

Yandex

MapReduce

Partitioner

World of MapReduce



Combiner

Partitioner

Comparator

World of MapReduce



Combiner

Partitioner

Comparator

World of MapReduce

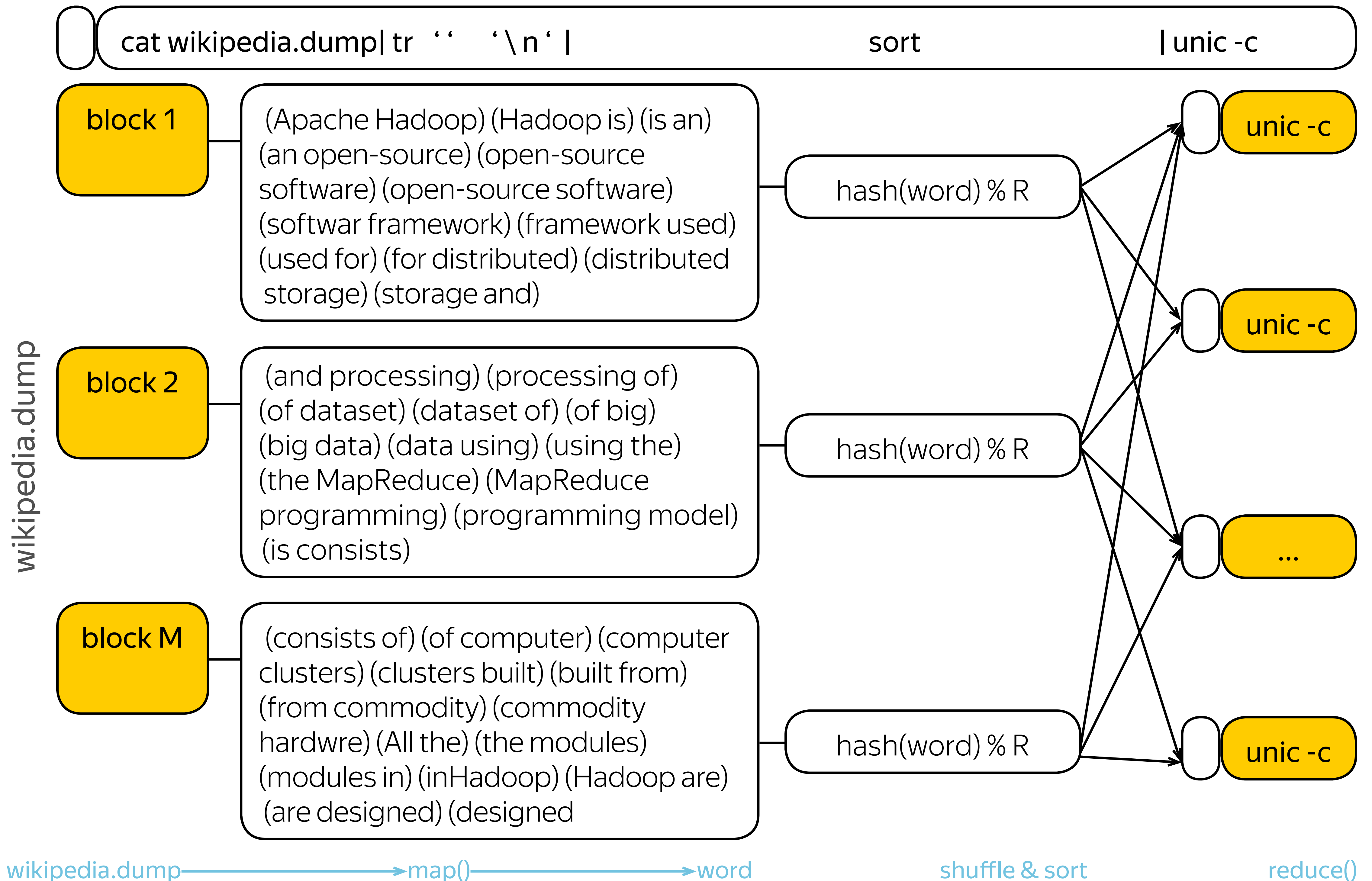


Combiner

Partitioner

Comparator

WordCount



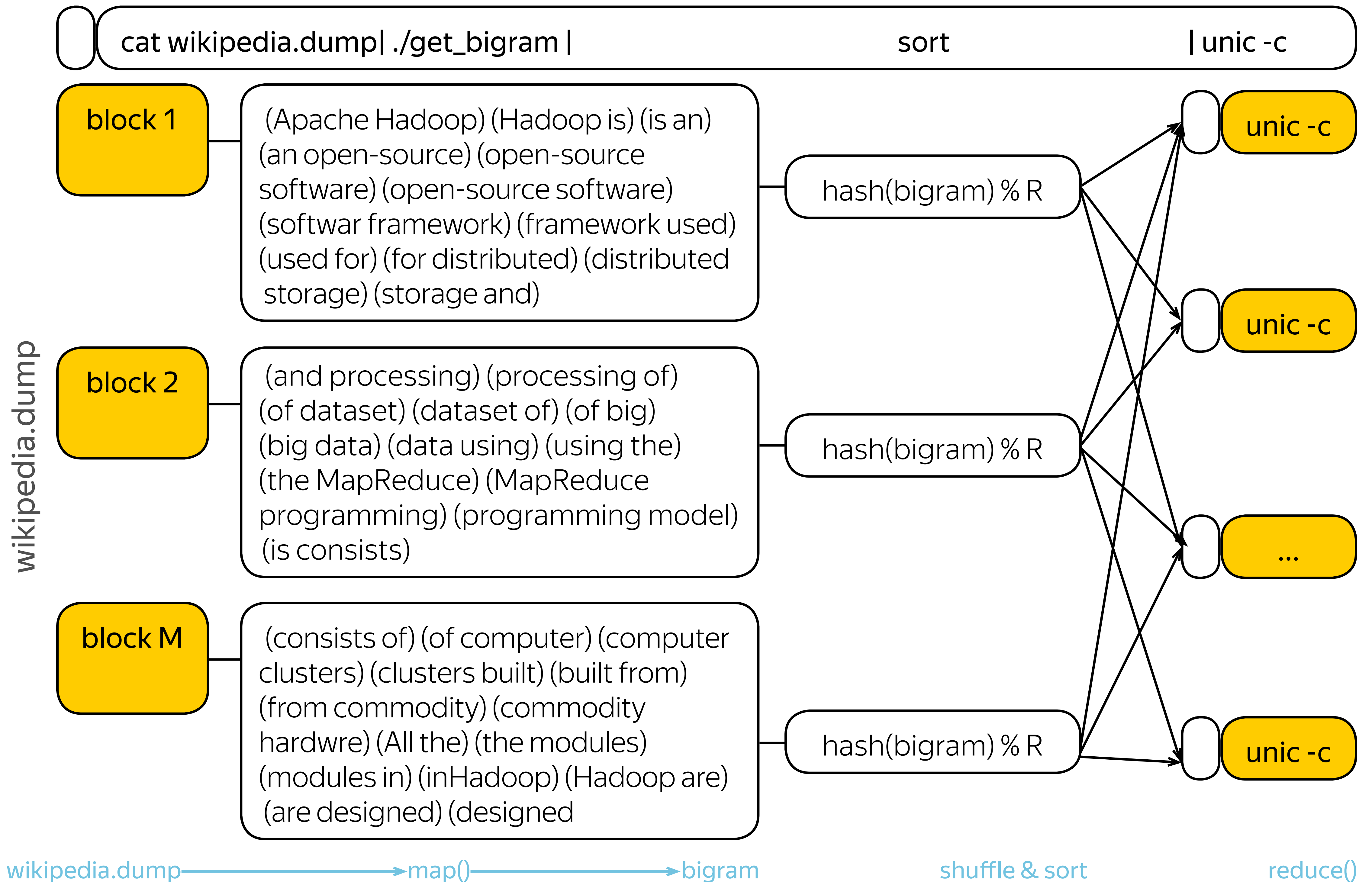
Collocations

natural English	unnatural English
the fast train	the quick train
fast food	quick food
a quick shower	a fast shower
a quick meal	a fast meal

