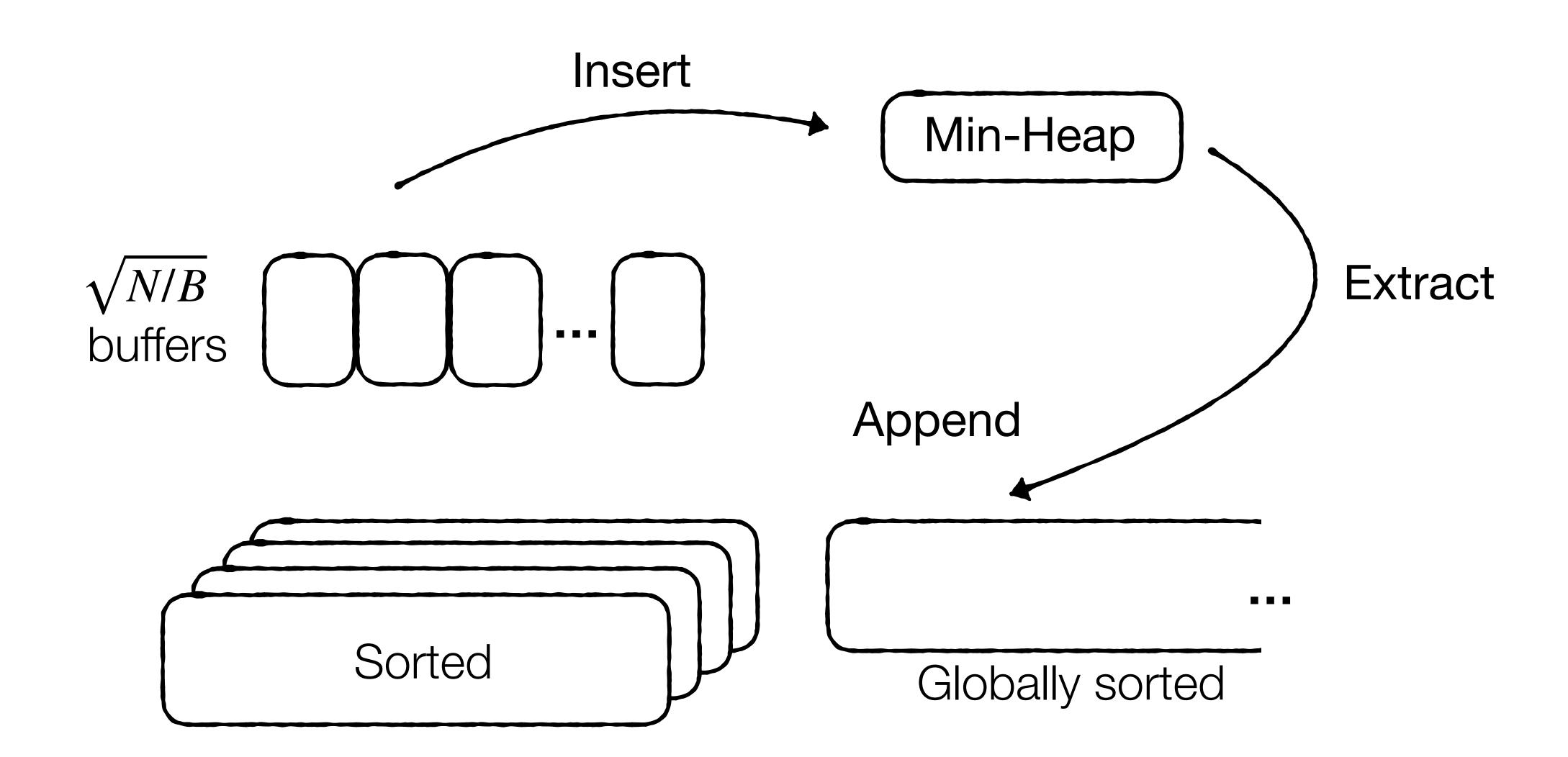
can do this with insert_and_extract



can do this with insert_and_extract: $O(log_2\sqrt{N/B})$ per entry



can do this with insert_and_extract: $O(log_2\sqrt{N/B})$ $O(N \cdot log_2\sqrt{N/B})$ overall



