

MapReduce

1997

Data

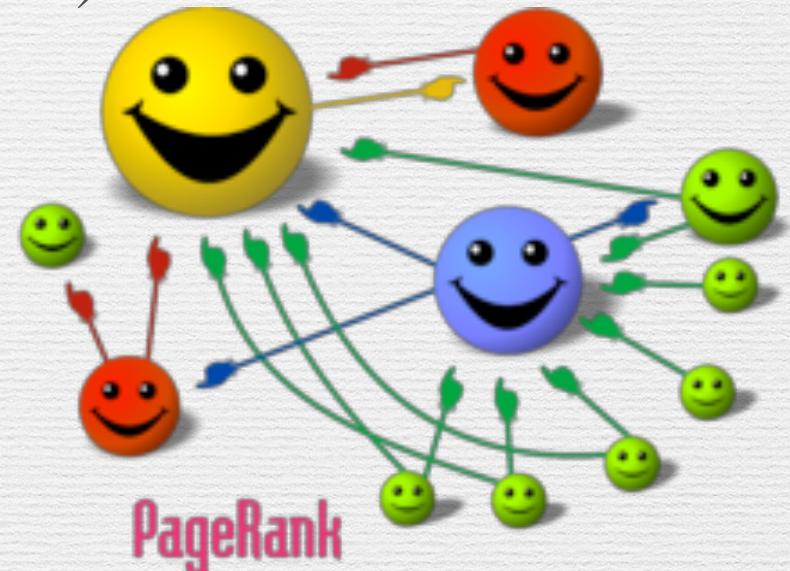
Problem

Idea: Algorithm

Webpages (URLs)

Find URLs

PageRank



Larry Page and Sergey Brin

Google!

BETA

Search the web using Google!

Special Searches
[Stanford Search](#)
[Linux Search](#)

Help!
[About Google](#)
[Company Info](#)
[Google! Logos](#)

Get Google! updates monthly:

Copyright ©1998 Google Inc.

2003

Data

Webpages (URLs)



