

Hash functions

- cryptographic hash functions

$$h: \Omega \longrightarrow \{0, \dots, n-1\} = [n]$$

$$\rightarrow i \in [n] ; \Pr[h(x) = i] = \frac{1}{n}$$

$$\{h(x_1), \dots, h(x_k)\}$$



→ "safety"

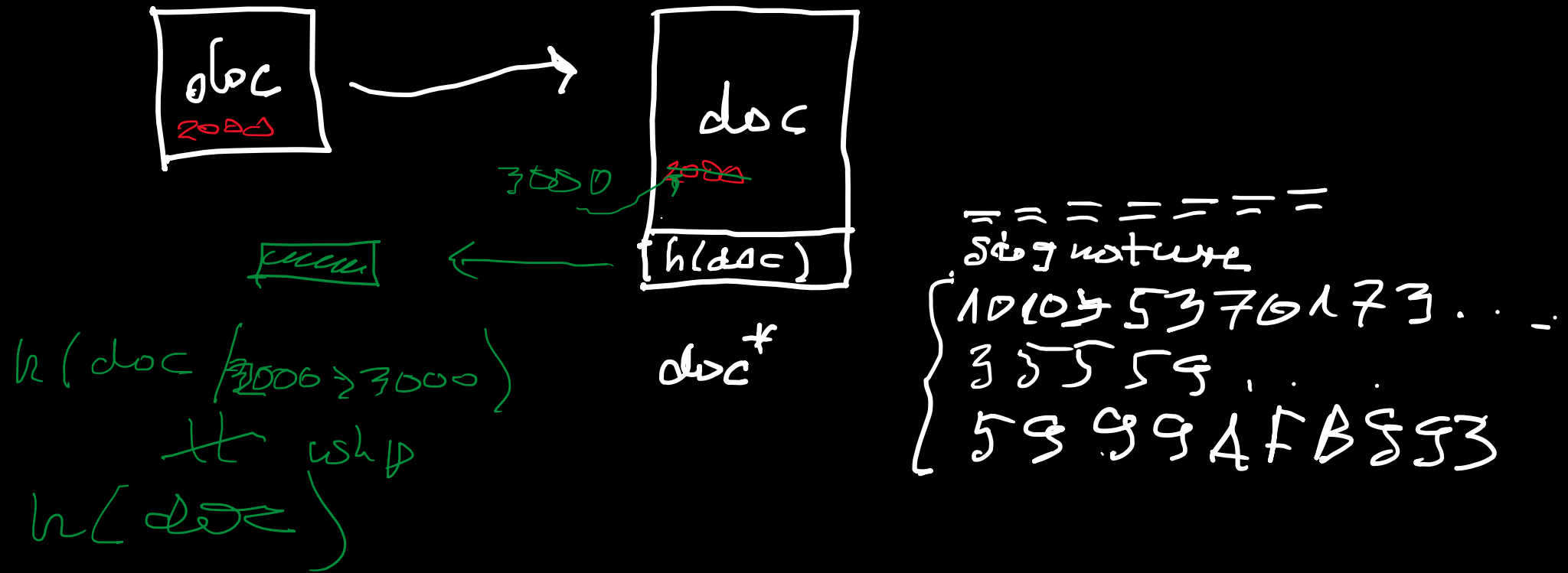
If we know $y = h(x)$

then it is very diff to find x^*

s.t. $y = h(x^*)$

random subset
of $[n]$ with t elem.

Example, Suppose doc is a document (txt)



Typical examples: MD5
SHA-X

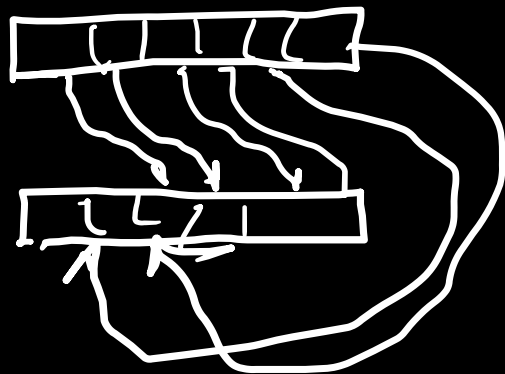
X = 128, 256, 512, 1024

non-crypt. functions:

k-independent for ~~any~~ large k

Popular: Murmur hash function

{ Mu \leftarrow multiplication
 r \leftarrow rotate (wiki)



rot. space.