Collocations

natural English	unnatural English	New and York
the fast train	the quick train	
fast food	quick food	
a quick shower	a fast shower	United and States
a quick meal	a fast meal	

WordCount



input

```
$ head -c 100 wikipedia_sample.txt  
12 Anarchism Anarchism is often defined as a political  
philosophy which holds the state to ...
```



Mapper (Python): bigram_mapper.py

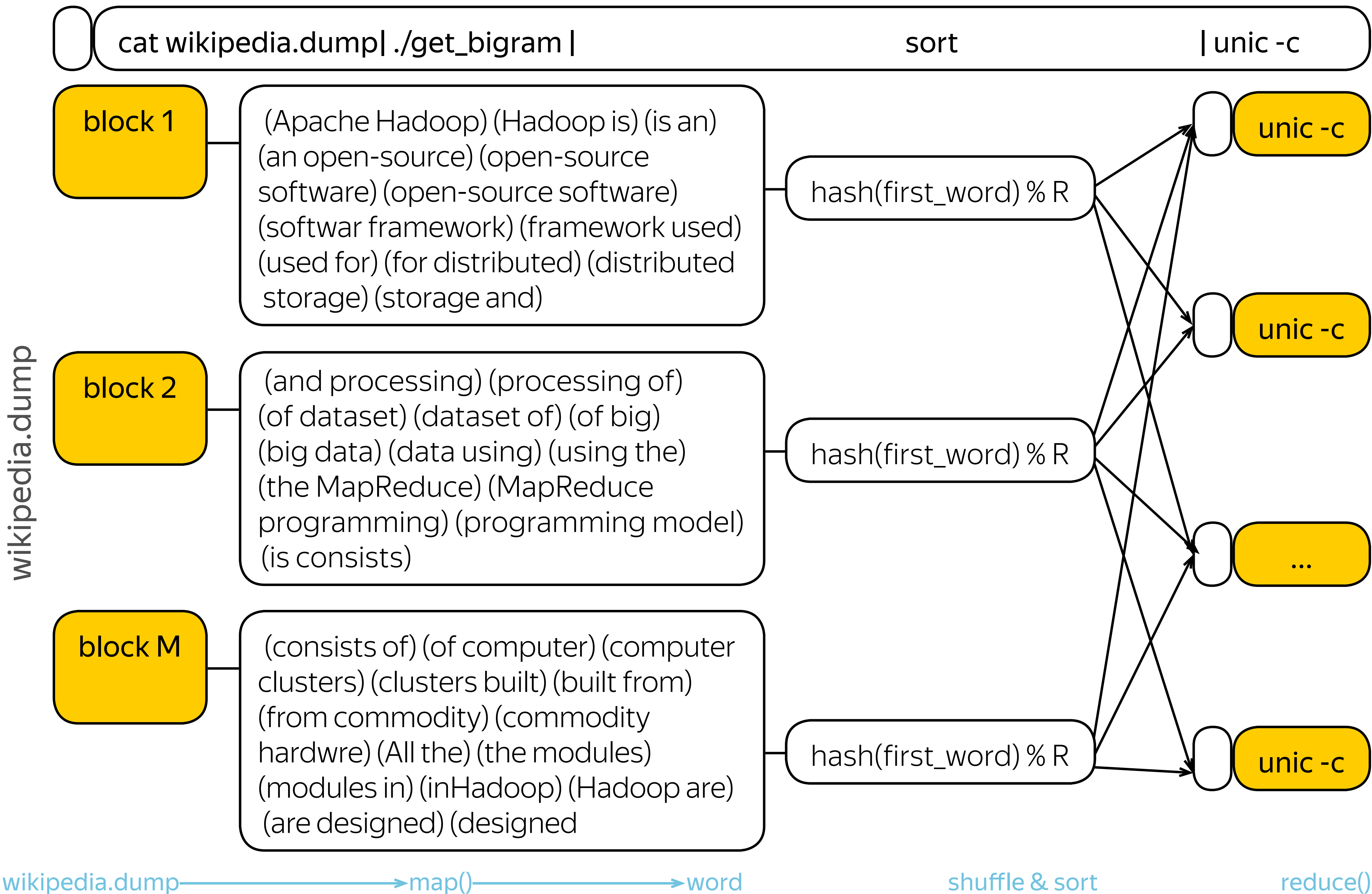
```
from __future__ import print_function  
import re  
import sys  
for line in sys.stdin:  
    article_id, content = line.split("\t", 1)  
    words = re.split("\W+", content)  
    for index in range(len(words) - 1):  
        print(words[index], words[index + 1], 1, sep="\t")
```



output

```
$ head word_count/part-00000  
Anarchism Anarchism 1  
Anarchism is 1  
is often 1
```

WordCount



Mapper (Python): inmemroy_bigram_reducer.py

```
from __future__ import print_function
from collections import Counter
import sys
```

