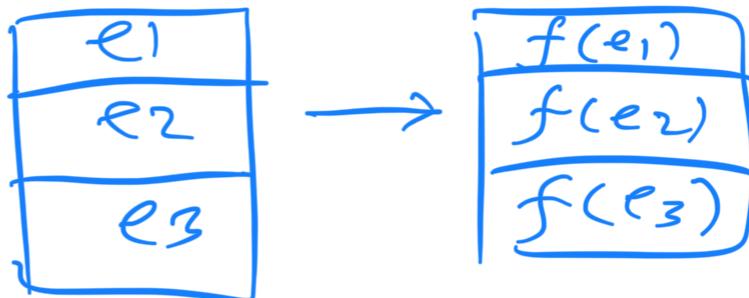


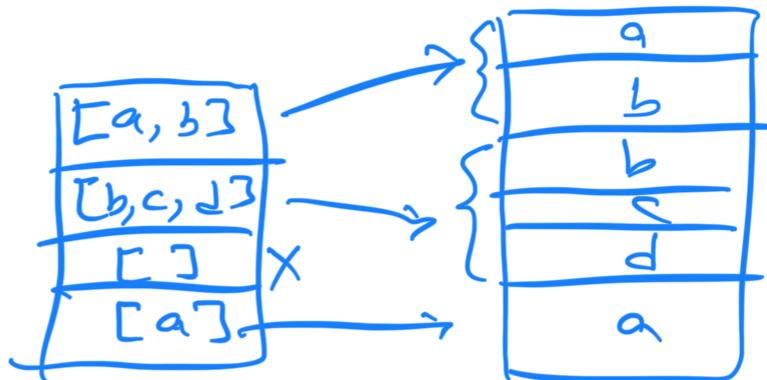
# Summary of mappers & Reducers

Mappers }      map()    1 - to - 1  
              flatMap()    1 - to - M  
              mapPartitions(): M - to - 1

map()