Scala

```
import scala.util.hashing { MurmurHash3 => MH3 }  
...
```

```
val seed = 195353 // unique
```

```
myHash(data, my seed) {
```

```
  val n = MH3.stringHash(data, my seedseed)
```

```
  n  
}
```

this must be gener.

indop. of the program

independence

```
val seeds = [3535, 35610, 4033]
```

random int values

```
myHashDel(x, i) { myHash(x, seeds[i]) }
```

Suppose $h: \mathbb{Q} \rightarrow [u]$
is good function.

$$\underbrace{h(x_1), h(x_2), \dots, h(x_k)}_{\substack{x_i \neq x_j \\ i \neq j}}$$

\approx random numbers from $\{0, \dots, u-1\}$
with distribution.

$\rightarrow x_1, x_2, x_2, \dots$ rand. numb for $[u]$
indep.

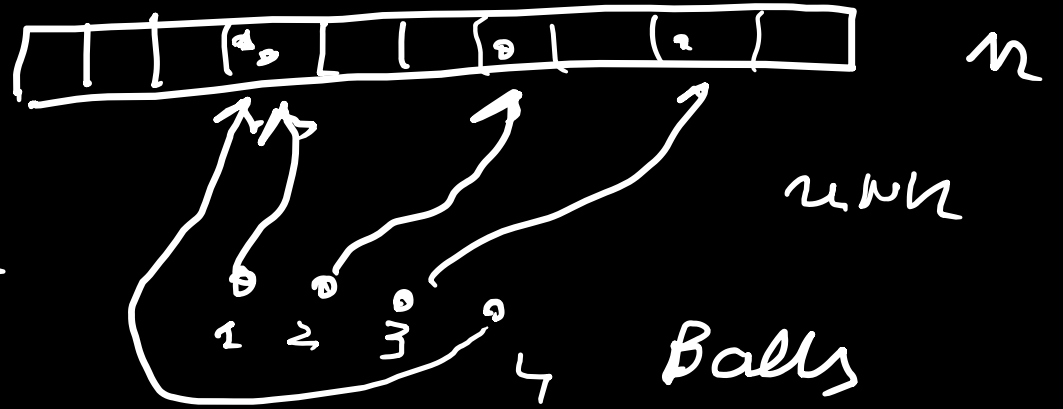
$$B = \min \{k : X_k \in \{X_1, \dots, X_{k-1}\}\}$$

$$1 \leq B \leq n+1$$

$$E[B] = \sqrt{\frac{\pi \cdot n}{2}} + \frac{3}{2} + O\left(\frac{1}{\sqrt{n}}\right)$$

$$\approx \sqrt{n}$$

Birthday Paradox



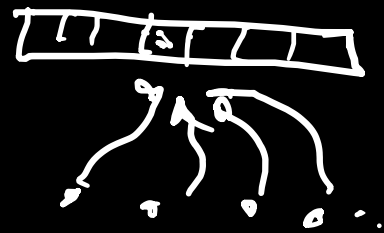
$$\sqrt{\frac{1^2}{2}} \approx \sqrt{\frac{3.14}{2}} \approx \sqrt{1.6}$$

$$\approx 1.3 \approx 1$$

Coupon Collector Problem

$$C = \min \{k : \{X_1, \dots, X_k\} = [n]\}$$

$$P[C = \infty] = 0$$



$$E[C] = n \cdot H_n, \quad H_n = \frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n}$$

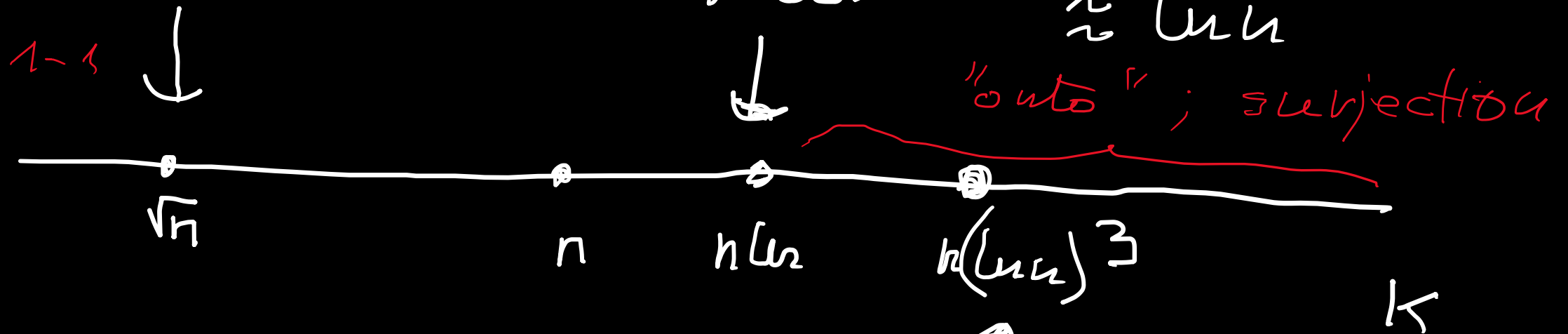
$$\approx n \ln n, \quad H_n = \ln n + \gamma + O\left(\frac{1}{n}\right)$$

BCH. year.

