

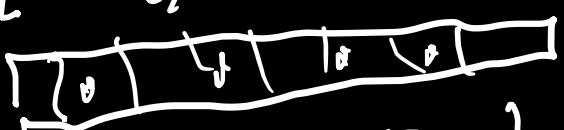
Funkcje hashujące

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

• kryptograficzne

MD5, SHA-...

→ $\{h(x_i) : i=1 \dots n\}$



[n]

funkc. rozdzielca:

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

→ Znanego $y = h(x)$ bardzo
trudno jest znaleźć x ,

• nie-kryptograficzne

MurMur hash

• k - niezależności
dla długości k

$\{M \leftarrow \text{multiply}$
 $\{R \leftarrow \text{rotate}$

MurmurHash

Scala:

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

.....

```
val seed = 1353451; (uint)
```

```
myHash(data) {
```

```
  val m = MH3.stringHash(data, seed)
```

```
  m
```

```
}
```

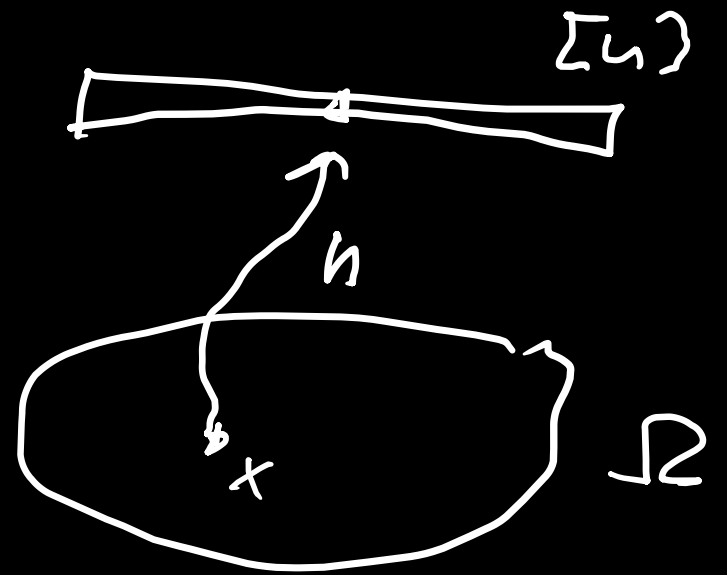
h_1, \dots, h_k : user-defined, fork, hash.

val = $[s_1, s_2, \dots, s_k]$ ← seed's

Define mesh functions,

$h \sim$ losowa niezależna
wybory $z [n]$

- X_1, X_2, X_3, \dots - ciąg niezależnych losowych
o wartości $z [n]$



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

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

$$\approx \sqrt{n}$$

← moment
pierwsz
kolizji

Birthday Paradox

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

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

$$E[C] = n H_n$$

$$\approx n \ln n$$

Coupon Collector Problem

$$H_n = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}$$

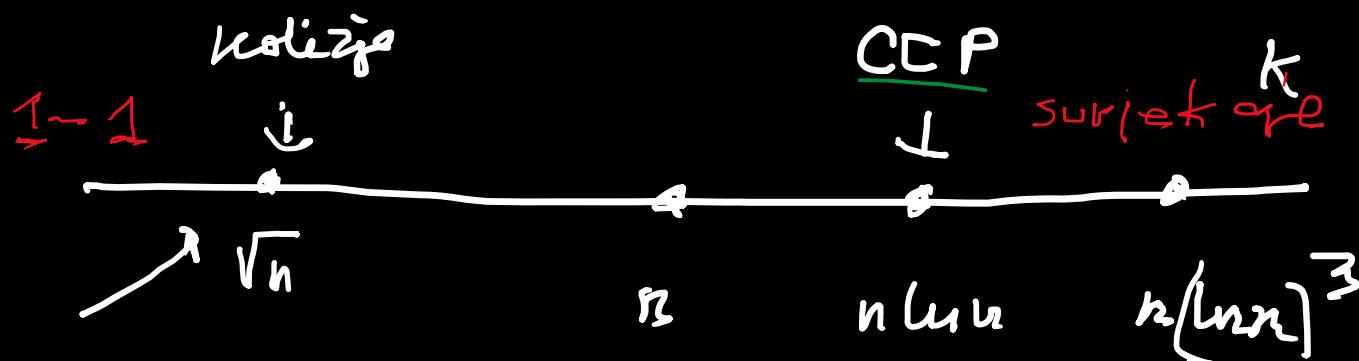
$$H_n \approx \ln n + \gamma + O\left(\frac{1}{n}\right)$$

