

Example: Two-Way External Merge Sort

3,4 6,2 9,4 8,7 5,6 3,1 7,4 6,1

Input file

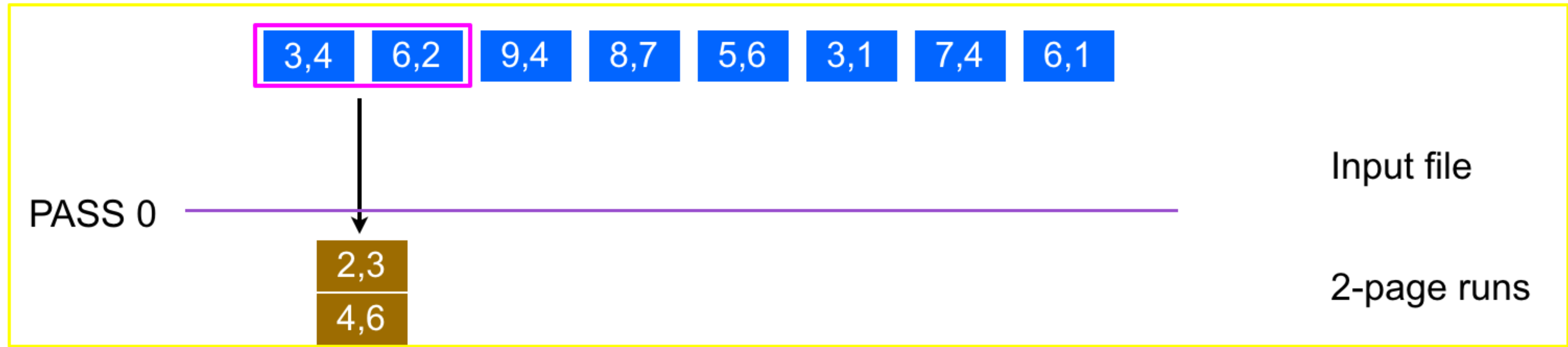
↑
page
(1) data bigger than mem memory
how to sort?

