

Recap

problem: Retrieving data efficiently

solution: Layout Strategies

- delimited

- fixed size

- directory

X

- For ~~Records~~ fields in a record

- For records on a page

problem: Retrieving specific data efficiently

solution: 1) Sort data

2) Sort + Tree index

3) Secondary Index

4) Hash table

problem: Sorting data

solution: ?

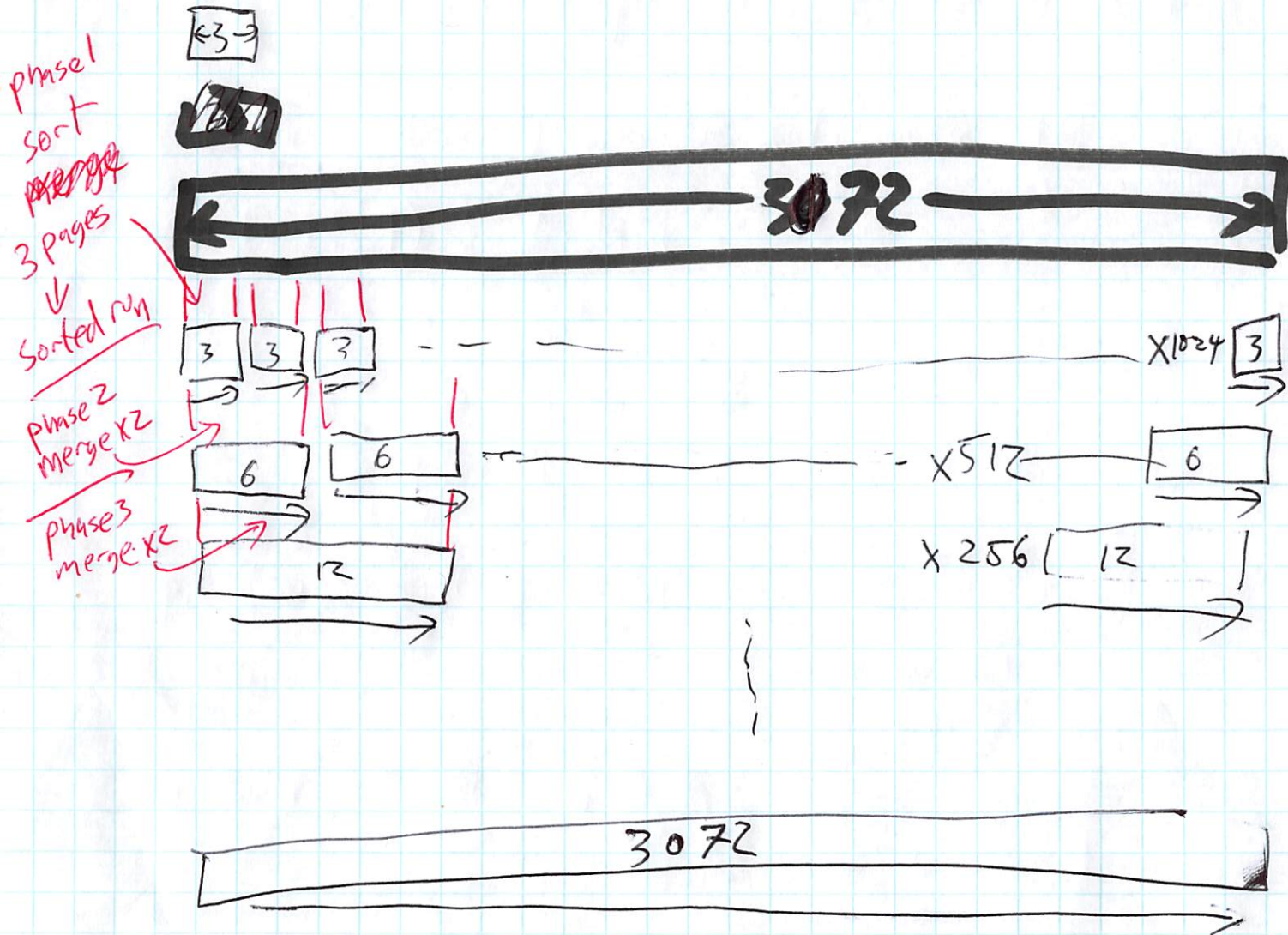
problem: Summarizing Data

solution: ?