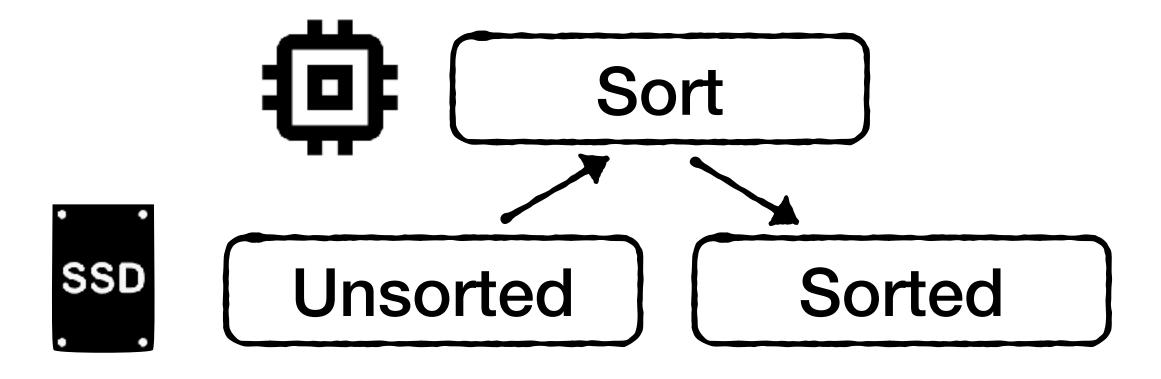
Analyzing CPU

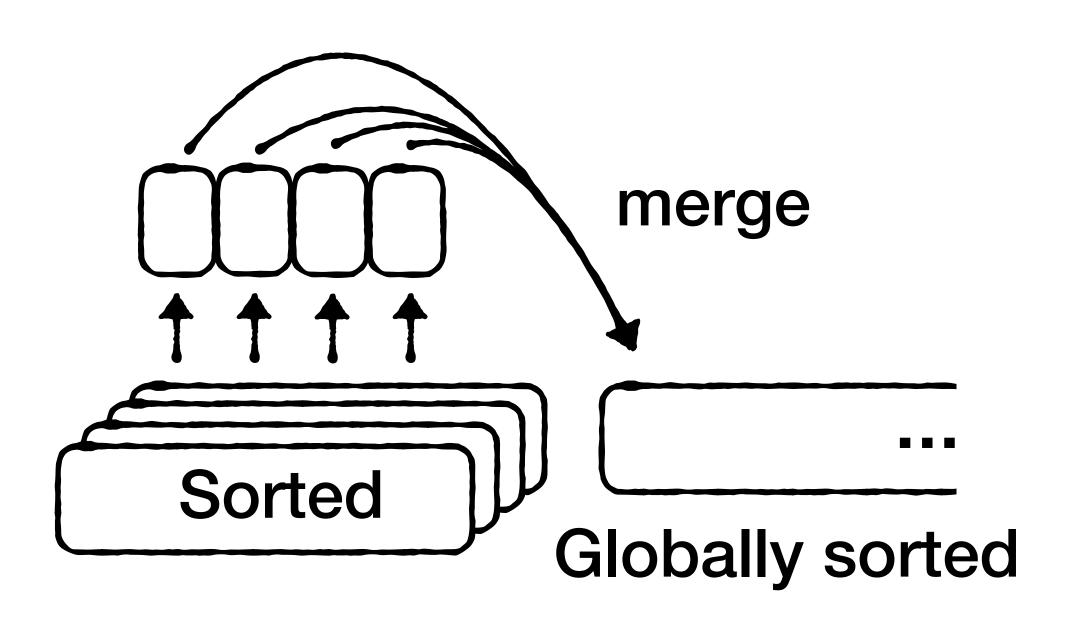
Partitioning Phase

 $O(N \cdot log_2M)$

Merging Phase

$$O(N \cdot log_2 \sqrt{N/B})$$



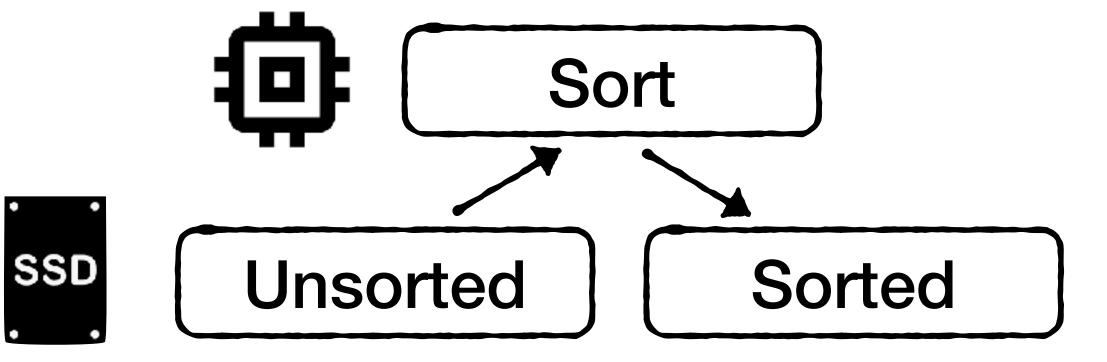


Analyzing CPU

Partitioning Phase

