

```
$> ls • /lib
      /build
      /data
```

```
/data/[filename].csv
```

```
$> echo "SELECT id * FROM filename;" |
```

```
Java -cp build: jsqlparser dubstep.Main -
```

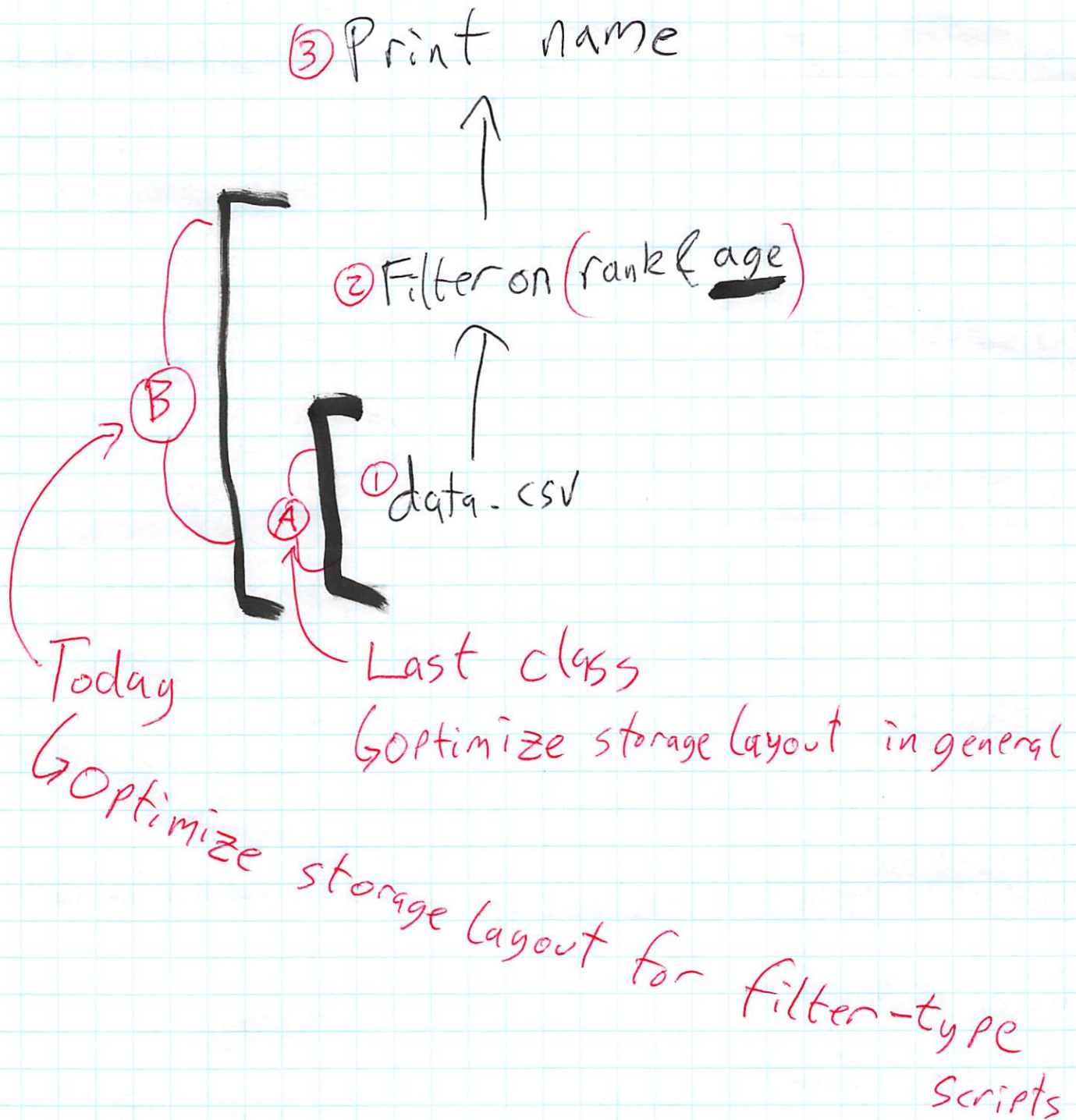
```
red | 23 | 59.265 | pickup
```

```
blue | 47 | 63.921 | dropout
```

```
⋮
```

Expected
output
for
checkpoint

Workflow for Script

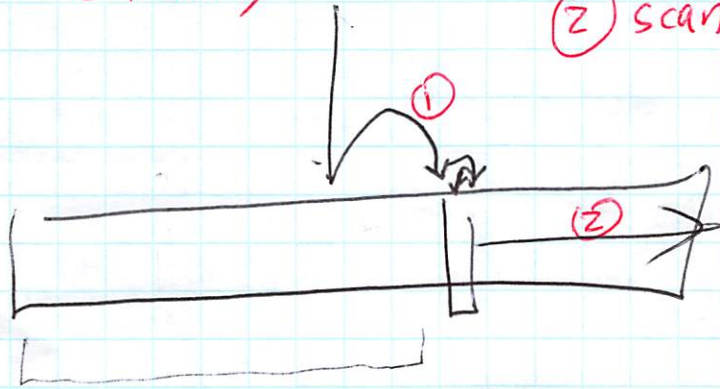


How do we make filter (age > 25) fast?

