**Slide 1**

# Advanced Database Systems
Spring 2026

## Q&A Session 2

THE UNIVERSITY of EDINBURGH

If you require this document in an alternative format, such as large print or a coloured background, please contact milos.nikolic@ed.ac.uk

---

**Slide 2**

## ABOUT THIS SESSION

Practice Worksheet 2 is now available on Learn

We will work through some questions during this session

This slide deck includes animations of discussed algorithms

External sorting

Join algorithms

**Main Purpose:** To reinforce and strengthen your understanding

---

**Slide 3**

## QUESTION 1

---

**Slide 4**

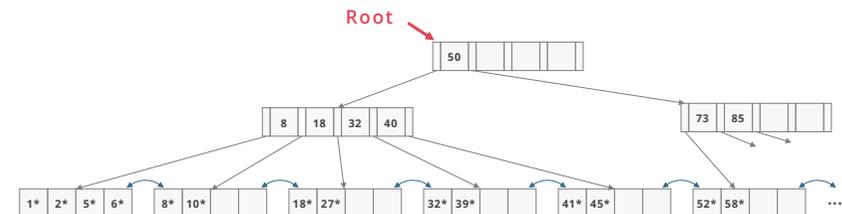## B+ TREE: SEARCH FOR 39
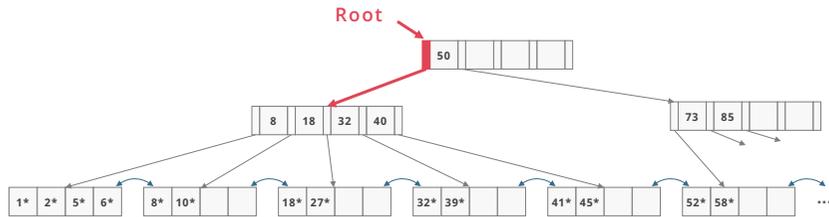
Find key = 39

Find split on each node

Follow pointer to next node

# B+ Tree: Search for 39

Find key = 39

Find split on each node

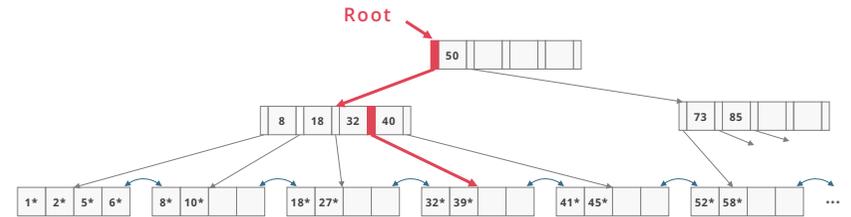Follow pointer to next node

Use binary search on each page

**Root**

| 50 | | | |

| 8 | 18 | 32 | 40 |

| 73 | 85 | | |

| 1* | 2* | 5* | 6* | | 8* | 10* | | 18* | 27* | | 32* | 39* | | 41* | 45* | | 52* | 58* | ...

# B+ Tree: Search for 39

Find key = 39

Find split on each node

Follow pointer to next node

Use binary search on each page

**Root**

| 50 | | | |

| 8 | 18 | 32 | 40 |

| 73 | 85 | | |

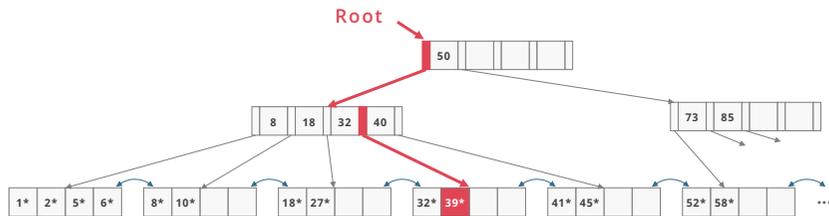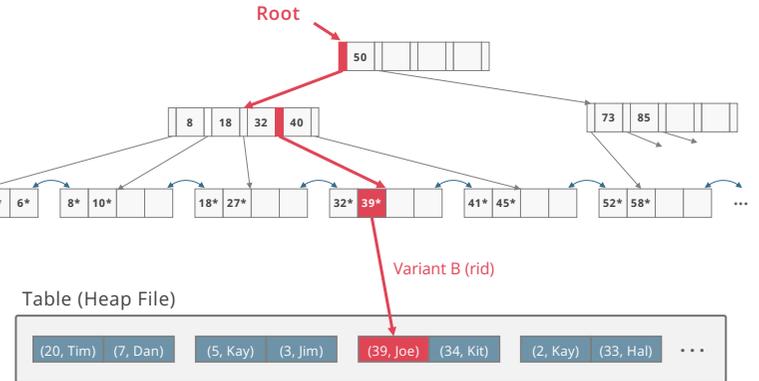| 1* | 2* | 5* | 6* | | 8* | 10* | | 18* | 27* | | 32* | 39* | | 41* | 45* | | 52* | 58* | ...

# B+ Tree: Search for 39

Find key = 39

Find split on each node

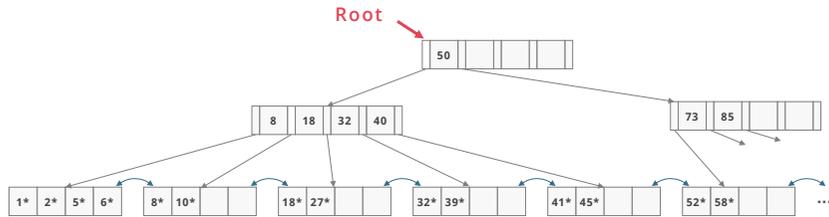Follow pointer to next node

Use binary search on each page

**Root**

| 50 | | | |

| 8 | 18 | 32 | 40 |

| 73 | 85 | | |

| 1* | 2* | 5* | 6* | | 8* | 10* | | 18* | 27* | | 32* | **39*** | | 41* | 45* | | 52* | 58* | ...

# B+ Tree: Search for 39

**Root**

| 50 | | | |

| 8 | 18 | 32 | 40 |

| 73 | 85 | | |

| 1* | 2* | 5* | 6* | | 8* | 10* | | 18* | 27* | | 32* | **39*** | | 41* | 45* | | 52* | 58* | ...

Variant B (rid)

Table (Heap File)

| (20, Tim) | (7, Dan) | | (5, Kay) | (3, Jim) | | (39, Joe) | (34, Kit) | | (2, Kay) | (33, Hal) | ...

# B+ Tree: Insert Entry 3*



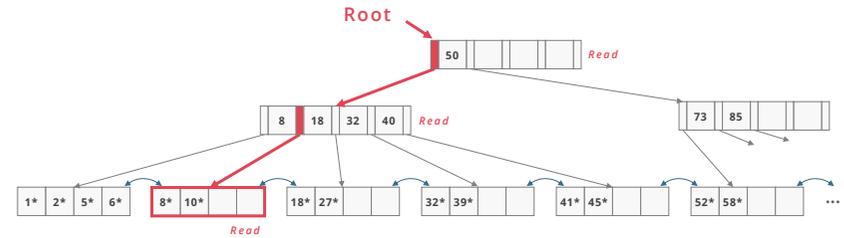**I/O Total (so far): 0**

13

# B+ Tree: Insert Entry 3*

Find the correct leaf node



**I/O Total (so far): 3**

14

# B+ Tree: Insert Entry 3*

Split leaf if not enough room: into two leaves with **d** and **d + 1** entries
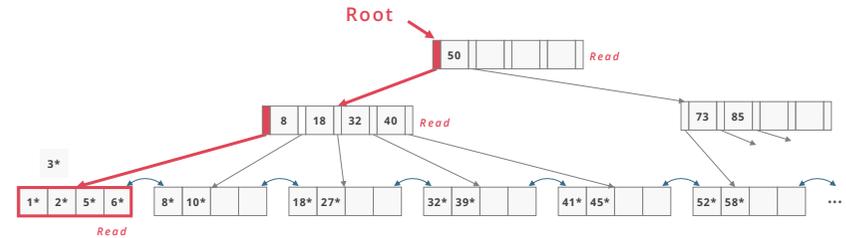