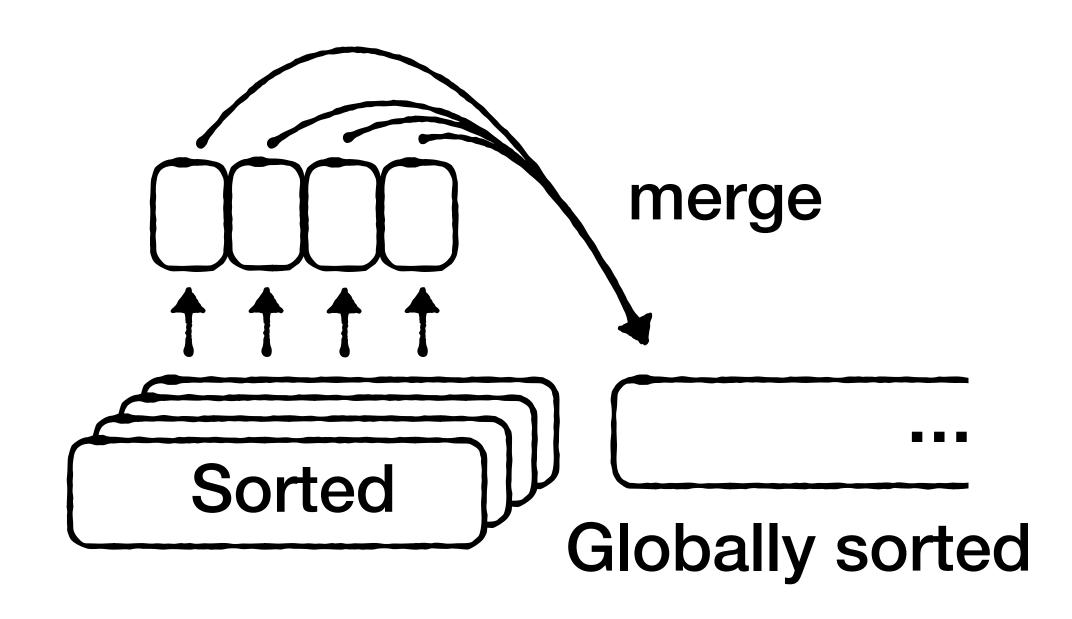
$$O(N \cdot log_2M)$$

$$= O(N \cdot log_2 \sqrt{N \cdot B})$$



Merging Phase

$$O(N \cdot log_2 \sqrt{N/B})$$



Analyzing CPU

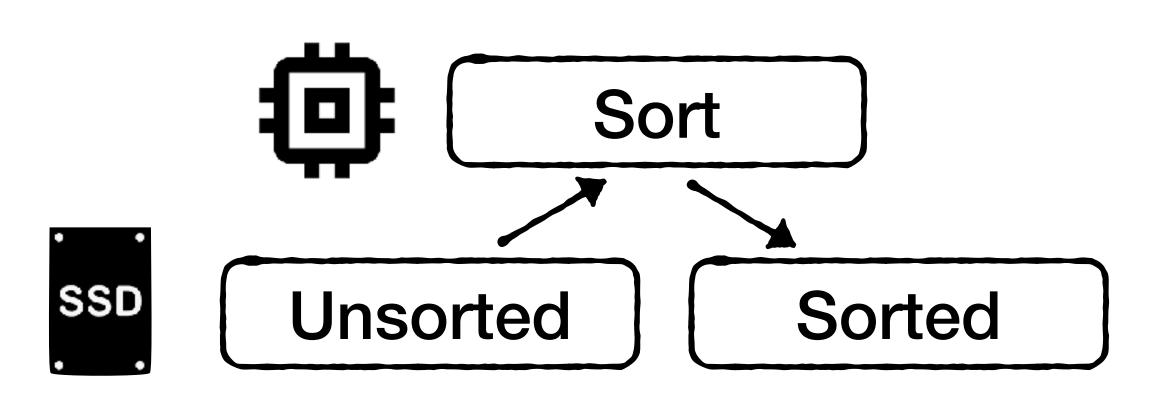
Partitioning Phase

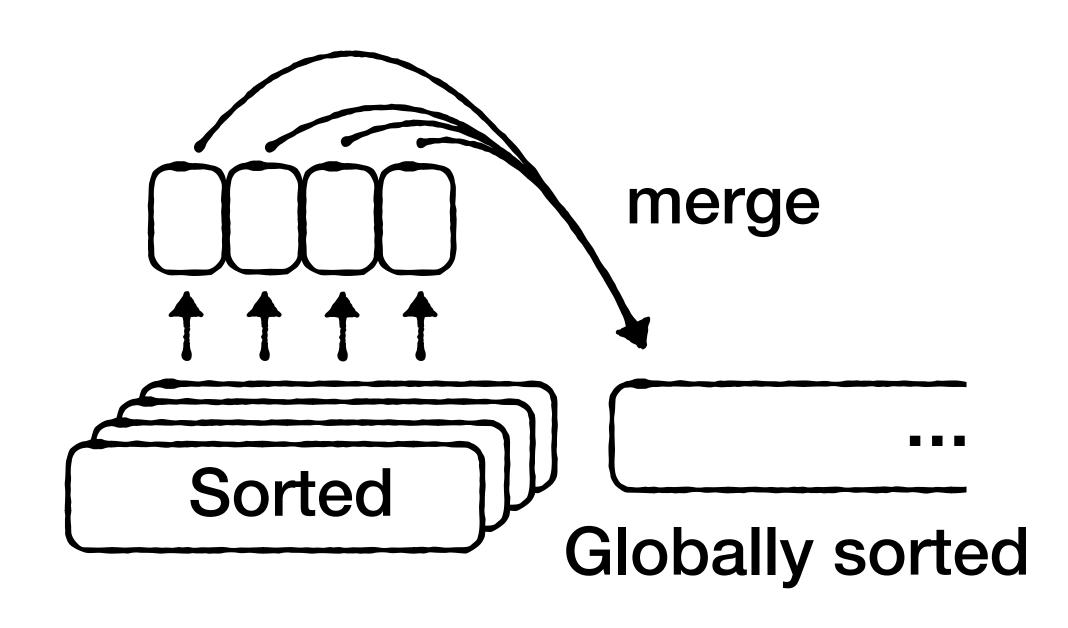
