

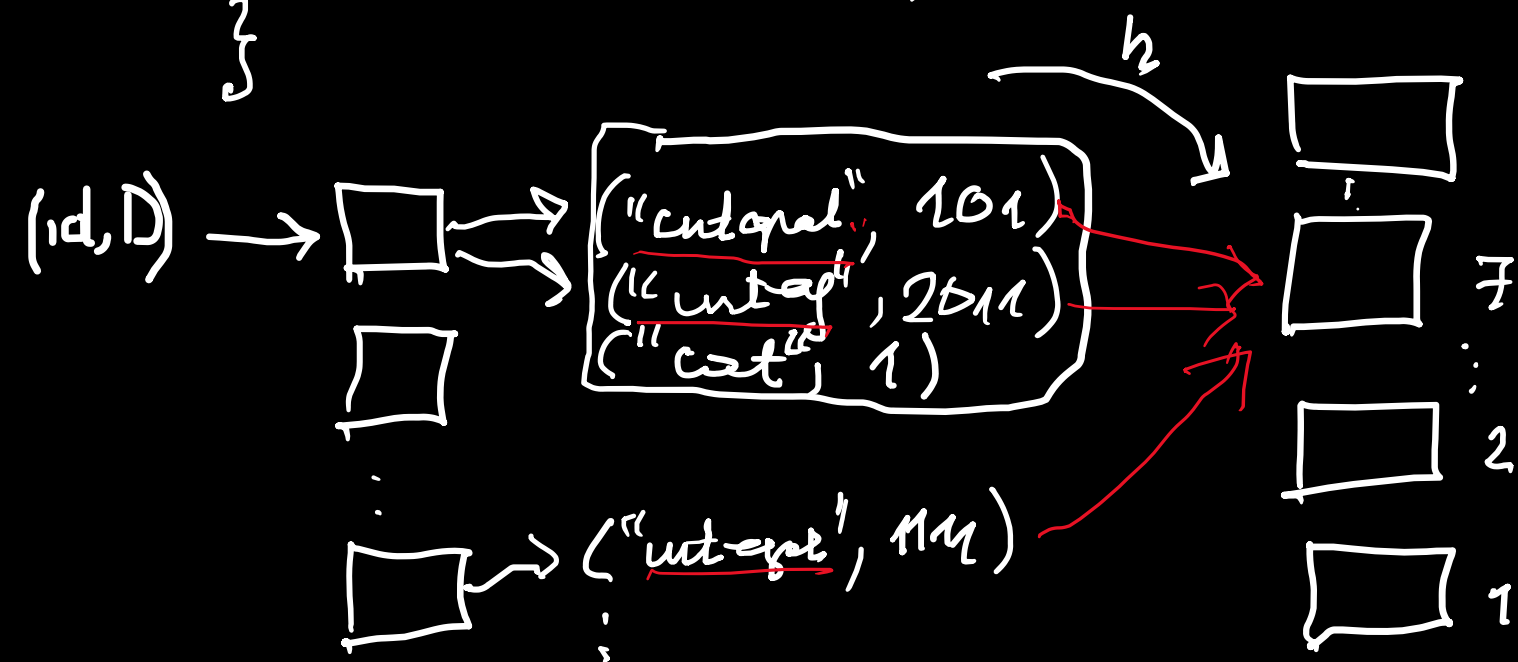
Example. Inverted Index

d_1, \dots, d_n : documents e.g. $n = 10^9$ all web pages

```
map (id, D) {  
  for all w in D do  
    emit(w, id);  
}
```

```
reducer (k, L) {  
  emit(k, L);  
}
```

identity mapp.