**I/O Total (so far): 3**

15

# B+ Tree: Insert Entry 3*

**Copy** up the middle key to inner node (since leaf nodes have data)



**I/O Total (so far): 7**

16

# B+ Tree: Insert Entry 3*

If inner node is full, split the inner node into two and **push** the middle key up



I/O Total (so far): 7

17

# B+ Tree: Insert Entry 3*

If inner node is full, split the inner node into two and **push** the middle key up



I/O Total (so far): 8

18

# B+ Tree: Insert Entry 3*

If inner node is full, split the inner node into two and **push** the middle key up



I/O Total (so far): 9

19

# B+ Tree: Insert Entry 3*

If inner node is full, split the inner node into two and **push** the middle key up



I/O Total (so far): 10

20

# QUESTION 2

# EXTENDIBLE HASHING

# INSERT ENTRY WITH HASH VALUE 68

# INSERT ENTRIES WITH HASH VALUES 17 & 69

# EXTERNAL SORTING

# GENERAL EXTERNAL MERGE SORT

**Goal**: sort records on 10 data pages using 4 buffer pages

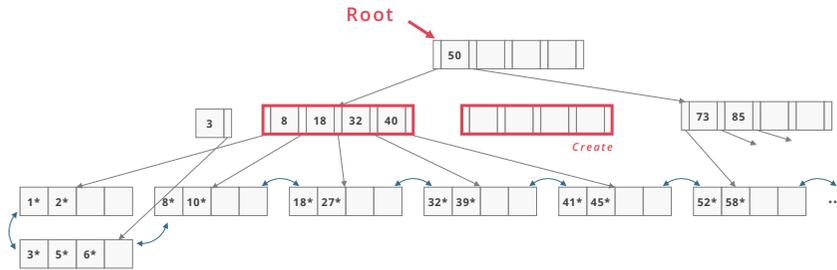| DISK | | MEMORY | DISK | |
|---|---|---|---|---|
| 6 | 8 | | 0 | 1 |
| 7 | 3 | | 2 | 3 |
| 9 | 1 | | 4 | 6 |
| 27 | 28 | | 7 | 8 |
| 12 | 17 | | 9 | 10 |
| 30 | 20 | | 11 | 12 |
| 10 | 4 | | 15 | 17 |
| 2 | 11 | | 20 | 25 |
| 15 | 25 | | 27 | 28 |
| 0 | 31 | | 30 | 31 |

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 1**

Read 4 pages into memory

| DISK | | MEMORY | |
|---|---|---|---|
| 6 | 8 | 6 | 8 |
| 7 | 3 | 7 | 3 |
| 9 | 1 | 9 | 1 |
| 27 | 28 | 27 | 28 |

**I/O Total (so far): 4**

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 1**

In-memory sort

| DISK | | MEMORY | |
|---|---|---|---|
| 6 | 8 | 1 | 3 |
| 7 | 3 | 6 | 7 |
| 9 | 1 | 8 | 9 |
| 27 | 28 | 27 | 28 |

**I/O Total (so far): 4**

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 1**

Write 4 sorted pages to disk

| DISK | | MEMORY | | DISK |
|---|---|---|---|---|
| 6 | 8 | 1 | 3 | 1 sorted run of 4 pages |
| 7 | 3 | 6 | 7 | 1  3  6  7  8  9  27 28 |
| 9 | 1 | 8 | 9 | |
| 27 | 28 | 27 | 28 | |

**I/O Total (so far): 8**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 2**

Read 4 pages into memory

| DISK | | MEMORY | | DISK |
|---|---|---|---|---|
| 12 | 17 | 12 | 17 | 1  3  6  7  8  9  27 28 |
| 30 | 20 | 30 | 20 | |
| 10 | 4 | 10 | 4 | |
| 2 | 11 | 2 | 11 | |

**I/O Total (so far): 12**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 2**

In-memory sort

| DISK | | MEMORY | | DISK |
|---|---|---|---|---|
| 12 | 17 | 2 | 4 | 1  3  6  7  8  9  27 28 |
| 30 | 20 | 10 | 11 | |
| 10 | 4 | 12 | 17 | |
| 2 | 11 | 20 | 30 | |

**I/O Total (so far): 12**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 2**

Write 4 sorted pages to disk

| DISK | | MEMORY | | DISK |
|---|---|---|---|---|
| 12 | 17 | 2 | 4 | 1  3  6  7  8  9  27 28 |
| 30 | 20 | 10 | 11 | 2  4  10 11 12 17 20 30 |
| 10 | 4 | 12 | 17 | |
| 2 | 11 | 20 | 30 | |

**I/O Total (so far): 16**

## Slide 33

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 3**

Read 2 pages into memory

| | | |
|---|---|---|
| **DISK** | **MEMORY** | **DISK** |

DISK: 15 25 / 0 31

MEMORY: 15 25 / 0 31 / *unused* / *unused*

DISK: 1 3 6 7 8 9 27 28 / 2 4 10 11 12 17 20 30

**I/O Total (so far): 18**

## Slide 34

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 3**

In-memory sort

DISK: 15 25 / 0 31

MEMORY: 0 15 / 25 31 / *unused* / *unused*

DISK: 1 3 6 7 8 9 27 28 / 2 4 10 11 12 17 20 30

| | | |
|---|---|---|
| **DISK** | **MEMORY** | **DISK** |

**I/O Total (so far): 18**

## Slide 35

# GENERAL EXTERNAL MERGE SORT

**Pass 0, Run 3**

Write 2 sorted pages to disk

DISK: 15 25 / 0 31

MEMORY: 0 15 / 25 31 / *unused* / *unused*

DISK: 1 3 6 7 8 9 27 28 / 2 4 10 11 12 17 20 30 / 0 15 25 31

| | | |
|---|---|---|
| **DISK** | **MEMORY** | **DISK** |

**I/O Total (so far): 20**

## Slide 36

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1: 1 3 6 7 8 9 27 28
Run 2: 2 4 10 11 12 17 20 30
Run 3: 0 15 25 31

MEMORY: *input* / *input* / *input* / *output*

| | | |
|---|---|---|
| **DISK** | **MEMORY** | **DISK** |

Reserve *B-1* input buffers and *1* output buffer. Load one page from each run at a time.

Store sorted results in output buffer. Write to disk when output buffer is full

# General External Merge Sort

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1: 1 3 | 6 7 | 8 9 | 27 28
Run 2: 2 4 | 10 11 | 12 17 | 20 30
Run 3: 0 15 | 25 31

Memory:
1 3
2 4
0 15
*empty*

DISK | MEMORY | DISK

**I/O Total (so far): 23**

37

---

# General External Merge Sort

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1: 6 7 | 8 9 | 27 28
Run 2: 10 11 | 12 17 | 20 30
Run 3: 25 31

Memory:
**1** 3
**2** 4
**0** 15
*empty*

DISK | MEMORY | DISK

**I/O Total (so far): 23**

38

---

# General External Merge Sort

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1: 6 7 | 8 9 | 27 28
Run 2: 10 11 | 12 17 | 20 30
Run 3: 25 31

Memory:
**1** 3
**2** 4
15
**0**

DISK | MEMORY | DISK

**I/O Total (so far): 23**

39

---

# General External Merge Sort

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1: 6 7 | 8 9 | 27 28
Run 2: 10 11 | 12 17 | 20 30
Run 3: 25 31

Memory:
**1** 3
**2** 4
**15**
0

DISK | MEMORY | DISK

**I/O Total (so far): 23**

40

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 6 7 8 9 27 28 | |
| Run 2 | 10 11 12 17 20 30 | |
| Run 3 | 25 31 | |

MEMORY:
3
2 4
15
0 1

DISK        MEMORY        DISK

I/O Total (so far): 23

41

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 6 7 8 9 27 28 | |
| Run 2 | 10 11 12 17 20 30 | |
| Run 3 | 25 31 | |