$$\gamma \approx 0.577$$

$$\approx \ln n$$

$$k \approx n(\ln n)^3$$

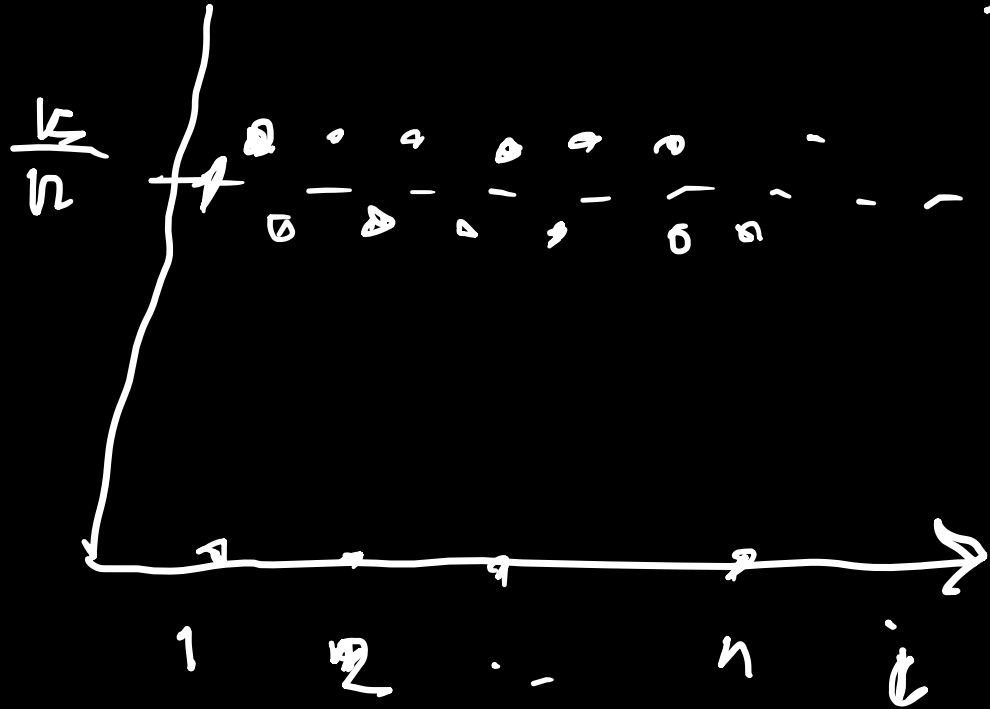
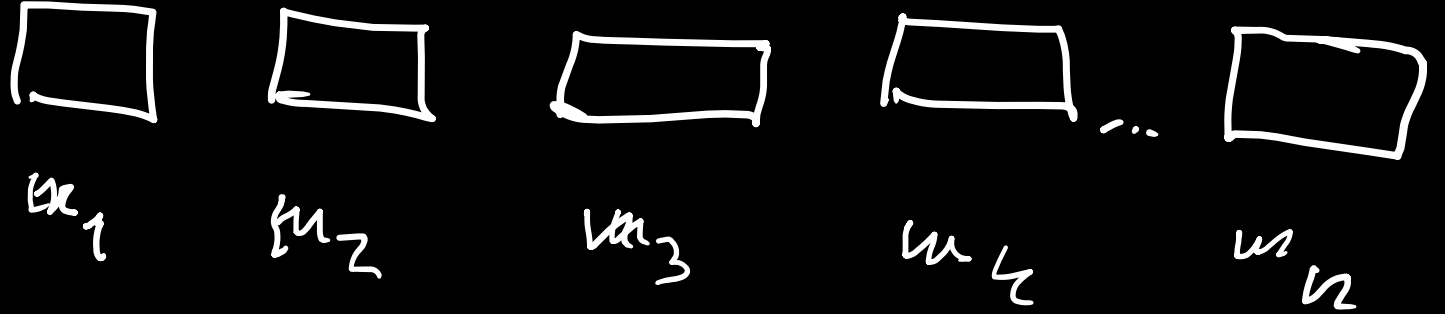
Ukazy są prawie jednakowe.
Wypowiedz : $\approx \frac{k}{n}$ elementów.



