# Big Data Management & Analytics

EXERCISE 4 – MAPREDUCE, SPARK

23rd of November 2015

**Sabrina Friedl**
LMU Munich

# 1. Matrix Multiplication with MapReduce

REVISION AND EXAMPLE

# MapReduce – Matrix Multiplication

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} \quad B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} \quad A \cdot B = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} \end{pmatrix}$$

Can be written as $\quad A = (I, J, V), B = (J, K, W) \quad$ where $\quad [0] = row, \quad [1] = column \quad$ and $\quad [2] = values$

**Steps**

- 1. Map $\qquad (i, j, a_{ij}) \longrightarrow (j, (A, i, a_{ij})) \qquad (j, k, b_{jk}) \longrightarrow (j, (B, k, b_{jk}))$

- 2. Join $\qquad (j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

- 3. Map $\qquad (j, [(A, i, a_{ij}), (B, k, b_{jk})]) \longrightarrow ((i, k), (a_{ij}b_{jk}))$

- 4. ReduceByKey $\quad ((i, k), [(a_{ij}b_{jk})]) \longrightarrow ((i, k), \sum(a_{ij}b_{jk}))$

# Matrix Multiplication - Example

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \qquad B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} = \begin{pmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{pmatrix} \qquad A \cdot B = C = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} = \begin{pmatrix} 58 & 64 \\ 139 & 154 \end{pmatrix}$$
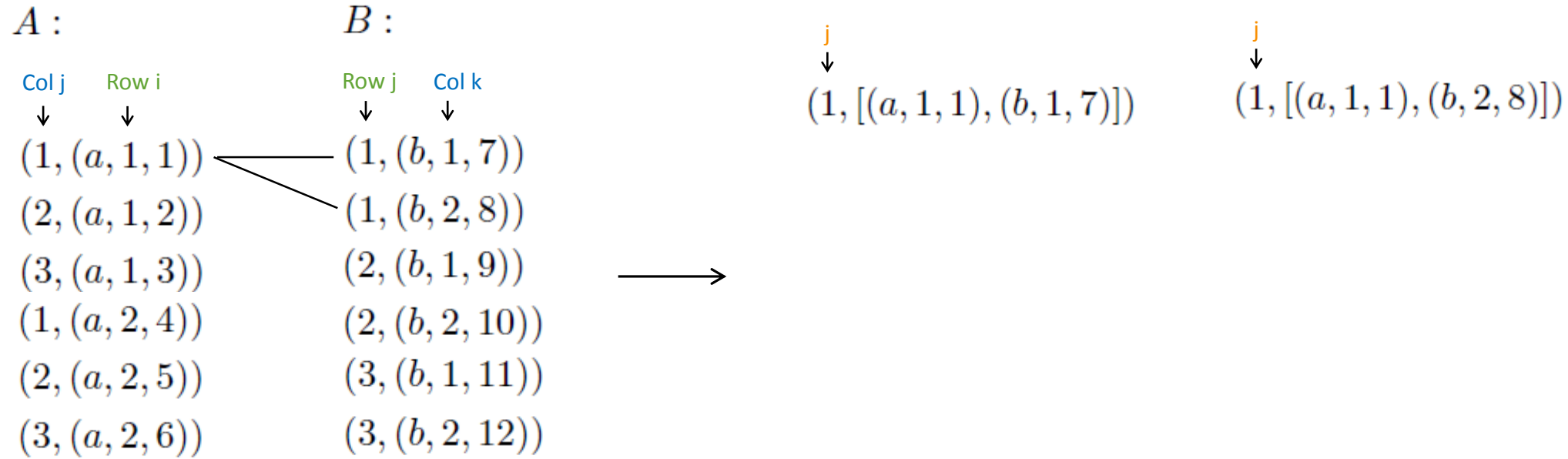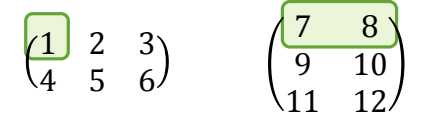
**1. Map:** $\qquad (i, j, a_{ij}) \longrightarrow (j, (A, i, a_{ij})), \qquad\qquad\qquad (j, k, b_{jk}) \longrightarrow (j, (B, k, b_{jk}))$

row col      col ID row               row col      row ID col

$$A: \quad (1,1,1) \longrightarrow (1,(a,1,1))$$
$$(1,2,2) \longrightarrow (2,(a,1,2))$$
$$(1,3,3) \longrightarrow (3,(a,1,3))$$
$$(2,1,4) \longrightarrow (1,(a,2,4))$$
$$(2,2,5) \longrightarrow (2,(a,2,5))$$
$$(2,3,6) \longrightarrow (3,(a,2,6))$$

$$B: \quad (1,1,7) \longrightarrow (1,(b,1,7))$$
$$(1,2,8) \longrightarrow (1,(b,2,8))$$
$$(2,1,9) \longrightarrow (2,(b,1,9))$$
$$(2,2,10) \longrightarrow (2,(b,2,10))$$
$$(3,1,11) \longrightarrow (3,(b,1,11))$$
$$(3,2,12) \longrightarrow (3,(b,2,12))$$