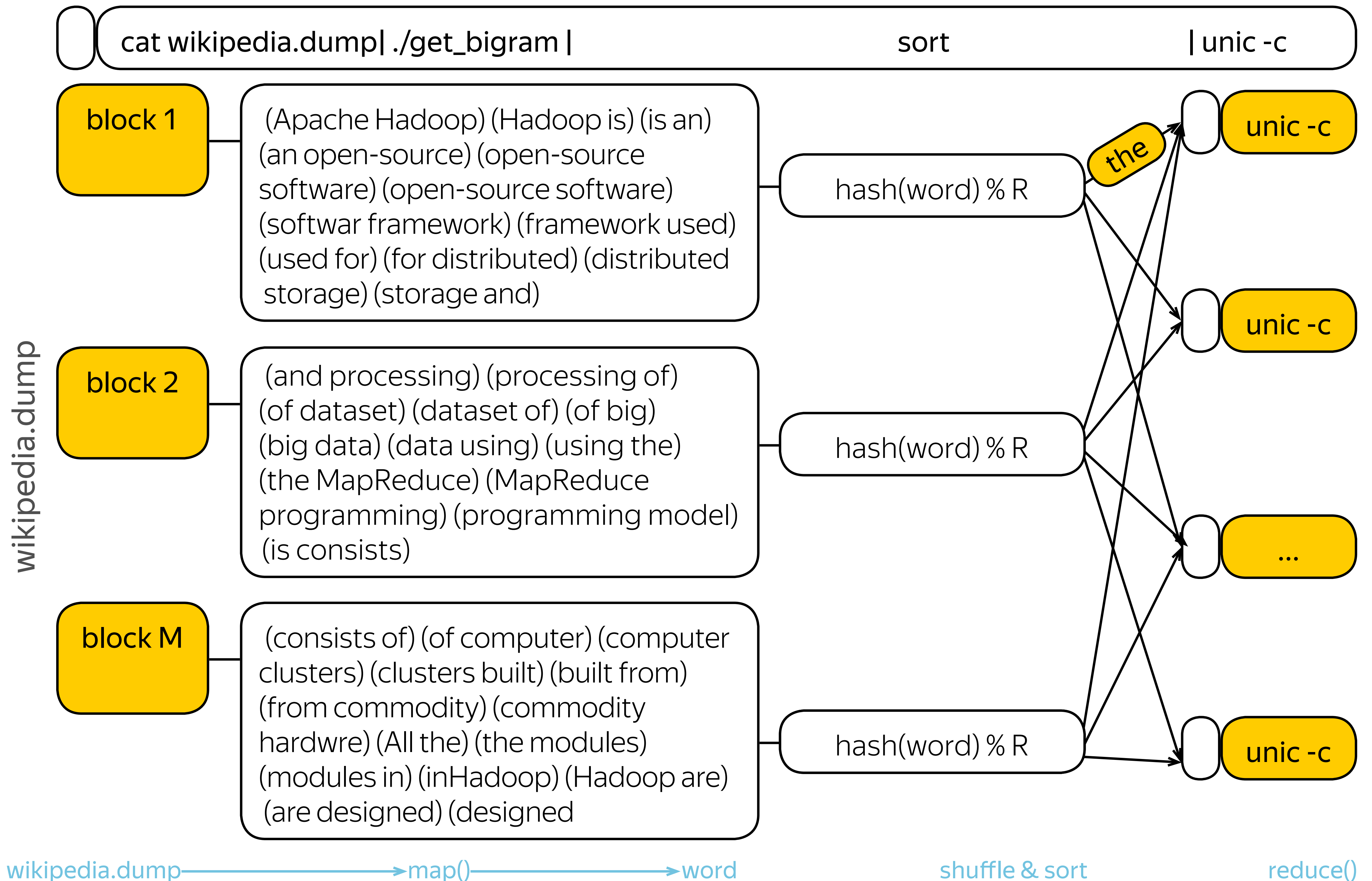
```
current_word = None
bigram_count = Counter()
```

```
for line in sys.stdin:
    first_word, second_word, counts = line.split("\t", 2)
    counts = Counter({second_word: int(counts)})
    if first_word == current_word:
        bigram_count += counts
    else:
        if current_word:
            for second_word, bigram_count in bigram_count.items():
                print(current_word, second_word,
```

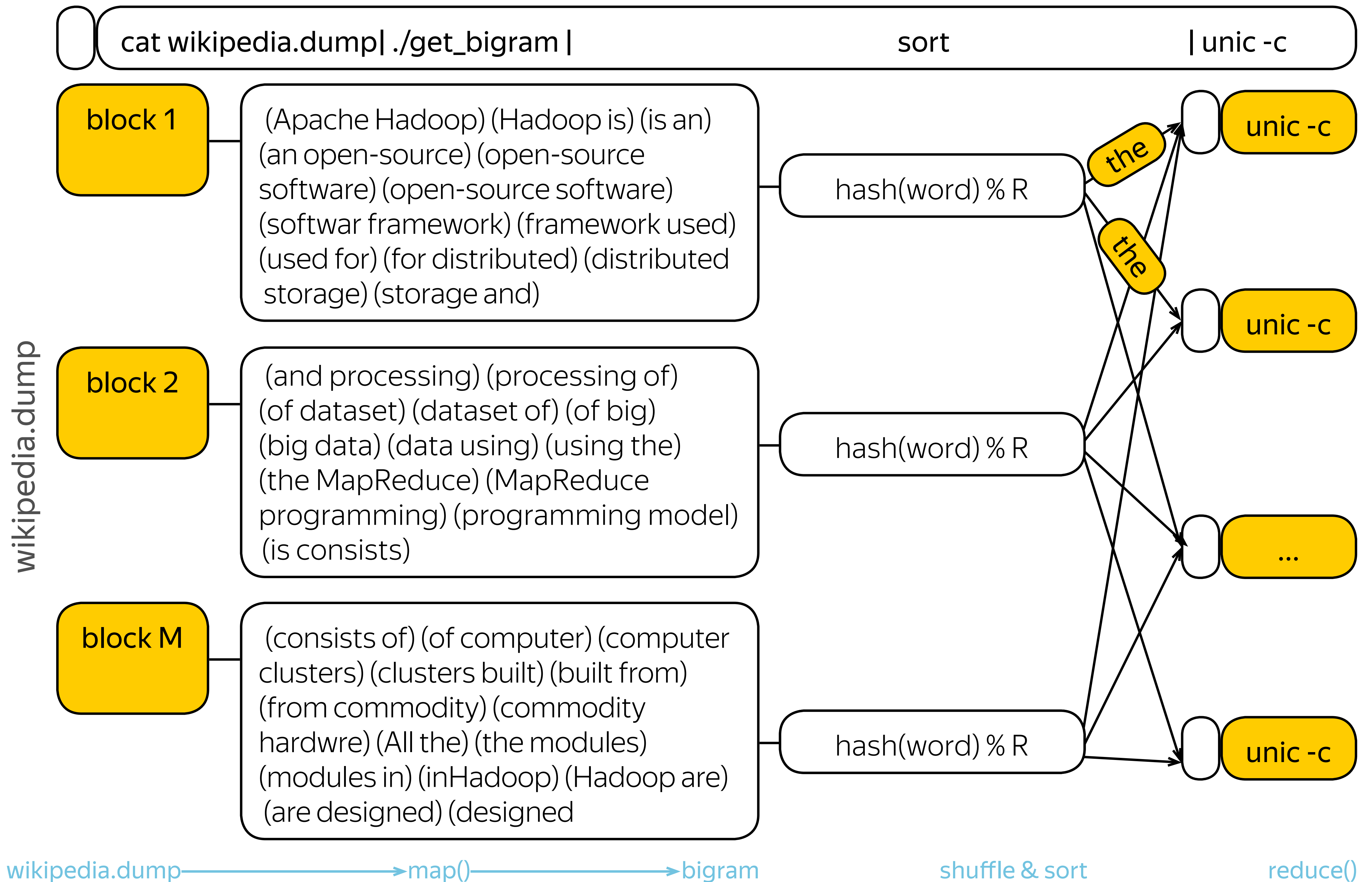
```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-files bigram_mapper.py,inmemory_bigram_reducer.py \  
-mapper 'python bigram_mapper.py' \  
-reducer 'python inmemory_bigram_reducer.py' \  
-numReduceTasks 5 \  
-input wikipedia_sample.txt \  
-output word_count \
```

```
$ grep $'^new\t' word_count/* | sort -nrk3,3 | head -4  
word_count/part-00001:new    york    112  
word_count/part-00001:new    world   15  
word_count/part-00001:new    jewrsey 12  
word_count/part-00001:new    constitution 12  
$ grep $'\tyork\t' word_count/* | sort -nrk3,3 | head  
word_count/part-00001:new    york    112  
word_count/part-00004:sergeant york    1
```


