# MapReduce
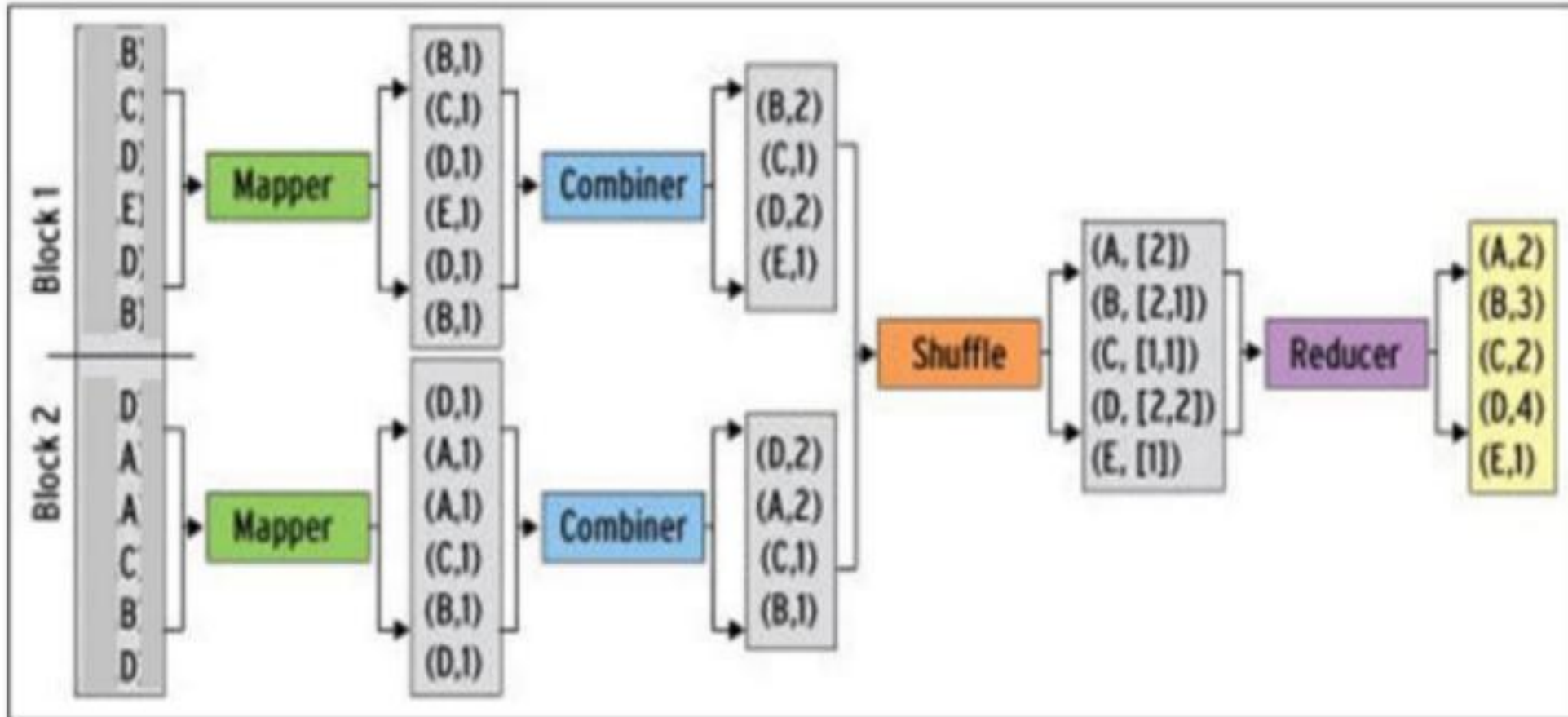
Lec 04

# MapReduce

# Word Count using MapReduce

```python
from mrjob.job import MRJob

class WordCount(MRJob):

    def mapper(self, _, line):
        for word in line.split():
            yield(word, 1)

    def combiner(self, word, counts):
        yield(word, sum(counts))

    def reducer(self, word, counts):
        yield(word, sum(counts))

if __name__ == '__main__':
    WordCount.run()
```

```
map(key, value):
// key: document name; value: text of
   the document

    for each word w in value:

        emit(w, 1)
```

```
reduce(key, values):
// key: a word; value:an array counts
    result = 0
    for each count v in values:
        result += v
    emit(key, result)
```