MEMORY:
3
2 4
15
empty

0 1

DISK        MEMORY        DISK

I/O Total (so far): 24

42

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 6 7 8 9 27 28 | |
| Run 2 | 10 11 12 17 20 30 | |
| Run 3 | 25 31 | |

MEMORY:
3
2 4
15
empty

0 1

DISK        MEMORY        DISK

I/O Total (so far): 24

43

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 6 7 8 9 27 28 | |
| Run 2 | 10 11 12 17 20 30 | |
| Run 3 | 25 31 | |

MEMORY:
3
4
15
2

0 1

DISK        MEMORY        DISK

I/O Total (so far): 24

44

## Slide 45

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 6 7 8 9 27 28
Run 2 | 10 11 12 17 20 30
Run 3 | 25 31

MEMORY: 3 / 4 / 15 / 2

0 1

DISK        MEMORY        DISK

I/O Total (so far): 24

---

## Slide 46

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 6 7 8 9 27 28
Run 2 | 10 11 12 17 20 30
Run 3 | 25 31

MEMORY: empty / 4 / 15 / 2 3

0 1

DISK        MEMORY        DISK

I/O Total (so far): 24

---

## Slide 47

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 6 7 8 9 27 28
Run 2 | 10 11 12 17 20 30
Run 3 | 25 31

MEMORY: 6 7 / 4 / 15 / empty

0 1 2 3

DISK        MEMORY        DISK

I/O Total (so far): 26

---

## Slide 48

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 8 9 27 28
Run 2 | 10 11 12 17 20 30
Run 3 | 25 31

MEMORY: 6 7 / 4 / 15 / empty

0 1 2 3

DISK        MEMORY        DISK

I/O Total (so far): 26

## Slide 49

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

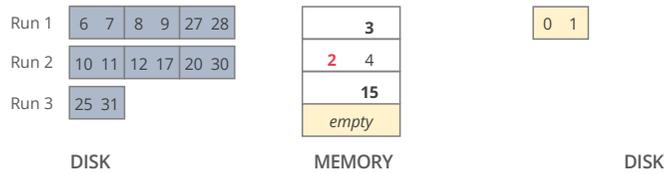Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 8 9 27 28
Run 2 | 10 11 12 17 20 30
Run 3 | 25 31

MEMORY: 6 7 / empty / 15 / 4

DISK: 0 1 2 3

**DISK**          **MEMORY**          **DISK**

**I/O Total (so far): 26**

## Slide 50

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 8 9 27 28
Run 2 | 10 11 12 17 20 30
Run 3 | 25 31

MEMORY: 6 7 / 10 11 / 15 / 4

DISK: 0 1 2 3

**DISK**          **MEMORY**          **DISK**

**I/O Total (so far): 27**

## Slide 51

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 8 9 27 28
Run 2 | 12 17 20 30
Run 3 | 25 31

MEMORY: 6 7 / 10 11 / 15 / 4

DISK: 0 1 2 3

**DISK**          **MEMORY**          **DISK**

**I/O Total (so far): 27**

## Slide 52

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1 | 8 9 27 28
Run 2 | 12 17 20 30
Run 3 | 25 31

MEMORY: 7 / 10 11 / 15 / 4 6

DISK: 0 1 2 3

**DISK**          **MEMORY**          **DISK**

**I/O Total (so far): 27**

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

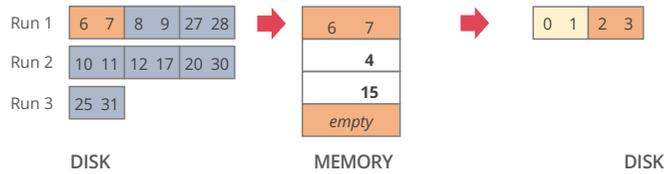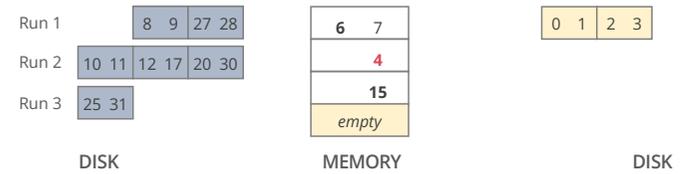Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 8  9  27  28 | |
| Run 2 | 12  17  20  30 | |
| Run 3 | 25  31 | |

MEMORY: 7 / 10  11 / 15 / empty

DISK: 0  1  2  3  4  6

**DISK**     **MEMORY**     **DISK**

**I/O Total (so far): 28**

53

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 8  9  27  28 | |
| Run 2 | 12  17  20  30 | |
| Run 3 | 25  31 | |

MEMORY: 7 / 10  11 / 15 / empty

DISK: 0  1  2  3  4  6

**DISK**     **MEMORY**     **DISK**

**I/O Total (so far): 28**

54

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 8  9  27  28 | |
| Run 2 | 12  17  20  30 | |
| Run 3 | 25  31 | |

MEMORY: empty / 10  11 / 15 / 7

DISK: 0  1  2  3  4  6

**DISK**     **MEMORY**     **DISK**

**I/O Total (so far): 28**

55

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

| | | |
|---|---|---|
| Run 1 | 8  9  27  28 | |
| Run 2 | 12  17  20  30 | |
| Run 3 | 25  31 | |

MEMORY: 8  9 / 10  11 / 15 / 7

DISK: 0  1  2  3  4  6

**DISK**     **MEMORY**     **DISK**

**I/O Total (so far): 29**

56

# GENERAL EXTERNAL MERGE SORT

Pass 1

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1    27 28
Run 2    12 17 20 30
Run 3    25 31

| 8 | 9 |
| 10 | 11 |
| 15 | |
| 7 | |

0 1 2 3 4 6

**DISK**        **MEMORY**        **DISK**

I/O Total (so far): 29

57

---

# GENERAL EXTERNAL MERGE SORT

Pass 1

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1    27 28
Run 2    12 17 20 30
Run 3    25 31

| | 9 |
| 10 | 11 |
| 15 | |
| 7 | 8 |

0 1 2 3 4 6

**DISK**        **MEMORY**        **DISK**

I/O Total (so far): 29

58

---

# GENERAL EXTERNAL MERGE SORT

Pass 1

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1    27 28
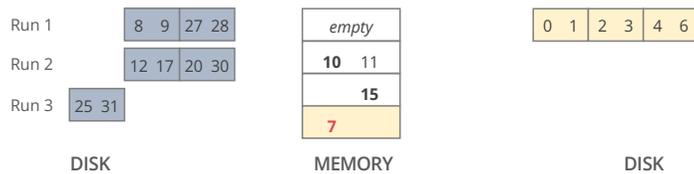Run 2    12 17 20 30
Run 3    25 31

| | 9 |
| 10 | 11 |
| 15 | |
| empty | |

0 1 2 3 4 6 7 8

**DISK**        **MEMORY**        **DISK**

I/O Total (so far): 30

59

---

# GENERAL EXTERNAL MERGE SORT

Pass 1

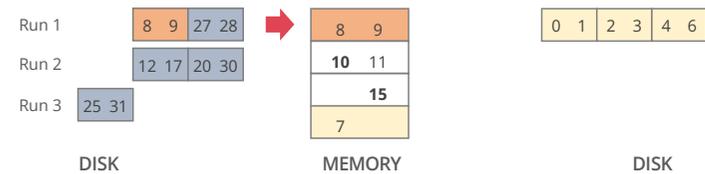Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1    27 28
Run 2    12 17 20 30
Run 3    25 31

| | 9 |
| 10 | 11 |
| 15 | |
| empty | |

0 1 2 3 4 6 7 8

**DISK**        **MEMORY**        **DISK**

I/O Total (so far): 30

60

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

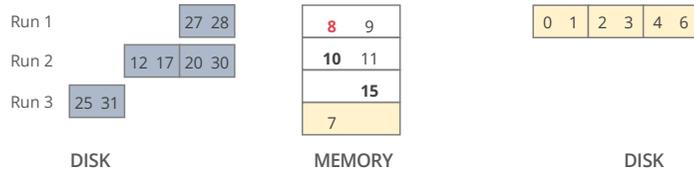Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2    12 17 20 30
Run 3    25 31

| 27 | 28 |
| | 11 |
| **15** | |
| *empty* | |

➡

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

DISK    MEMORY    DISK

**I/O Total (so far): 32**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
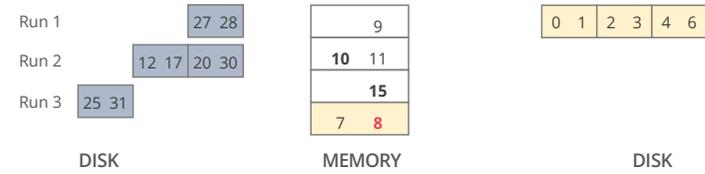Run 2    12 17 20 30
Run 3    25 31

| 27 | 28 |
| **11** | |
| **15** | |
| *empty* | |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