Problem

Find
Webpage



Idea: Algorithm

PageRank

System

MapReduce



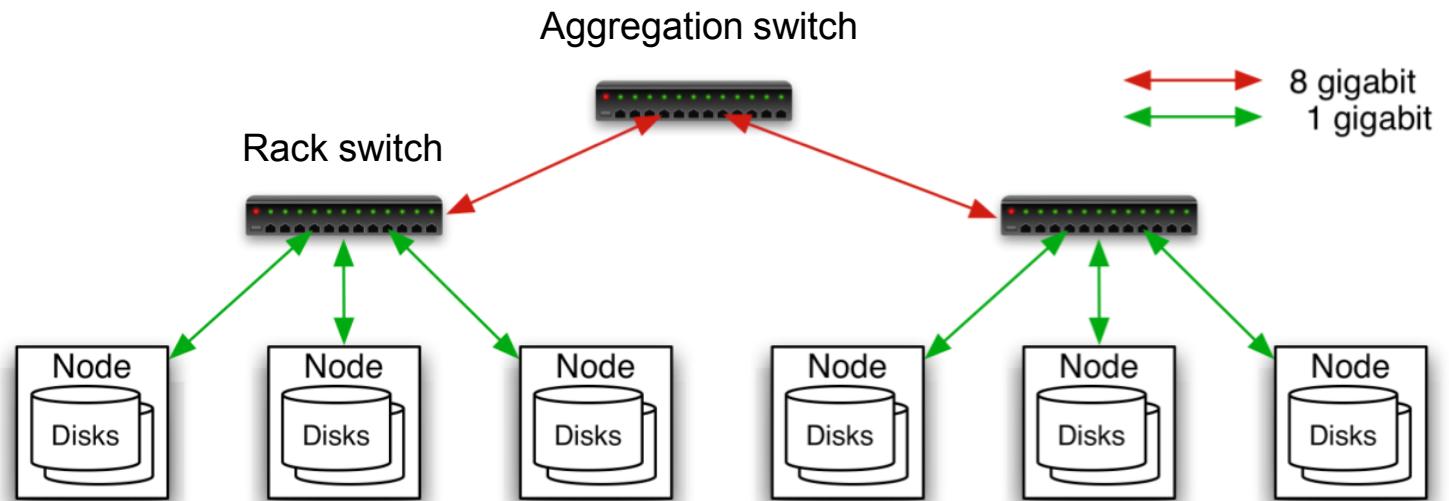
Google

Sanjay Ghemawat and Jeffrey Dean

System

2005

MapReduce System + Open source



Doug Cutting and Mike Cafarella



now used by Yahoo
Facebook Amazon

Hardware

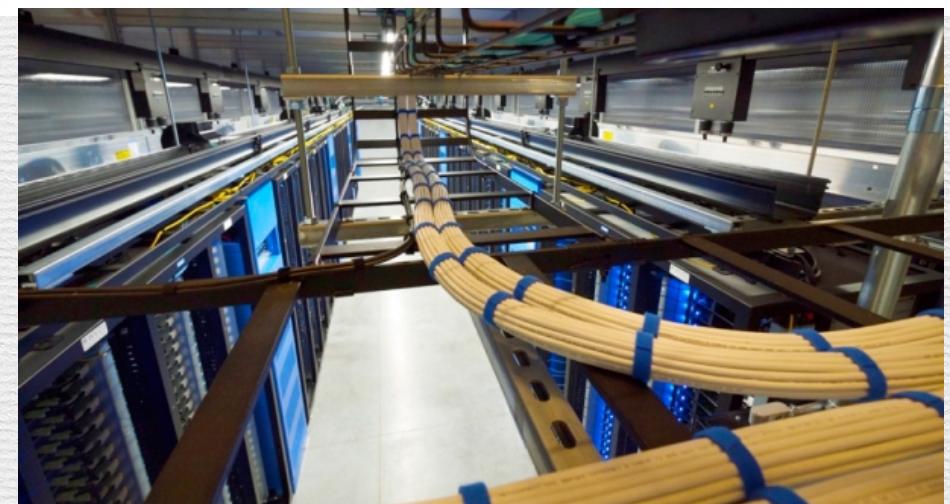
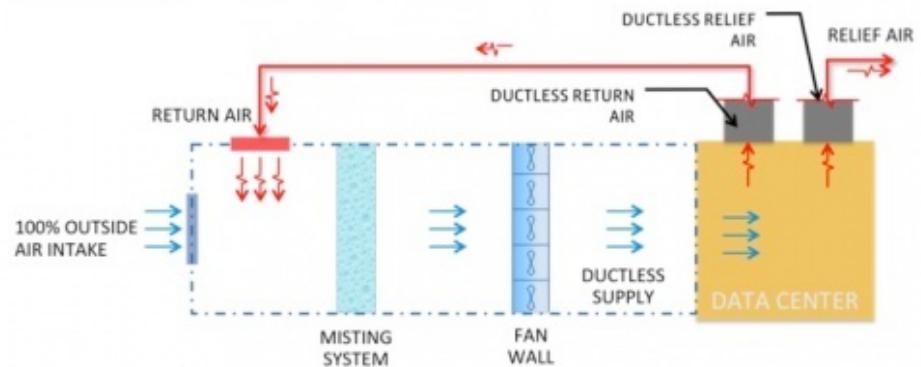
2011

Data Center Design + Open source

https://www.youtube.com/watch?v=_r97qdyQtIk



5.3.2 Airflow Overview



facebook®



OPEN
Compute Project



data center



cooling system



storage



power plant

Programming: **MapReduce, Spark ...**

Programming: Python, Java, ...

