# Pair RDD

- Spark gives you more transformations if the RDD type is pair RDD.

- Spark can shuffle elements based on the first element of the pair

- Called Key-Value pair RDD

  - First element in pair is called "key"

  - Second element in pair is called "value"

# Create Pair RDD

```
pairs = sc.range(1, 100).map(lambda n: (n, 1))
```

```
val pairs = sc.range(1, 100).map(_ -> 1)
```

```
JavaRDD<Integer> numbers =
  sc.parallelize(Arrays.asList(1, 2, 3, 4, 5));

JavaPairRDD<Integer, Integer> pairs =
  rdd.mapToPair(n -> new Tuple2(n, 1))
```

- In Java we use Scala's Tuple2
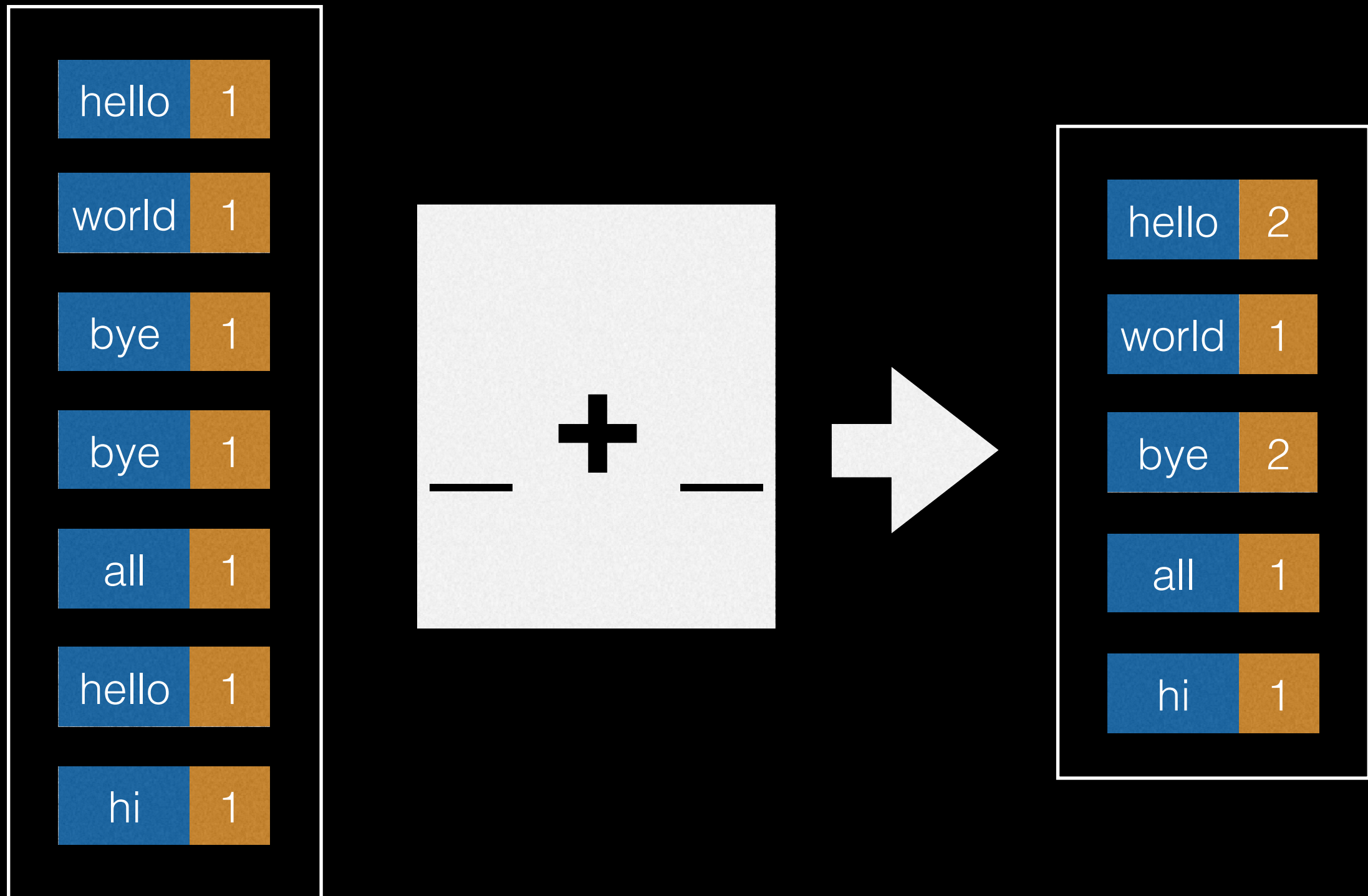
# Transformations

