

# Lecture 21:

Data into  
your program

1

getting data from files

2

map + reduce =

map reduce

↳ Files

↳ Networked data

3

extract-combine (?)

# Combining map and reduce

```
fun map (f, t) =
```

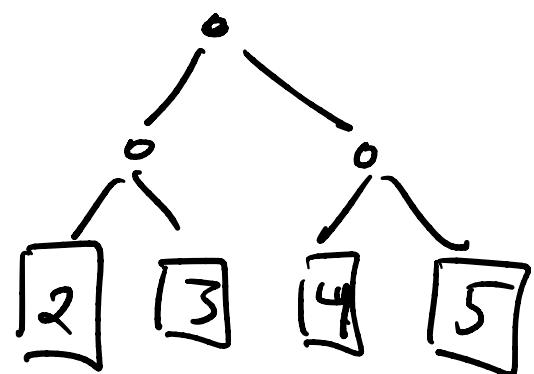
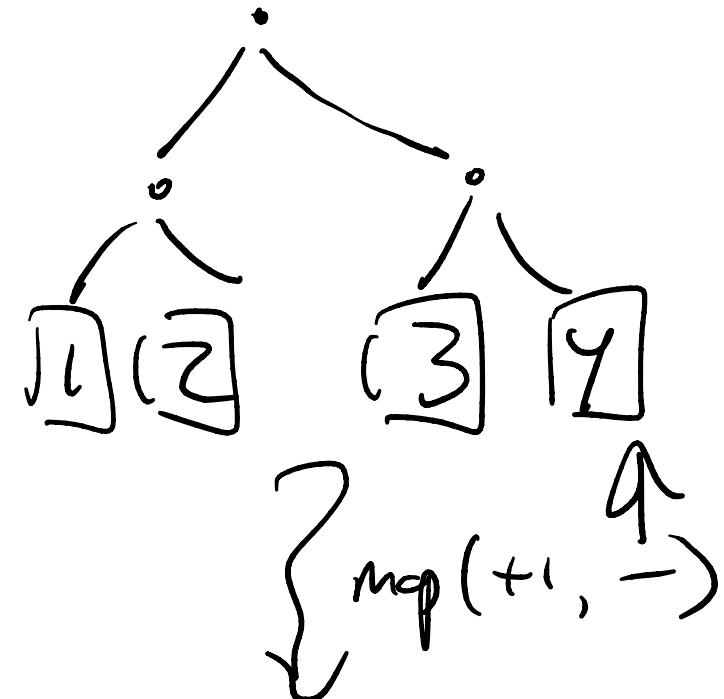
```
case + of
```

Empty  $\rightarrow$  Empty

| Leaf  $x \rightarrow$  Leaf ( $f x$ )

| Node ( $l, r$ )  $\Rightarrow$

Node (map ( $f, l$ ),  
map ( $f, r$ ))



fun reduce( $\text{a}^+, \text{e}^0, \epsilon$ ) =

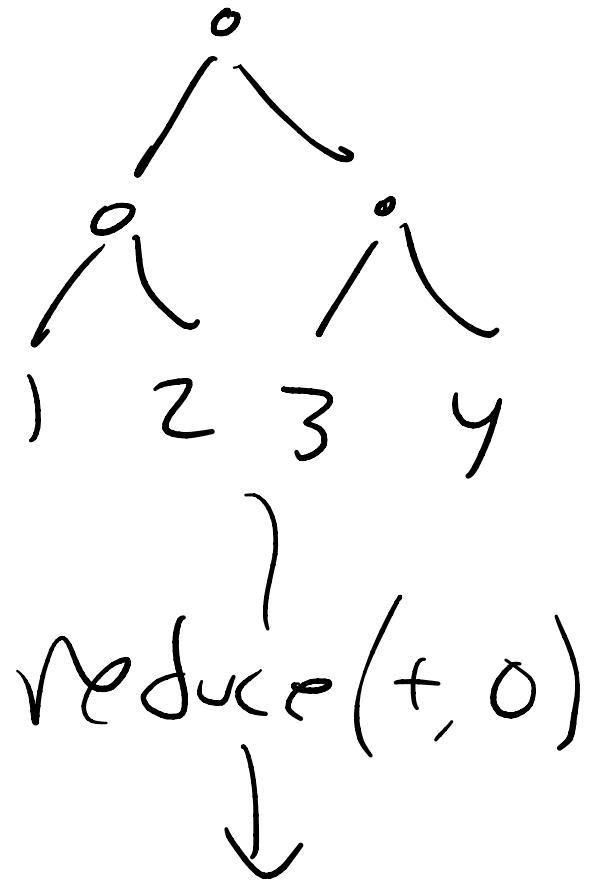
case + of

Empty  $\Rightarrow e$

| Leaf  $x \Rightarrow x$

| Node  $(l, r) \Rightarrow$

$+ n(\text{reduce}(n, e, l),$   
 $\leqslant \text{reduce}(n, e, r))$



$$\begin{aligned} & 1 \underline{+} 2 \underline{+} 3 \times 4 \\ & = 10 \end{aligned}$$