laptop



fan



hard drive



power cord

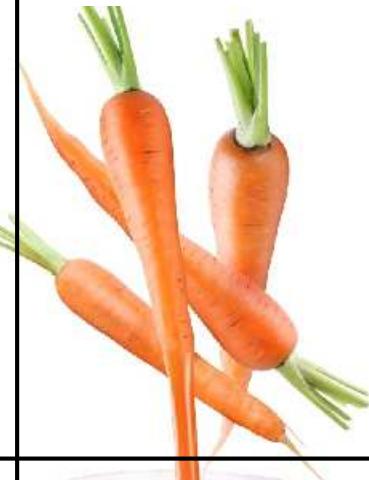
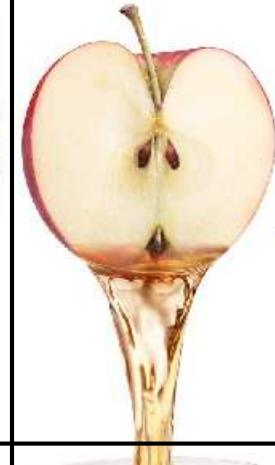
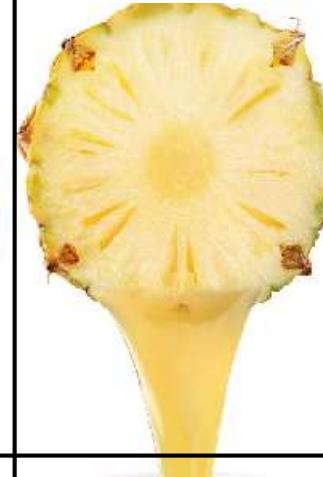
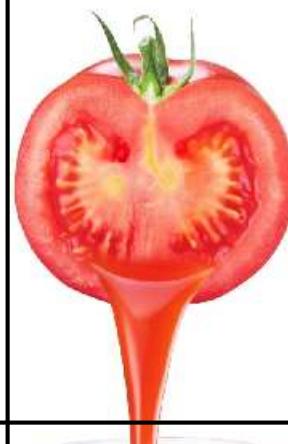


MapReduce

- MapReduce
 - Framework for parallel computing
 - Programmers get simple API
 - Don't have to worry about handling
 - parallelization
 - data distribution
 - load balancing
 - fault tolerance
- Allows one to process huge amounts of data (terabytes and petabytes) on thousands of processors

Mappers

Input:
one data
record



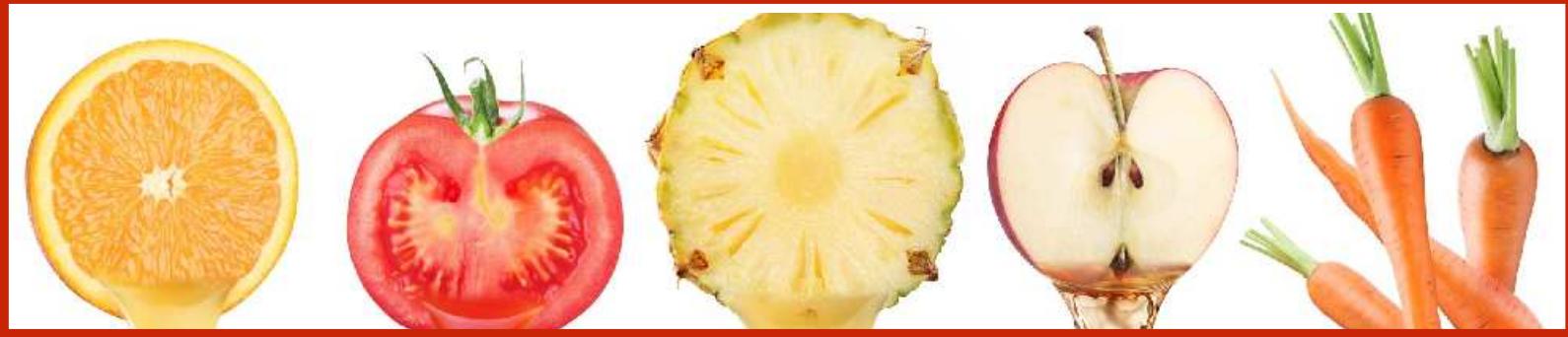
Output:
one/more
results



Reducer

a list of data records

Input:



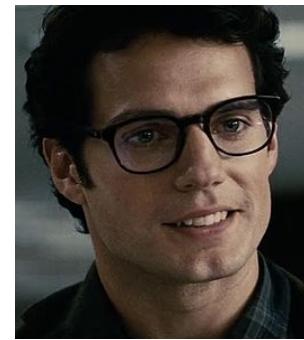
Output:

one/more
results



Mappers

Input:
one data
record



Output:
one/more
results



Reducer

Input:

a collection of data records

Output:
one (or multiple)
result(s)