WordCount



WordCount



new york

united states

south korea

kota kinabalu

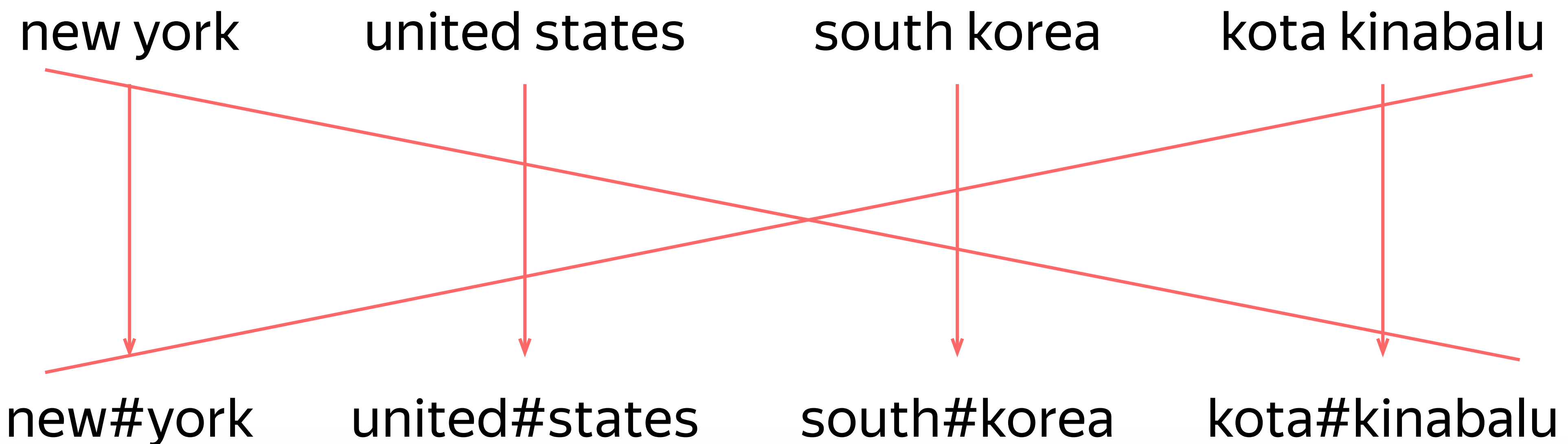


new#york

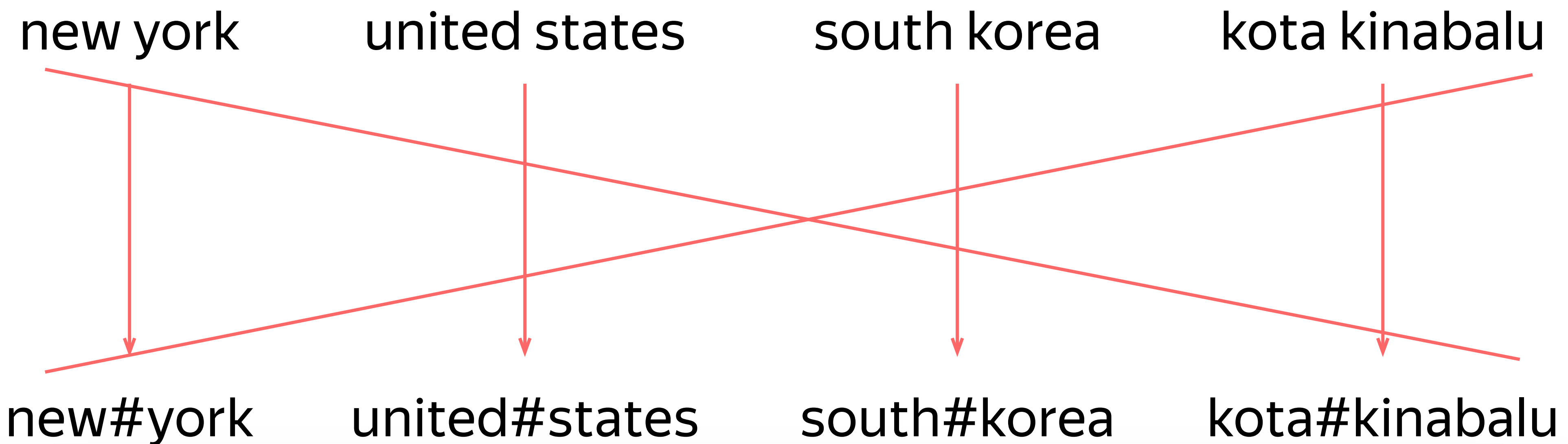
united#states

south#korea

kota#kinabalu



```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
→ -D stream.num.map.output.key.fields=2 \  
-files bigram_mapper.py,bigram_reducer.py \  
-mapper 'python bigram_mapper.py' \  
-reducer 'python bigram_reducer.py' \  
-numReduceTasks 5 \  
-input wikipedia_sample.txt \  
-output word_count \
```