(2) Optimizations

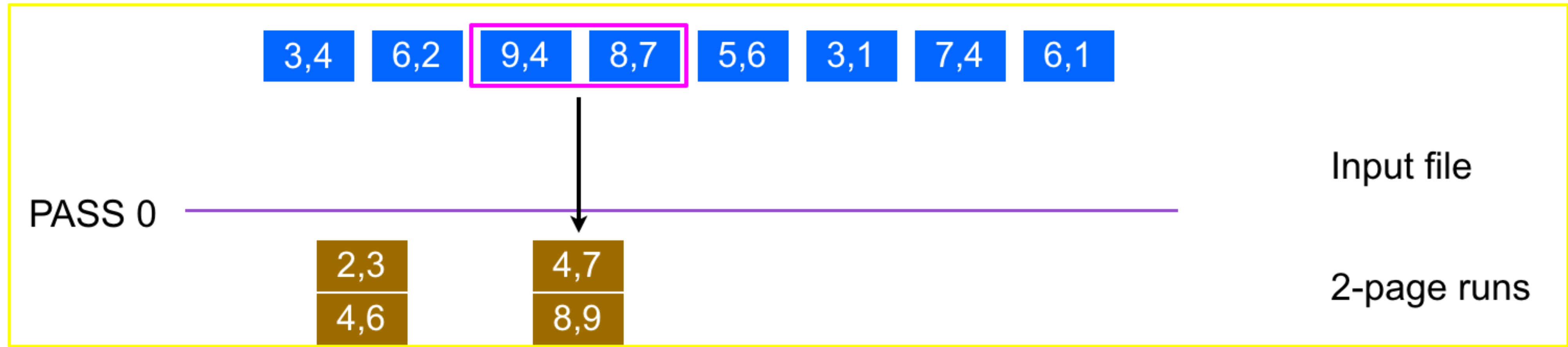
(3) Transition to bigger group of files

single algo $\xrightarrow{\text{EM 5}}$ multiple algo

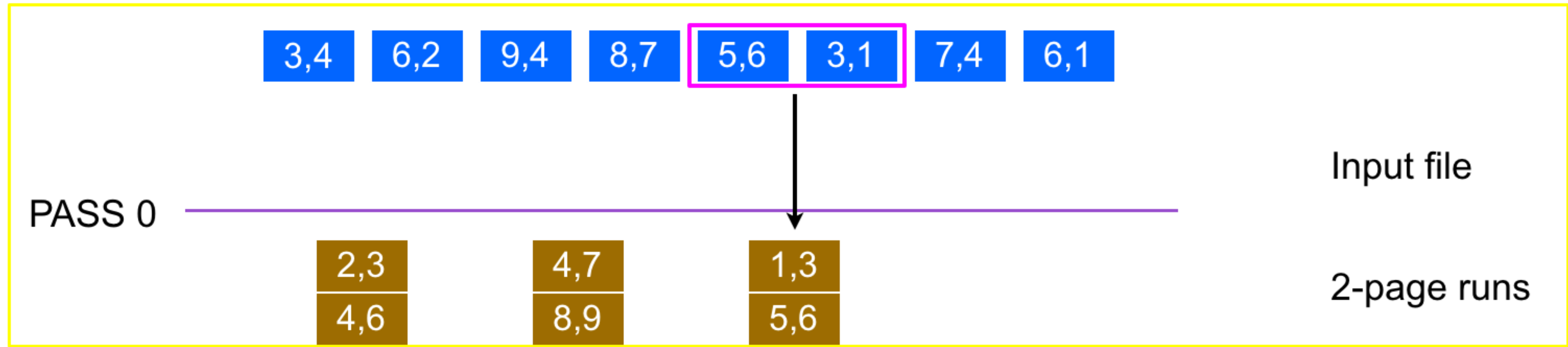
Run Generation: First Run



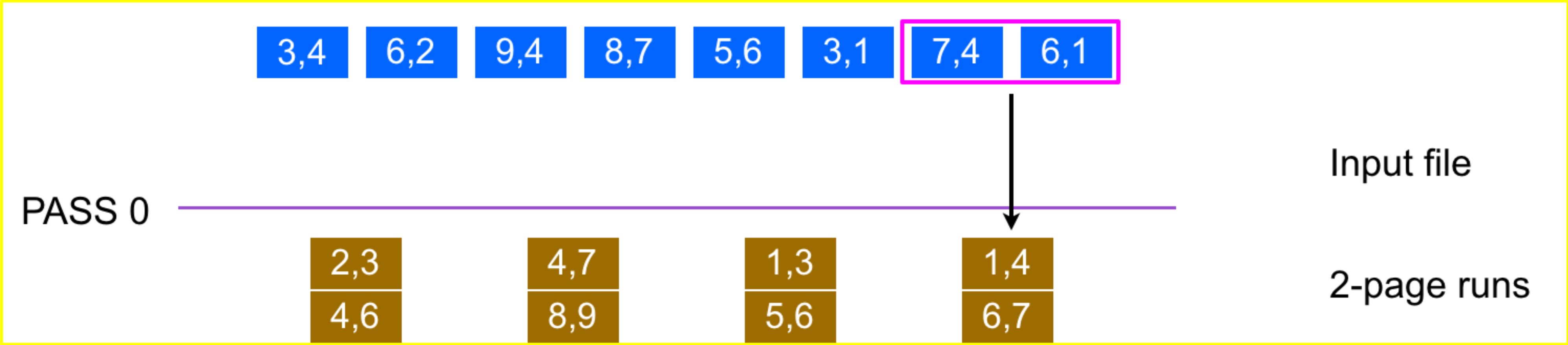
Run Generation: Second Run



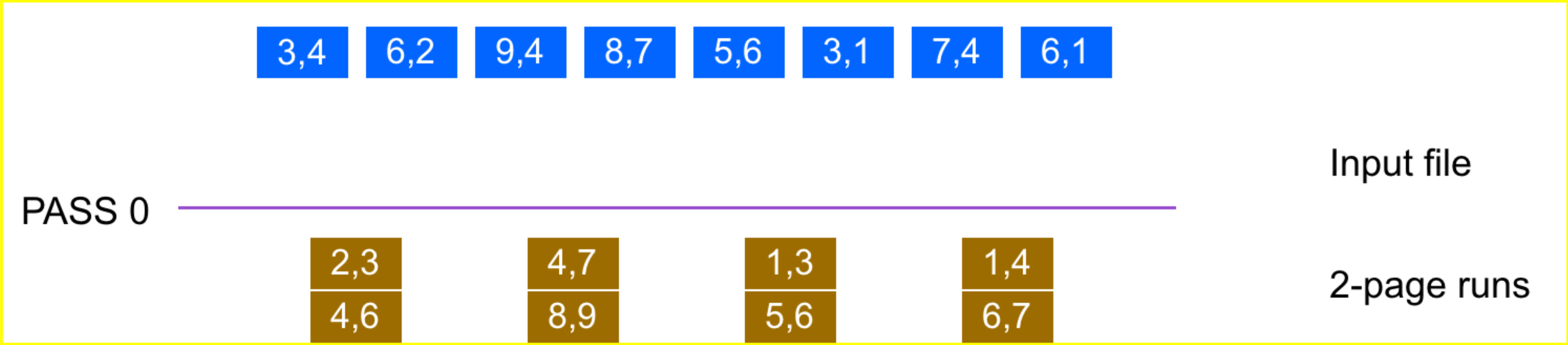
Run Generation: Third Run



Run Generation: Fourth Run



Run Generation Done.