DISK    MEMORY    DISK

**I/O Total (so far): 32**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
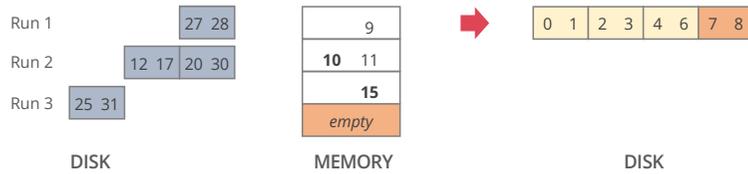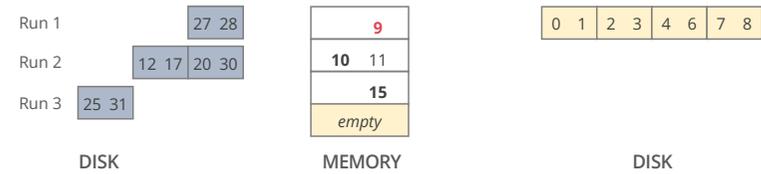Run 2    12 17 20 30
Run 3    25 31

| 27 | 28 |
| *empty* | |
| **15** | |
| **11** | |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

DISK    MEMORY    DISK

**I/O Total (so far): 32**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
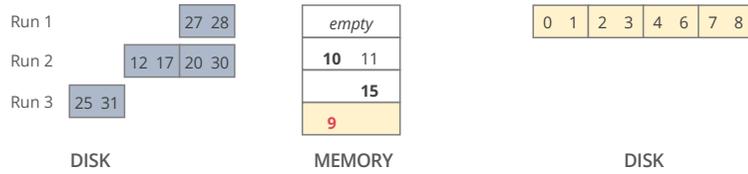Run 2    12 17 20 30
Run 3    25 31

| 27 | 28 |
| 12 | 17 |
| **15** | |
| 11 | |

➡

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

DISK    MEMORY    DISK

**I/O Total (so far): 33**

# GENERAL EXTERNAL MERGE SORT

Pass 1

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2    20 30
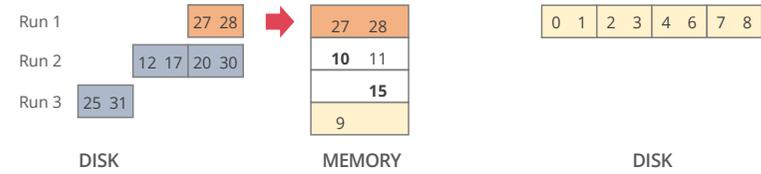Run 3    25 31

| 27 | 28 |
| 12 | 17 |
| **15** | |
| 11 | |

DISK          MEMORY          DISK

0 1 2 3 4 6 7 8 9 10

I/O Total (so far): 33

---

# GENERAL EXTERNAL MERGE SORT

Pass 1

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2    20 30
Run 3    25 31

| 27 | 28 |
| 17 | |
| **15** | |
| 11 | 12 |

DISK          MEMORY          DISK

0 1 2 3 4 6 7 8 9 10

I/O Total (so far): 33

---

# GENERAL EXTERNAL MERGE SORT

Pass 1

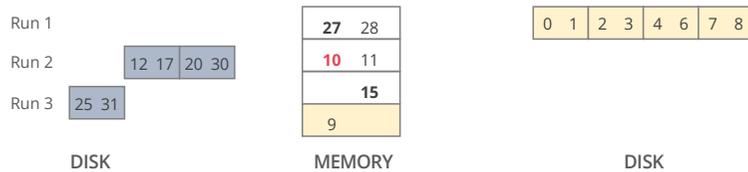Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2    20 30
Run 3    25 31

| 27 | 28 |
| 17 | |
| **15** | |
| *empty* | |

DISK          MEMORY          DISK

0 1 2 3 4 6 7 8 9 10
11 12

I/O Total (so far): 34

---

# GENERAL EXTERNAL MERGE SORT

Pass 1

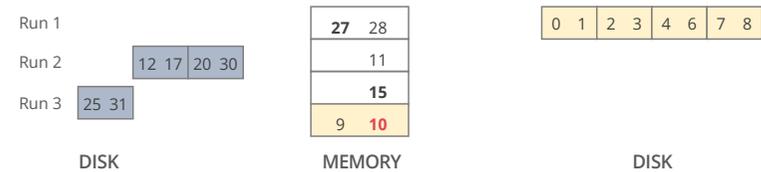Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2    20 30
Run 3    25 31

| 27 | 28 |
| **17** | |
| 15 | |
| *empty* | |

DISK          MEMORY          DISK

0 1 2 3 4 6 7 8 9 10
11 12

I/O Total (so far): 34

# GENERAL EXTERNAL MERGE SORT

### Pass 1

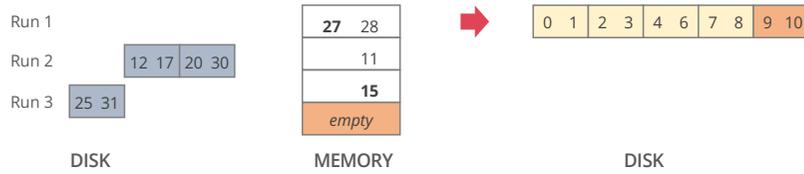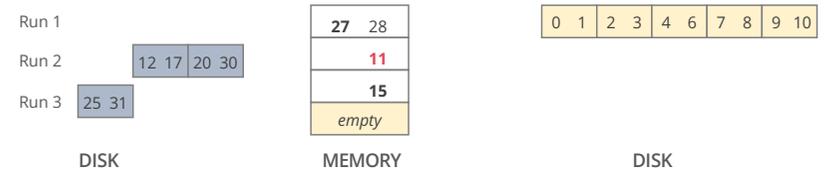Read 3 sorted runs into memory, write 1 sorted run to disk

| | |
|---|---|
| Run 1 | |
| Run 2 | 20 30 |
| Run 3 | 25 31 |

| **27** | 28 |
|---|---|
| **17** | |
| *empty* | |
| **15** | |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 11 | 12 | | | | | | | | |

**DISK**　　　　**MEMORY**　　　　**DISK**

**I/O Total (so far): 34**

---

# GENERAL EXTERNAL MERGE SORT

### Pass 1

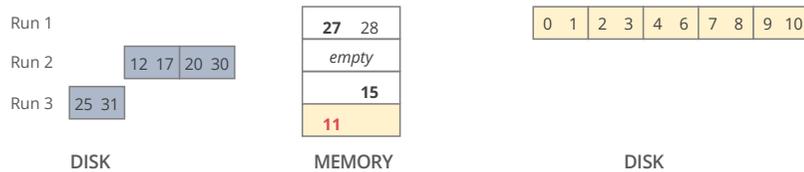Read 3 sorted runs into memory, write 1 sorted run to disk

| | |
|---|---|
| Run 1 | |
| Run 2 | 20 30 |
| Run 3 | 25 31 |

| **27** | 28 |
|---|---|
| **17** | |
| 25 | 31 |
| 15 | |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 11 | 12 | | | | | | | | |