map



map



map



map



map



map



reduce



reduce





map



map



map



map



map



map



reduce

4



reduce

2



map



map



map



map



map



map

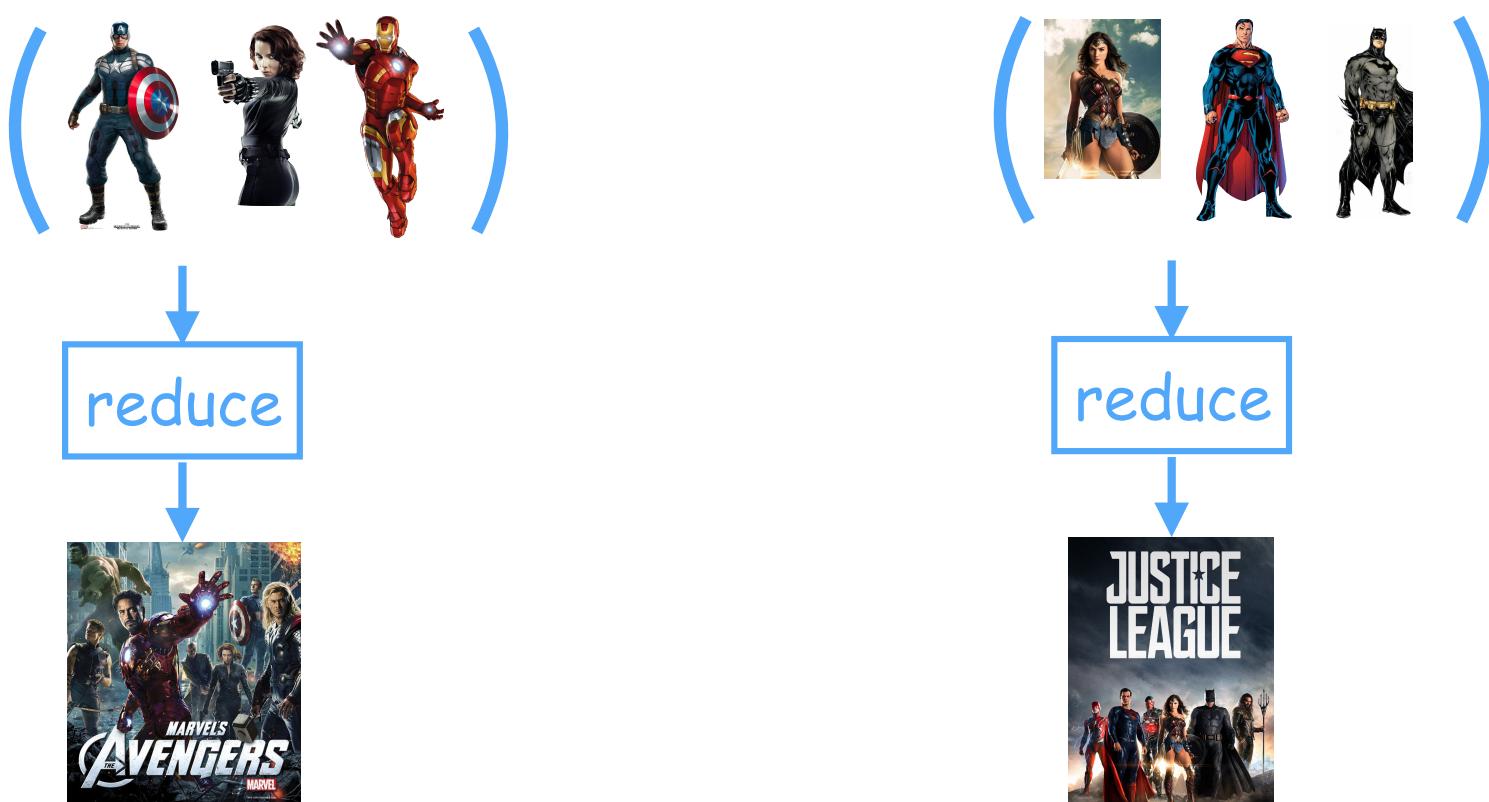
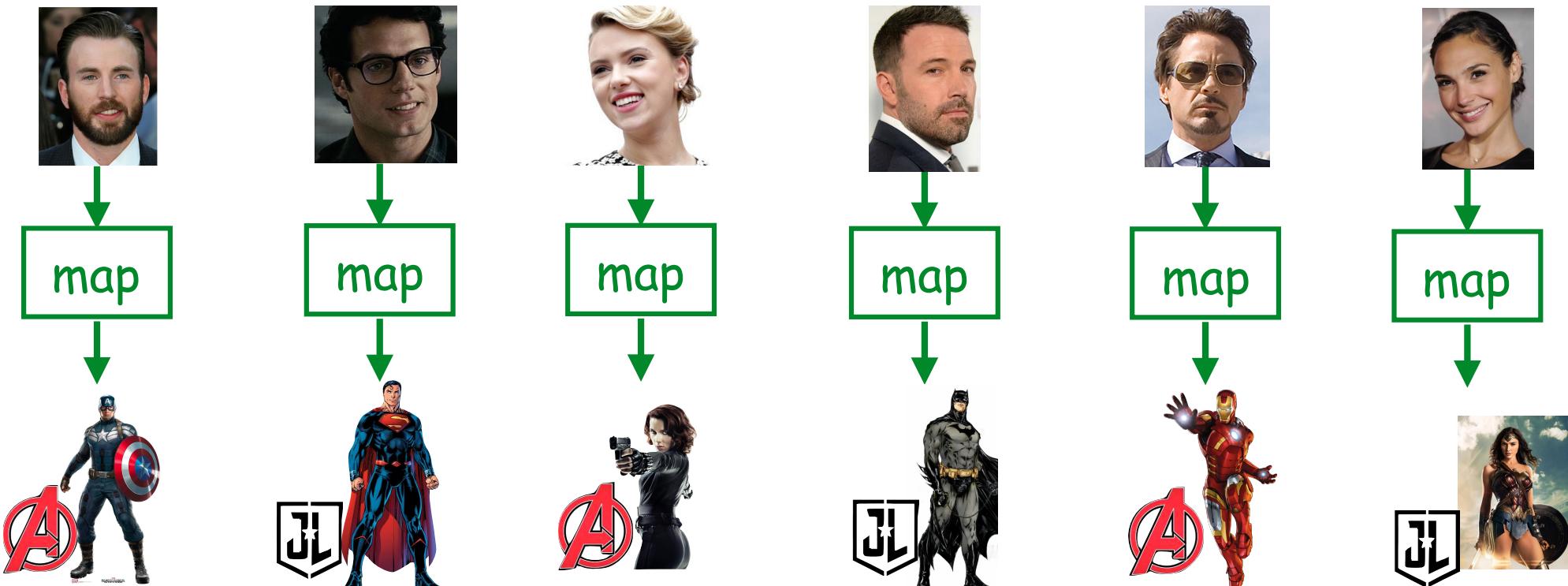