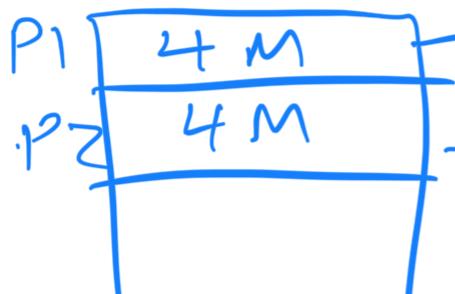


flatMap



mapPartitions

source RDD



+ target



# of elements: 80,000,000,000

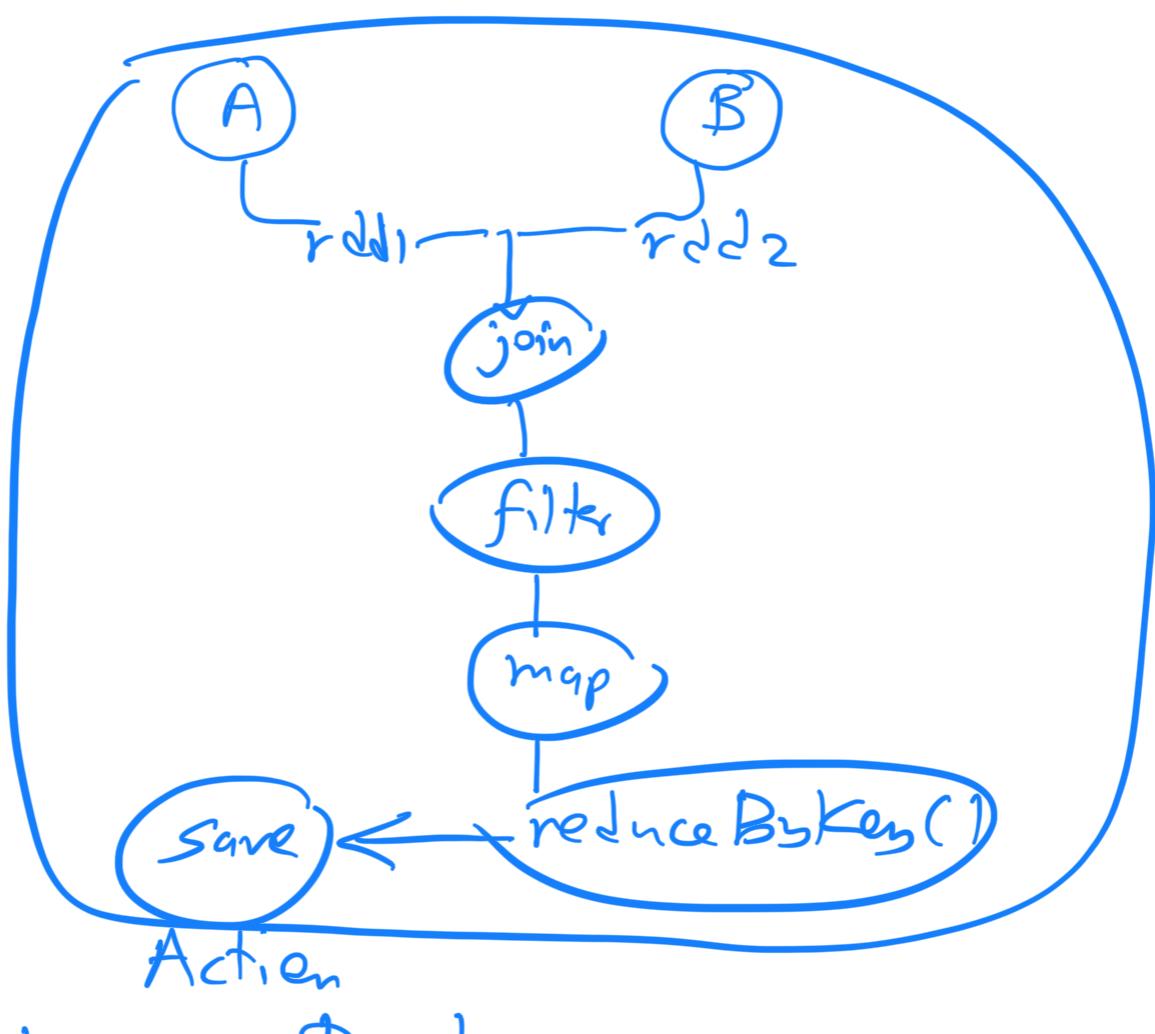
$$N = 20,000$$

each Part will have  $\approx 4000,000$

`target = source.mapPartitions(func)`

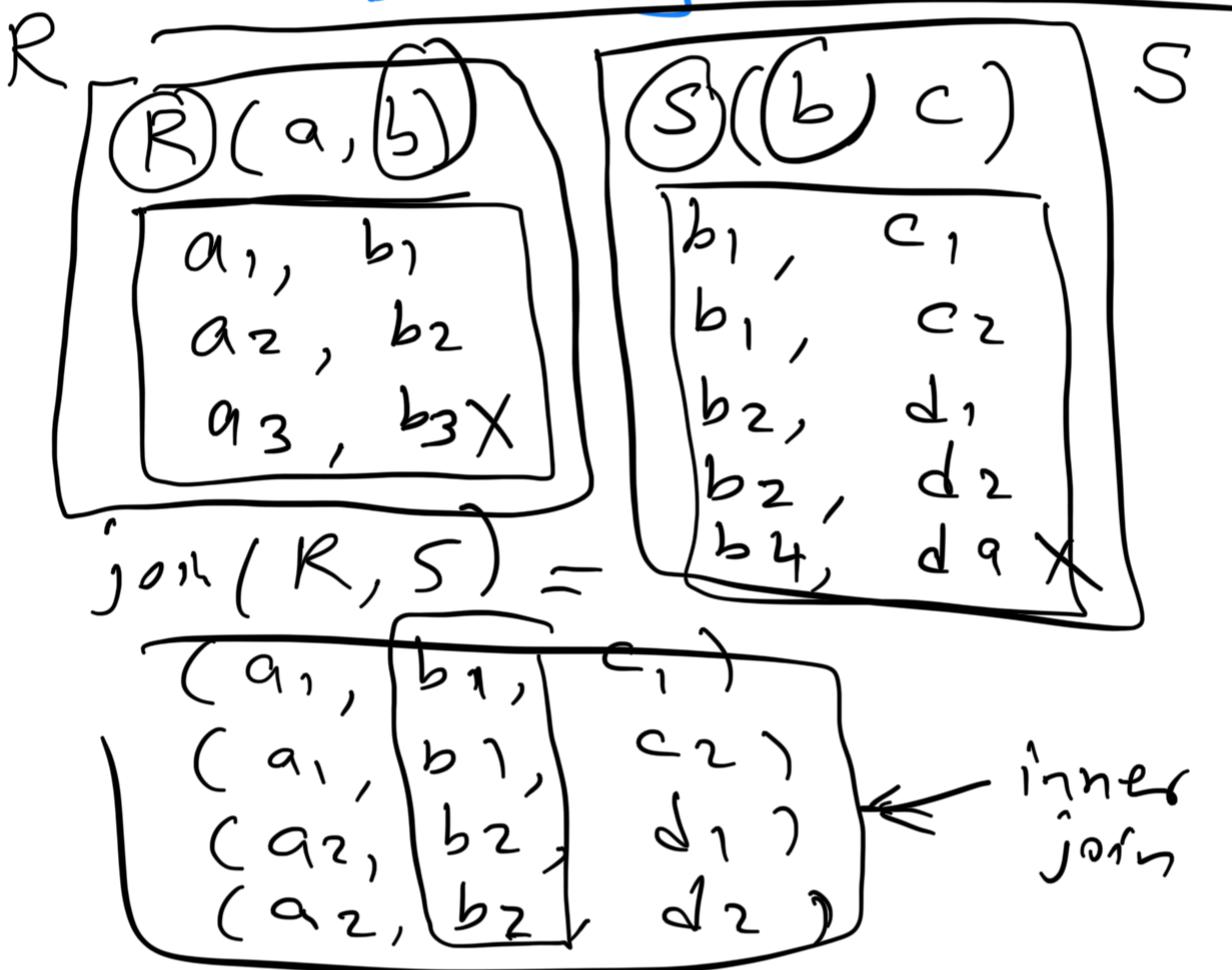