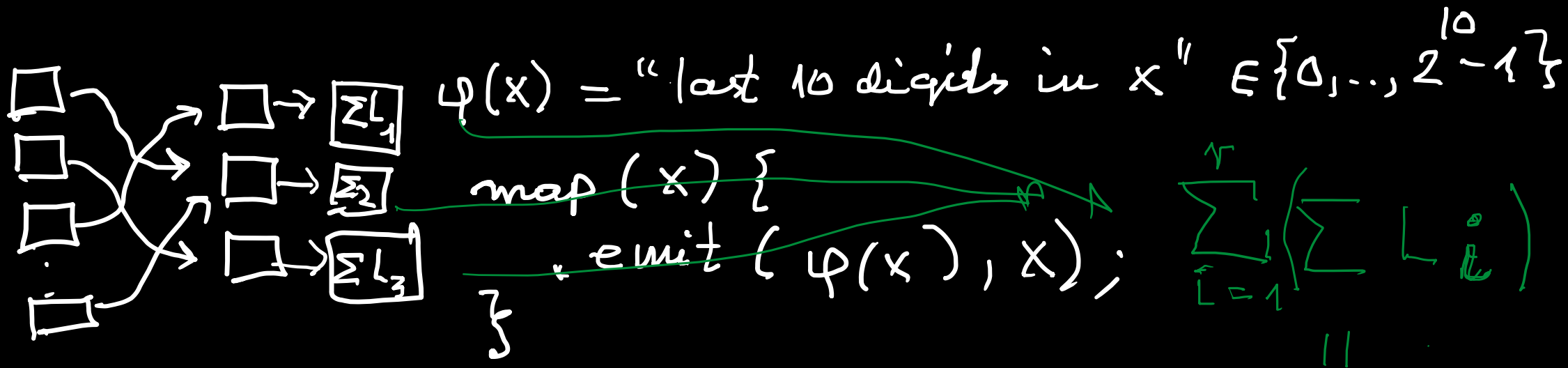


$h(\text{"int-epal"}) = 7$
what reducer do:
 $L = [(\text{"int-epal"}, 101), (\text{"table"}, 2011), \dots]$

Aggregates of data:

$\vec{x} = (x_1, \dots, x_N) \leftarrow$ doubles

- sum: $\sum x_i$; • min(\vec{x}) ; • max(\vec{x})

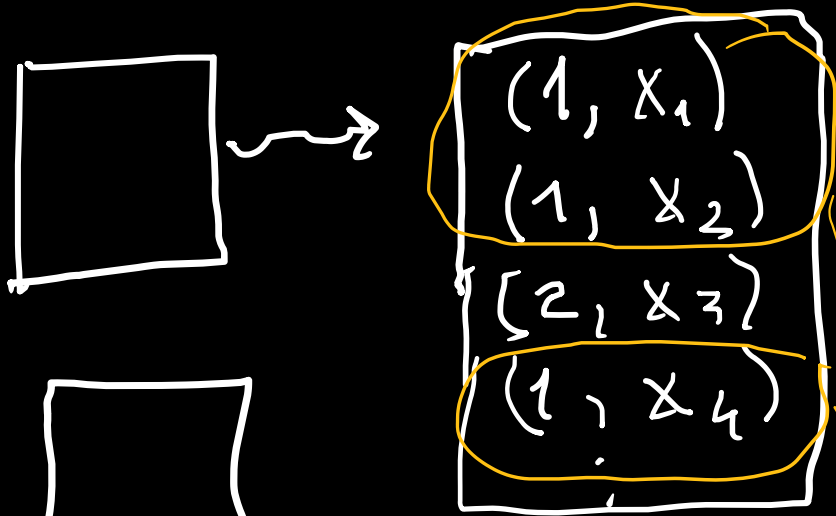


reduce(k, L) $\{ \text{emit}(1, \sum L); \}$

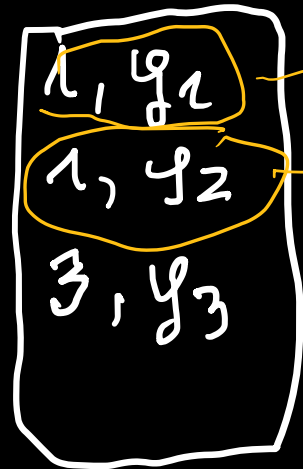
$\{ \text{emit}(1, \text{min}(L)); \}$...