**2. Join:** $(j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \qquad \begin{pmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{pmatrix}$$

$A:$          $B:$

Col j    Row i        Row j    Col k

$(1, (a, 1, 1))$      $(1, (b, 1, 7))$

$(2, (a, 1, 2))$      $(1, (b, 2, 8))$

$(3, (a, 1, 3))$      $(2, (b, 1, 9))$

$(1, (a, 2, 4))$      $(2, (b, 2, 10))$

$(2, (a, 2, 5))$      $(3, (b, 1, 11))$

$(3, (a, 2, 6))$      $(3, (b, 2, 12))$

"Join over j"

$j$

$(1, [(a, 1, 1), (b, 1, 7)])$

$j$

$(1, [(a, 1, 1), (b, 2, 8)])$

**2. Join:** $(j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \qquad \begin{pmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{pmatrix}$$

$A:$       $B:$

Col j    Row i       Row j    Col k

$(1, (a, 1, 1))$     $(1, (b, 1, 7))$

$(2, (a, 1, 2))$     $(1, (b, 2, 8))$

$(3, (a, 1, 3))$     $(2, (b, 1, 9))$

$(1, (a, 2, 4))$     $(2, (b, 2, 10))$

$(2, (a, 2, 5))$     $(3, (b, 1, 11))$

$(3, (a, 2, 6))$     $(3, (b, 2, 12))$

$\longrightarrow$

j

$(1, [(a, 1, 1), (b, 1, 7)])$

$(1, [(a, 2, 4), (b, 1, 7)])$

j

$(1, [(a, 1, 1), (b, 2, 8)])$

$(1, [(a, 2, 4), (b, 2, 8)])$

"Join over j"

**2. Join:** $(j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$  $\begin{pmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{pmatrix}$