with
limited memory

3072 pages of data

3 pages of memory



can be as big as mem
create runs of size 3

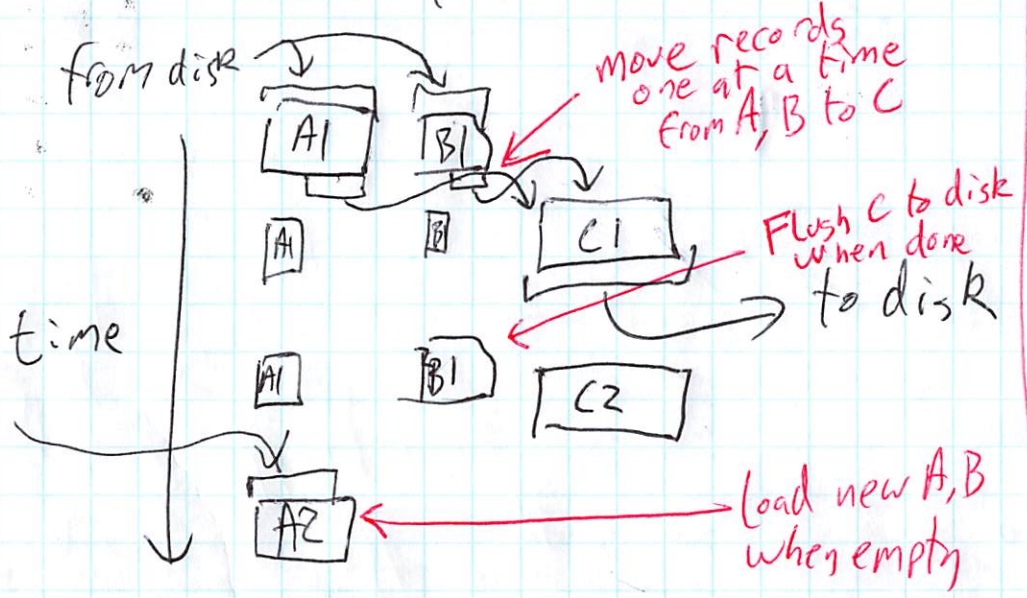
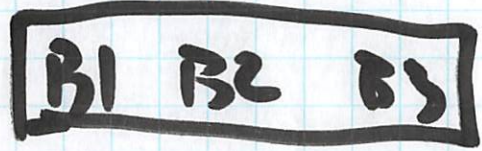
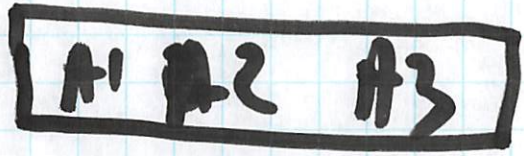
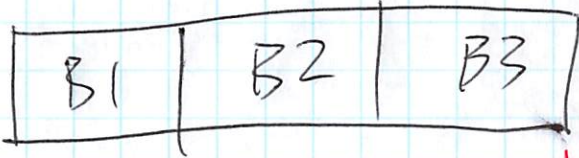
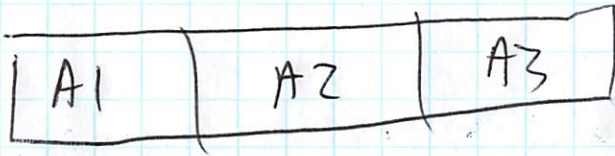
Phase 1

Phase 2 merge

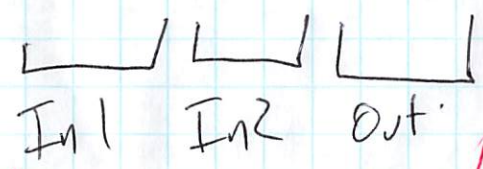
given k →
Phase 3 merge

merge k-1 sorted runs together

Phase II merge



With $R+1$ pages
merge k at a time



only 3 pages required

Cost analysis

Phase 1 : $2^{\# \text{pages}}$ IOs

Phases 2-11 : $2^{\# \text{pages}}$ IOs

After phase

1
2
3
N

this big sorted runs

3
6
12
 $3 \cdot 2^{N-1}$ }

$$\# \text{pages} \leq 3 \cdot 2^{N-1}$$

$$\frac{\# \text{pages}}{3} \leq 2^{N-1}$$

$$\log_2 \left(\frac{\# \text{pages}}{3} \right) \leq N-1$$

$$N \geq 1 + \log \left(\frac{\# \text{pages}}{3} \right)$$

$$\frac{\text{Total cost}}{(1 + \log \left(\frac{\# \text{pages}}{3} \right)) \cdot \# \text{pages}} \cdot 2$$