$\sum (L_1 || L_2 || \dots || L_N)$
 + is commutative
 and associative

MAPERS

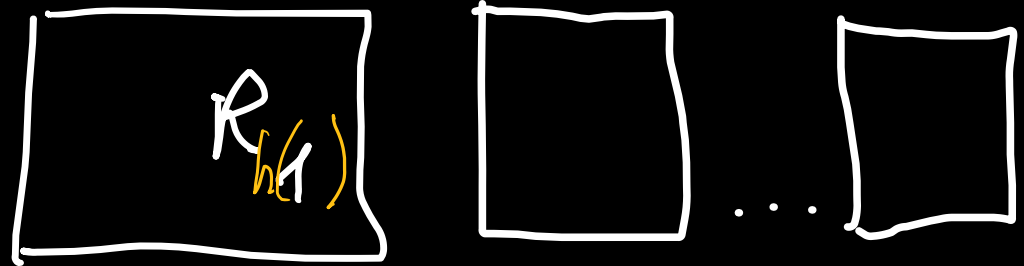


$(1, x_1 + x_2 + x_4)$



$(1, y_1 + y_2)$

reducers



Composers : used after job of mapper before sending data to reducer.

mapper

compose (k, L) {
 emit (k, ΣL);
 }

when we can use it (composer)?

reduce (k, L) {
 emit (k, $\theta(L)$);
 }

$$\theta(x_1, \dots, x_n) = \begin{cases} \Sigma x_i \\ \text{min } x_i \\ \text{max } x_i \end{cases}$$

$$\theta(\theta(L_1), \theta(L_2), \dots, \theta(L_m)) =$$

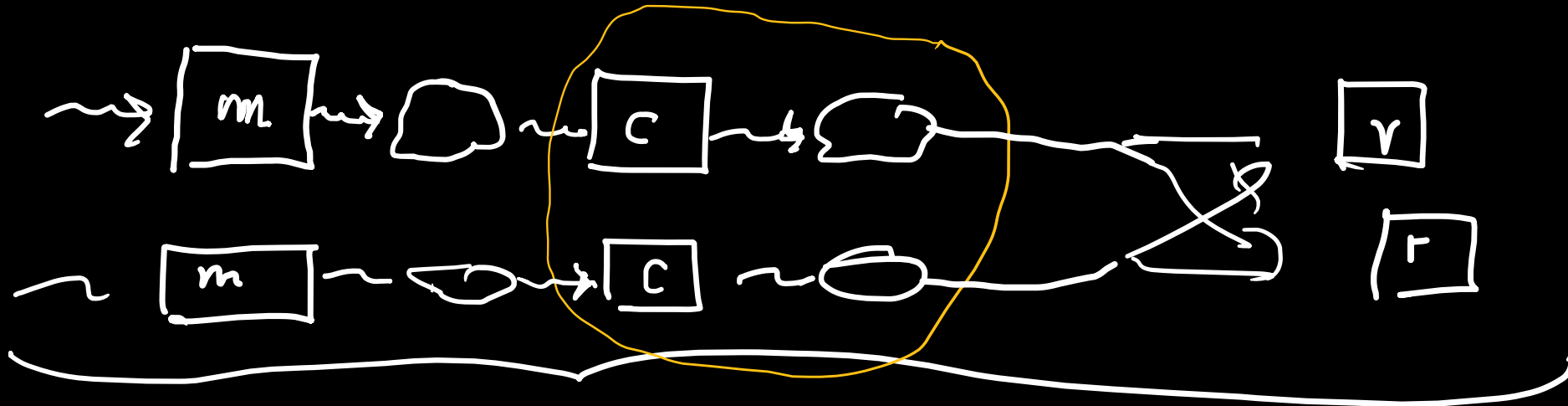
$$\pi \leftarrow \text{perm. of } 1, \dots, m. \quad \theta(L_{\pi(1)} \parallel \dots \parallel L_{\pi(m)})$$