Comp. Coll. Problem

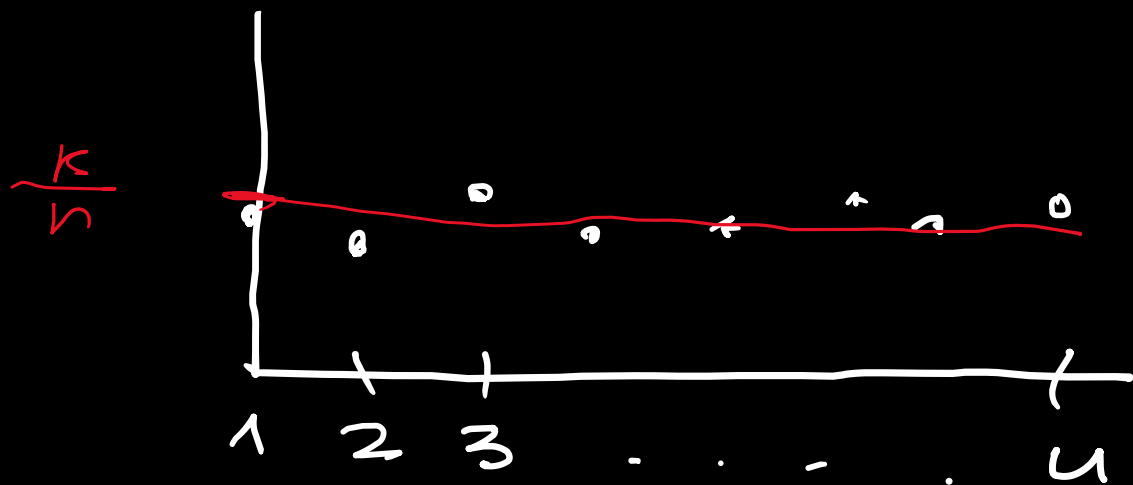
$$\gamma = 0.5771\dots$$

$$\approx \ln n$$



$$x_{11} \dots x_{kk} \ll \sqrt{n}$$

uniform distribution



uniformity of distribution -

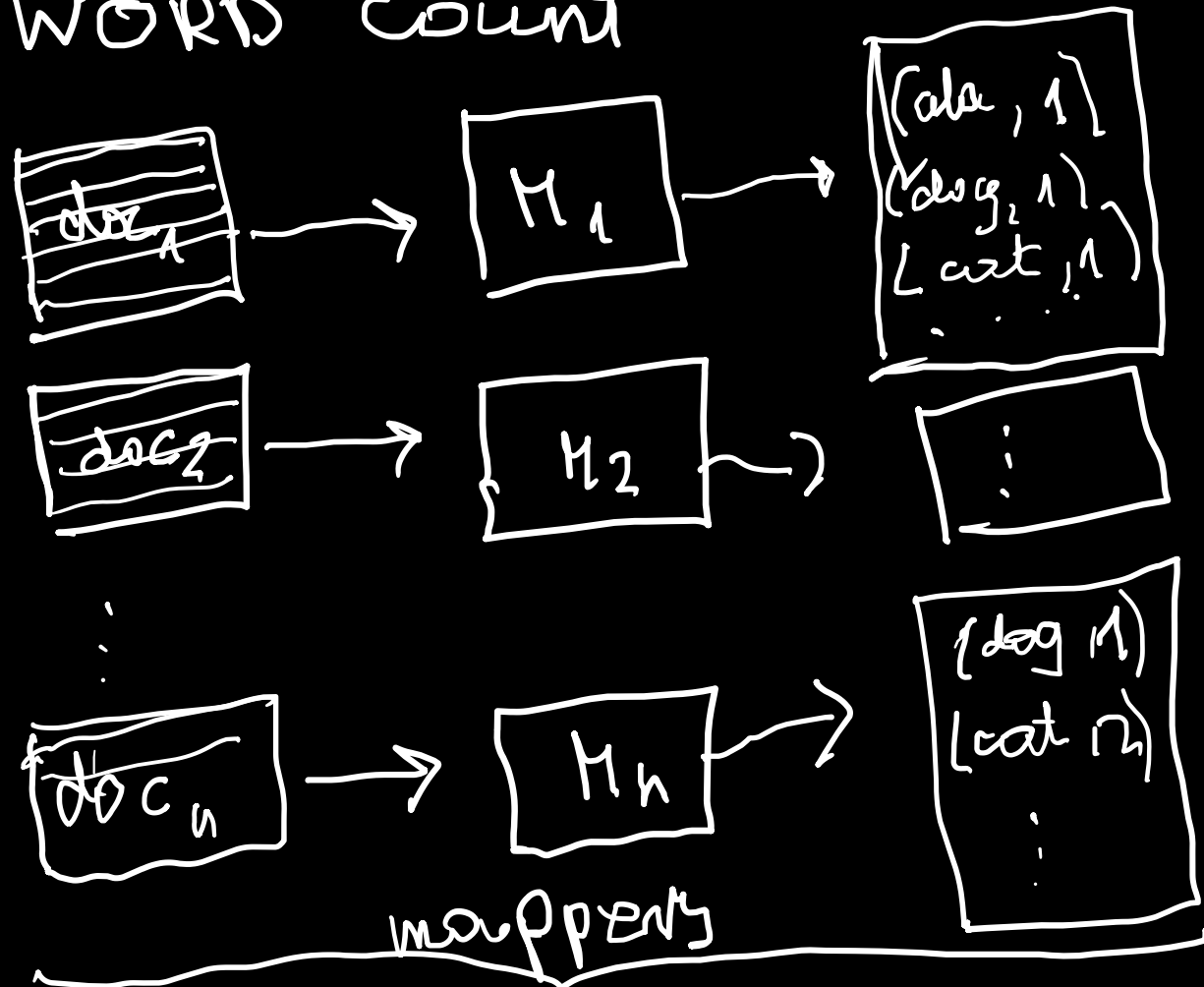
$$k > n(uu)^3$$