$A:$

Col j    Row i
↓      ↓

$(1, (a, 1, 1))$

$(2, (a, 1, 2))$

$(3, (a, 1, 3))$

$(1, (a, 2, 4))$

$(2, (a, 2, 5))$

$(3, (a, 2, 6))$

$B:$

Row j    Col k
↓      ↓

$(1, (b, 1, 7))$

$(1, (b, 2, 8))$

$(2, (b, 1, 9))$

$(2, (b, 2, 10))$

$(3, (b, 1, 11))$

$(3, (b, 2, 12))$

$\longrightarrow$

j
↓

$(1, [(a, 1, 1), (b, 1, 7)])$

$(1, [(a, 2, 4), (b, 1, 7)])$

$(2, [(a, 1, 2), (b, 1, 9)])$

j
↓

$(1, [(a, 1, 1), (b, 2, 8)])$
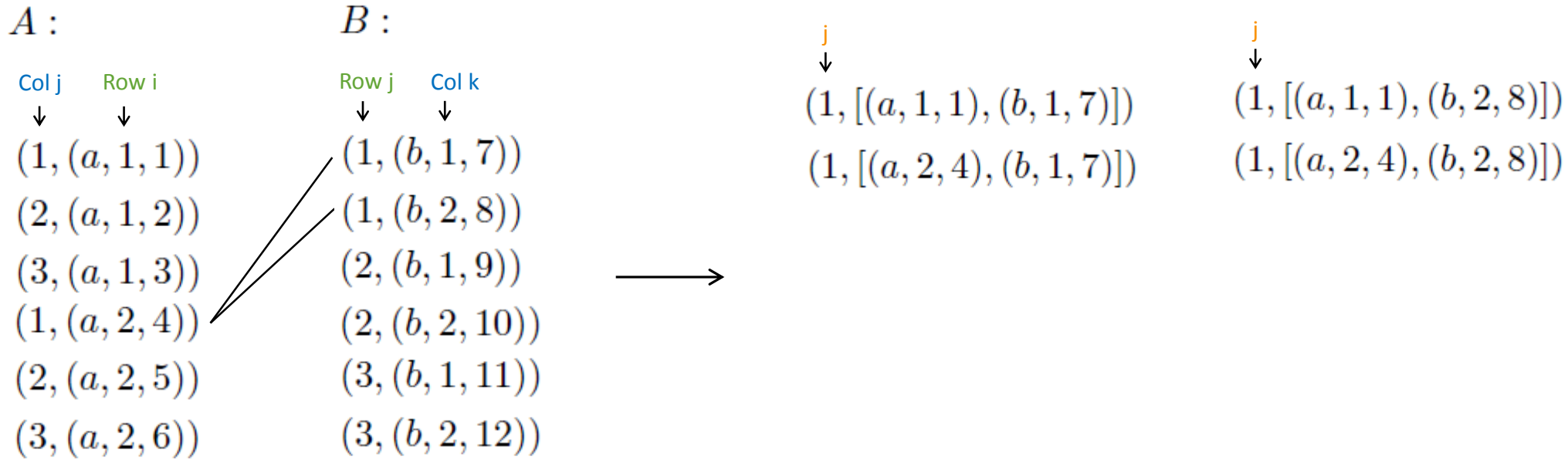
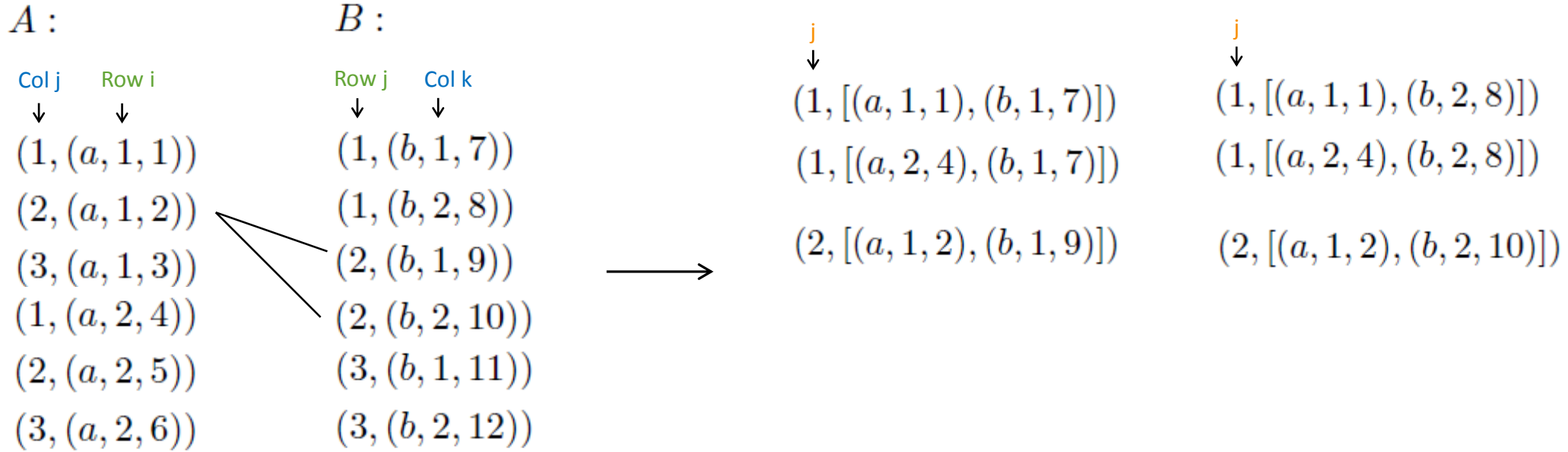$(1, [(a, 2, 4), (b, 2, 8)])$

$(2, [(a, 1, 2), (b, 2, 10)])$

"Join over j"

**2. Join:** $(j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

$\begin{pmatrix} 1 & \boxed{2} & 3 \\ 4 & \boxed{5} & 6 \end{pmatrix}$    $\begin{pmatrix} 7 & 8 \\ \boxed{9 \quad 10} \\ 11 & 12 \end{pmatrix}$