---

## Directed Acyclic Graph



" + . . . "

## Lazy Joining



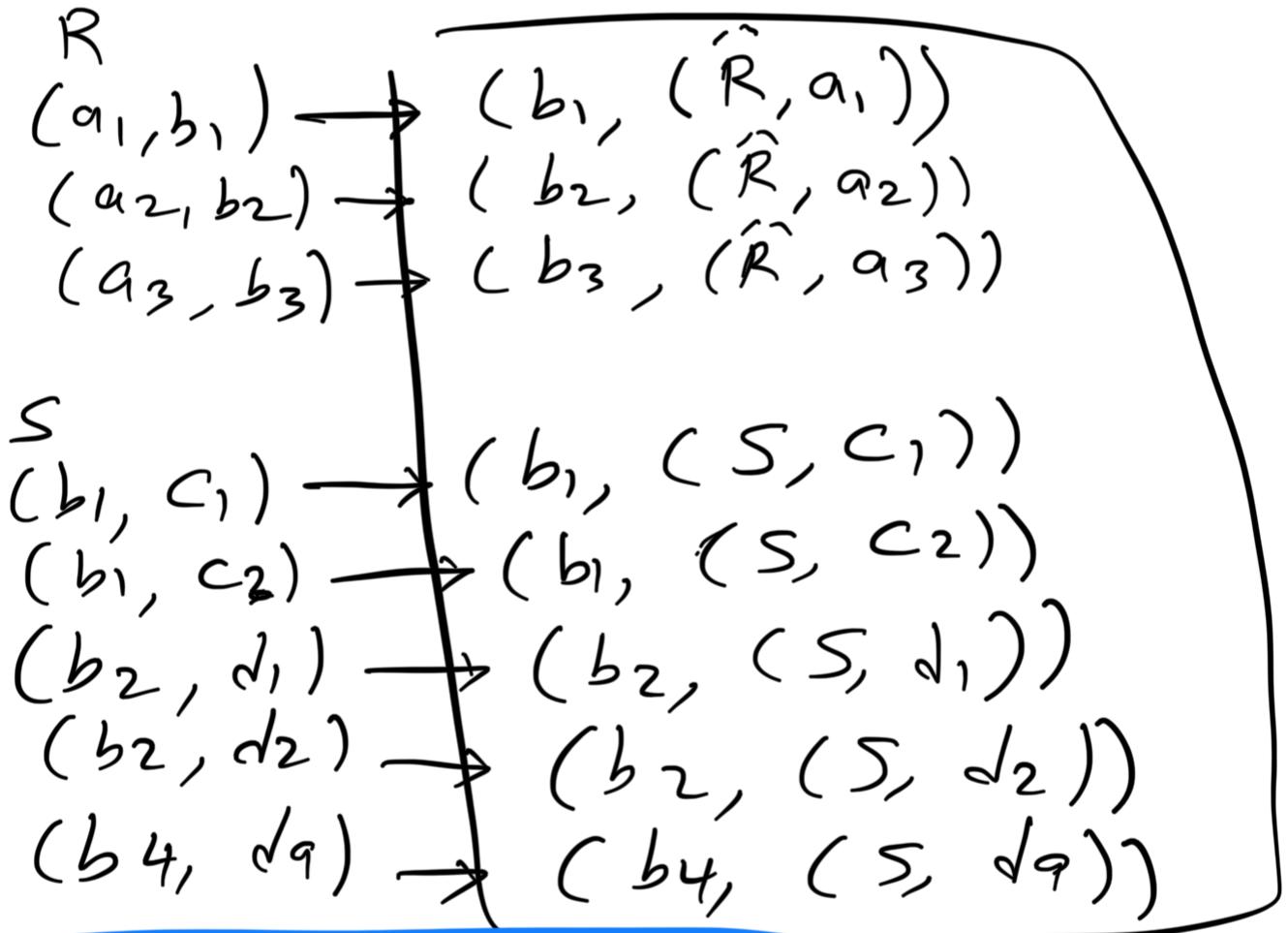
① map  $R$  record number, ignore actual record

```
map(k, v) {
    tokens = v.split(",")
    key = tokens[1]
    value = tokens[0]
    emit(key, ("R", value))
}
```