Prop. wersja

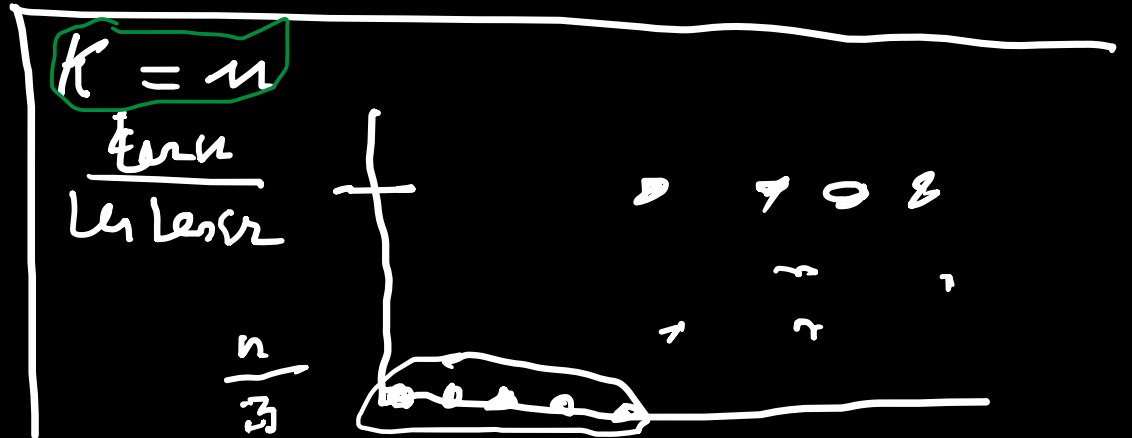
średnio pustych celów $\approx \frac{n}{e}$