$A:$    $B:$

Col j    Row i    Row j    Col k
↓    ↓    ↓    ↓

$(1, (a, 1, 1))$    $(1, (b, 1, 7))$

$(2, (a, 1, 2))$    $(1, (b, 2, 8))$

$(3, (a, 1, 3))$    $(2, (b, 1, 9))$

$(1, (a, 2, 4))$    $(2, (b, 2, 10))$

$(2, (a, 2, 5))$    $(3, (b, 1, 11))$

$(3, (a, 2, 6))$    $(3, (b, 2, 12))$

$\longrightarrow$

j
↓

$(1, [(a, 1, 1), (b, 1, 7)])$

$(1, [(a, 2, 4), (b, 1, 7)])$

$(2, [(a, 1, 2), (b, 1, 9)])$

$(2, [(a, 2, 5), (b, 1, 9)])$

j
↓

$(1, [(a, 1, 1), (b, 2, 8)])$

$(1, [(a, 2, 4), (b, 2, 8)])$

$(2, [(a, 1, 2), (b, 2, 10)])$
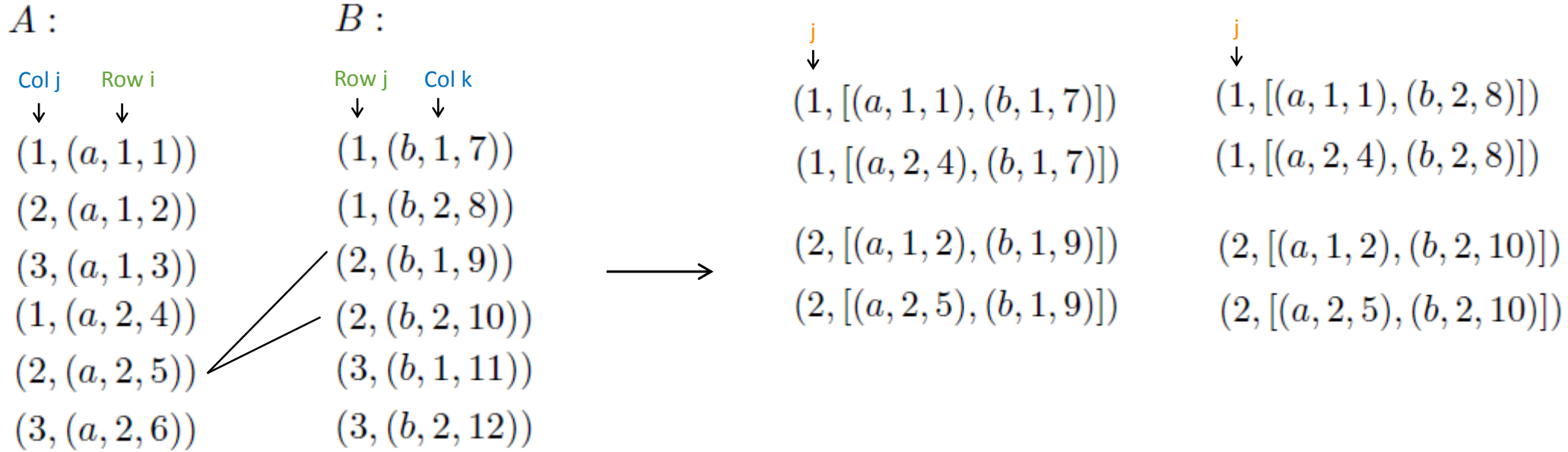
$(2, [(a, 2, 5), (b, 2, 10)])$

"Join over j"

**2. Join:** $(j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

$\begin{pmatrix} 1 & 2 & \boxed{3} \\ 4 & 5 & 6 \end{pmatrix}$ $\quad \begin{pmatrix} 7 & 8 \\ 9 & 10 \\ \boxed{11} & \boxed{12} \end{pmatrix}$