# Average Temperatures

```python
from mrjob.job import MRJob
class AvgTemperature(MRJob):
    def mapper(self, _, line):
        month, temperature = line.split()
        yield (month, (int(temperature),1))
```

```python
def _reducer_combiner(self, month, temperatures):
    sum, count = 0, 0
    for tmp, c in temperatures:
        sum = sum + tmp
        count += c
    avg = sum/count
    return (month, (avg, count))
```

```python
    def combiner(self, month, temperatures):
        # (May, (28 Degrees, 8 item average)
        yield self._reducer_combiner(month, temperatures)

    def reducer(self, month, temperatures):
        month, (avg, count) = self._reducer_combiner(month, temperatures)
        # (May, 28 Degrees)
        yield (month, avg)
```

Output

```
!python avg2.py temp.txt

"Apr"    12.666666666666666
"Feb"    2.0
"Jan"    -4.25
"Jun"    22.0
"Mar"    2.0
"May"    14.5
```

# Movie Rating

USER ID | MOVIE ID | RATING | TIMESTAMP

| | | | |
|---|---|---|---|
| 196 242 3 | 881250949 | | 3,1 |
| 186 302 3 | 891717742 | | 3,1 |
| 196 377 1 | 878887116 | Map | 1,1 |
| 244 51 2 | 880606923 | | 2,1 |
| 166 346 1 | 886397596 | | 1,1 |
| 186 474 4 | 884182806 | | 4,1 |
| 186 265 2 | 881171488 | | 2,1 |

Shuffle & Sort

1 -> 1, 1
2 -> 1, 1
3 -> 1, 1
4 -> 1

Reduce

1, 2
2, 2
3, 2
4, 1

```python
from mrjob.job import MRJob
from mrjob.step import MRStep

class RatingBreakdown(MRJob):
    def steps(self):
        return [
            MRStep(
                        mapper=self.mapper_get_ratings,
                        reducer=self.reducer_count_ratings
            )
        ]
    def mapper_get_ratings(self,_,line):
        (userId,movieID,rating,timestamp)=line.split('\t')
        yield rating,1


    def reducer_count_ratings(self,key,values):
        yield key,sum(values)



if __name__=="__main__":
    RatingBreakdown.run()
```

```python
from mrjob.job import MRJob
from mrjob.step import MRStep
import re


class MRMostUsedWord(MRJob):
    def mapper_get_words(self, _, line):
        # yield each word in the line
        for word in line.split(''):
            yield (word.lower(), 1)


    def combiner_count_words(self, word, counts):
        # sum the words we've seen so far
        yield (word, sum(counts))
```

```python
def reducer_count_words(self, word, counts):
    # send all (num_occurrences, word) pairs to the same reducer.
    # num_occurrences is so we can easily use Python's max() function.
    yield None, (sum(counts), word)
    # discard the key; it is just None


def reducer_find_max_word(self, _, word_count_pairs):
    # each item of word_count_pairs is (count, word),
    # so yielding one results in key=counts, value=word
    yield max(word_count_pairs)


def steps(self):
    return [
        MRStep(
            mapper=self.mapper_get_words,
            combiner=self.combiner_count_words,
            reducer=self.reducer_count_words),
        MRStep(
            reducer=self.reducer_find_max_word)
    ]

if __name__ == '__main__':
    MRMostUsedWord.run()
```

# Analysis of Weather Dataset

- ## Data from NCDC(National Climatic Data Center): A large volume of log data collected by weather sensors: e.g. temperature

- ## Data format
  - *Line-oriented ASCII format with many elements*
  - *We focus on the temperature element*
  - *Data files are organized by date and weather station*

```
        Year                              Temperature

0067011990999991950051507004...9999999N9+00001+99999999999...
0043011990999991950051512004...9999999N9+00221+99999999999...
0043011990999991950051518004...9999999N9-00111+99999999999...
0043012650999991949032412004...0500001N9+01111+99999999999...
0043012650999991949032418004...0500001N9+00781+99999999999...
```

**Contents of data files**

```
% ls raw/1990 | head
010010-99999-1990.gz
010014-99999-1990.gz
010015-99999-1990.gz
010016-99999-1990.gz
010017-99999-1990.gz
010030-99999-1990.gz
010040-99999-1990.gz
010080-99999-1990.gz
010100-99999-1990.gz
010150-99999-1990.gz
```