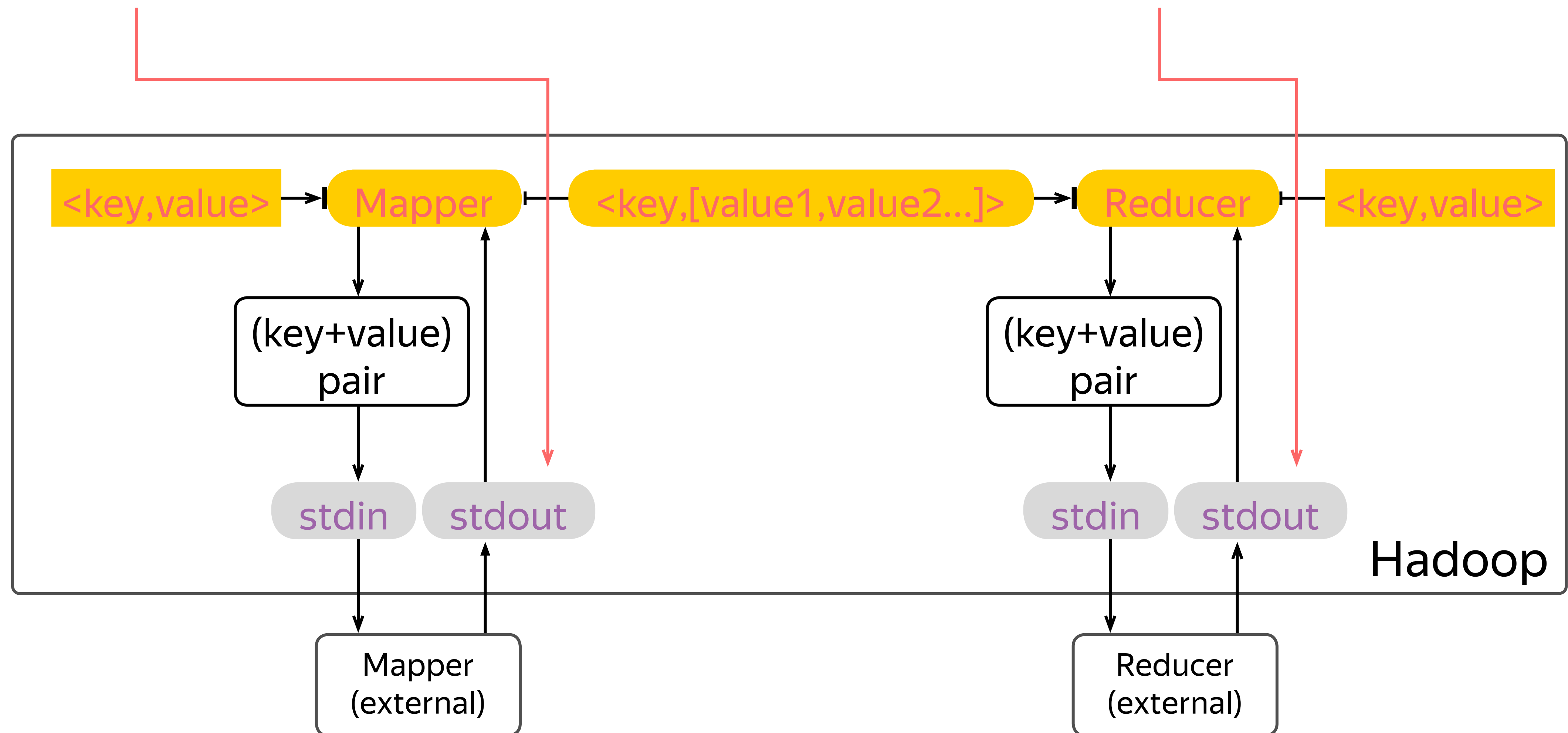


```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
→ -D stream.num.map.output.key.fields=2 \  
-files bigram_mapper.py,bigram_reducer.py \  
-mapper 'python bigram_mapper.py' \  
-reducer 'python bigram_reducer.py' \  
-numReduceTasks 5 \  
-input wikipedia_sample.txt \  
-output word_count
```

```
$ grep $'^new\t' word_count/* | sort -nrk3,3 | head -4  
word_count/part-00001:new york 112  
word_count/part-00001:new world 15  
word_count/part-00002:new constitution 12  
word_count/part-00000:new jersey 12
```

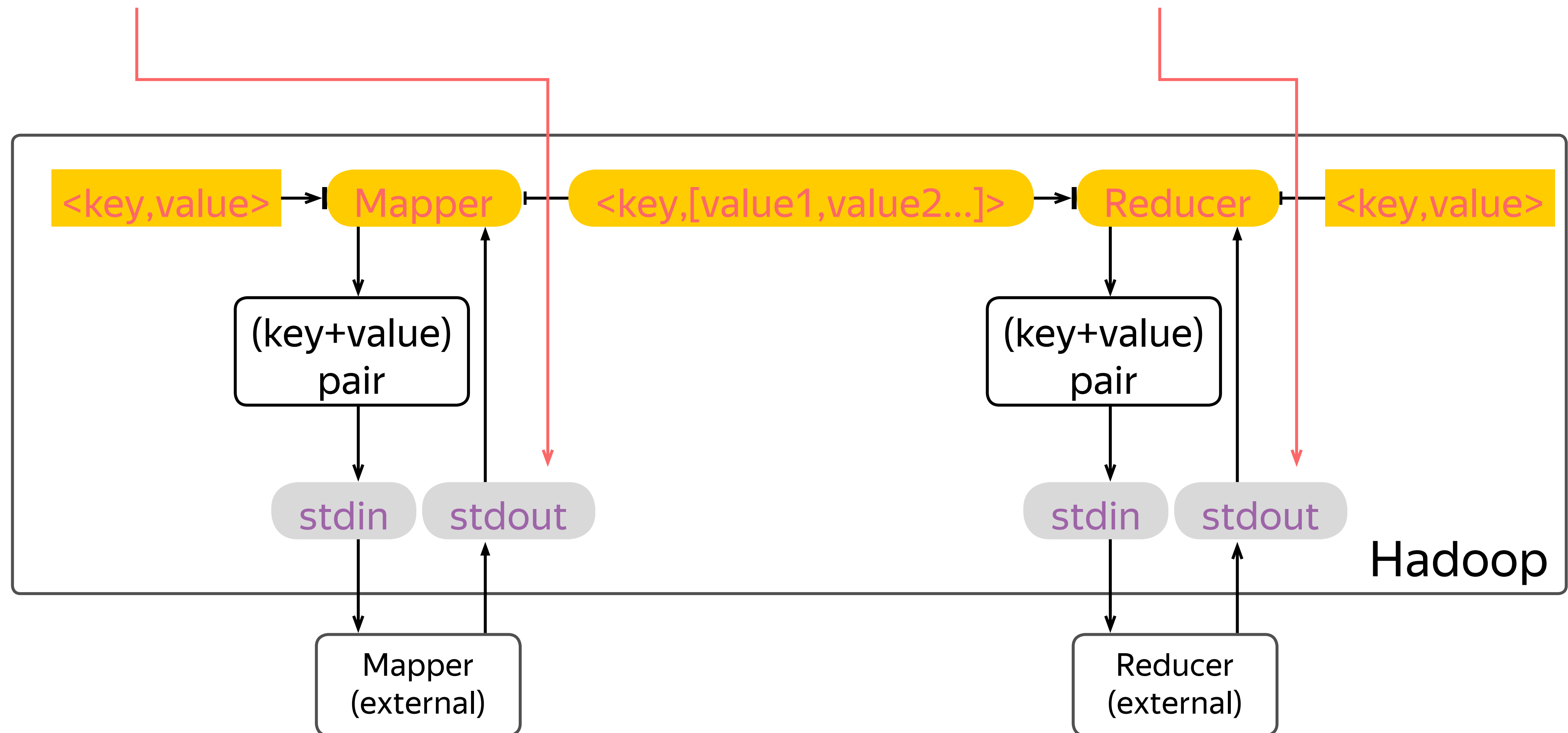
-D stream.map.output.field.separator=.
-D stream.num.map.output.key.fields=1

-D stream.reduce.output.field.separator=.
-D stream.num.reduce.output.key.fields=2



-D stream.map.output.field.separator=.
-D stream.num.map.output.key.fields=1

-D stream.reduce.output.field.separator=.
-D stream.num.reduce.output.key.fields=2