$A:$

Col j    Row i
$\downarrow$     $\downarrow$

$(1, (a, 1, 1))$

$(2, (a, 1, 2))$

$(3, (a, 1, 3))$

$(1, (a, 2, 4))$

$(2, (a, 2, 5))$

$(3, (a, 2, 6))$

$B:$

Row j    Col k
$\downarrow$     $\downarrow$

$(1, (b, 1, 7))$

$(1, (b, 2, 8))$

$(2, (b, 1, 9))$

$(2, (b, 2, 10))$

$(3, (b, 1, 11))$

$(3, (b, 2, 12))$

$\longrightarrow$

j
$\downarrow$

$(1, [(a, 1, 1), (b, 1, 7)])$

$(1, [(a, 2, 4), (b, 1, 7)])$

$(2, [(a, 1, 2), (b, 1, 9)])$

$(2, [(a, 2, 5), (b, 1, 9)])$

$(3, [(a, 1, 3), (b, 1, 11)])$

j
$\downarrow$

$(1, [(a, 1, 1), (b, 2, 8)])$

$(1, [(a, 2, 4), (b, 2, 8)])$

$(2, [(a, 1, 2), (b, 2, 10)])$

$(2, [(a, 2, 5), (b, 2, 10)])$

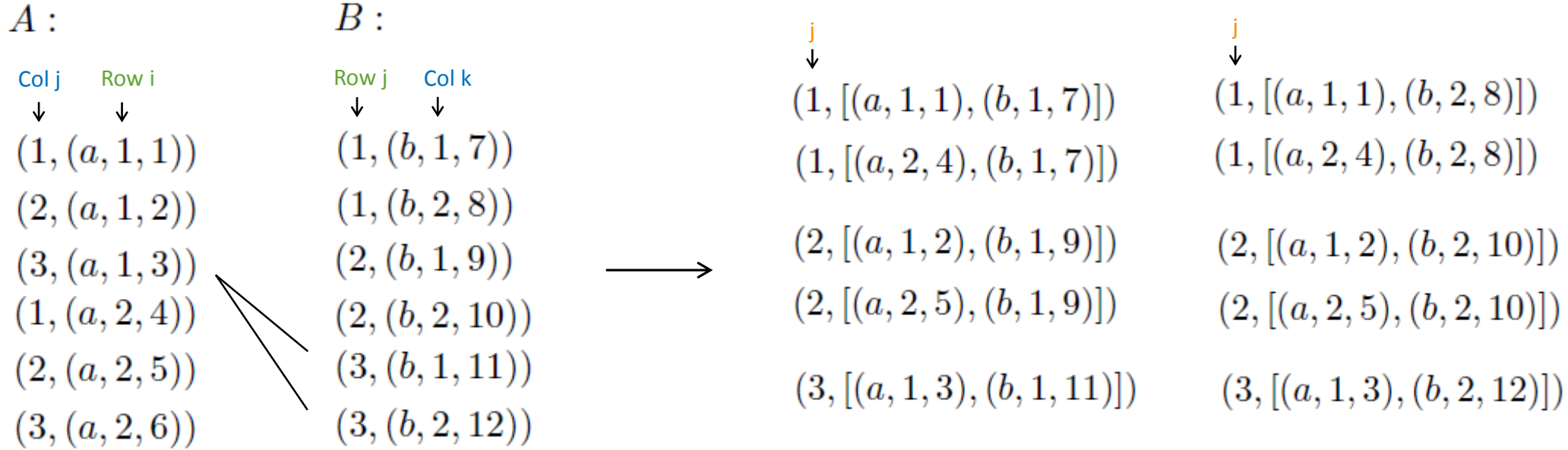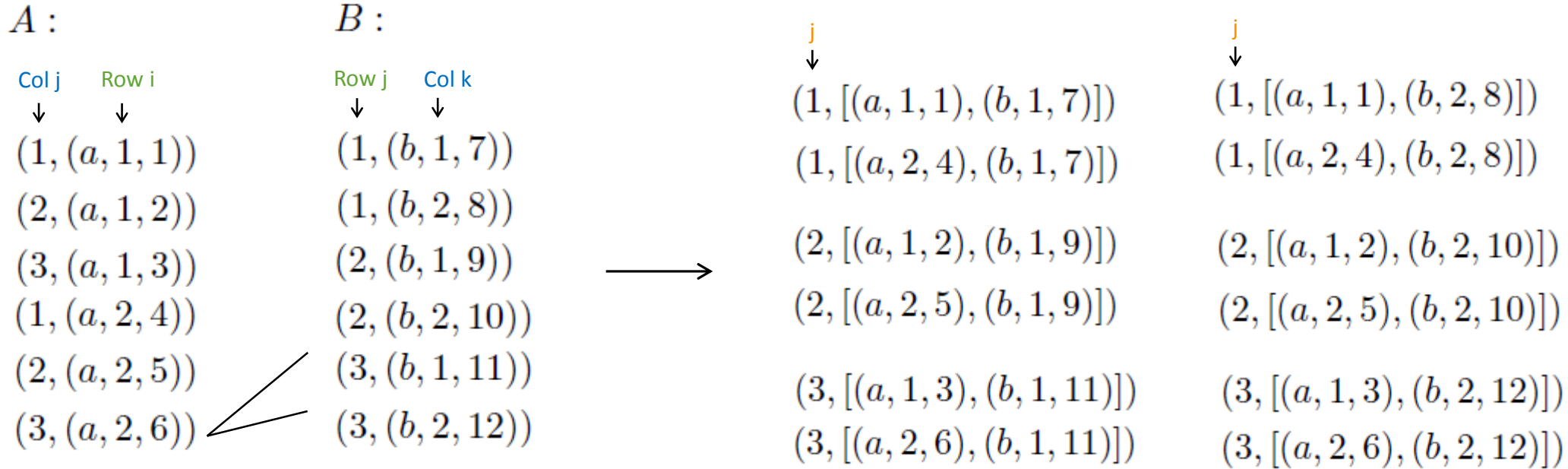$(3, [(a, 1, 3), (b, 2, 12)])$

"Join over j"