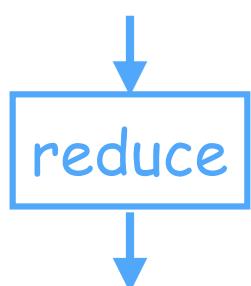
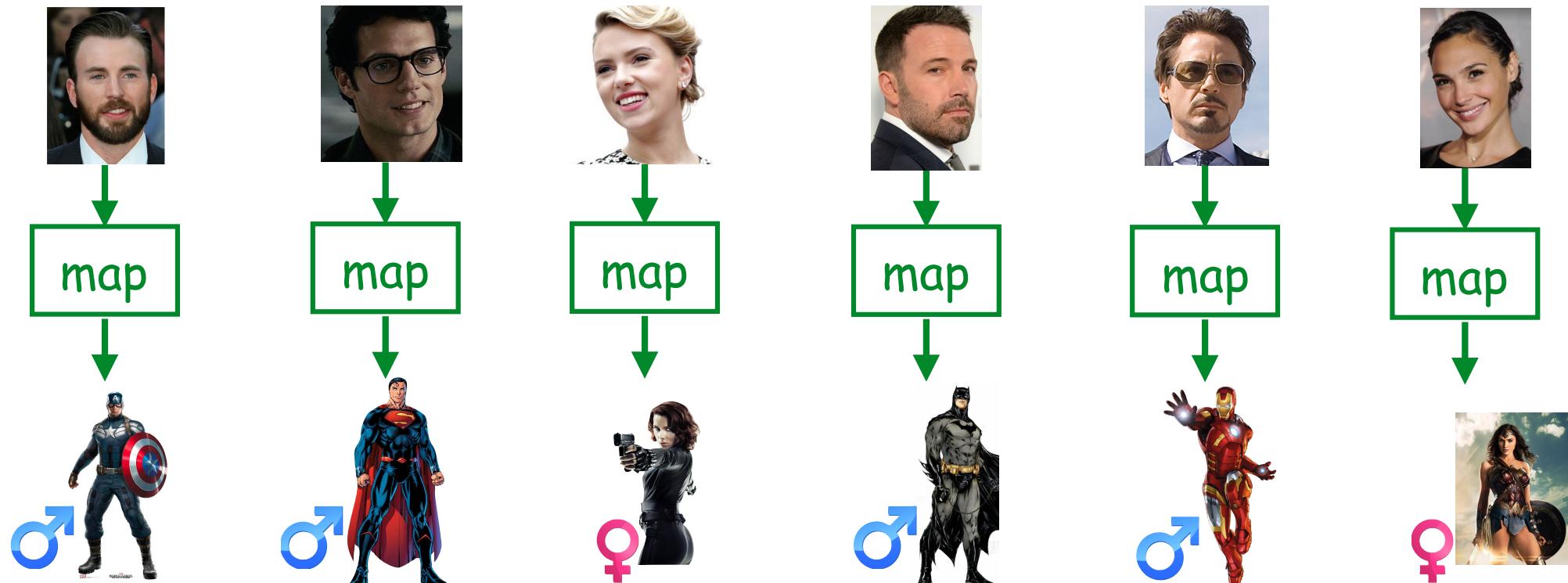
reduce