**DISK**　　　　**MEMORY**　　　　**DISK**

**I/O Total (so far): 35**

---

# GENERAL EXTERNAL MERGE SORT

### Pass 1

Read 3 sorted runs into memory, write 1 sorted run to disk

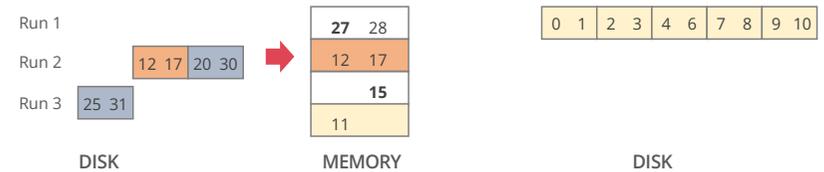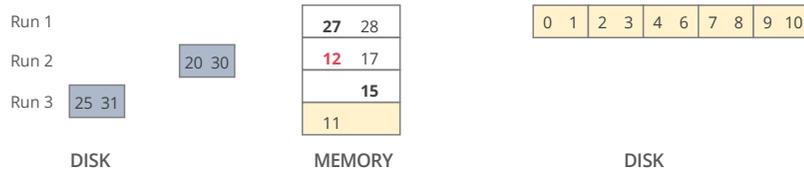| | |
|---|---|
| Run 1 | |
| Run 2 | 20 30 |
| Run 3 | |

| **27** | 28 |
|---|---|
| | **17** |
| **25** | 31 |
| 15 | |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 11 | 12 | | | | | | | | |

**DISK**　　　　**MEMORY**　　　　**DISK**

**I/O Total (so far): 35**

---

# GENERAL EXTERNAL MERGE SORT

### Pass 1

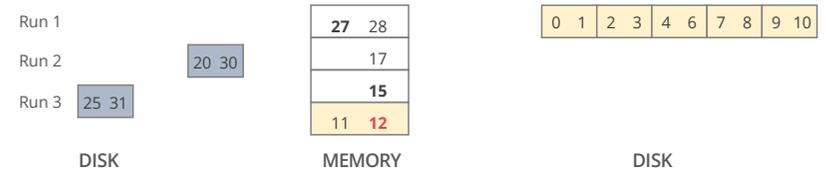Read 3 sorted runs into memory, write 1 sorted run to disk

| | |
|---|---|
| Run 1 | |
| Run 2 | 20 30 |
| Run 3 | |

| **27** | 28 |
|---|---|
| *empty* | |
| **25** | 31 |
| 15 | **17** |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 11 | 12 | | | | | | | | |

**DISK**　　　　**MEMORY**　　　　**DISK**

**I/O Total (so far): 35**

# GENERAL EXTERNAL MERGE SORT

## Pass 1

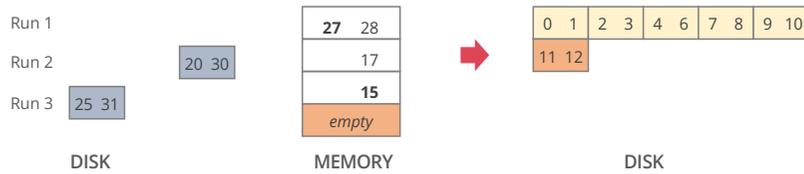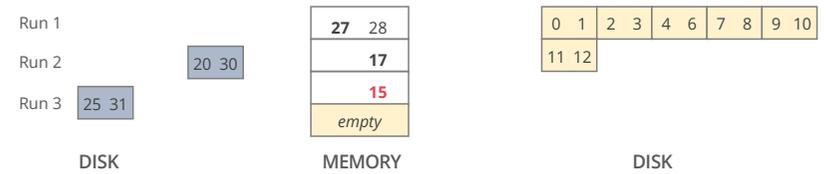Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2
Run 3

| 27 | 28 |
| 30 |
| 31 |
| 20 | 25 |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 |

DISK          MEMORY          DISK

I/O Total (so far): 37

---

# GENERAL EXTERNAL MERGE SORT

## Pass 1

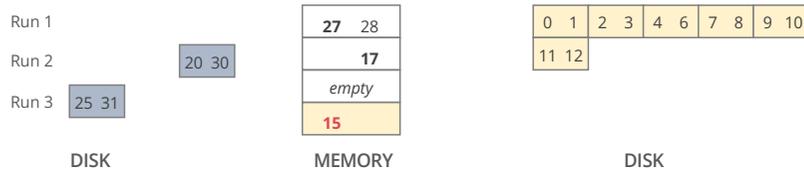Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2
Run 3

| 27 | 28 |
| 30 |
| 31 |
| empty |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 |

DISK          MEMORY          DISK

I/O Total (so far): 38

---

# GENERAL EXTERNAL MERGE SORT

## Pass 1

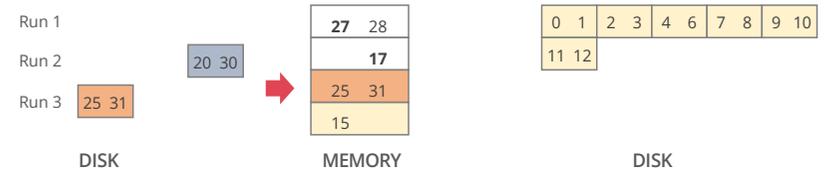Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2
Run 3

| 27 | 28 |
| 30 |
| 31 |
| empty |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 |

DISK          MEMORY          DISK

I/O Total (so far): 38

---

# GENERAL EXTERNAL MERGE SORT

## Pass 1

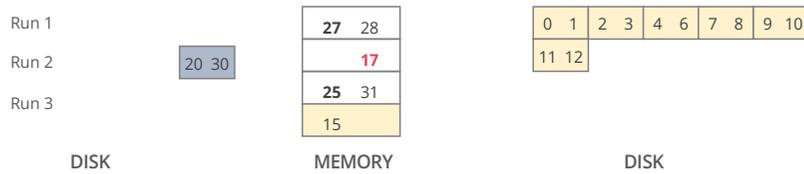Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2
Run 3

| 28 |
| 30 |
| 31 |
| 27 |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 |

DISK          MEMORY          DISK

I/O Total (so far): 38

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

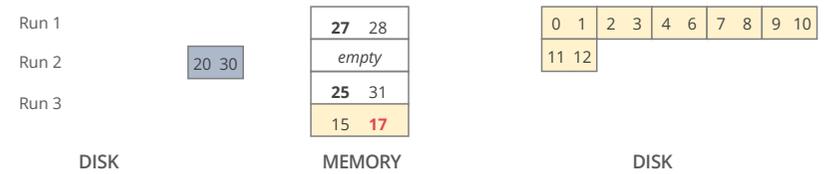Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1  
Run 2  
Run 3

| 28 |
| 30 |
| 31 |
| 27 |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 |

DISK         MEMORY         DISK

**I/O Total (so far): 38**

85

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk
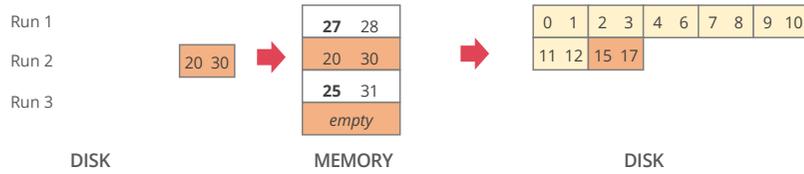
Run 1  
Run 2  
Run 3

| empty |
| 30 |
| 31 |
| 27   28 |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 |

DISK         MEMORY         DISK

**I/O Total (so far): 38**

86

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk
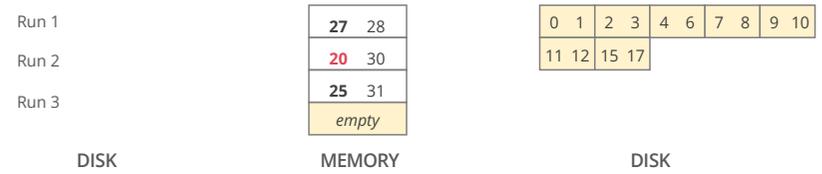
Run 1  
Run 2  
Run 3

| empty |
| 30 |
| 31 |
| empty |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 |