m_i = number of balls in i th urn

$$m_1 + \dots + m_k = k$$

mean number of balls in levels : $\frac{k}{u}$

MAP - REDUCE MODEL OF COMPUTATION, WORD COUNT

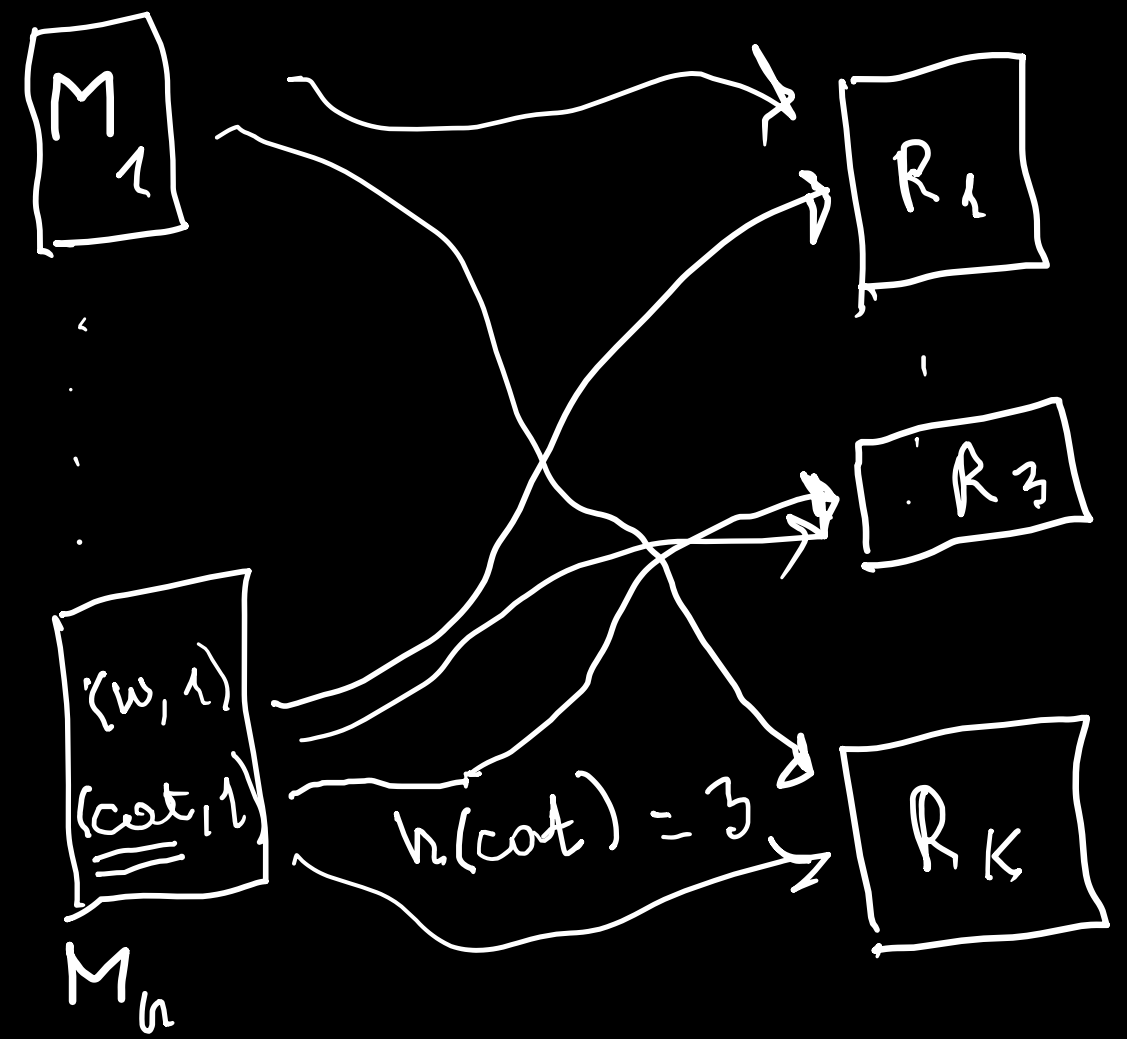


STEP 1

```
map(L) {  
  for each  $w \in L$  do {  
    emit (w, 1)  
  }  
}
```