$$k > n(n+1)^3$$



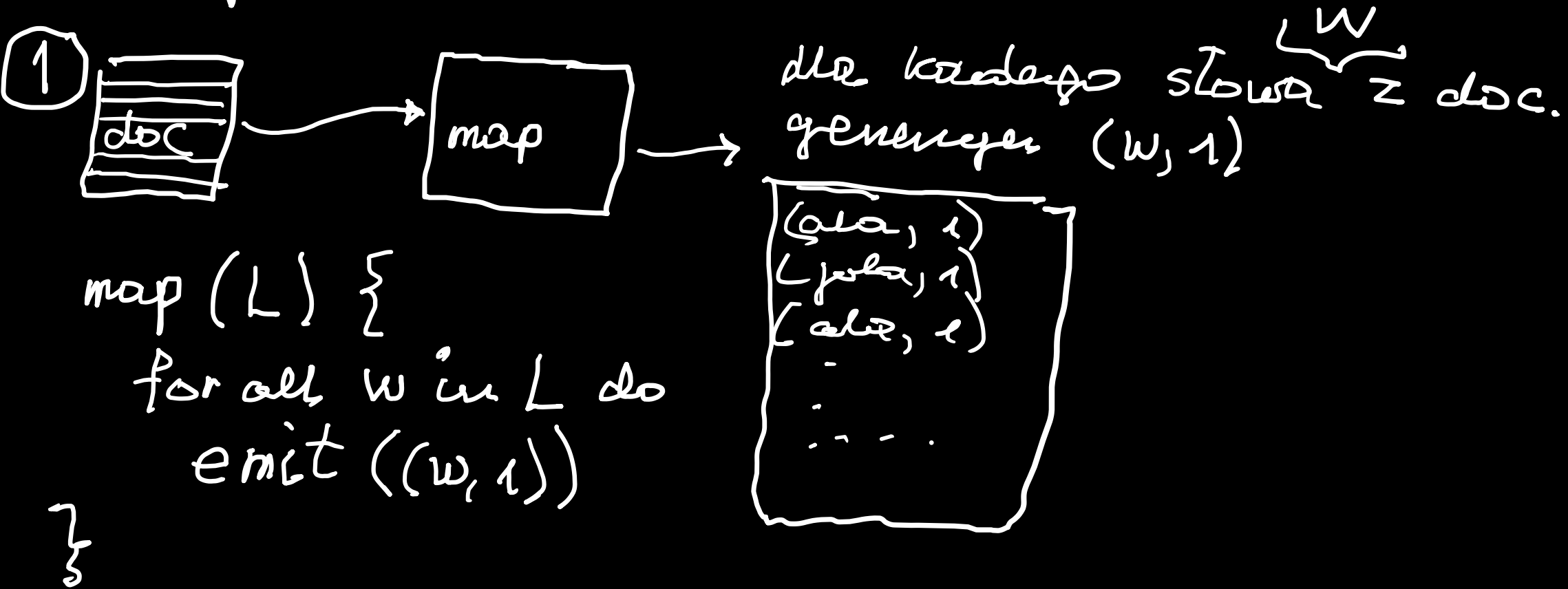
← mussn bliz
des pedeeost

to nie jet "hyresialne"