Total cost

$$O(N \cdot log_2 \sqrt{N \cdot B})$$
 + $O(N \cdot log_2 \sqrt{N/B})$

$$O(N \cdot log_2 \sqrt{N/B})$$

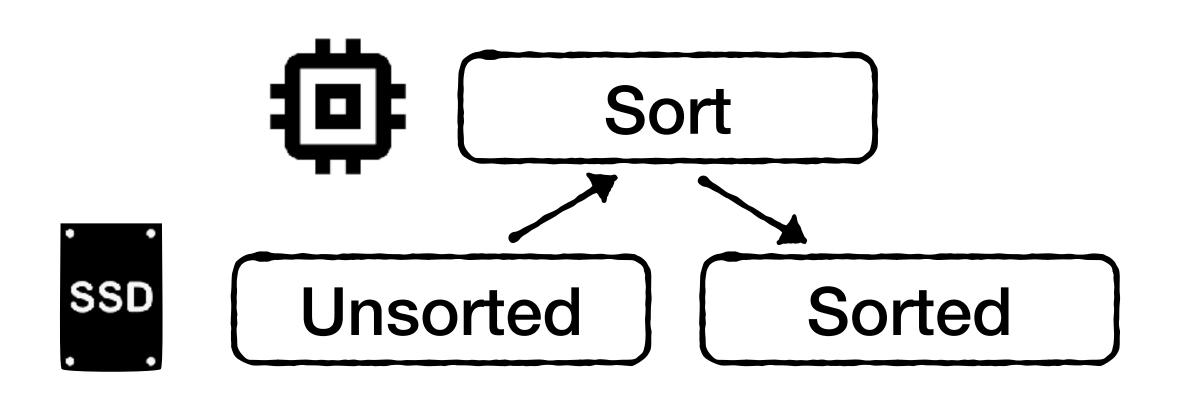
$$O(N \cdot log_2N)$$

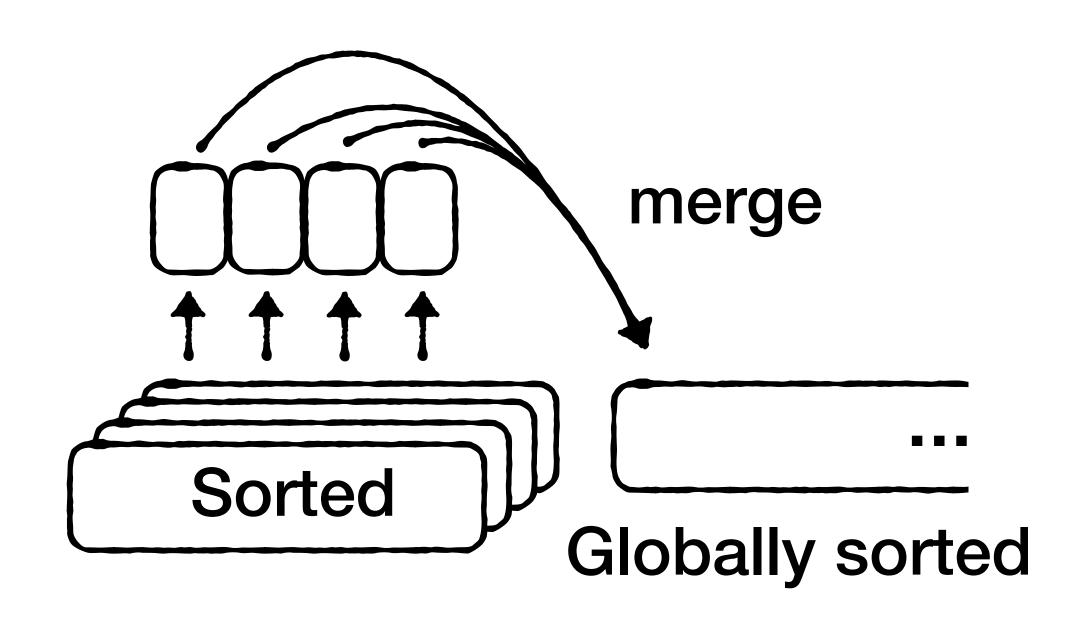




Same as in-memory algorithms:)