**2. Join:** $(j, (A, i, a_{ij})) \bowtie (j, (B, k, b_{jk})) \longrightarrow (j, [(A, i, a_{ij}), (B, k, b_{jk})])$

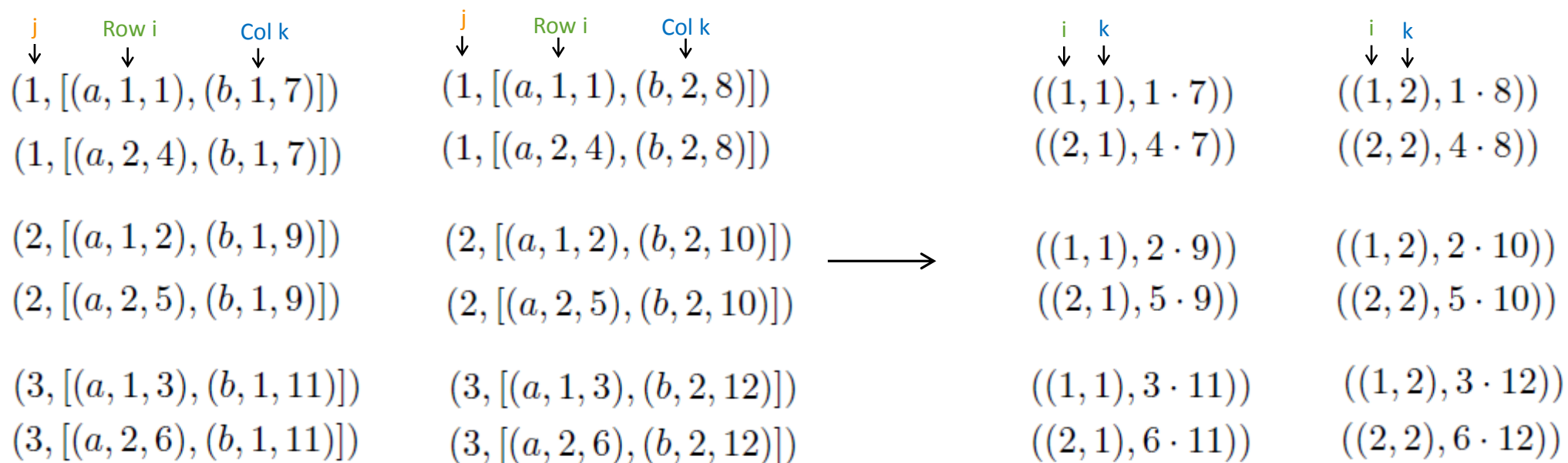$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} \qquad \begin{pmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{pmatrix}$$

$A:$

Col j  Row i

$(1, (a, 1, 1))$
$(2, (a, 1, 2))$
$(3, (a, 1, 3))$
$(1, (a, 2, 4))$
$(2, (a, 2, 5))$
$(3, (a, 2, 6))$

$B:$

Row j  Col k

$(1, (b, 1, 7))$
$(1, (b, 2, 8))$
$(2, (b, 1, 9))$
$(2, (b, 2, 10))$
$(3, (b, 1, 11))$
$(3, (b, 2, 12))$