IPv4 address (example): 69.89.31.226

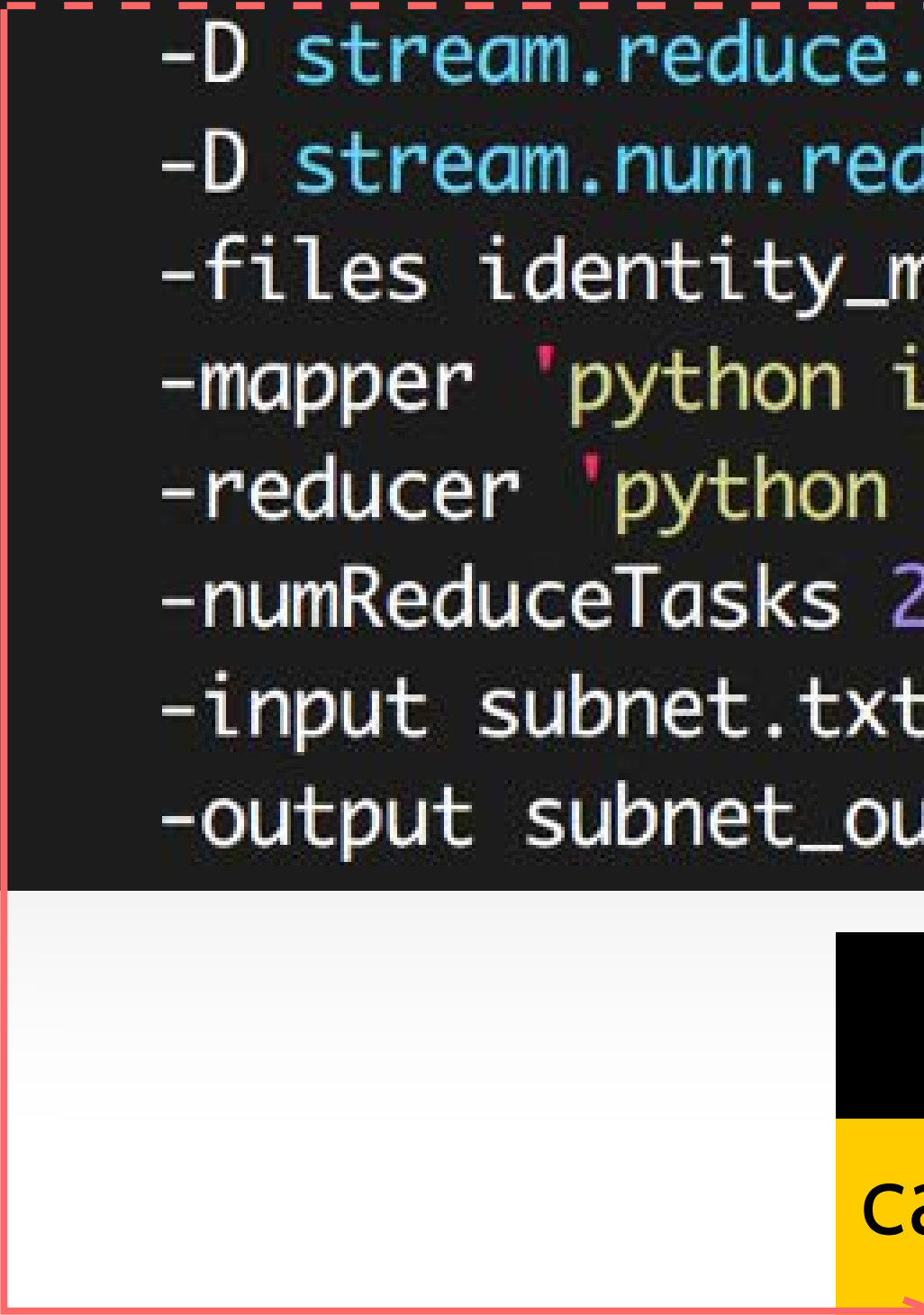
input

```
$ cat subnet.txt  
1.2.3.4  
2.3.4.5  
3.4.5.6  
4.5.6.7
```

```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-D stream.map.output.field.separator=. \  
-D stream.num.map.output.key.fields=1 \  
-D stream.reduce.output.field.separator=. \  
-D stream.num.reduce.output.key.fields=2 \  
-files identity_mr.py \  
-mapper 'python identity_mr.py' \  
-reducer 'python identity_mr.py' \  
-numReduceTasks 2 \  
-input subnet.txt \  
-output subnet_out
```

output

```
cat subnet_out/*  
2 3.4 5  
4 5.6 7  
1 2.3 4  
3 4.5 6
```




input

```
$ cat subnet.txt  
1.2.3.4  
2.3.4.5  
3.4.5.6  
4.5.6.7
```

```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-D stream.map.output.field.separator=. \  
-D stream.num.map.output.key.fields=1 \  
-D stream.reduce.output.field.separator=. \  
-D stream.num.reduce.output.key.fields=2 \  
-files identity_mr.py \  
-mapper 'python identity_mr.py' \  
-reducer 'python identity_mr.py' \  
-numReduceTasks 2 \  
-input subnet.txt \  
-output subnet_out
```

output

```
cat subnet_out/*  
2 3.4 5  
4 5.6 7  
1 2.3 4  
3 4.5 6
```



input

```
$ cat subnet.txt  
1a.2.3.4  
2a.3.4.5  
3b.4.5.6  
4b.5.6.7
```

```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-D map.output.key.field.separator=. \  
-D mapreduce.partition.keypartitioner.options=-k1.2,1.2 \  
-files identity_mr.py \  
-mapper 'python identity_mr.py' \  
-reducer 'python identity_mr.py' \  
-numReduceTasks 2 \  
-input subnet.txt \  
-output subnet_out \  
-partitioner org.apache.hadoop.mapred.lib.KeyFieldBasedPartitioner
```

output

```
cat subnet_out/*  
3b.4.5.6  
4b.5.6.7  
1a.2.3.4  
2a.3.4.5
```

```
ls -lth subnet_out  
total 8.0K  
0 Apr 1 14:59 _SUCCESS  
20 Apr 1 14:59 part-00001  
20 Apr 1 14:59 part-00000
```

input

```
$ cat subnet.txt  
1a.2.3.4  
2a.3.4.5  
3b.4.5.6  
4b.5.6.7
```

```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-D map.output.key.field.separator=. \  
-D mapreduce.partition.keypartitioner.options=-k1.2,1.2 \  
-files identity_mr.py \  
-mapper 'python identity_mr.py' \  
-reducer 'python identity_mr.py' \  
-numReduceTasks 2 \  
-input subnet.txt \  
-output subnet_out \  
-partitioner org.apache.hadoop.mapred.lib.KeyFieldBasedPartitioner
```

output

```
cat subnet_out/*  
3b.4.5.6  
4b.5.6.7  
1a.2.3.4  
2a.3.4.5
```

```
ls -lth subnet_out  
total 8.0K  
0 Apr 1 14:59 _SUCCESS  
20 Apr 1 14:59 part-00001  
20 Apr 1 14:59 part-00000
```


input

```
$ cat subnet.txt  
1a.2.3.4  
2a.3.4.5  
3b.4.5.6  
4b.5.6.7
```

```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-D map.output.key.field.separator=. \  
-D mapreduce.partition.keypartitioner.options=-k1.2,1.2 \  
-files identity_mr.py \  
-mapper 'python identity_mr.py' \  
-reducer 'python identity_mr.py' \  
-numReduceTasks 2 \  
-input subnet.txt \  
-output subnet_out \  
-partitioner org.apache.hadoop.mapred.lib.KeyFieldBasedPartitioner
```

output

```
cat subnet_out/*  
3b.4.5.6  
4b.5.6.7  
1a.2.3.4  
2a.3.4.5
```

```
ls -lth subnet_out  
total 8.0K  
0 Apr 1 14:59 _SUCCESS  
20 Apr 1 14:59 part-00001  
20 Apr 1 14:59 part-00000
```

input

```
$ cat subnet.txt  
1a.2.3.4  
2a.3.4.5  
3b.4.5.6  
4b.5.6.7
```

man sort | grep KEYDEF




```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
-D map.output.key.field.separator=. \  
-D mapreduce.partition.keypartitioner.options=-k1.2,1.2 \  
-files identity_mr.py \  
-mapper 'python identity_mr.py' \  
-reducer 'python identity_mr.py' \  
-numReduceTasks 2 \  
-input subnet.txt \  
-output subnet_out \  
-partitioner org.apache.hadoop.mapred.lib.KeyFieldBasedPartitioner
```