$$\theta(x_1, \dots, x_5) = x_1 + \dots + x_5$$

$$\begin{aligned} \theta(\theta(x_5, x_6), \theta(x_1, x_2, x_3, x_4)) &= \\ = \theta(x_5 + x_6, x_1 + x_2 + \dots + x_4) &= \\ = (x_5 + x_6) + (x_1 + x_2 + x_3 + x_4) & \text{ or } \\ = x_1 + x_2 + x_3 + x_4 + x_5 + x_6 \end{aligned}$$

⊕ is commutative, associative

Use composer if operation
is commutative and associative



composer
phase

x_1, \dots, x_n ; Σ, \min, \max

map (x) {
 emit($\varphi(x), (x, x, x)$);
}

compose (k, L) {

 emit($k, \Sigma \pi_1(L), \min(\pi_2(L)), \max(\pi_3(L))$);
}

reducer (k, L) {
 -- 11 --
}

$L = [(x_1, x_1, x_1), (x_2, x_2, x_2), \dots, (x_k, x_k, x_k)]$

$L = [(s_1, m_1, H_1), (s_2, m_2, H_2), \dots]$

Mean value of $x_1 \dots x_n$.

$$\frac{x_1 + \dots + x_n}{n}$$

map(x) { emit($\psi(x), (1, x)$); }

compose(k, L) { $L = [(1, x_1), (1, x_2), \dots, (1, x_n)]$

emit($k, (\sum \pi_1(L), \sum \pi_2(L))$); }

\downarrow
($3, x_1 + x_2 + x_3$)

reduce(k, L) {

$R_1: \begin{bmatrix} n_1 & s_1 \\ n_2 & s_2 \\ \vdots & \vdots \end{bmatrix}$



$$\frac{s_1 + \dots + s_n}{n_1 + \dots + n_n}$$

}

VARIANCE

$$\text{var}(\vec{x}) = \frac{1}{n} \sum_{i=1}^n \left(x_i - \underbrace{\frac{x_1 + \dots + x_n}{n}}_{\mu} \right)^2 =$$

$$n = |\vec{x}|$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2) = \frac{1}{n} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + \mu^2 \sum_{i=1}^n 1 \right)$$

$$= \frac{1}{n} \left(\sum x_i^2 - 2\mu \cdot (\mu \cdot n) + \mu^2 \cdot n \right) =$$

$$= \frac{1}{n} \sum x_i^2 - \mu^2 = \frac{1}{n} \sum x_i^2 - \left(\frac{\sum x_i}{n} \right)^2$$

$$\boxed{\text{var}(X) = E(X^2) - [E(X)]^2}$$

map (x) { emit($\psi(x)$, (1, x, x²)); }

reduce (k, L) {

emit (1, $\sum \pi_1(L)$, $\sum \pi_2(L)$, $\sum \pi_3(L)$)

k:

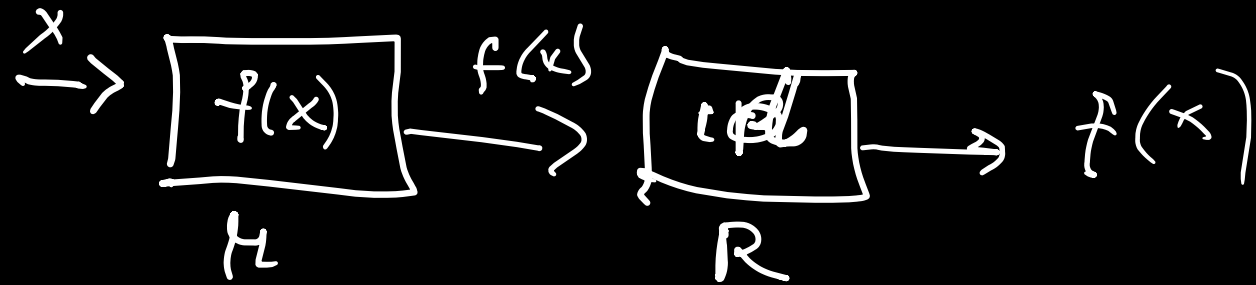
L = [(1, x₁, x₁²), (1, x₂, x₂²), ...]