- reduceByKey

- groupByKey
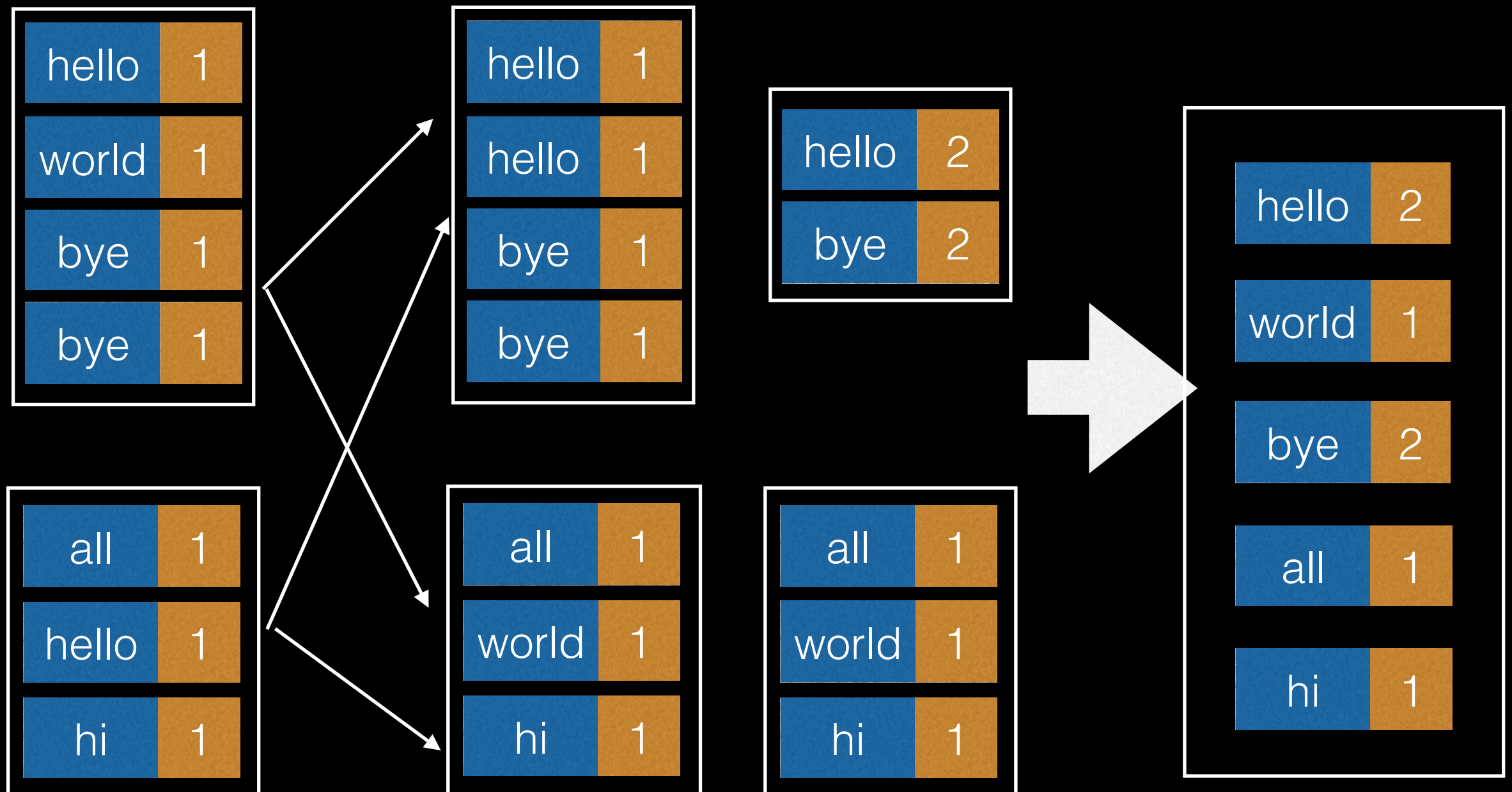
- mapValues

- keys

- values

# reduceByKey

- Combines values that have same key