DISK         MEMORY         DISK

**I/O Total (so far): 39**

87

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

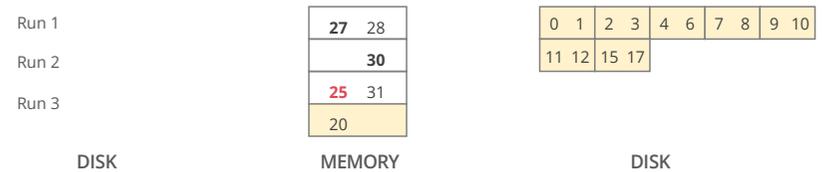Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1  
Run 2  
Run 3

| empty |
| 30 |
| 31 |
| empty |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 |

DISK         MEMORY         DISK

**I/O Total (so far): 39**

88

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
Run 2
Run 3

| empty |
| empty |
| **31** |
| **30** |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 |

DISK          MEMORY          DISK

**I/O Total (so far): 39**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

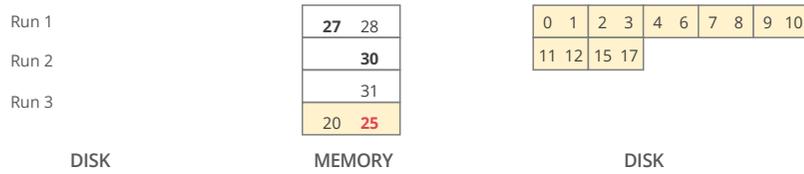Read 3 sorted runs into memory, write 1 sorted run to disk
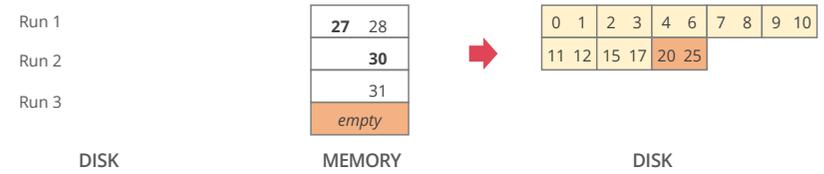
Run 1
Run 2
Run 3

| empty |
| empty |
| **31** |
| 30 |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 |

DISK          MEMORY          DISK

**I/O Total (so far): 39**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

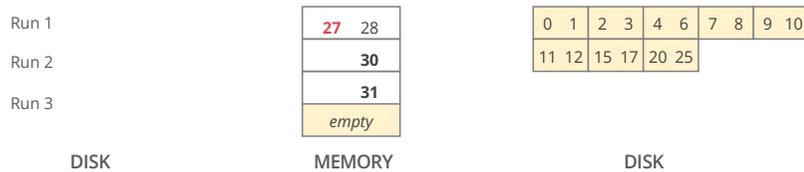Run 1
Run 2
Run 3

| empty |
| empty |
| empty |
| 30 | **31** |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 |

DISK          MEMORY          DISK

**I/O Total (so far): 39**

---

# GENERAL EXTERNAL MERGE SORT

**Pass 1**

Read 3 sorted runs into memory, write 1 sorted run to disk

Run 1
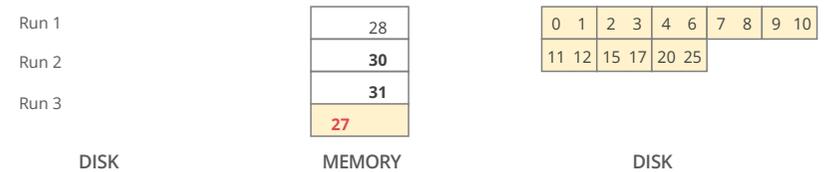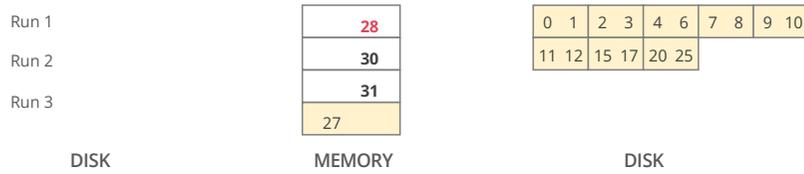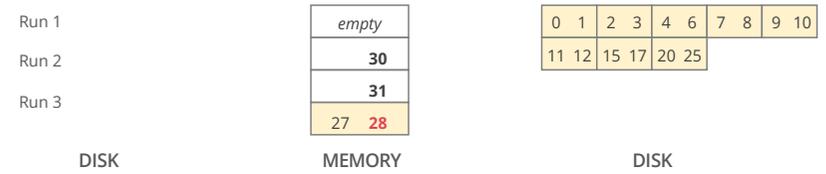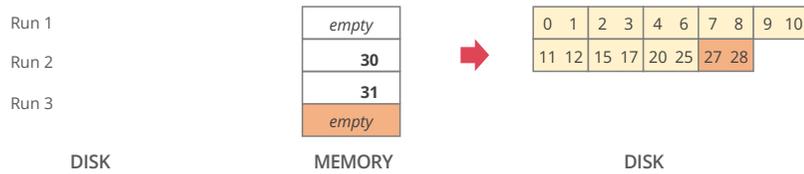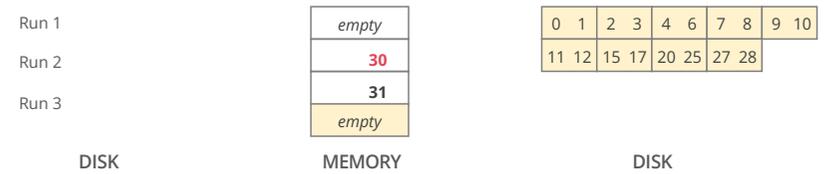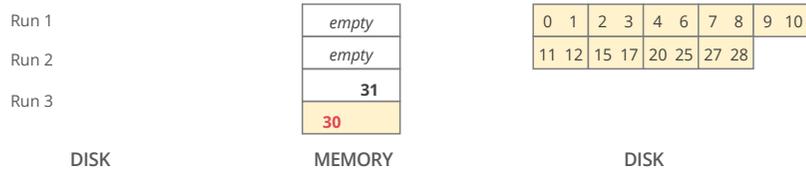Run 2
Run 3

| empty |
| empty |
| empty |
| empty |

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |

| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 | 30 | 31 |

1 sorted run of 10 pages

DISK          MEMORY          DISK

**I/O Total: 40**

## SANITY CHECK

N = 10, B = 4

Cost = 2N * (1 + ⌈ $\log_{B-1}$(⌈N/B⌉) ⌉)

= 2 * 10 * (1 + ⌈ $\log_3$(2.5) ⌉)

= 20 * (1 + 1)

= 40 I/Os  ✓

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 11 | 12 | 15 | 17 | 20 | 25 | 27 | 28 | 30 | 31 |

1 sorted run of 10 pages

**I/O Total: 40**

---

# JOINS

---

## JOINS

Bit of notation:

*M* = number of pages in *R*

*m* = number of records in *R*  (the **cardinality** of *R*)

*N* = number of pages in *S*

*n* = number of records in *S*

R(id, …)   S(id, …)

*M* pages
*m* tuples

*N* pages
*n* tuples

We typically exclude the final write's I/O cost

Don't add the cost of writing the joined output to disk

We might decide to stream it to the next operator instead of materializing results!

---

## SIMPLE NESTED LOOPS JOIN (SNLJ)

Direct translation of the definition of join into code

To perform $R \bowtie_\theta S$, take each tuple in R and scan through S to find the matching S-tuples!

```
foreach tuple r ∈ R:
  foreach tuple s ∈ S:
    if θ(r,s): output r joined with s
```

## Slide 97

# SIMPLE NESTED LOOPS JOIN (SNLJ)



First iteration of outer loop…

Compare with all tuples in *S*
… and add matches to result

*M* pages
*m* tuples
*R*

*N* pages
*n* tuples
*S*

## Slide 98

# SIMPLE NESTED LOOPS JOIN (SNLJ)



Second iteration of outer loop…

Compare with all tuples in *S*
… and add matches to result

*M* pages
*m* tuples
*R*

*N* pages
*n* tuples
*S*

## Slide 99

# SIMPLE NESTED LOOPS JOIN (SNLJ)



$$\text{Cost} = M + m \cdot N$$

Flipping the order of *R* and *S*
can change the I/O cost

Which relation to use as outer?