$$O(N \cdot log_2N)$$

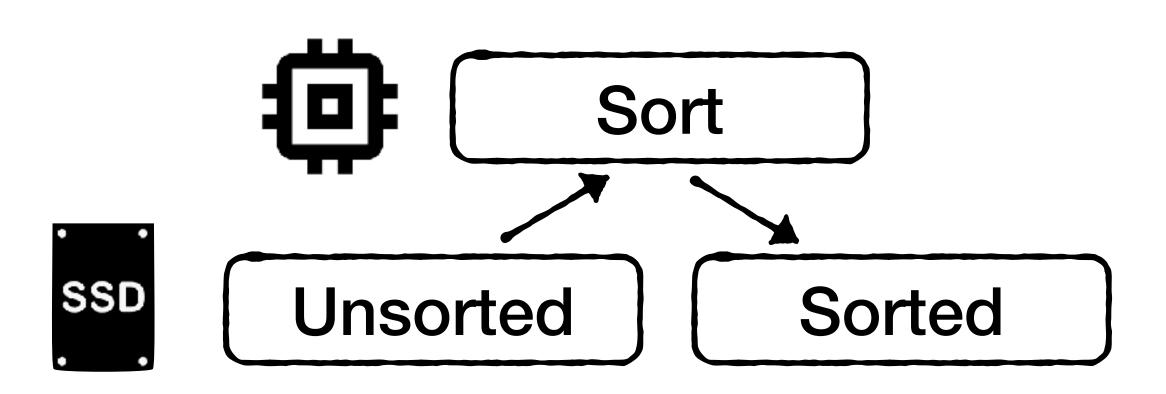


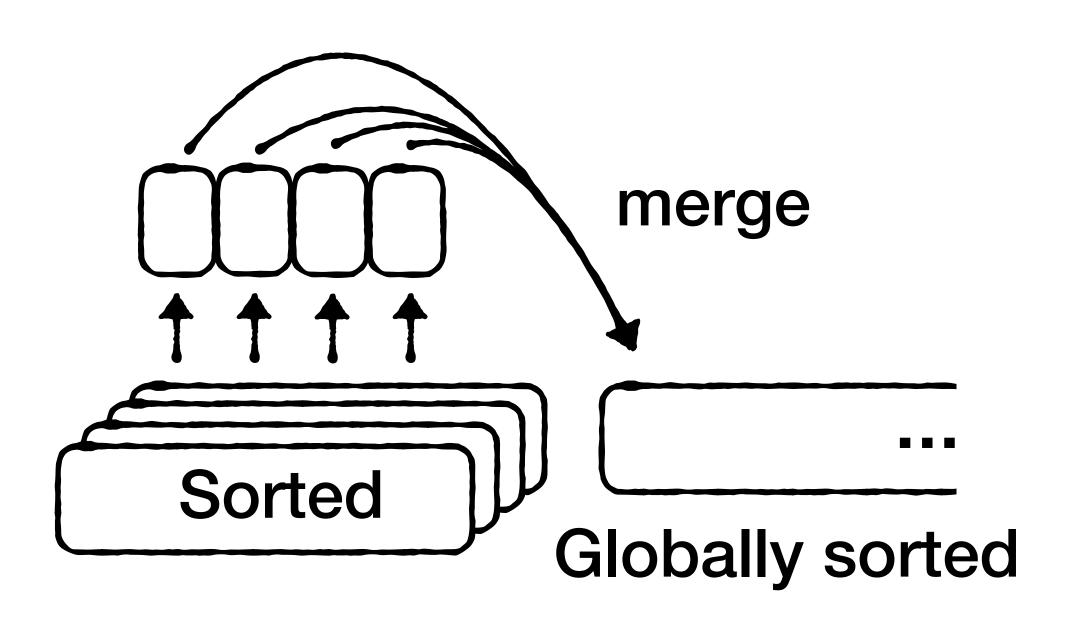


Overall costs

 $O(N \cdot log_2N)$ CPU

 $O(N/B \cdot log_{M/B}(N/B))$ //O

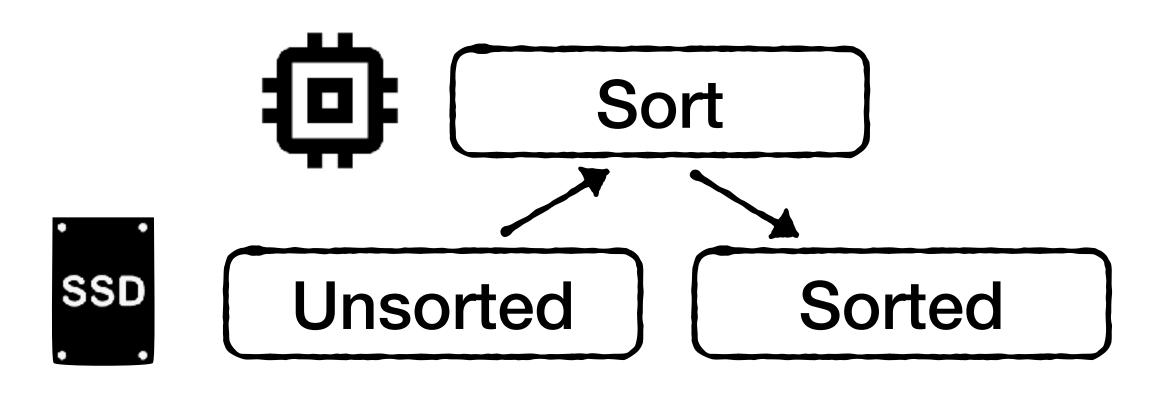