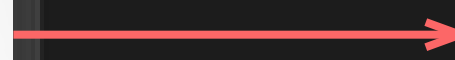
output

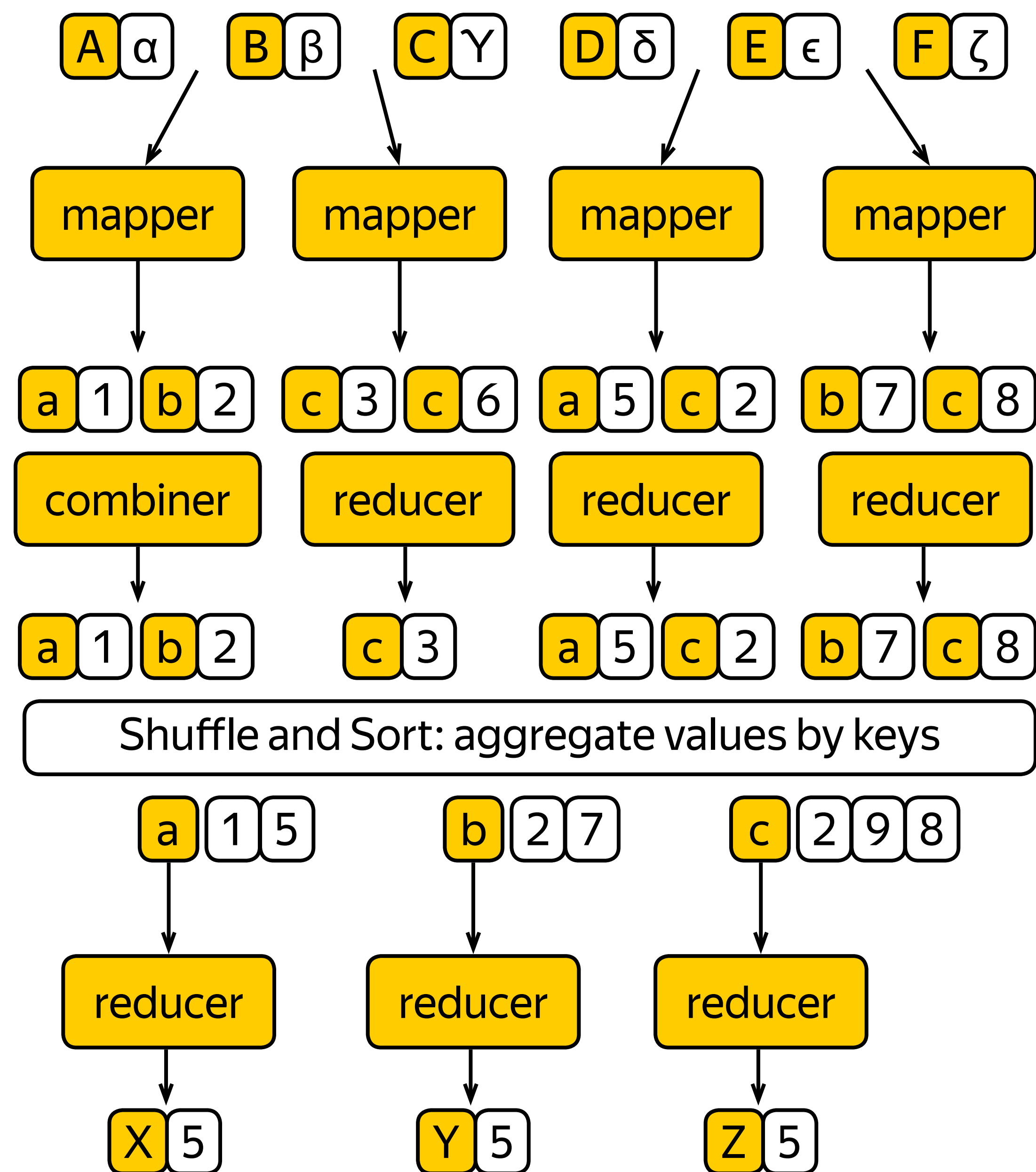
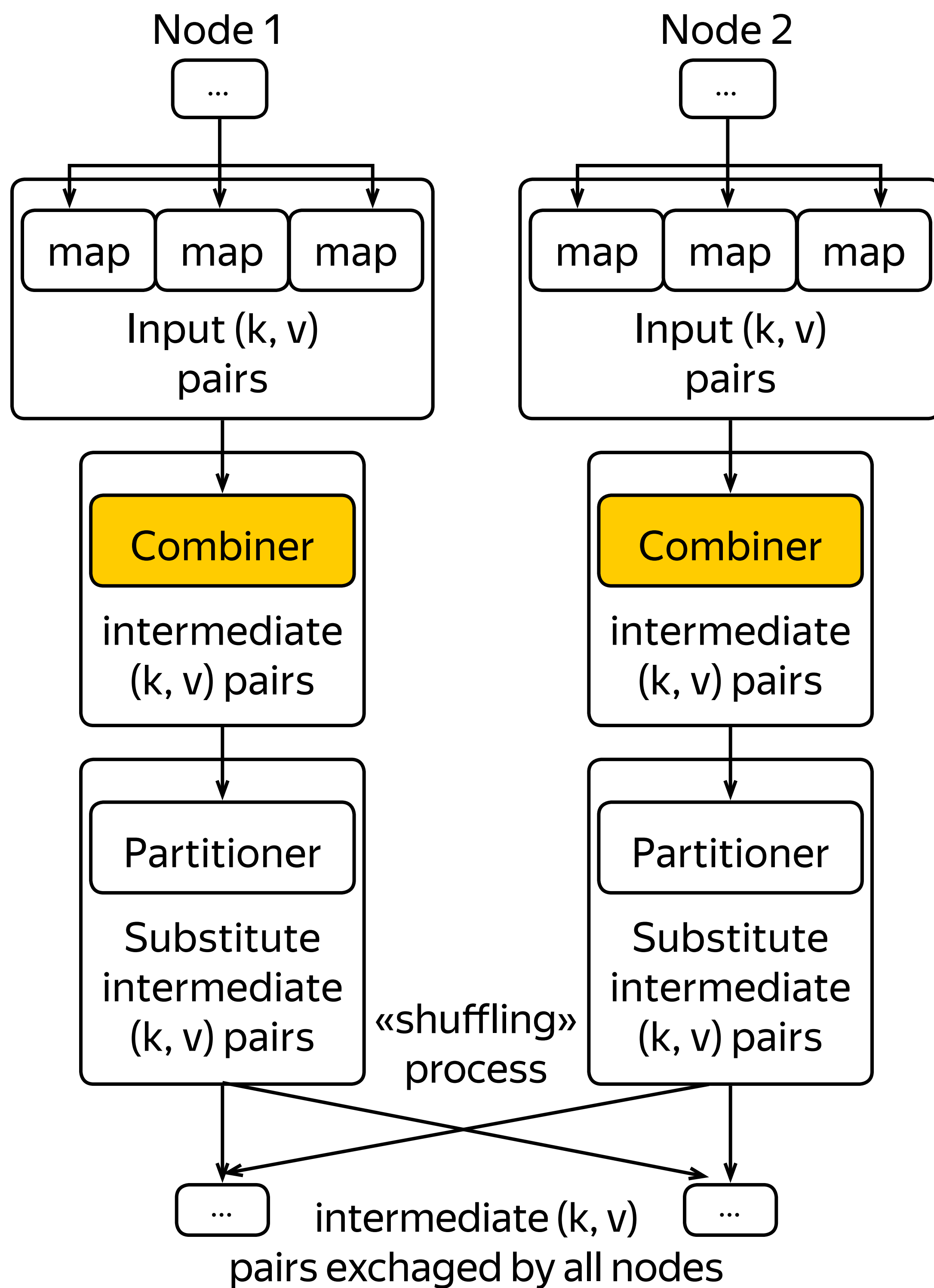
```
cat subnet_out/*  
3b.4.5.6  
4b.5.6.7  
1a.2.3.4  
2a.3.4.5
```