- Runs reduce on each group alone

- Returns RDD consisting of each key and the reduced value for that key.
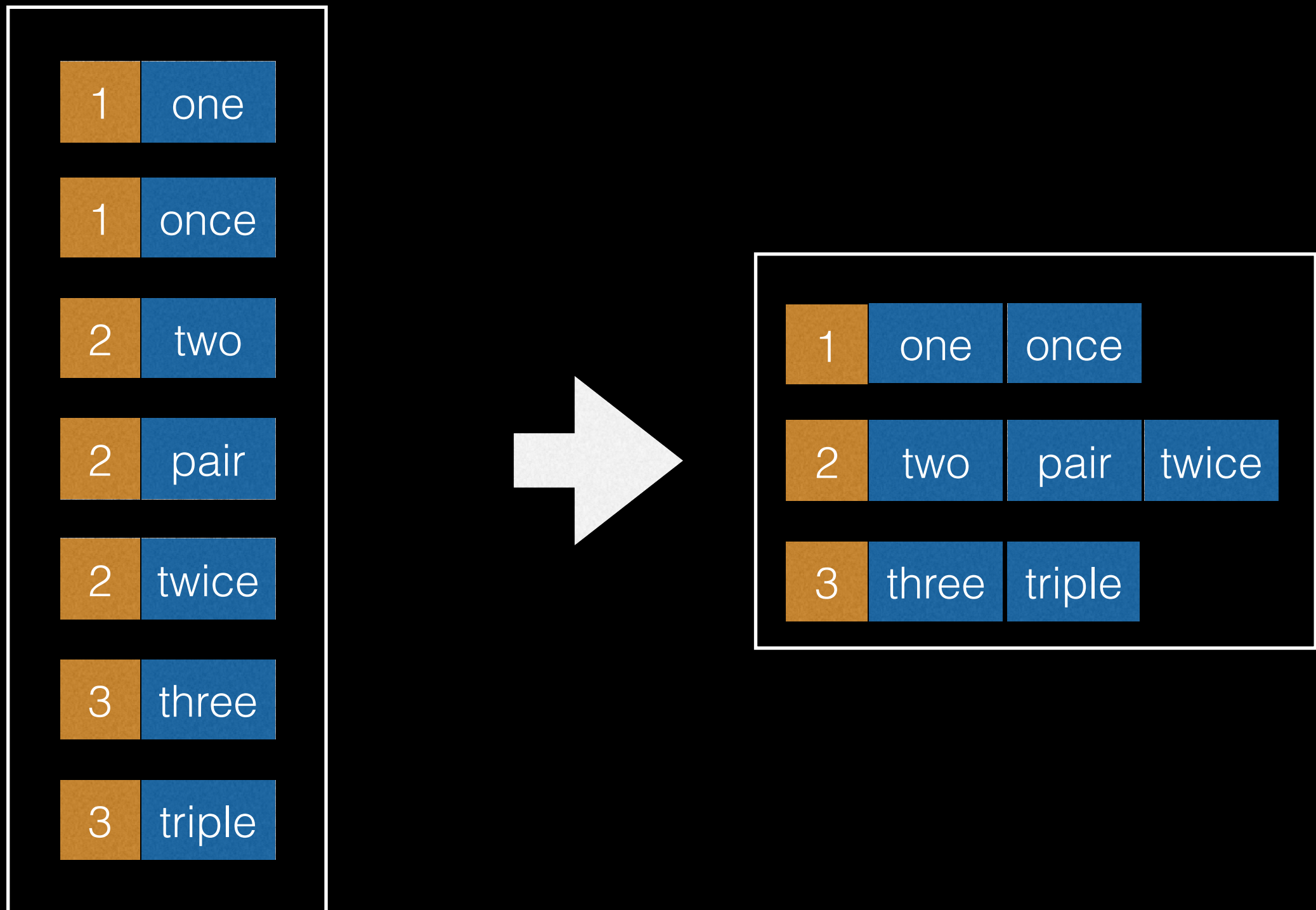
# reduceByKey