Deforestation =

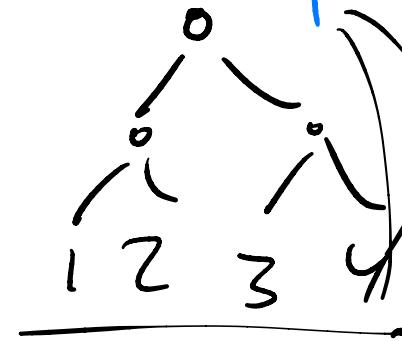
Eliminating intermediate trees

$O(n)$

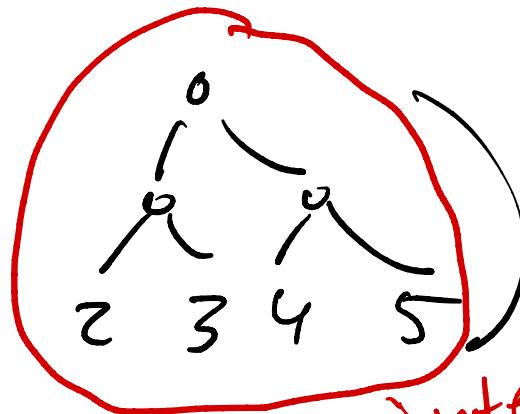
work  $O(n)$

reduce(+, 0,

map( $f(x) \Rightarrow x+1$ ,



$\mapsto$  reduce(+, 0



intermediate

$$\mapsto (2+3)+(4+5)$$

uses  $O(n)$   
Space

$$= 14$$

Spec  $\frac{\text{Id-act}}{\text{mapreduce}(f, e, n, t) = \text{reduce}(n, e, \text{map}(f, t))}$   $\Rightarrow$

fun mapreduce(f, e, n, t) =

(case t of

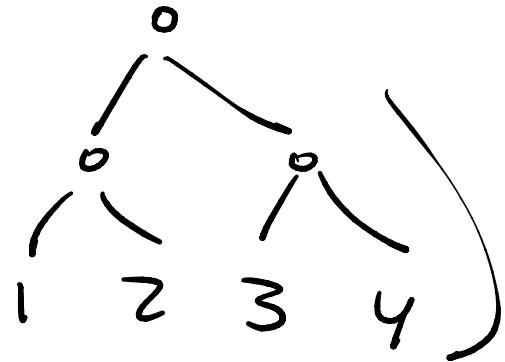
Empty,  $\Rightarrow e$

| Leaf x  $\Rightarrow f(x)$

| Node(l, r)  $\Rightarrow n(\text{mapreduce}(f, e, n, l),$

$\text{mapreduce}(f, e, n, r))$

Map reduce ( $f_n x \Rightarrow x+1$ ,  $O$ ,  $+$ ,



$\mapsto (1+1) + (2+1) + (3+1) + (4+1)$

$\mapsto 2 + 3 + 4 + 5$

$\mapsto 14$

don't  
make  
the  
 $O(n)$   
tree

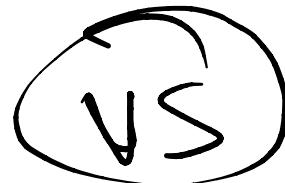
Imagine that data

1  
2  
3  
4  
...

is in  
a file

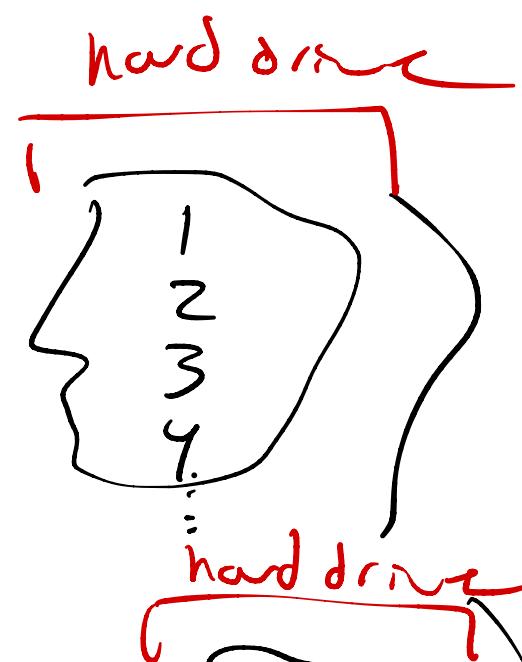
rather than a tree

mapreduce ( +, 0, +,



reduce ( +, 0,

map ( +,



never load  
this tree/scg  
into RAM

Signature MAP-REDUCE =

sig

type 'a mapreducible

Val mapreduce:

$$('a \rightarrow 'b)$$

\* 'b

$$\times ('b * 'b \rightarrow 'b)$$

\* 'a mapreducible

$$\rightarrow 'b$$

f

e

n

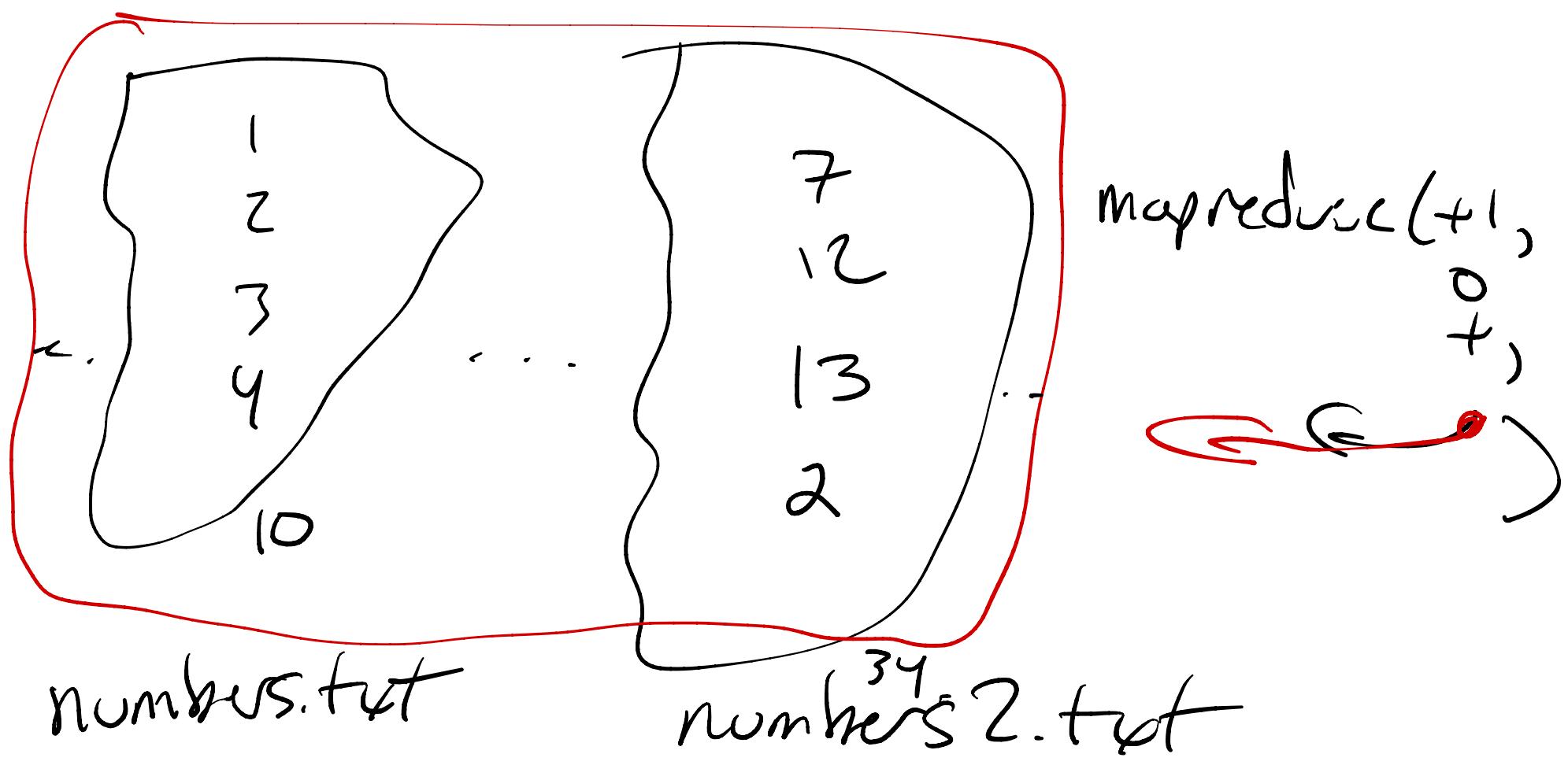
"tree"