*M* pages
*m* tuples
*R*

*N* pages
*n* tuples
*S*

## Slide 100

# PAGE NESTED LOOPS JOIN (PNLJ)

Can we do better?

We scan *S* for every tuple in *R*,

… but we had to load an entire page of *R* into memory to get that tuple!

Instead of finding the tuples in *S* that match a tuple in *R*,

… do the check for all tuples in a page in *R* at once

# PAGE NESTED LOOPS JOIN (PNLJ)

SNLJ

```
foreach tuple r ∈ R:
  foreach tuple s ∈ S:
    if θ(r,s): output r joined with s
```

101

# PAGE NESTED LOOPS JOIN (PNLJ)

SNLJ (but with page fetches written out explicitly)

```
foreach page Pᵣ ∈ R:
  foreach tuple r ∈ Pᵣ:
    foreach page Pₛ ∈ S:
      foreach tuple s ∈ Pₛ:
        if θ(r,s): output r joined with s
```

102

# PAGE NESTED LOOPS JOIN (PNLJ)

PNLJ

```
foreach page Pᵣ ∈ R:
  foreach page Pₛ ∈ S:          flipped loops
    foreach tuple r ∈ Pᵣ:
      foreach tuple s ∈ Pₛ:
        if θ(r,s): output r joined with s
```

103

# PAGE NESTED LOOPS JOIN (PNLJ)



First iteration of outer loop…

Compare with all tuples in *S*
… and add matches to result

*M* pages
*m* tuples          *R*

*N* pages
*n* tuples          *S*

104

# PAGE NESTED LOOPS JOIN (PNLJ)



Second iteration of outer loop...

Compare with all tuples in *S*
... and add matches to result

*M* pages
*m* tuples  *R*

*N* pages
*n* tuples  *S*

---

# PAGE NESTED LOOPS JOIN (PNLJ)



$\text{Cost} = M + M \cdot N$

*M* pages
*m* tuples  *R*

*N* pages
*n* tuples  *S*

---

# BLOCK NESTED LOOPS JOIN (BNLJ)

Can we do even better?

We only use three page of memory for PNLJ

... (one buffer for *R*, one buffer for *S*, one output buffer),

... but we usually have more memory!

Instead of fetching one page of *R* at a time,

.. why not fetch as many pages of *R* as we can fit (*B - 2* pages)!

---

# BLOCK NESTED LOOPS JOIN (PNLJ)

PNLJ

```
foreach page P_R ∈ R:
  foreach page P_S ∈ S:
    foreach tuple r ∈ P_R:
      foreach tuple s ∈ P_S:
        if θ(r,s): output r joined with s
```

# BLOCK NESTED LOOPS JOIN (PNLJ)

BNLJ

```
foreach block C_R of B-2 pages ∈ R:
  foreach page P_S ∈ S:
    foreach tuple r ∈ C_R:
      foreach tuple s ∈ P_S:
        if θ(r,s): output r joined with s
```
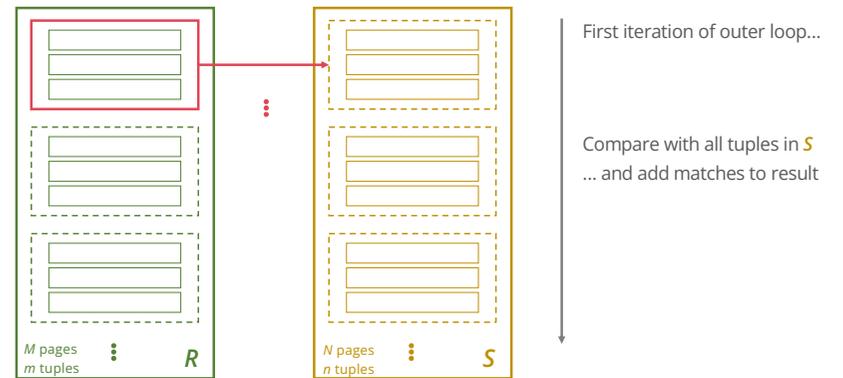
109



$B = 4$

First iteration of outer loop…

Compare with all tuples in $S$
… and add matches to result

*M* pages
*m* tuples
*R*

*N* pages
*n* tuples
*S*

110



$B = 4$

First iteration of outer loop…

Compare with all tuples in $S$
… and add matches to result

*M* pages
*m* tuples
*R*

*N* pages
*n* tuples
*S*

111



$B = 4$

First iteration of outer loop…

Compare with all tuples in $S$
… and add matches to result

*M* pages
*m* tuples
*R*

*N* pages
*n* tuples
*S*

112

**113**



**114**



**115**



**116**

## Slide 117

$B = 4$

$M$ pages
$m$ tuples
$R$

$N$ pages
$n$ tuples
$S$

Second iteration of outer loop…

Compare with all tuples in $S$
… and add matches to result

## Slide 118

$B = 4$

$M$ pages
$m$ tuples
$R$

$N$ pages
$n$ tuples
$S$

Cost

$= M + (\text{\# blocks in } R) \cdot N$

$= M + \lceil M / \text{chunk size} \rceil \cdot N$

$= M + \lceil M / B - 2 \rceil \cdot N$

## Slide 119

# INDEX NESTED LOOPS JOIN (INLJ)

A join is essentially

```
foreach tuple r ∈ R:
  foreach tuple s ∈ S that satisfies θ(r,s):
    output r joined with s
```

## Slide 120

# INDEX NESTED LOOPS JOIN (INLJ)

An **index** on $S$ allows us to do the inner loop efficiently!

```
foreach tuple r ∈ R:
  foreach tuple s ∈ S that satisfies θ(r,s):
  (found using the index)
    output r joined with s
```

# INDEX NESTED LOOPS JOIN (INLJ)

Cost = $M$ + $m$ * cost to find matching $S$ tuples

$M$ from scanning through $R$

Cost to find matching $S$ tuples via **tree index**

Variant **A**: cost to traverse root to leaf + read all the leaves with matching tuples

Variants **B** or **C**: cost of retrieving RIDs (similar to Variant **A**) + cost to fetch actual records

1 I/O per **page** if clustered, 1 I/O per **matching tuple** if not

# INDEX NESTED LOOPS JOIN (INLJ)

Cost = $M$ + $m$ * cost to find matching $S$ tuples

$M$ from scanning through $R$

Cost to find matching $S$ tuples via **hash index**

1-2 I/Os to reach the target bucket

Then scan pages in that bucket

May stop early if search key values are unique and we found a match!

# INDEX NESTED LOOPS JOIN (INLJ)

Cost = $M$ + $m$ * cost to find matching $S$ tuples

$M$ from scanning through $R$

If we have **no index**

Then the only way to search for matching $S$ tuples is by scanning all of $S$ ⇒ SNLJ

Cost to find matching $S$ tuples is then $N$, giving us the formula for SNLJ cost

# INDEX NESTED LOOPS JOIN (INLJ)

# INDEX NESTED LOOPS JOIN (INLJ)

| 43 | ... | Match
| 5 | ... |
| 11 | ... |

R

Output

# INDEX NESTED LOOPS JOIN (INLJ)

| 43 | ... | Match
| 5 | ... |
| 11 | ... |

R

Output

125

126

# INDEX NESTED LOOPS JOIN (INLJ)

| 43 | ... |
| 5 | ... | Not a match
| 11 | ... |

R

Output

# INDEX NESTED LOOPS JOIN (INLJ)

| 43 | ... |
| 5 | ... |
| 11 | ... | Match

R

Output

127

128

# INDEX NESTED LOOPS JOIN (INLJ)



Match

*R*

*Output*

---

# SORT MERGE JOIN (SMJ)

What if we process the data a bit before we join things together?