**List of data files**

# MapReduce Design of NCDC Example

## ■ Map phase

- – *Text input format of the dataset files*
    - ■ Key: offset of the line (unnecessary)
    - ■ Value: each line of the files
- – *Pull out the year and the temperature*
    - ■ The map phase is simply data preparation phase
    - ■ Drop bad records(filtering)

```
0067011990999991950051507004...9999999N9+00001+99999999999...
0043011990999991950051512004...9999999N9+00221+99999999999...
0043011990999991950051518004...9999999N9-00111+99999999999...
0043012650999991949032412004...0500001N9+01111+99999999999...
0043012650999991949032418004...0500001N9+00781+99999999999...
```

**Input File**

### Input of Map Function (key, value)

```
(0,   0067011990999991950051507004...9999999N9+00001+99999999999...)
(106, 0043011990999991950051512004...9999999N9+00221+99999999999...)
(212, 0043011990999991950051518004...9999999N9-00111+99999999999...)
(318, 0043012650999991949032412004...0500001N9+01111+99999999999...)
(424, 0043012650999991949032418004...0500001N9+00781+99999999999...)
```

**Map**

### Output of Map Function (key, value)
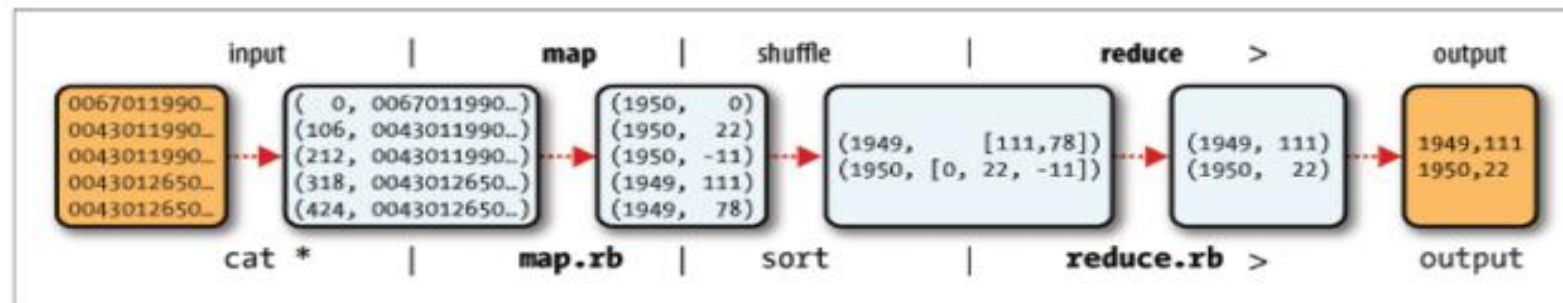
```
(1950, 0)
(1950, 22)
(1950, -11)
(1949, 111)
(1949, 78)
```

# MapReduce Design of NCDC Example

The output from the map function is processed by MapReduce framework

```
(1950, 0)
(1950, 22)
(1950, -11)
(1949, 111)
(1949, 78)
```

**Sort and Group By**
→

```
(1949, [111, 78])
(1950, [0, 22, -11])
```

- Reduce function iterates through the list and pick up the maximum value

```
(1949, [111, 78])
(1950, [0, 22, -11])
```

**Reduce**
→

```
(1949, 111)
(1950, 22)
```

| input | | map | | shuffle | | reduce | > | output |
|---|---|---|---|---|---|---|---|---|
| 0067011990…<br>0043011990…<br>0043011990…<br>0043012650…<br>0043012650… | | ( 0, 0067011990…)<br>(106, 0043011990…)<br>(212, 0043011990…)<br>(318, 0043012650…)<br>(424, 0043012650…) | | (1950, 0)<br>(1950, 22)<br>(1950, -11)<br>(1949, 111)<br>(1949, 78) | | (1949, [111,78])<br>(1950, [0, 22, -11]) → (1949, 111)<br>(1950, 22) | | 1949,111<br>1950,22 |
| cat * | | map.rb | | sort | | reduce.rb > | | output |

# Reference

- https://buildmedia.readthedocs.org/media/pdf/mrjob/latest/mrjob.pdf