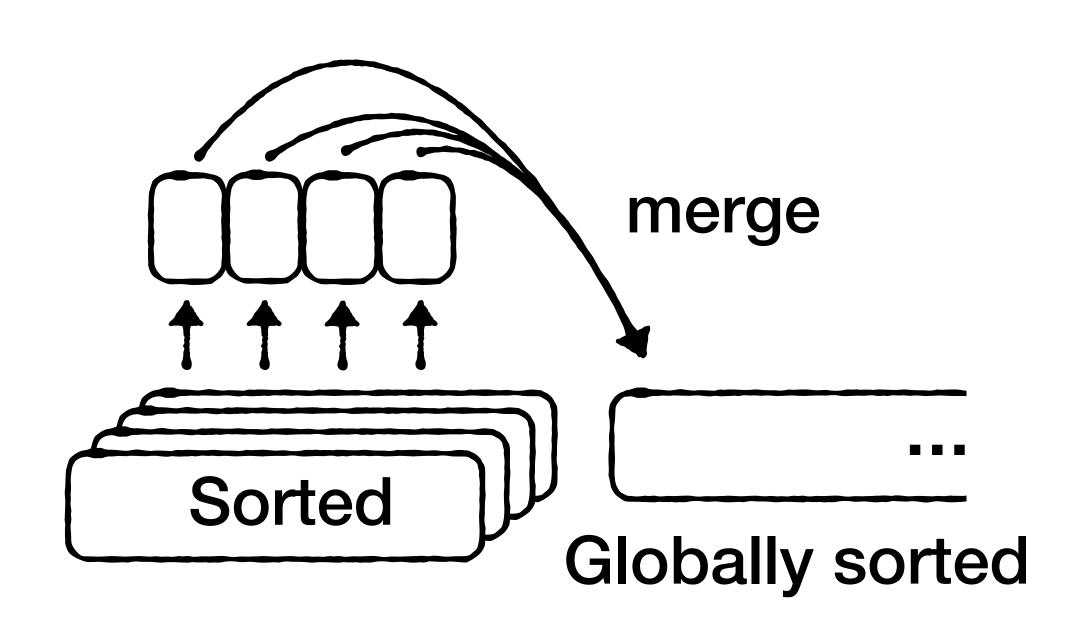


Overall costs

$$O(N \cdot log_2N)$$
 CPU

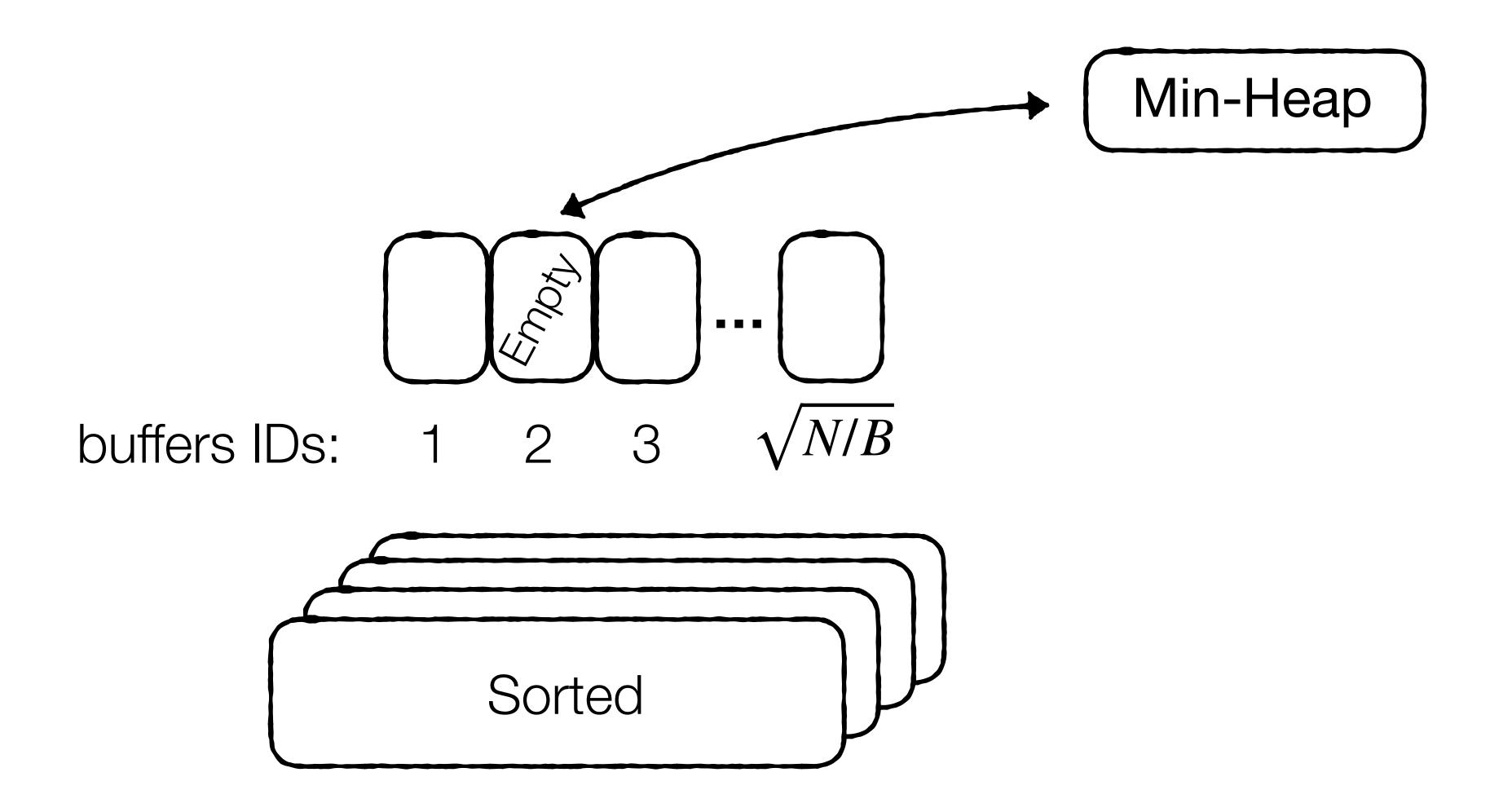
$$O(N/B \cdot log_{M/B}(N/B))$$
 I/O or $O(N/B)$ when $M > \sqrt{N \cdot B}$