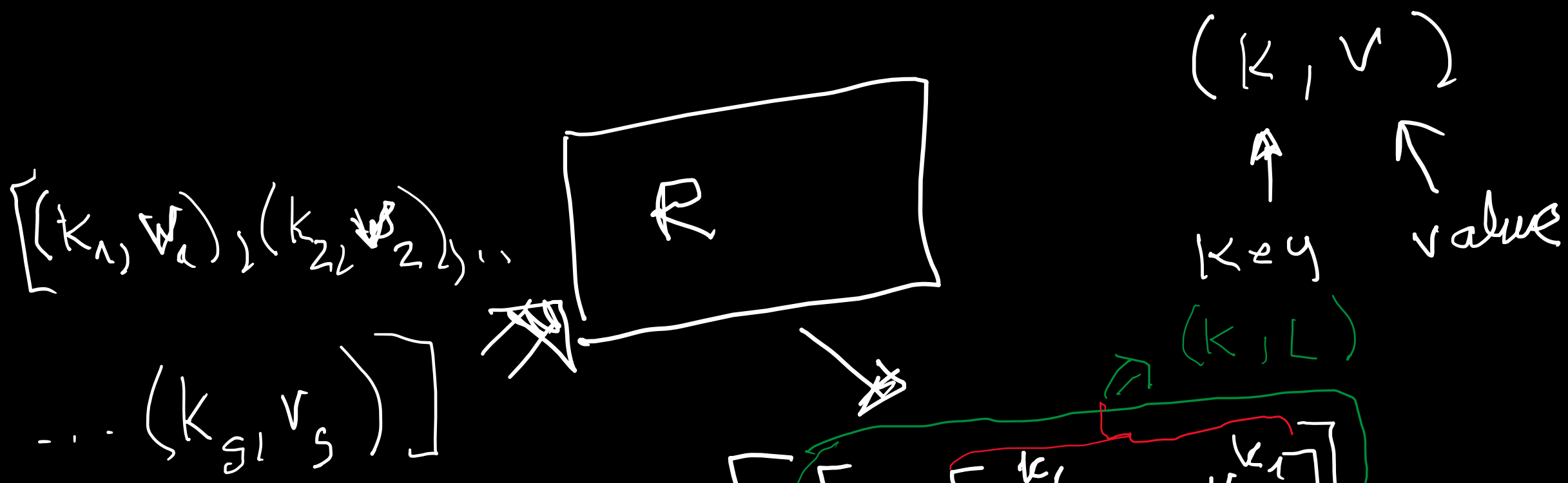
Hadoop
Apache Spark

STEP 2



$$h: \Omega \rightarrow \{1, \dots, K\}$$

fixed hash function



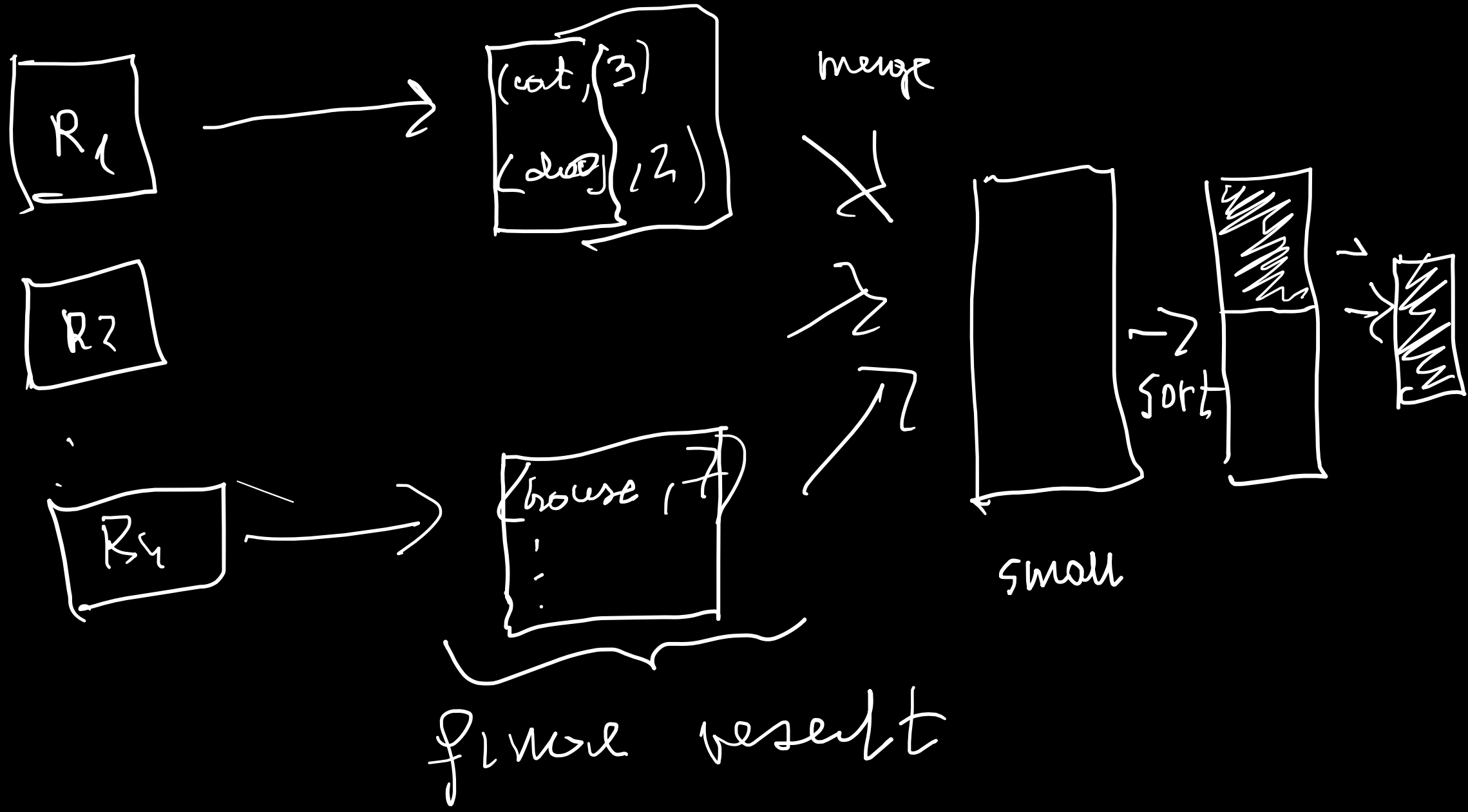