end

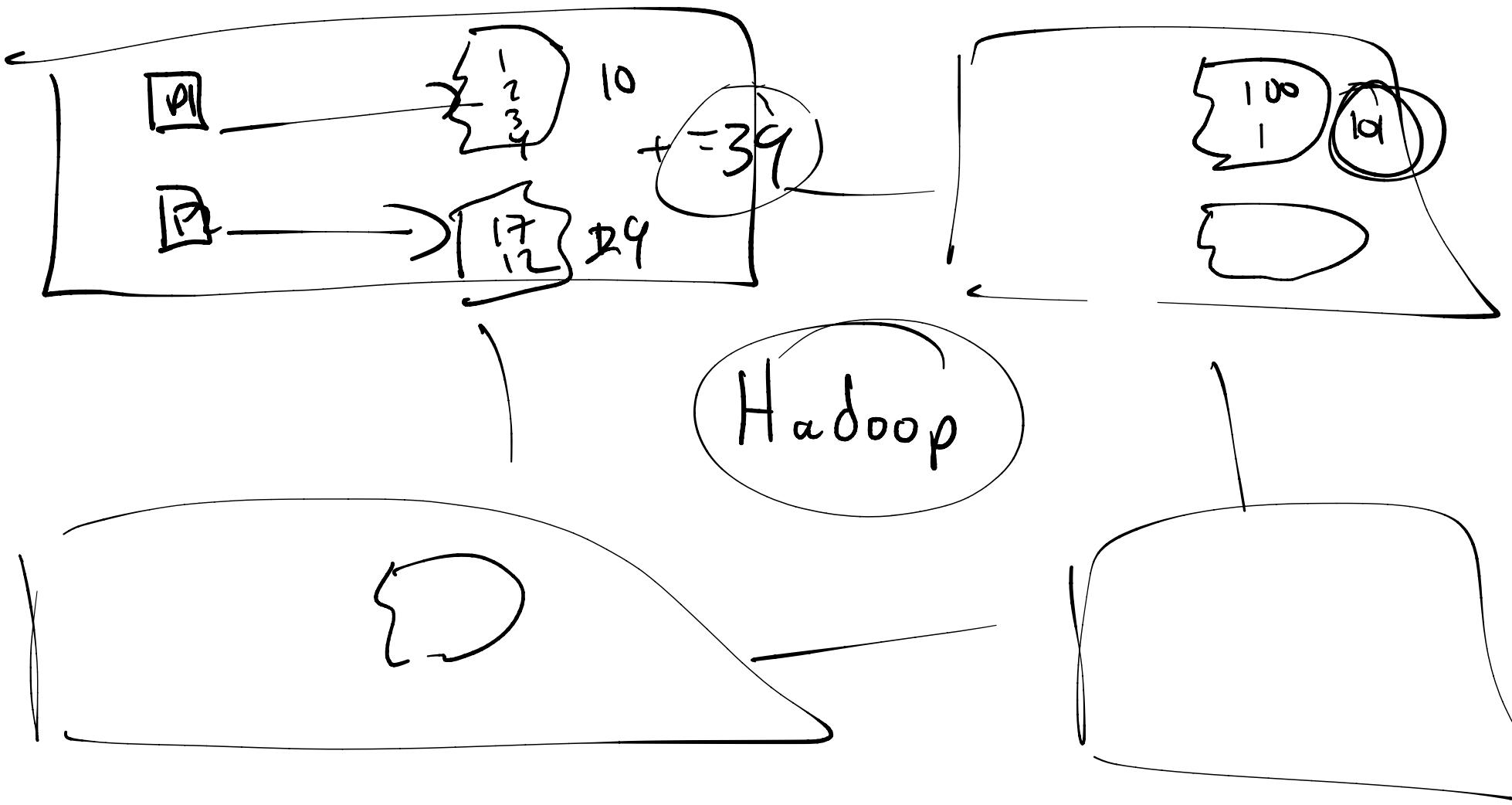
## Examples

- ① Sequences with Mapreduce =  
Reduce after map  
  
not deforested
- ② files with a parsing function
- ③ many files
- ④ many files or many computers

Lab 10



- Idea:
- ① compute the  $\sum$  of each file
  - ② add summaries for each file
- $= 44$



Idea:

- ① each computer does the summary for its files
- ② add/summarize all the computers

frequencies counting via

map/reduce



Naïve Bayes