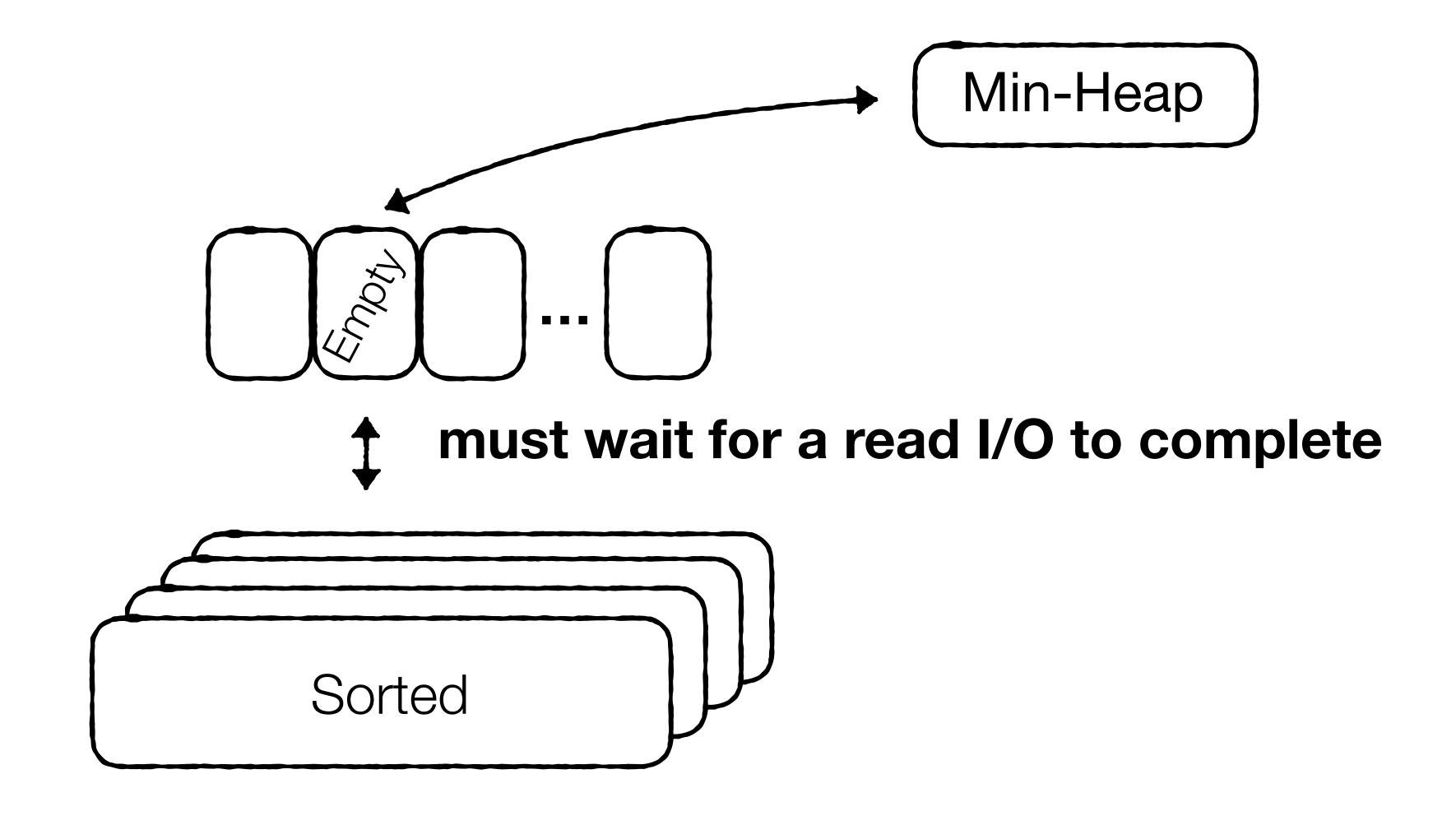


And now: TA office hours

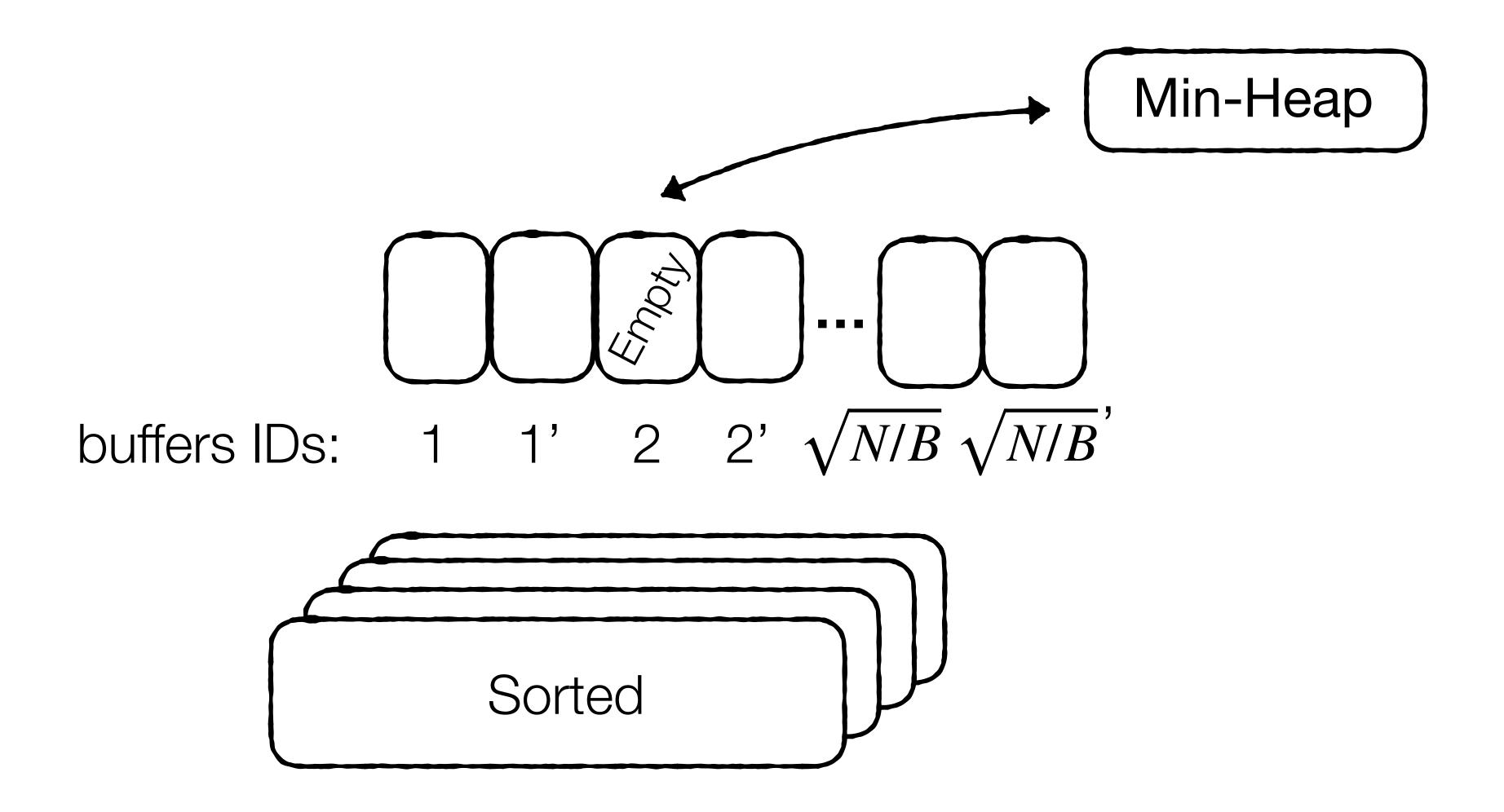
Suppose we need next min entry from buffer 2 but it is empty.