sorting with respect to keys

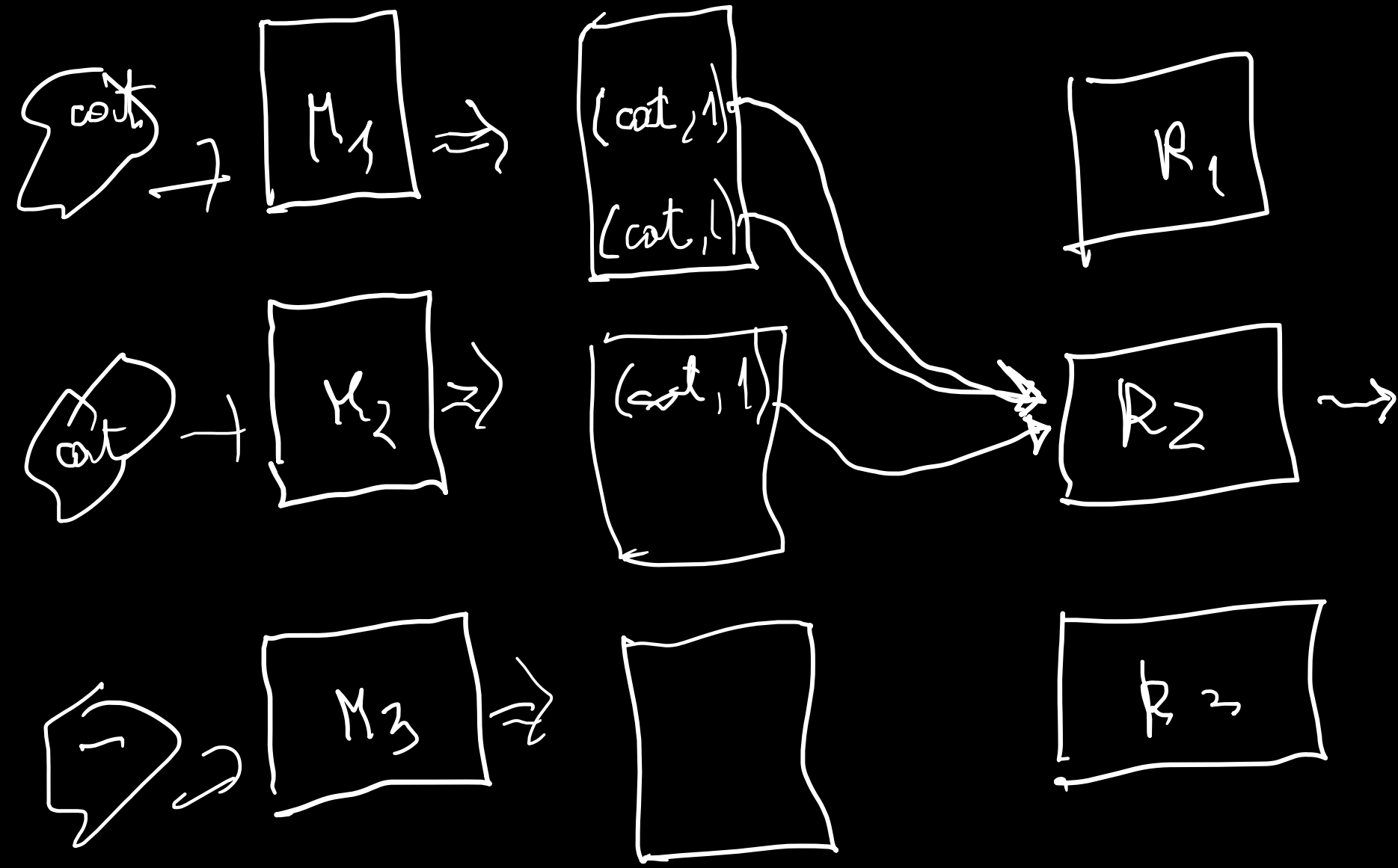
grouping $[(cat, 1), (cat, 1), (cat, 1), (dog, 1), (dog, 1)]$
 $[(cat, [1, 1, 1]), (dog, [1, 1])]$