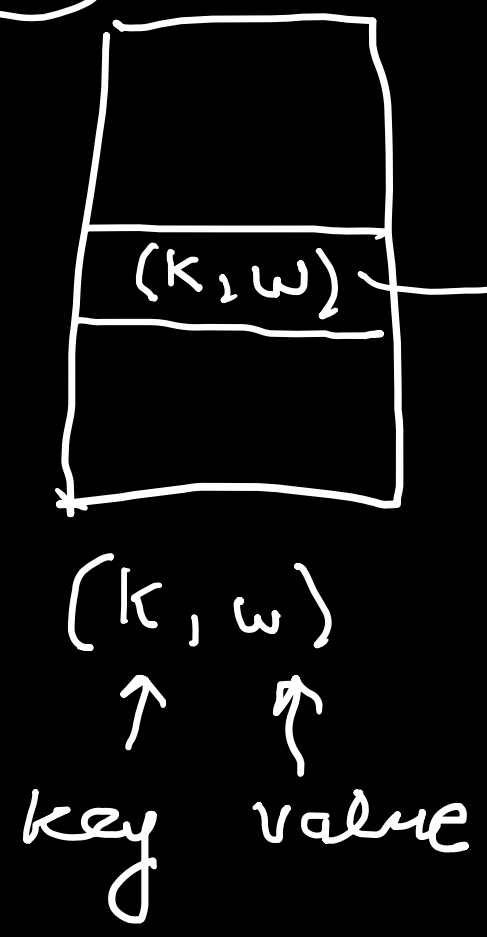


MODEL MAP-REDUCE

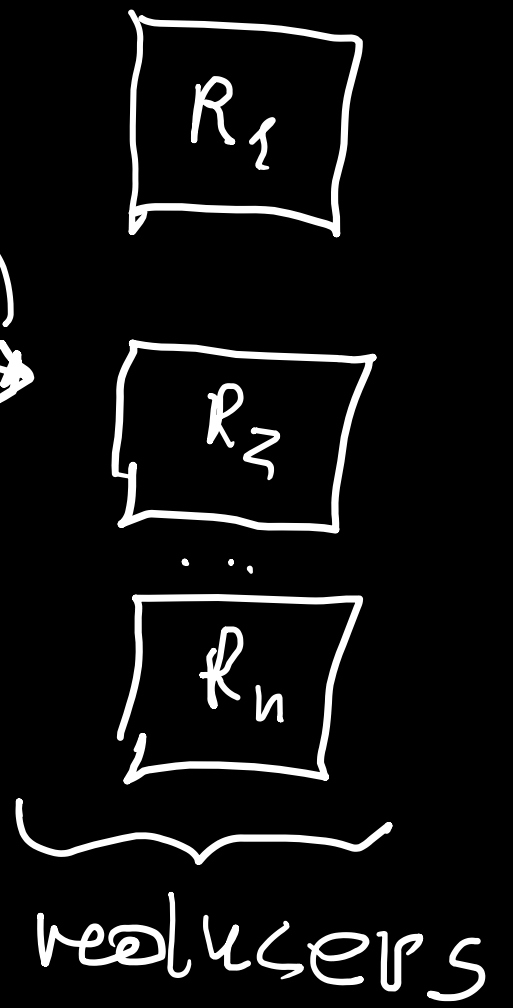
Przykład: word count



2

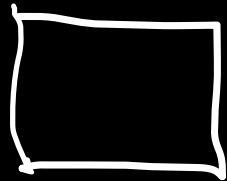


$h(k)$
 (k, w)



$h: \Omega \rightarrow \{1, \dots, n\}$
FIX

③ reducer