Idea: Sort data by age!

① binary search to find age = 25

② scan + print all later records



Question

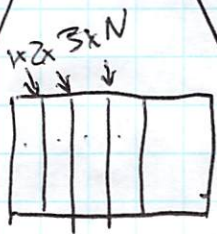
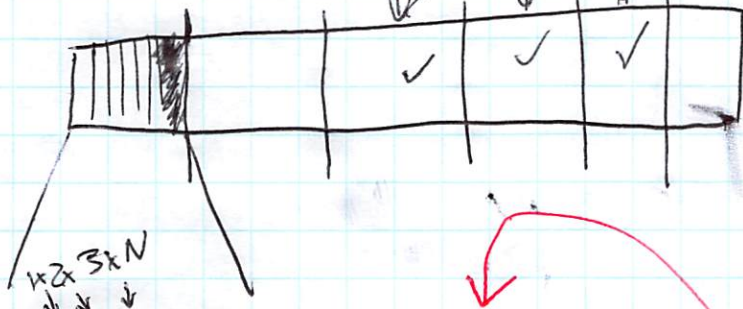
Big blob of bits

Where does the middle record start?

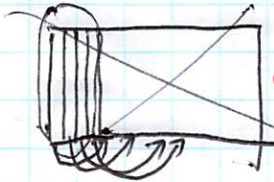
Use fixed length records

Create fixed size "chunks" of records

still need to do binary search within a chunk

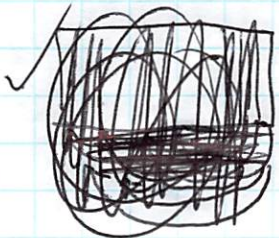


Option 1: fixed size records

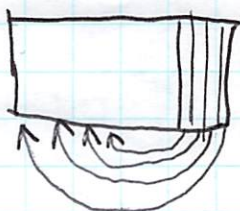


Option 2: Header w directory of record positions

(not generally used)
bonus Q: Why?



~~Option 3: Footer w directory of record positions (actually used... why?)~~

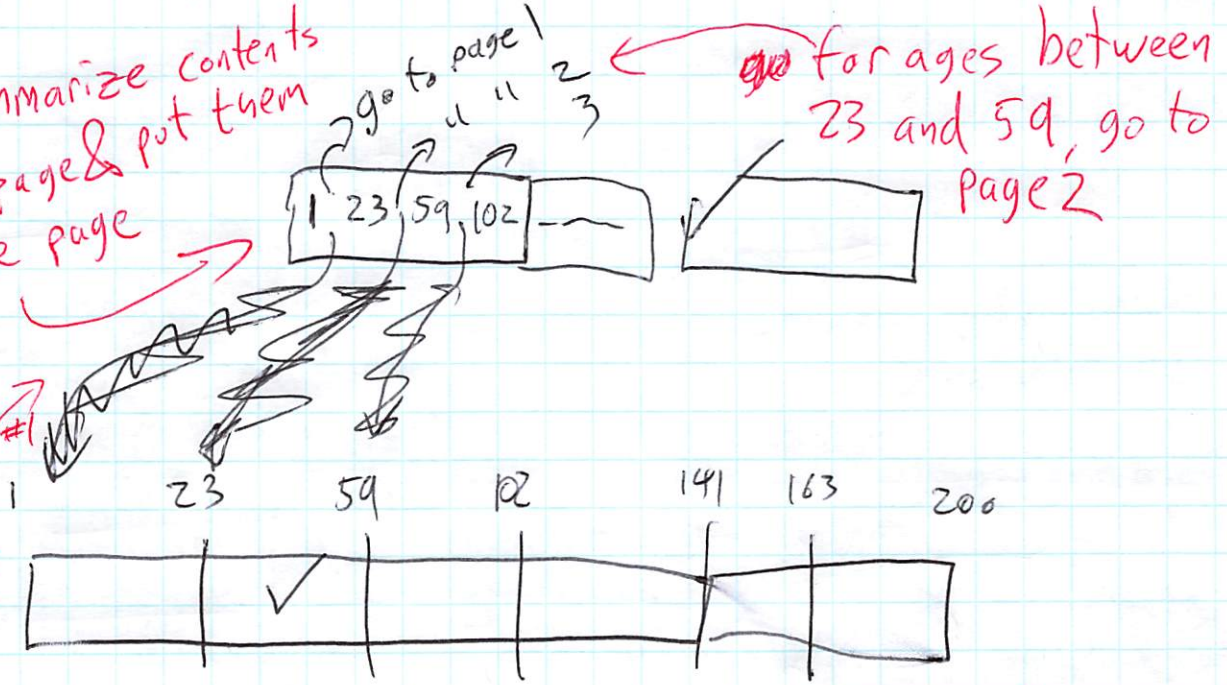


Free Bonus: Chunk = disk page (~4kb)
= cache line (64b)