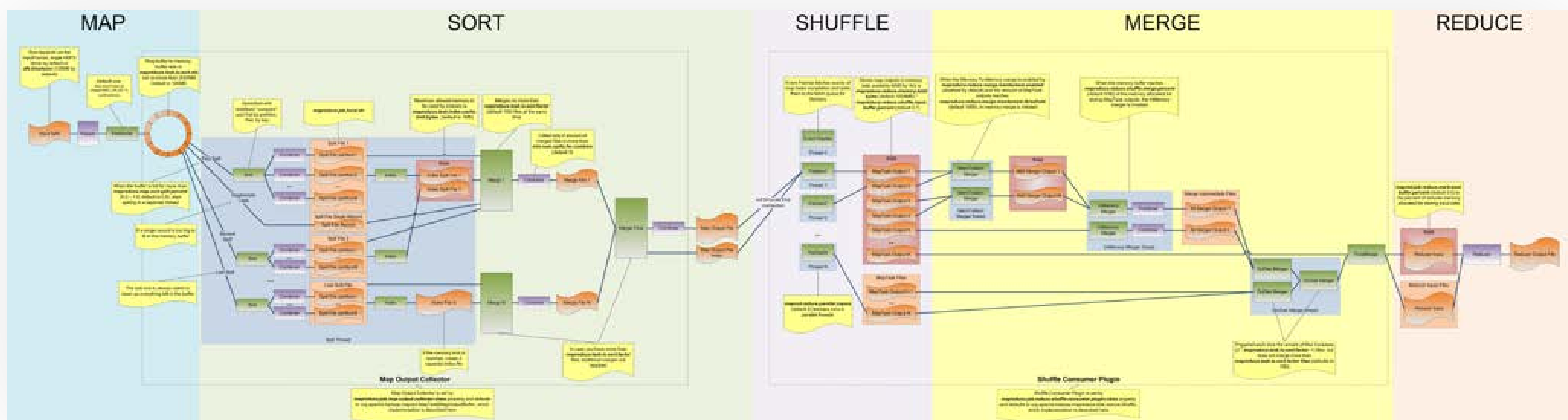
```
ls -lth subnet_out  
total 8.0K  
0 Apr 1 14:59 _SUCCESS  
20 Apr 1 14:59 part-00001  
20 Apr 1 14:59 part-00000
```

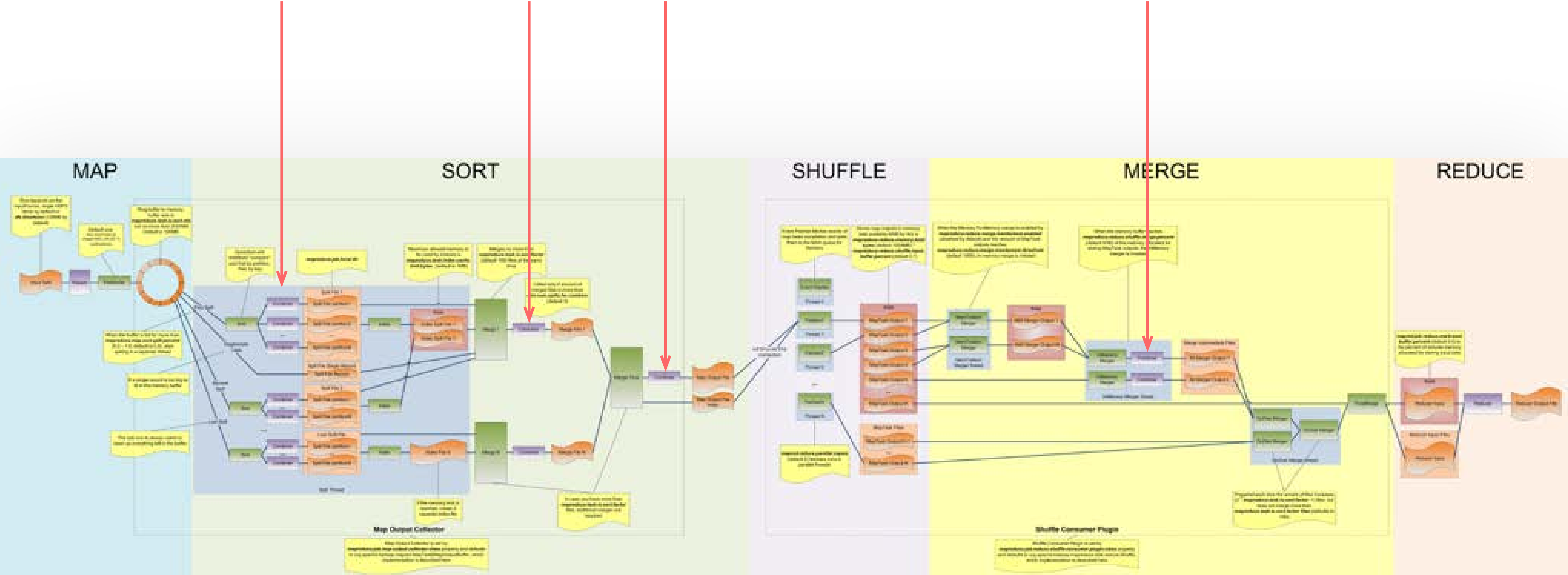


```
yarn --config $HADOOP_EMPTY_CONFIG jar $HADOOP_STREAMING_JAR \  
  -D map.output.key.field.separator=. \  
  -D mapreduce.partition.keypartitioner.options=-k1.2,1.2 \  
  -files identity_mr.py \  
  -mapper 'python identity_mr.py' \  
  -reducer 'python identity_mr.py' \  
  -numReduceTasks 2 \  
  -input subnet.txt \  
  -output subnet_out \  
  -partitioner org.apache.hadoop.mapred.lib.KeyFieldBasedPartitioner
```









Summary

Summary

- › You know how to use Partitioner in MapReduce application

Summary

- › You know **how to use** Partitioner in MapReduce application
- › You know **how to count** bigrams with Hadoop MapReduce and **distribute load** over reducers

BigDATAteam