$\longrightarrow$

j

$(1, [(a, 1, 1), (b, 1, 7)])$
$(1, [(a, 2, 4), (b, 1, 7)])$

$(2, [(a, 1, 2), (b, 1, 9)])$
$(2, [(a, 2, 5), (b, 1, 9)])$

$(3, [(a, 1, 3), (b, 1, 11)])$
$(3, [(a, 2, 6), (b, 1, 11)])$

j

$(1, [(a, 1, 1), (b, 2, 8)])$
$(1, [(a, 2, 4), (b, 2, 8)])$

$(2, [(a, 1, 2), (b, 2, 10)])$
$(2, [(a, 2, 5), (b, 2, 10)])$

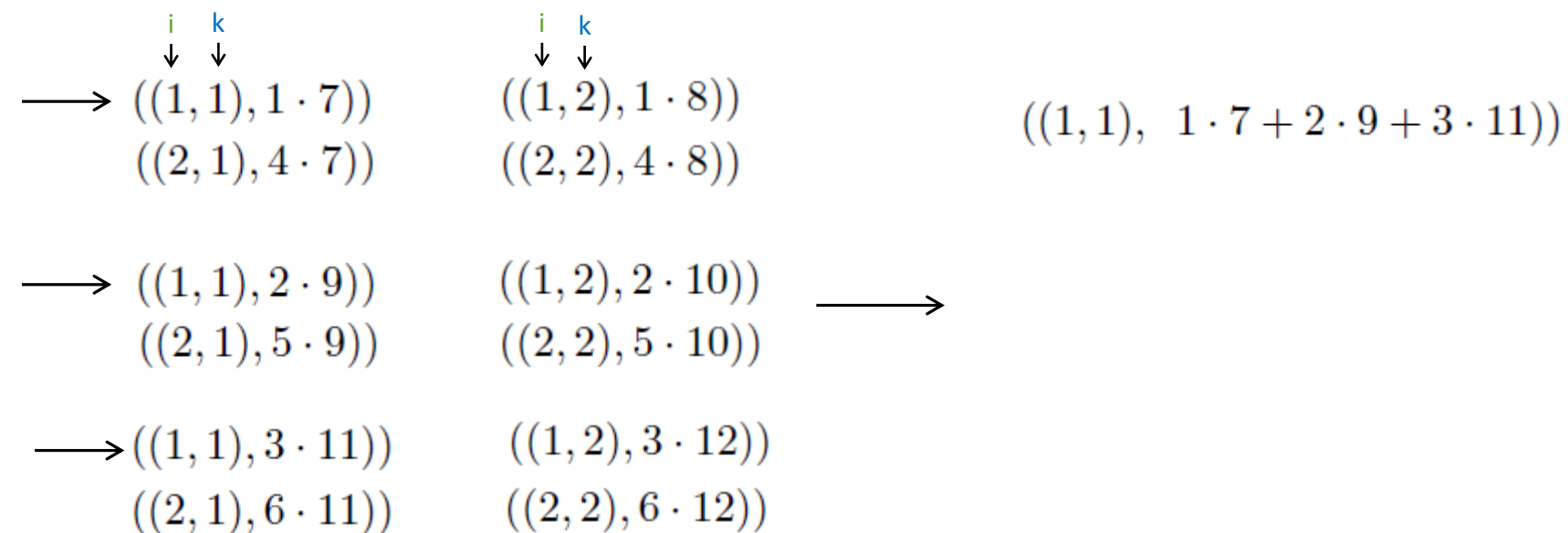$(3, [(a, 1, 3), (b, 2, 12)])$
$(3, [(a, 2, 6), (b, 2, 12)])$

"Join over j"

Number of key-value pairs: $i \cdot j \cdot k$

**3. Map:** $(j, [(A, i, a_{ij}), (B, k, b_{jk})]) \longrightarrow ((i, k), (a_{ij}b_{jk}))$

| j | Row i | Col k |
|---|---|---|

$(1, [(a, 1, 1), (b, 1, 7)])$

$(1, [(a, 2, 4), (b, 1, 7)])$

$(2, [(a, 1, 2), (b, 1, 9)])$

$(2, [(a, 2, 5), (b, 1, 9)])$

$(3, [(a, 1, 3), (b, 1, 11)])$

$(3, [(a, 2, 6), (b, 1, 11)])$

| j | Row i | Col k |
|---|---|---|

$(1, [(a, 1, 1), (b, 2, 8)])$

$(1, [(a, 2, 4), (b, 2, 8)])$

$(2, [(a, 1, 2), (b, 2, 10)])$

$(2, [(a, 2, 5), (b, 2, 10)])$

$(3, [(a, 1, 3), (b, 2, 12)])$

$(3, [(a, 2, 6), (b, 2, 12)])$

$\longrightarrow$

| i | k |
|---|---|

$((1, 1), 1 \cdot 7))$

$((2, 1), 4 \cdot 7))$

$((1, 1), 2 \cdot 9))$

$((2, 1), 5 \cdot 9))$

$((1, 1), 3 \cdot 11))$

$((2, 1), 6 \cdot 11))$

| i | k |
|---|---|

$((1, 2), 1 \cdot 8))$

$((2, 2), 4 \cdot 8))$

$((1, 2), 2 \cdot 10))$

$((2, 2), 5 \cdot 10))$

$((1, 2), 3 \cdot 12))$

$((2, 2), 6 \cdot 12))$

## 4. ReduceByKey: $(lambda\ x, y : x + y)$

$$((1,1), 1 \cdot 7))$$
$$((2,1), 4 \cdot 7))$$

$$((1,2), 1 \cdot 8))$$
$$((2,2), 4 \cdot 8))$$

$$((1,1),\ 1 \cdot 7 + 2 \cdot 9 + 3 \cdot 11))$$

$$((1,1), 2 \cdot 9))$$
$$((2,1), 5 \cdot 9))$$

$$((1,2), 2 \cdot 10))$$
$$((2,2), 5 \cdot 10))$$

$$((1,1), 3 \cdot 11))$$
$$((2,1), 6 \cdot 11))$$

$$((1,2), 3 \cdot 12))$$
$$((2,2), 6 \cdot 12))$$

## 4. ReduceByKey: $(lambda\ x, y : x + y)$

i   k                i   k

$((1,1), 1 \cdot 7)) \longrightarrow ((1,2), 1 \cdot 8))$
$((2,1), 4 \cdot 7)) \qquad ((2,2), 4 \cdot 8))$

$((1,1), 2 \cdot 9)) \longrightarrow ((1,2), 2 \cdot 10))$
$((2,1), 5 \cdot 9)) \qquad ((2,2), 5 \cdot 10)) \longrightarrow$