# reduceByKey

# Word Count

- Let's implement word count program

- 1. We need to read a text file

- 2. Generate RDD of words

- 3. Count occurrences of each word

# groupByKey

- Groups values that have same key

- Returns RDD consisting of each key and list of values.
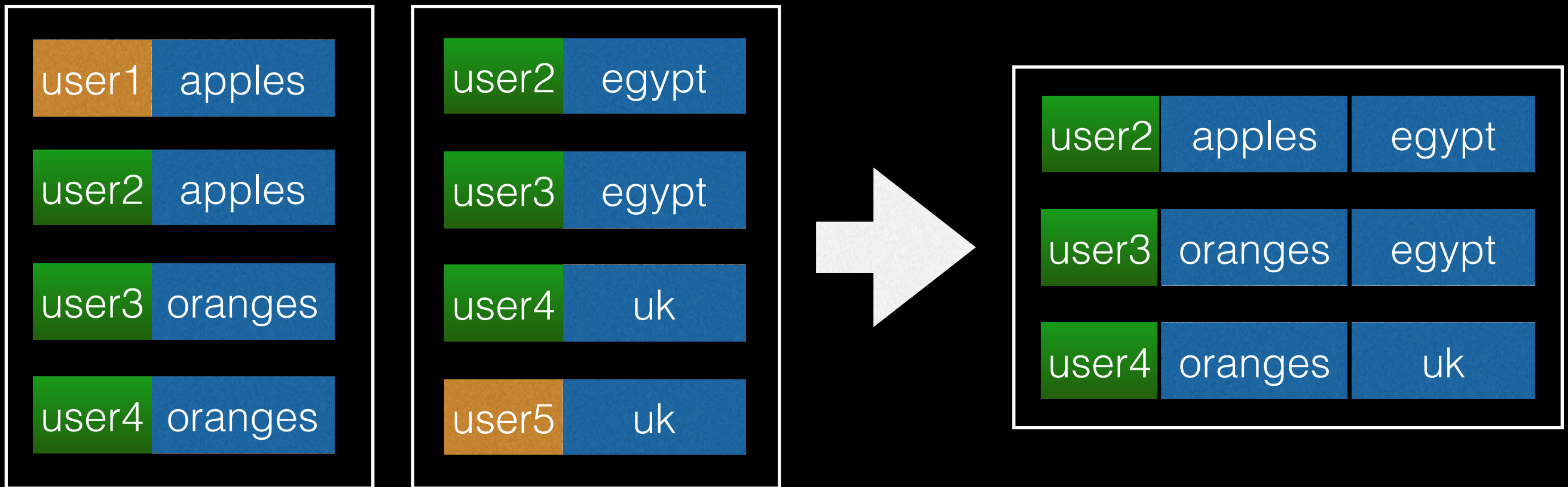
# groupByKey