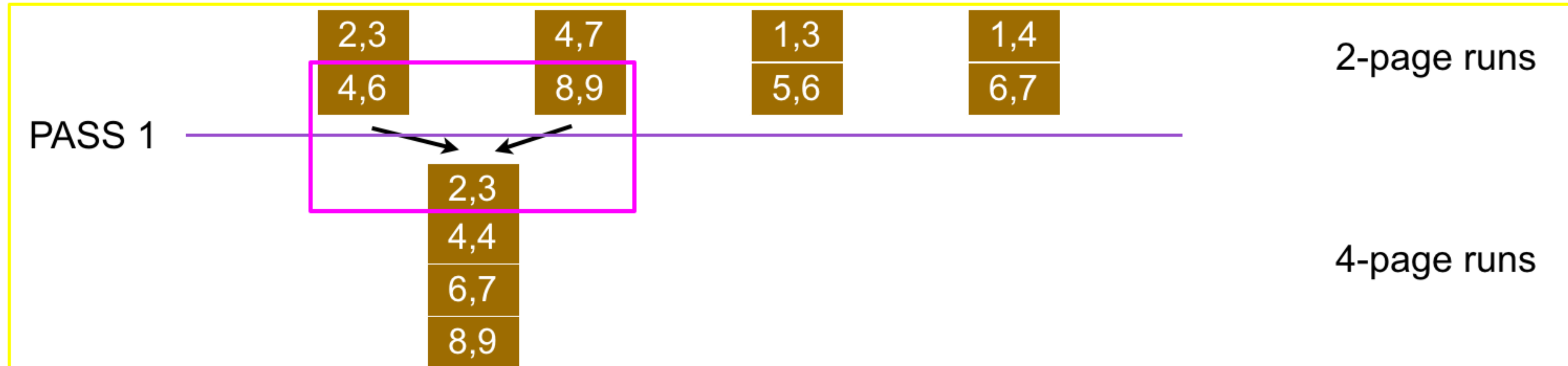


PASS 1: First Merge

3,4 6,2 9,4 8,7 5,6 3,1 7,4 6,1

Input file

PASS 0



Fen-ih
 $F = 2$

$m = \# \text{ of } m \text{ pages}$

$$F = m - 1$$

$$F \in [2; m-1]$$

PASS 1: Second Merge

3,4 6,2 9,4 8,7 5,6 3,1 7,4 6,1

Input file

PASS 0

2,3	4,7	1,3	1,4
4,6	8,9	5,6	6,7

2-page runs

PASS 1

2,3	1,1
4,4	3,4
6,7	5,6
8,9	6,7

4-page runs

PASS 1 Done.

3,4 6,2 9,4 8,7 5,6 3,1 7,4 6,1

Input file

PASS 0

2,3	4,7	1,3	1,4
4,6	8,9	5,6	6,7

2-page runs

PASS 1

2,3	1,1
4,4	3,4
6,7	5,6
8,9	6,7

4-page runs

PASS 2: First (and only) Merge

3,4 6,2 9,4 8,7 5,6 3,1 7,4 6,1

Input file

PASS 0

2,3	4,7	1,3	1,4
4,6	8,9	5,6	6,7

2-page runs

PASS 1

2,3	1,1
4,4	3,4
6,7	5,6
8,9	6,7

4-page runs

PASS 2

1,1
2,3
3,4
4,4
5,6
6,6
7,7
8,9

8-page runs

PASS 2: First (and only) Merge Done.

3,4 6,2 9,4 8,7 5,6 3,1 7,4 6,1

Input file

PASS 0

2,3 4,7 1,3 1,4
4,6 8,9 5,6 6,7

2-page runs

PASS 1

2,3 1,1
4,4 3,4
6,7 5,6
8,9 6,7

4-page runs

PASS 2

1,1
2,3
3,4
4,4
5,6
6,6
7,7
8,9

8-page runs

External Merge Sort

R

F := fan-in of the merge-phase

m := # pages available for a run (typically = available main memory)

ExternalSorting(R):

Heap<Run> runs;

//heap of runs to consider

While R not empty:

//Pass 0:

Run run = runGenerate(R, m);

//read m pages of input from R and sort

runs.add(run);

//add reference to this run to heap

size of the run

small \rightarrow highest priority

External Merge Sort

R

F := fan-in of the merge-phase

m := # pages available for a run (typically = available main memory)

ExternalSorting(R):

Heap<Run> runs;

While R not empty:

 Run run = runGenerate(R, m);

 runs.add(run);

While runs.size() > 1:

 List<Run> inputs;

 inputs = runs.popK(F);

 Run run = mergeRuns(inputs);

 runs.add(run);

//heap of runs to consider

//Pass 0:

 //read m pages of input from R and sort

 //add reference to this run to heap

//Passes 1 and following:

 //list of inputs to merge (in a single merge)

 //remove next F inputs from the heap

 //merge runs into one output run

 //add reference to merged run to runs

Very simple version of it!