$((1,1), 3 \cdot 11)) \longrightarrow ((1,2), 3 \cdot 12))$
$((2,1), 6 \cdot 11)) \qquad ((2,2), 6 \cdot 12))$

$((1,1),\ 1 \cdot 7 + 2 \cdot 9 + 3 \cdot 11)) \qquad ((1,2),\ 1 \cdot 8 + 2 \cdot 10 + 3 \cdot 12)$

**4. ReduceByKey:** $(lambda\ x, y : x + y)$

i  k
↓  ↓
$((1,1), 1 \cdot 7))$      i  k
↓  ↓
$((1,2), 1 \cdot 8))$

⟶ $((2,1), 4 \cdot 7))$      $((2,2), 4 \cdot 8))$

$((1,1), 2 \cdot 9))$      $((1,2), 2 \cdot 10))$

⟶ $((2,1), 5 \cdot 9))$      $((2,2), 5 \cdot 10))$

$((1,1), 3 \cdot 11))$      $((1,2), 3 \cdot 12))$

⟶ $((2,1), 6 \cdot 11))$      $((2,2), 6 \cdot 12))$

$((1,1),\ 1 \cdot 7 + 2 \cdot 9 + 3 \cdot 11))$      $((1,2),\ 1 \cdot 8 + 2 \cdot 10 + 3 \cdot 12)$

$((2,1),\ 4 \cdot 7 + 5 \cdot 9 + 6 \cdot 11))$

## 4. ReduceByKey: $(lambda\ x, y : x + y)$

$$((1,1), 1 \cdot 7)) \qquad ((1,2), 1 \cdot 8))$$
$$((2,1), 4 \cdot 7)) \longrightarrow ((2,2), 4 \cdot 8))$$

$$((1,1), 2 \cdot 9)) \qquad ((1,2), 2 \cdot 10))$$
$$((2,1), 5 \cdot 9)) \longrightarrow ((2,2), 5 \cdot 10)) \qquad \longrightarrow$$

$$((1,1), 3 \cdot 11)) \qquad ((1,2), 3 \cdot 12))$$
$$((2,1), 6 \cdot 11)) \longrightarrow ((2,2), 6 \cdot 12))$$

$$((1,1),\ 1 \cdot 7 + 2 \cdot 9 + 3 \cdot 11)) \qquad ((1,2),\ 1 \cdot 8 + 2 \cdot 10 + 3 \cdot 12)$$

$$((2,1),\ 4 \cdot 7 + 5 \cdot 9 + 6 \cdot 11)) \qquad ((2,2),\ 4 \cdot 8 + 5 \cdot 10 + 6 \cdot 12)$$

## 4. ReduceByKey: $(lambda\ x, y : x + y)$

$$((1,1), 1 \cdot 7))$$
$$((2,1), 4 \cdot 7))$$

$$((1,2), 1 \cdot 8))$$
$$((2,2), 4 \cdot 8))$$

$$((1,1),\ 1 \cdot 7 + 2 \cdot 9 + 3 \cdot 11))$$
$$((1,2),\ 1 \cdot 8 + 2 \cdot 10 + 3 \cdot 12)$$

$$((1,1), 2 \cdot 9))$$
$$((2,1), 5 \cdot 9))$$

$$((1,2), 2 \cdot 10))$$
$$((2,2), 5 \cdot 10))$$

$$((2,1),\ 4 \cdot 7 + 5 \cdot 9 + 6 \cdot 11))$$
$$((2,2),\ 4 \cdot 8 + 5 \cdot 10 + 6 \cdot 12)$$

$$((1,1), 3 \cdot 11))$$
$$((2,1), 6 \cdot 11))$$

$$((1,2), 3 \cdot 12))$$
$$((2,2), 6 \cdot 12))$$

$$\longrightarrow$$

$$C = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} = \begin{pmatrix} 58 & 64 \\ 139 & 154 \end{pmatrix}$$

Number of elements:   $i \cdot k$

# 2. KMeans with MapReduce

Revision

# MapReduce - KMeans

Randomly initialize k centers:

$$\mu^{(0)} = \mu_1^{(0)}, \ldots, \mu_k^{(0)}$$

**Classify:** Assign each point $j \in \{1, \ldots, m\}$ to nearest centre:

$$z^j \leftarrow \arg\min_i \|\mu_i - x^j\|_2^2$$

**Map**

**Recenter:** $\mu_i$ becomes centroid of its points:

$$\mu_i^{(t+1)} \leftarrow \arg\min_\mu \sum_{j:z^j=i} \|\mu - x^j\|_2^2$$

**Reduce**

(a) Initialization

(b) First Iteration

(c) Convergence