$[(k_1, a_1)] [(k_2, a_2)] \dots [(k_n, a_n)]$

sort online

$[[k_1, [a_{1,1} \dots a_{1,n}]] , \dots , [k_r, [a_{r,1} \dots a_{r,n}]]]$

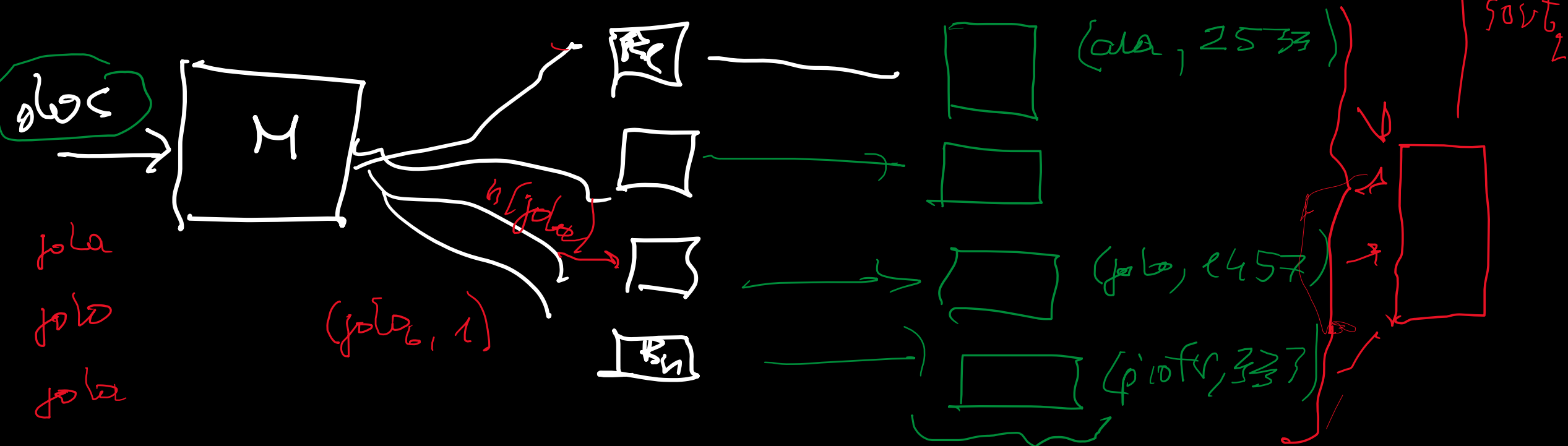
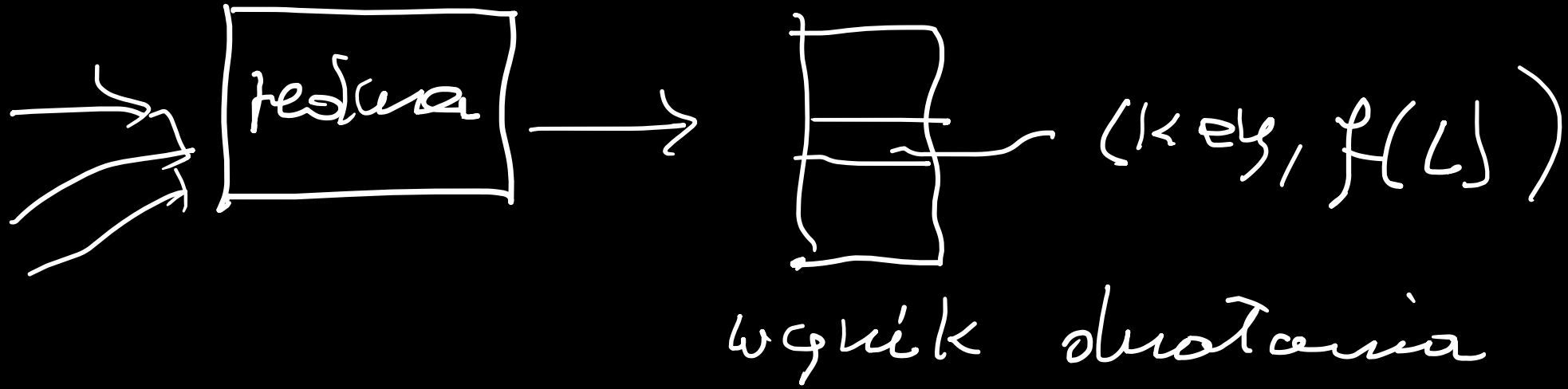
$[a_{1,1}, [1, 1, 1, \dots, 1]]$