Programmer's job with reducer

!!
o o
[reduce (key, L) }
emit (key, length(L))
}

(cat, [1, 1, 1])





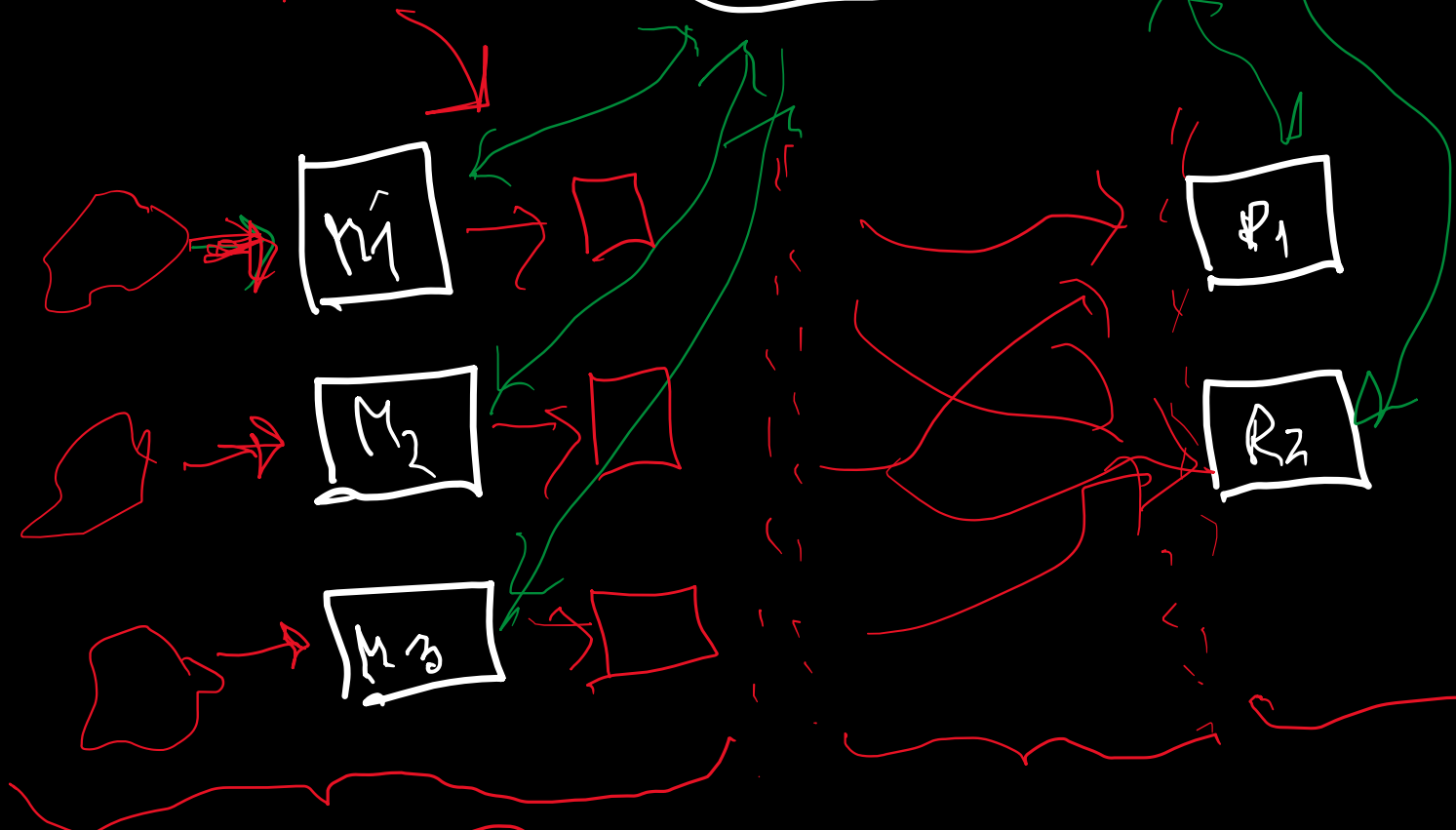
$$i = h(\text{cat}) = 2$$

(cat, 3) !!!
 o n o

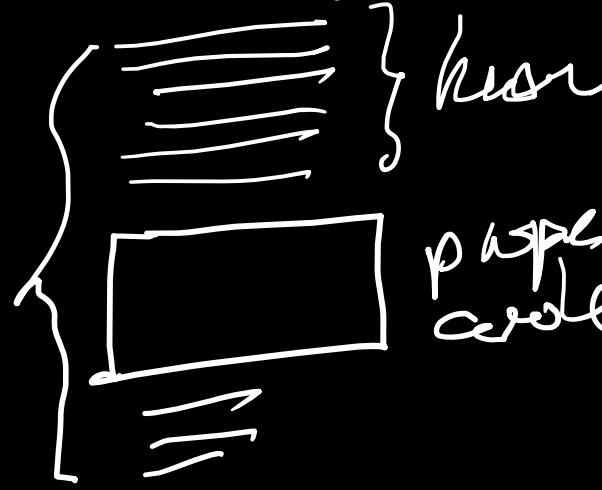
MAP-
 REDUCE

map

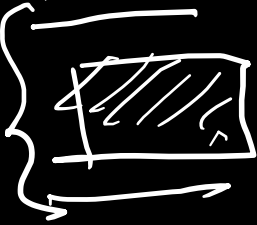
CONTROLLER



Bootstrap



Scala