# mapValues

- Applies a function to each value of RDD

- The key stays the same

- Better than "map", as it allows Spark to know that key is the same (No need to reshuffle)
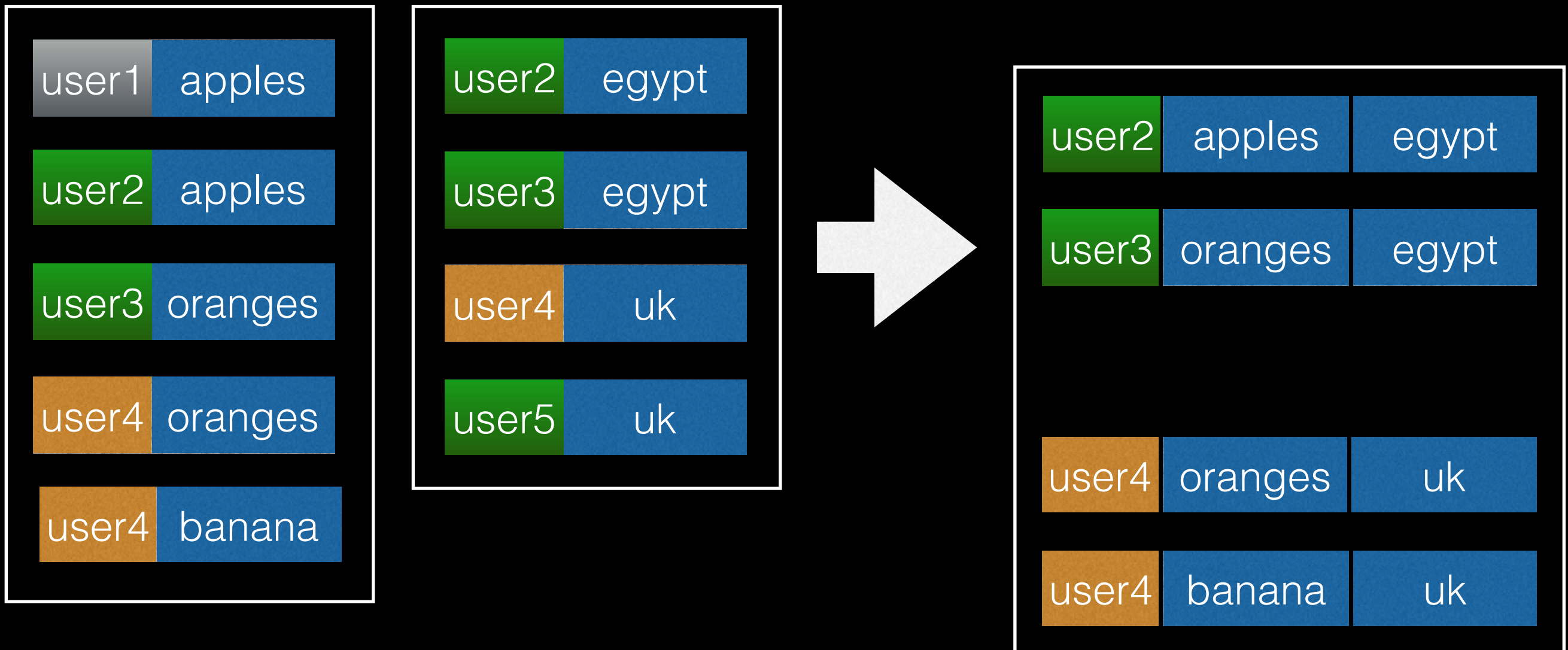
# Transformations