mean (x), var (x)

EXERCISE: calc.: $\frac{1}{n} \sum (x_i - \mu)^3$ { central
3-th moment

x → ($\psi(x)$, (1, x, x², x³))

What can be calculated?



Set operations : \cup, \cap, \setminus

A, B - too big sets

$x \in A \rightsquigarrow ("A", x)$
 $x \in B \rightsquigarrow ("B", x)$ } input for mapper

```
map (X, x) {  
  emit (x, X);  
}
```

SUM

```
reduce (x, L) {  
  emit (1, x);  
}
```

INTERG.

```
reduce (x, L) {  
  if (|L| = 2) then emit (1, x);  
}
```

DIFF:

```
reduce (x, L) { if (L = ["A"]) then emit (1, x) }
```

reducer

(X, L)

L = ["A"]

L = ["B"]

L = ["A", "B"]

L = ["B", "A"]

SELECT a_1, \dots, a_n FROM R

WHERE $\theta(a_1, \dots, a_n)$

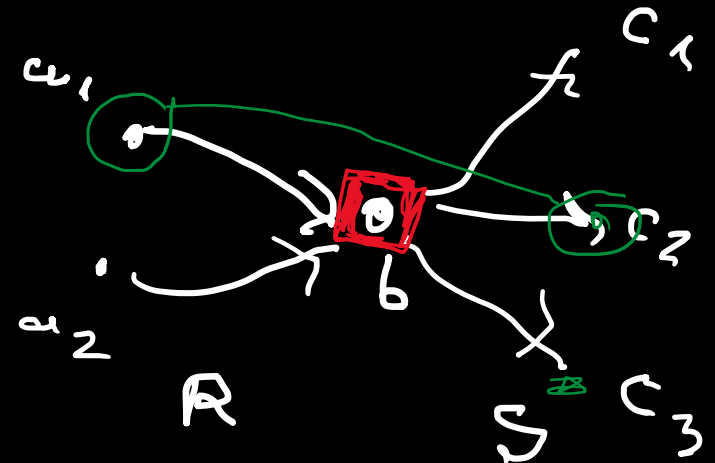
map(a_1, \dots, a_n) { emit ($\theta(a_1, \dots, a_n)$, (a_1, \dots, a_n)) }

JOIN: R(a, b), S(b, c)

SELECT R.a, S.c FROM R, S

WHERE R.b = S.b

INPUT: R(a, b) \rightsquigarrow ("R", a, b)
S(b, c) \rightsquigarrow ("S", b, c)



$$\begin{aligned} (R, a, b) &\longrightarrow (b, (R, a)) \\ ("S", b, c) &\longrightarrow (b, (S, c)) \end{aligned}$$

```
map (X, x, y) {
  if (x == "R") { emit (y, ("R", x)); }
  else { emit (x, ("S", y)); }
}
```

```
reducer (k, L) {
  (L1, L2) ← split(L)
  ↘ ↙
  starts from R starts from S
```

for all $(R, x) \in L_1$, for all $(S, y) \in L_2$ emit $(1, (x, y))$; }

$$\begin{aligned} L &= [(R, x_1), (S, y_1), \\ &\quad (R, x_2), (R, x_3), \\ &\quad (S, y_2), \dots] \\ &\quad \downarrow \\ &[(R, x_1), (R, x_2), \dots (R, x_k)] \\ &\quad [(S, y_1), \dots (S, y_e)] \end{aligned}$$