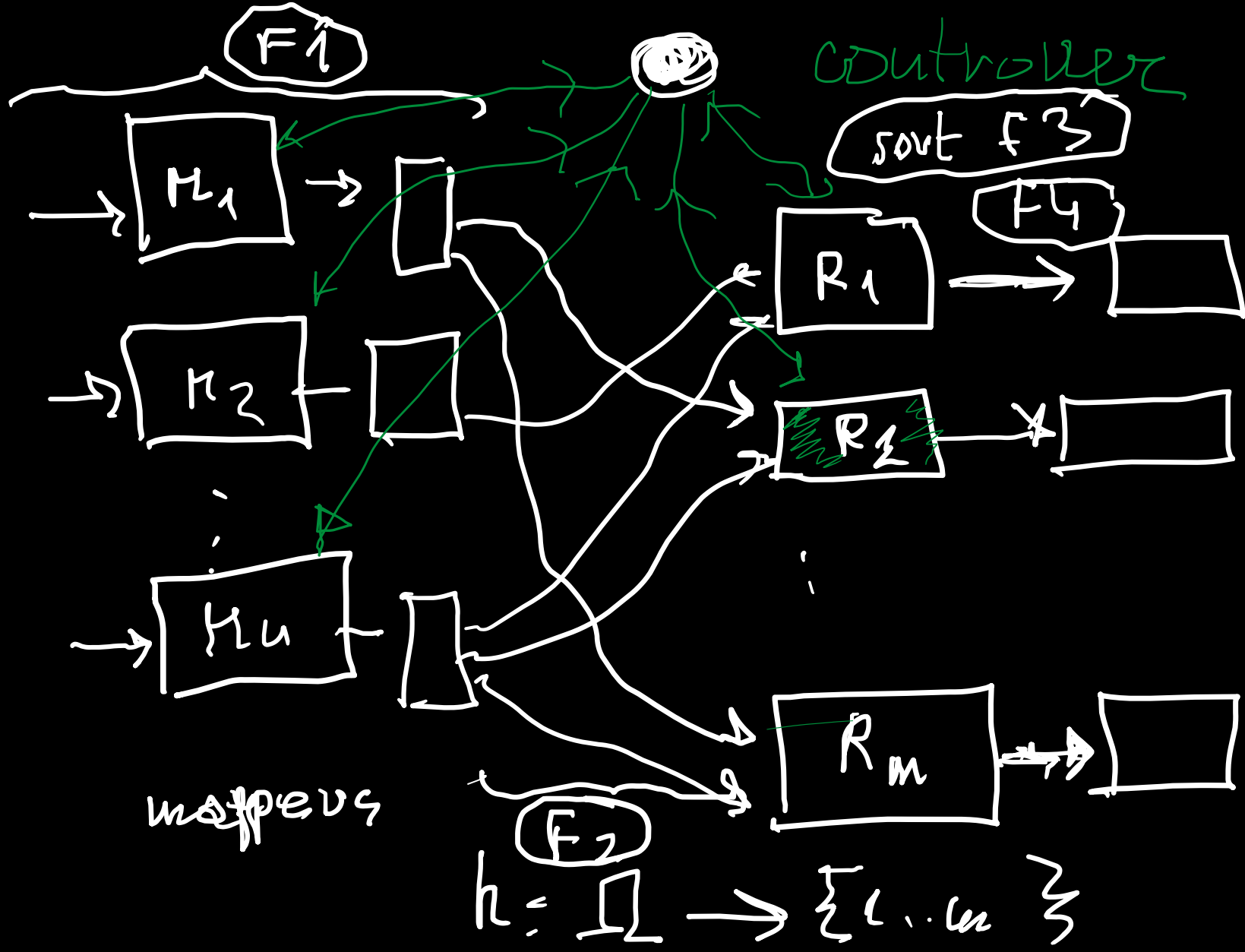
reduce (key, L) {

$m = f(L);$
emit (k, m)

word - count

reduce (key, L) {
emit (key, length(L))
}





```

map (data) {
  ...
  (k1, w1) ... L
}
reduce (key, L) {
  emit (key, f(L))
}

```

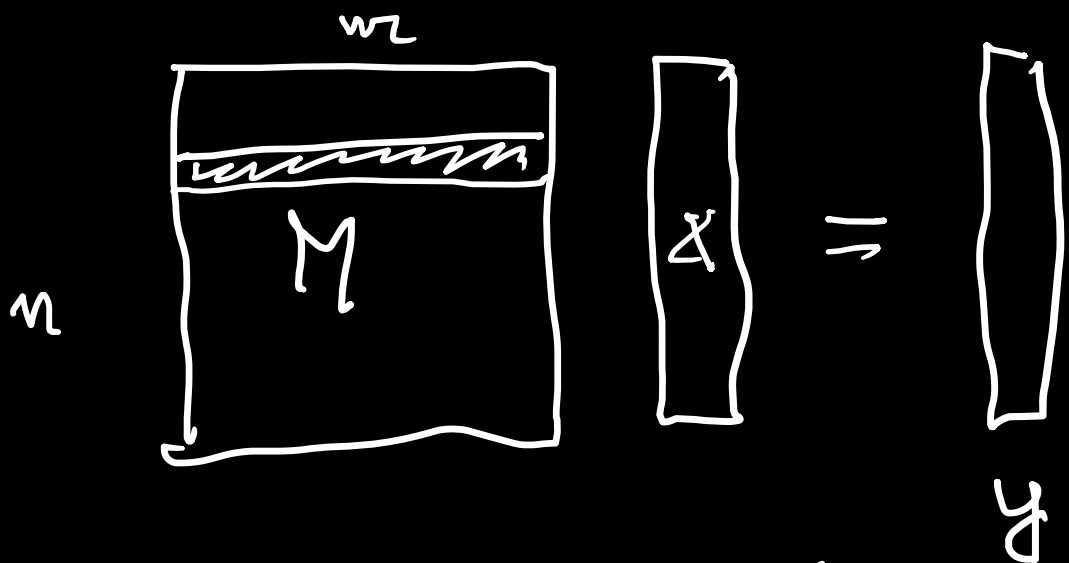
{ Hadoop &
Apache Spark

NIEZAWODNOŚĆ: duplikacja

(P)