Problem: Binary search $\rightarrow \log N$ pages loaded

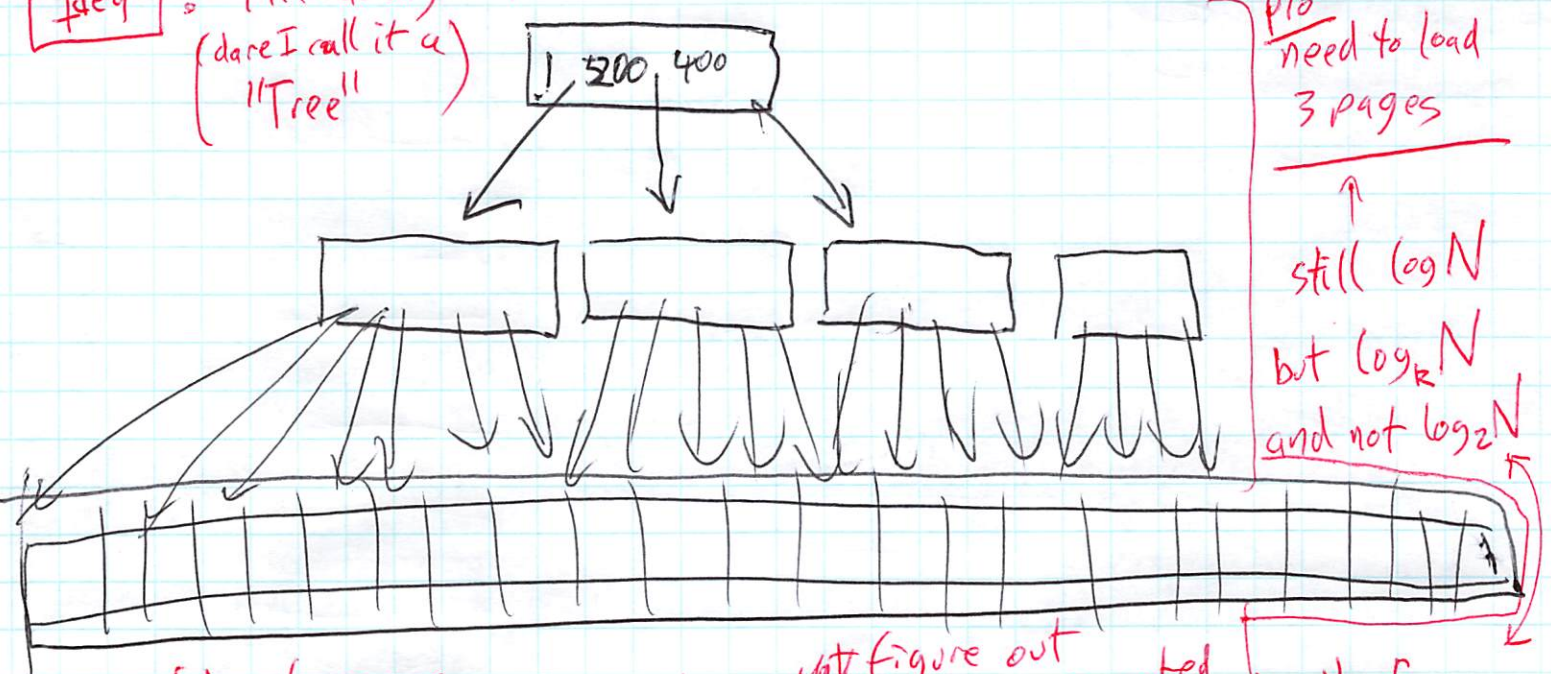
Idea: Summarize contents of each page & put them all on one page

Pro: only need #1 to load 2 pages = #2



Problem: what if you need more than one page of summaries?

Idea: Hierarchy of summaries \rightarrow Summarize the summaries (dare I call it a "Tree")



In the language of DBs

Index: Summaries to help you figure out where your data lives
 ISAM Index: The "Tree Index" we just invented
 Index page: A "summary" page or "chunk"
 Data page: # of page on "chunk" of records

$k = \#$ of page refs stored on one summary page

Generalizations to other filtering conditions

$age = X$

↳ binary search

$age > X$ ($age \geq X$)

↳ binary search
for page 1

↳ scan remaining

$age < X$

↳ scan up to last
record

$X < age < Y$