- subtractByKey
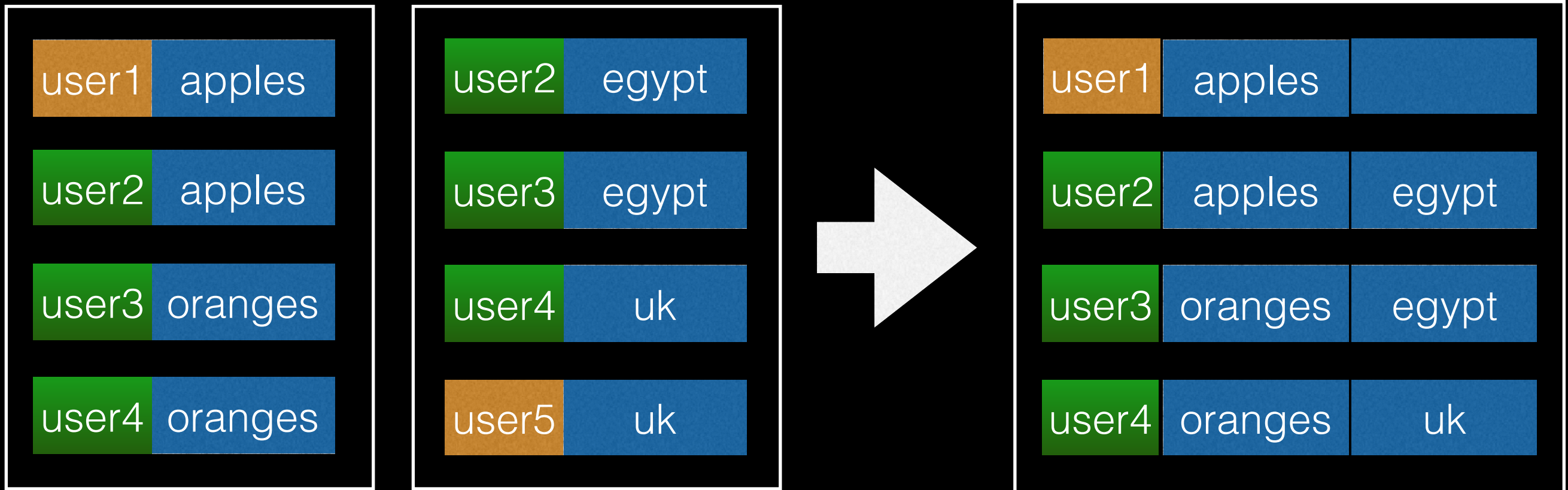
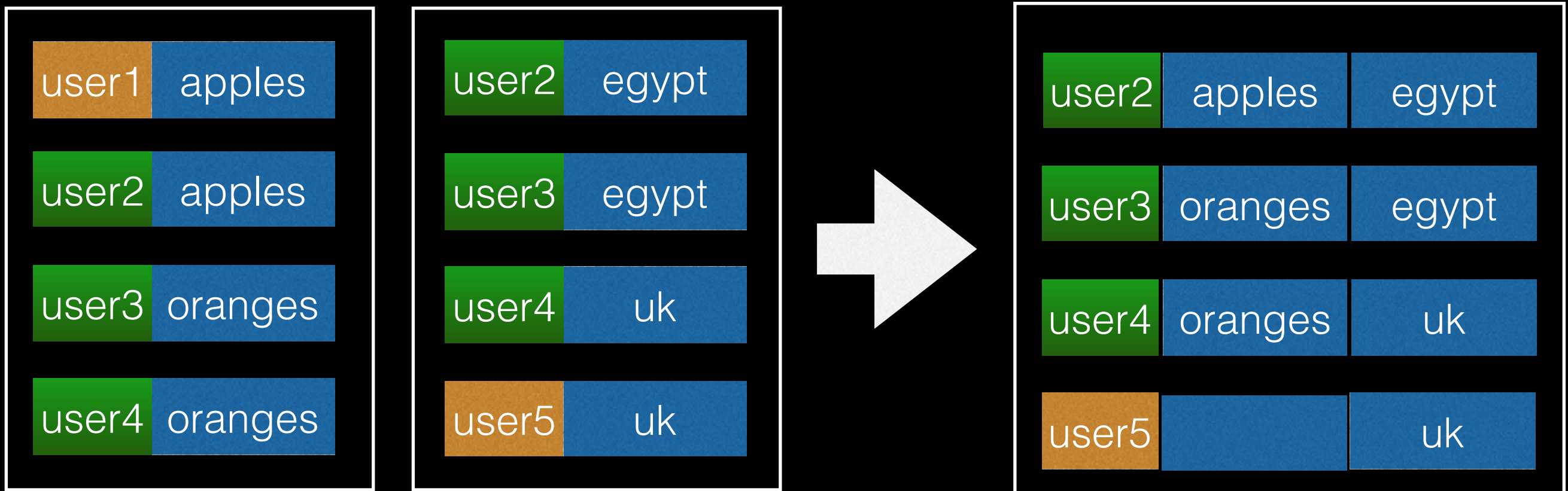- join