Mnożenie macierzy

przez
 wektor



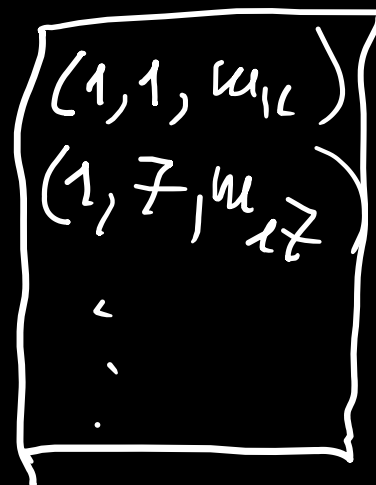
$$y_i = \sum_{j=1}^m m_{ij} x_j$$

$X \leftarrow$ wiersze $Y \leftarrow$ kolumny

$M \leftarrow$ całość

repr. macierzy: (i, j, m_{ij})

(można założyć, że $m_{ij} \neq 0$)



map (i, j, m) {

emit $(i, m \cdot x_j)$

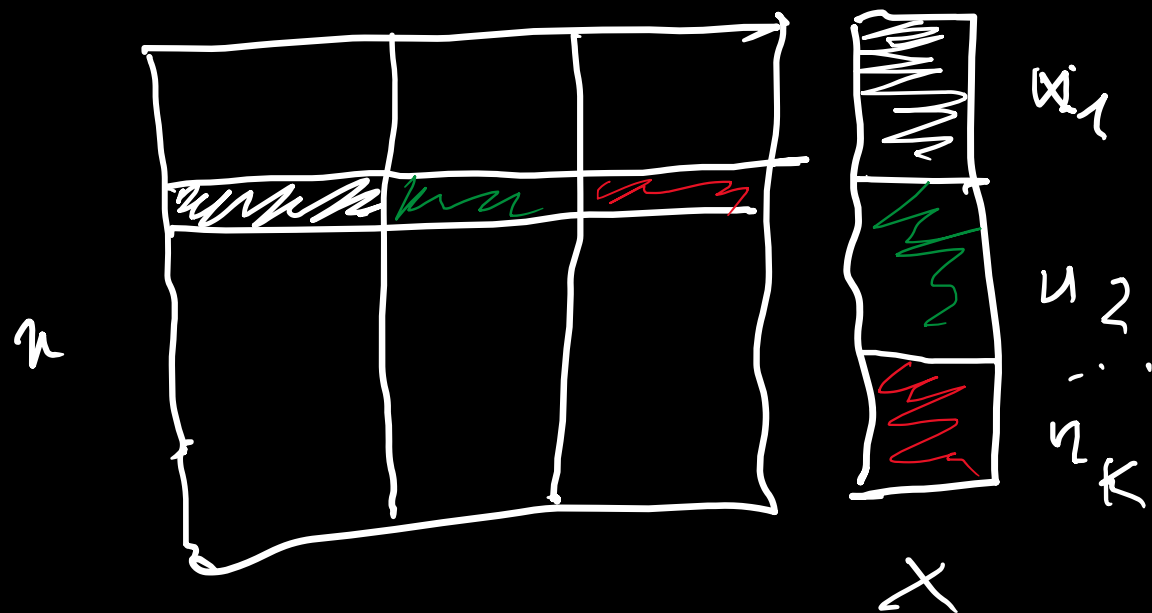
reduce (i, L) {

emit $(i, \sum L)$

$$y_i = \sum_j w_{ij} x_j$$

$[1, [w_{1,3} \cdot x_3, w_{1,5} \cdot x_5, \dots]]$

Za duży wektor?



$$u_1 + u_2 + \dots + u_k = n$$

jest to analogicznie $M-R$,