② map  $S$

```
map(k, v) {
    tokens = v.split(",")
}
```

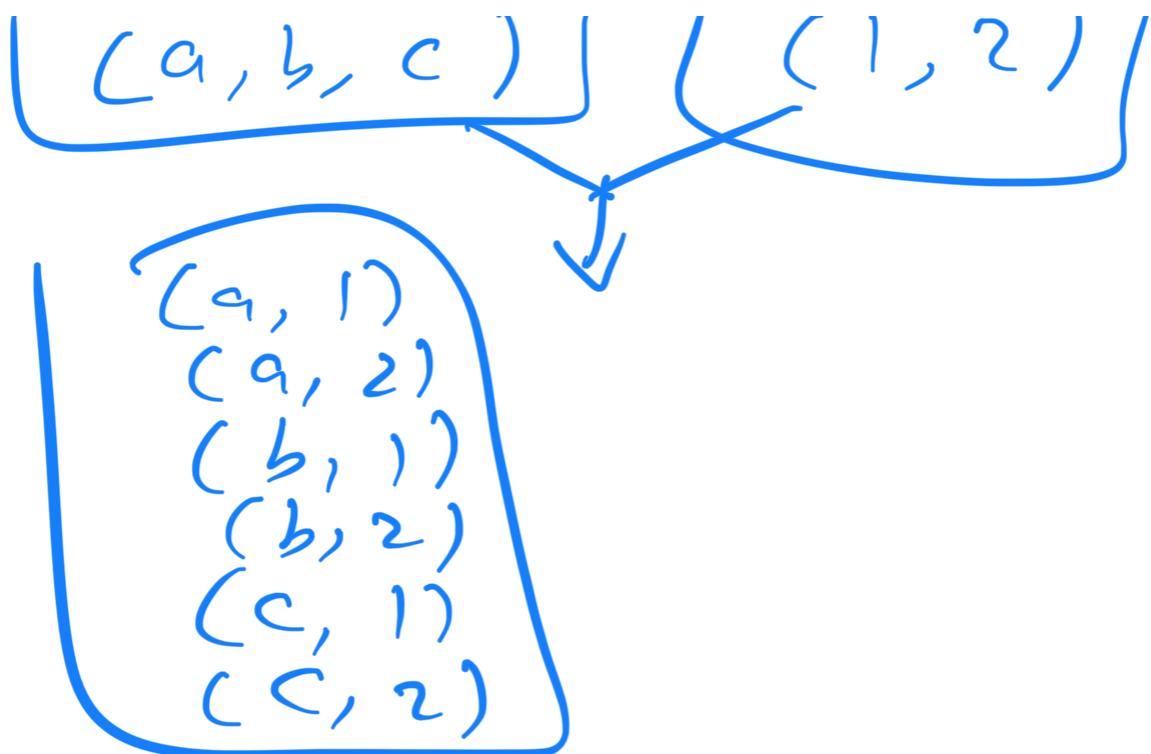
$\left\{ \begin{array}{l} \text{key} = \text{tokens}[0] \\ \text{value} = \text{tokens}[1] \\ \text{emit}(\text{key}, \underline{\text{"5"}}, \text{value}) \end{array} \right.$



$(b_1, [(R, a_1), (S, c_1), (S, c_2)])$   
 $(b_2, [(R, a_2), (S, d_1), (S, d_2)])$   
 $(b_3, [(R, a_3)])$   
 $(b_4, [(S, d_9)])$

drop

Cartesian Product



```

reduce(K, values) {
    # K = actual key
    # values = [ (R, v1) (R, v2), ... , (S, u1), (S, u2), ... ]
    RList = []
    SList = []
    for pair in values {
        relation = pair[0]
        value = pair[1]
        if (relation == 'R') {
            RList.append(value)
        } else {
            SList.append(value)
        }
    }
    return RList
}

```

# ~~new - 01~~

~~#~~ Rlist = (a,)

~~#~~ Slist = (c<sub>1</sub>, c<sub>2</sub>)

IF (len(R) == 0) {return}

if (len(S) == 0) {return}

~~#~~ Len(R) > 0 & len(S) > 0

~~#~~ find join by Cartesian Product

for (a in Rlist) {

    for (c in Slist) {

        emit (a, K, c)

}

}

# } end of reduce()