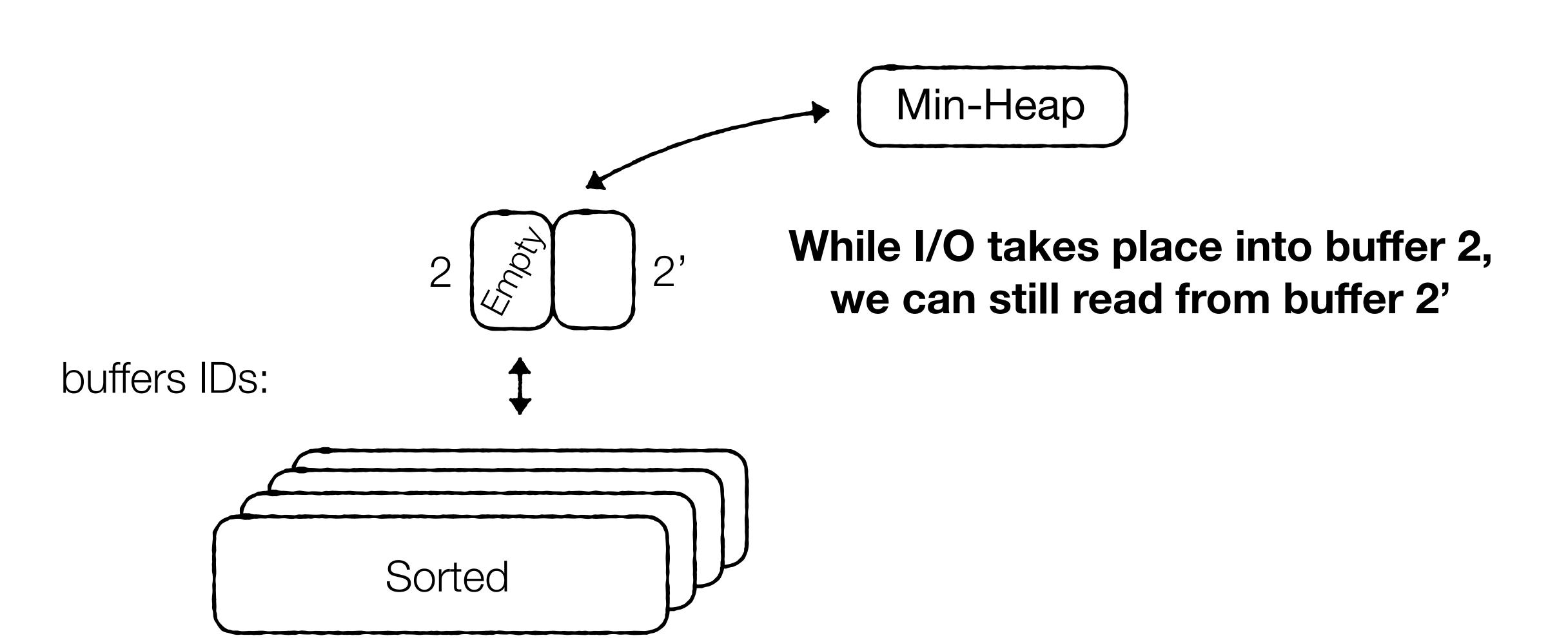
Suppose we need next min entry from buffer 2 but it is empty.



Double buffering: load one additional buffer preemptively for each partition before the first buffer empties



Double buffering: load one additional buffer preemptively for each partition before the first buffer empties



Larger though fewer buffers: more groups, so potentially more iterations, but each I/O reads more data