For example, sort both relations first! Then we can join them efficiently

In some cases, we might even have one of the relations already sorted on the right key, and then we don't even have to spend time sorting it!

---

# SORT MERGE JOIN (SMJ)

First step: **sort** both *R* and *S* (with external sorting)

Second step: **merge** matching tuples from *R* and *S* together

We do this efficiently by moving iterators over sorted *R* and sorted *S* in lockstep: move the iterator with the smaller key

We know that this key is smaller than all remaining key values in the other relation, so we're completely done joining that tuple!

---

# SORT MERGE JOIN (SMJ)

First step: **sort** both *R* and *S* (with external sorting)

Second step: **merge** matching tuples from *R* and *S* together

Need a bit more care than this:
we might have multiple tuples in *R* matching with multiple tuples in *S*

**Mark** the first matching tuple in *S*

Match tuples with the first matching tuple in *R*,

Then **reset** the iterator to the mark

… so we can go through the tuples in *S* again for the second matching tuple in *R*

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
| --- | --- |
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
| --- | --- |
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

133

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
| --- | --- |
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
| --- | --- |
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

134

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
| --- | --- |
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
| --- | --- |
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

135

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
| --- | --- |
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
| --- | --- |
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

136

## SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
  advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |

137

## SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
  advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |

138

## SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
  advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

139

## SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
  advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

140

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

141

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

142

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

143

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|---|---|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|---|---|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|---|---|---|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

144

## Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

145

## Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |

146

## Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |

147

## Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |

148

## Slide 149

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |

149

## Slide 150

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |

150

## Slide 151

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |

151

## Slide 152

# Sort Merge Join (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|----------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |

152

# SORT MERGE JOIN (SMJ)

```
while not done:
    while (r < s): advance r
    while (r > s): advance s

    mark s    // same start of "block"
    while (r == s):
        // outer loop over r
        while (r == s):
            // inner loop over s
            output r joined with s
            advance s
        reset s to mark
        advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |

153

# SORT MERGE JOIN (SMJ)

```
while not done:
    while (r < s): advance r
    while (r > s): advance s

    mark s    // same start of "block"
    while (r == s):
        // outer loop over r
        while (r == s):
            // inner loop over s
            output r joined with s
            advance s
        reset s to mark
        advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |

154

# SORT MERGE JOIN (SMJ)

```
while not done:
    while (r < s): advance r
    while (r > s): advance s

    mark s    // same start of "block"
    while (r == s):
        // outer loop over r
        while (r == s):
            // inner loop over s
            output r joined with s
            advance s
        reset s to mark
        advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

155

# SORT MERGE JOIN (SMJ)

```
while not done:
    while (r < s): advance r
    while (r > s): advance s

    mark s    // same start of "block"
    while (r == s):
        // outer loop over r
        while (r == s):
            // inner loop over s
            output r joined with s
            advance s
        reset s to mark
        advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

156

## Slide 157

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

157

## Slide 158

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

158

## Slide 159

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

159

## Slide 160

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

160

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

161

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

162

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

163

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s   // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

164

# SORT MERGE JOIN (SMJ)

```
while not done:
  while (r < s): advance r
  while (r > s): advance s

  mark s    // same start of "block"
  while (r == s):
    // outer loop over r
    while (r == s):
      // inner loop over s
      output r joined with s
      advance s
    reset s to mark
    advance r
```

| sid | sname |
|-----|-------|
| 22 | John |
| 28 | Smith |
| 31 | Gold |
| 31 | Alvaro |
| 44 | McDonald |
| 57 | Gupta |

| sid | bid |
|-----|-----|
| 28 | 103 |
| 28 | 104 |
| 31 | 101 |
| 31 | 102 |
| 42 | 142 |
| 58 | 107 |

| sid | sname | bid |
|-----|-------|-----|
| 28 | Smith | 103 |
| 28 | Smith | 104 |
| 31 | Gold | 101 |
| 31 | Gold | 102 |
| 31 | Alvaro | 101 |
| 31 | Alvaro | 102 |

165

---

# SORT MERGE JOIN (SMJ)

Total cost:

Cost of sorting $R$

Cost of sorting $S$

Cost of merging: $M + N$

Only one pass (if we assume there aren't a lot of duplicates)

166

---

# SORT MERGE JOIN (SMJ)

We can be sometimes smarter about SMJ

**Observation:**

We make **one pass** through each sorted relation (assuming no duplicate values in $R$)

⇒ We do **not need** the sorted relations to be materialized!

**Optimisation:**

In the final merge pass of sorting both relations, instead of writing the sorted relations to disk, we can stream them into the second part of SMJ!

167

---

# SORT MERGE JOIN (SMJ)



168

# SORT MERGE JOIN (SMJ)



Sort Phase         Merge Phase

We have to be able to fit the input buffers of the last merge pass of sorting *R* and sorting *S* in memory, as well as have one output buffer for joined tuples

Need: (# runs in last merge pass for *R*) + (# runs in last merge pass for *S*) ≤ *B − 1*

Reduces I/O cost by $2 \cdot (M + N)$!

---

# GRACE HASH JOIN

Similar idea as SMJ, but let's build some hash tables instead

Two passes: partition the data, then build a hash table and probe it

Partition *R* and *S* into *B - 1* partitions (like in external hashing) using same hash function

All the tuples in *R* matching a tuple in *S* must be in the same partition

⇒ We can consider each partition independently

---

# GRACE HASH JOIN

Similar idea as SMJ, but let's build some hash tables instead

Two passes: partition the data, then build a hash table and probe it

Then, build an in-memory hash table for a partition of *R*

We can use this in-memory hash table to find all the tuples in *R* that match a tuple in *S*

Stream in tuples of *S*, probe the hash table, output matching tuples

---

# GRACE HASH JOIN: PARTITION

## GRACE HASH JOIN: PARTITION



## GRACE HASH JOIN: PARTITION



## GRACE HASH JOIN: PARTITION



## GRACE HASH JOIN: PARTITION

# GRACE HASH JOIN: PARTITION



181

# GRACE HASH JOIN: PARTITION



182

# GRACE HASH JOIN: PARTITION



183

# GRACE HASH JOIN: PARTITION



184

## GRACE HASH JOIN

We need partitions of *R* (but not *S*!) to fit in *B - 2* pages

    1 page reserved for streaming *S* partition

    1 page reserved for streaming output

What if partitions of *R* are too big?

    If *S* is smaller, do $S \bowtie_\theta R$ instead

    Recursively partition! Make sure that for any partition of *R* you recursively partition, the matching *S* partition is also recursively partitioned!

185

## GRACE HASH JOIN: BUILD & PROBE

186

## GRACE HASH JOIN: BUILD & PROBE

187

## GRACE HASH JOIN: BUILD & PROBE

2 matches

188

189



190



191



192

GRACE HASH JOIN: **BUILD** & PROBE

Hash table *B-2* buffers

In buffer   Out buffer

3 matches

GRACE HASH JOIN: BUILD & **PROBE**

Hash table *B-2* buffers

In buffer   Out buffer

GRACE HASH JOIN: BUILD & PROBE

Hash table *B-2* buffers

In buffer   Out buffer

No match

GRACE HASH JOIN: BUILD & PROBE

Hash table *B-2* buffers

In buffer   Out buffer

No match

193

194

195

196

## GRACE HASH JOIN: BUILD & PROBE

**197**

Hash table *B-2* buffers

In buffer    Out buffer

1 match

---

## GRACE HASH JOIN: BUILD & PROBE

**198**

Hash table *B-2* buffers

In buffer    Out buffer

---

## GRACE HASH JOIN: BUILD & PROBE

**199**

Hash table *B-2* buffers

In buffer    Out buffer

2 matches

---

## GRACE HASH JOIN: BUILD & PROBE

**200**

Hash table *B-2* buffers

...

In buffer    Out buffer