- rightOuterJoin

- leftOuterJoin
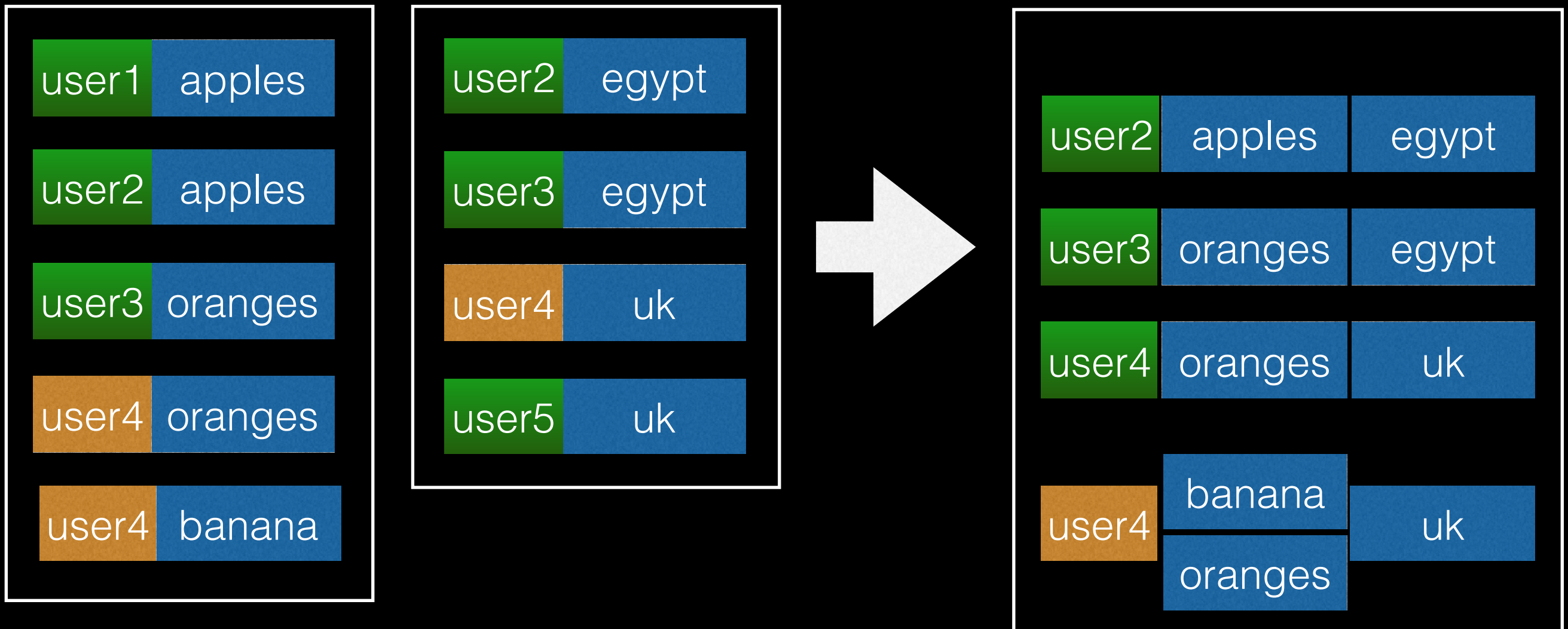
- cogroup

# join

# join

# leftOuterJoin

# rightOuterJoin

# cogroup

# Actions

- countByKey

- collectAsMap

- lookup(key)

# Page Rank

- Named after Larry Page

- assign a rank for each document based how many other documents links to it (and rank of those documents)

- Can also measure influence of users on social networks