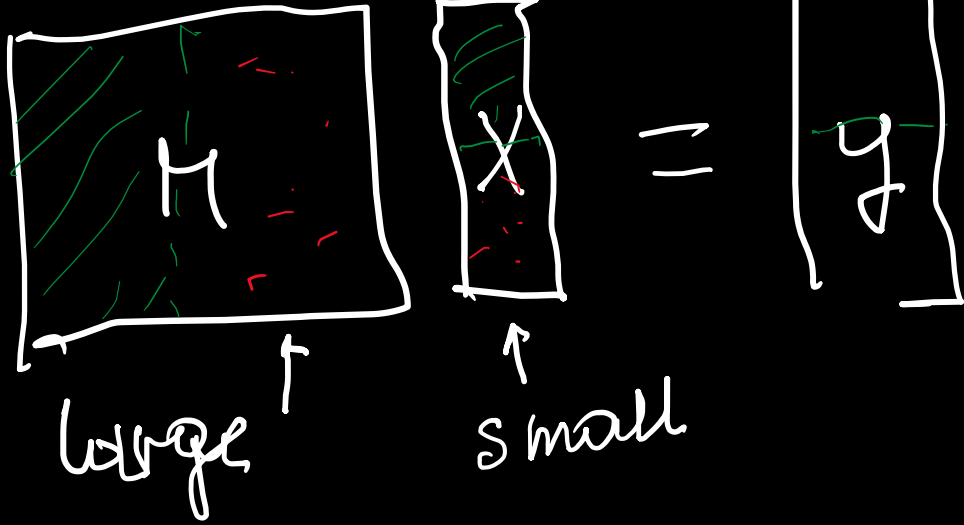
P_1

P_2
framework
shuffle

P_3
sort,
group

job with
redlect

Example



$$y_i = \sum_{j=1}^n m_{ij} \cdot x_j$$

what reducer gets

$$[3, [m_{3,7} \cdot x_7, m_{3,9} \cdot x_9, \dots]]$$

Input: (i, j, m_{ij})

map (i, j, m_{ij}) {
 emit $(i, m_{ij} \cdot x_j)$
 }

reduce (i, L) {
 emit $(i, \sum L)$
 }