IOs

have k pages of mem

After phase

1

2

3

N

this big sorted runs

k

$k \cdot (k-1)$

$k \cdot (k-1)^2$

$k \cdot (k-1)^{N-1}$

$$\# \text{pages} \leq k \cdot (k-1)^{N-1}$$

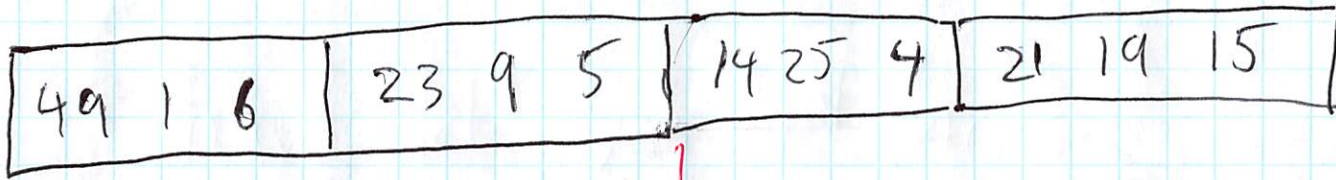
$$N \geq 1 + \log_{(k-1)} \left(\frac{\# \text{pages}}{k} \right)$$



$$\left(1 + \log_{(k-1)} \left(\frac{\# \text{pages}}{k} \right) \right) \cdot \# \text{pages} \cdot 2$$

IOs

$$\approx \log_k (\# \text{pages}) \cdot \# \text{pages} \cdot 2$$



#1 Load R-1 pages then sort



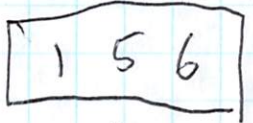
#2

move lowest records to output page



#3

Load next data page in



#3

flush output page to disk (save highest record)

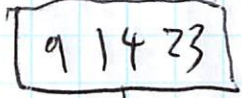
highest record 6



greater than 6

#5

find re-sort and copy lowest records greater than highest record in output page to continue sorted run



highest = 23

#6 keep flushing output to disk



#7 when can not create any more output pages, start a new sorted run

Fold (default) (merge (current_{value}, record_{value}))

Sum = ~~fold~~ fold(0) (current + record)

<u>record</u>		<u>current (sum)</u>	<u>current (count)</u>
1	+	0 $\downarrow +1$	0 $\downarrow +1$
10	+	1 $\downarrow +10$	1 $\downarrow +1$
19	+	11 $\downarrow +19$	2 $\downarrow +1$
13	+	30 $\downarrow +13$	3 $\downarrow +1$
4	+	43 $\downarrow +4$	4 $\downarrow +1$
2	+	47 $\downarrow +2$	5 $\downarrow +1$
		49	6 $\downarrow +1$

$$\text{Count} = \text{fold}(0)(\text{current} + 1)$$

$$\text{MIN} = \text{fold}(+\infty)(\min(\text{current}, \text{record}))$$

$$\text{MAX} = \text{fold}(-\infty)(\max(\text{current}, \text{record}))$$

$$\text{Average} = \text{fold}(\langle \text{count}: 0, \text{sum}: 0 \rangle) \left(\langle \text{count}: \text{current.count} + 1, \text{sum}: \text{current.sum} + \text{record} \rangle \right)$$

with $\frac{\text{current.sum}}{\text{current.count}}$

1 10 19 13 4 2 9

Stdder
~~variance~~
Stdder

$$\text{fold}(\langle \text{sum}: 0, \text{sumsq}: 0 \rangle) \left(\begin{array}{l} \text{sum} : \text{current.sum} + \text{record} \\ \text{sumsq} : \text{current.sumsq} + \text{record}^2 \\ \text{count} : \end{array} \right)$$

with $\frac{(\text{current.sum})^2 - (\text{current.sumsq})}{N^2}$ $R \rightarrow \frac{\left(\left(\frac{\sum x}{N} \right)^2 - \frac{\sum x^2}{N} \right)}{N^2}$

Aggregate

SUM

COUNT

AVERAGE

STDDEV

MIN

MAX

MEDIAN

COUNT DISTINCT

Alg. vs Hol.

~~Algebraic~~

"

"

"

"

"

Holistic

"

Intermediate

```
SELECT class, Average(grade)
FROM Grades
GROUP BY class
```

Grades (Student, Grade, class)

* Classes = { }

for record in grades:

classes[record.class] += [record.grade]