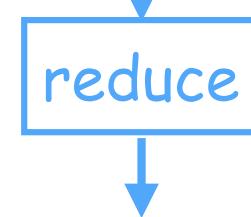
reduce



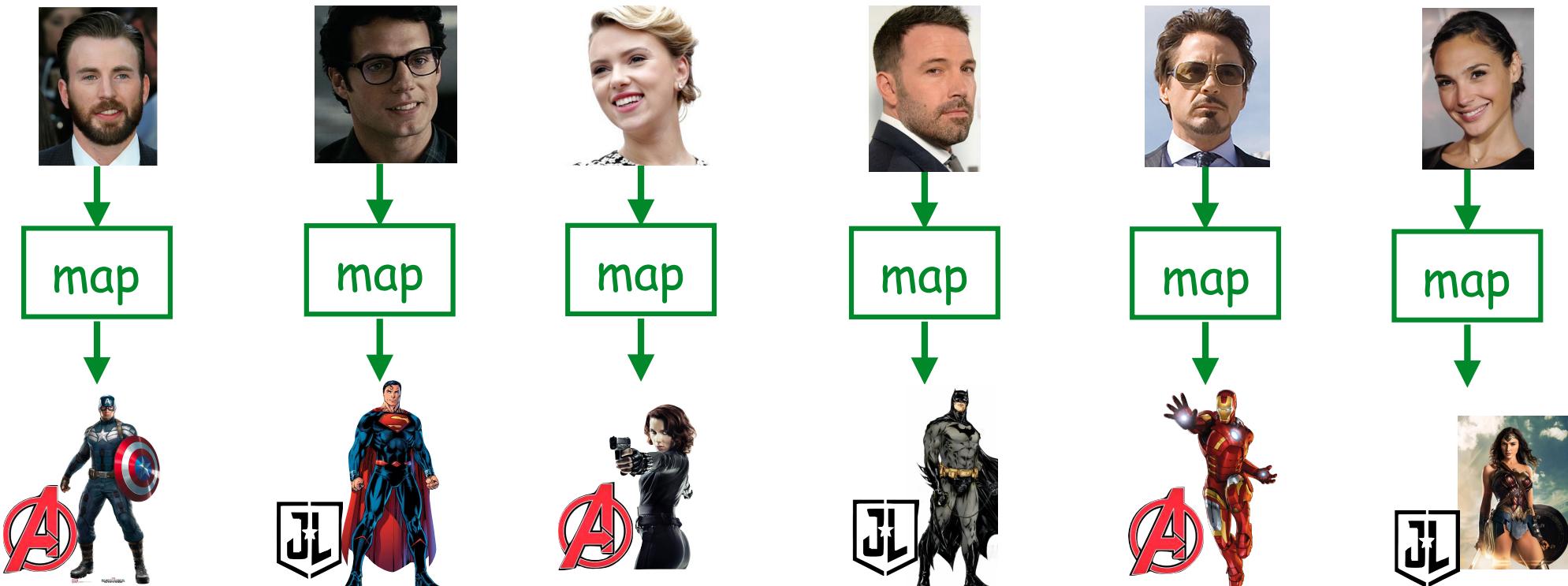




4



2

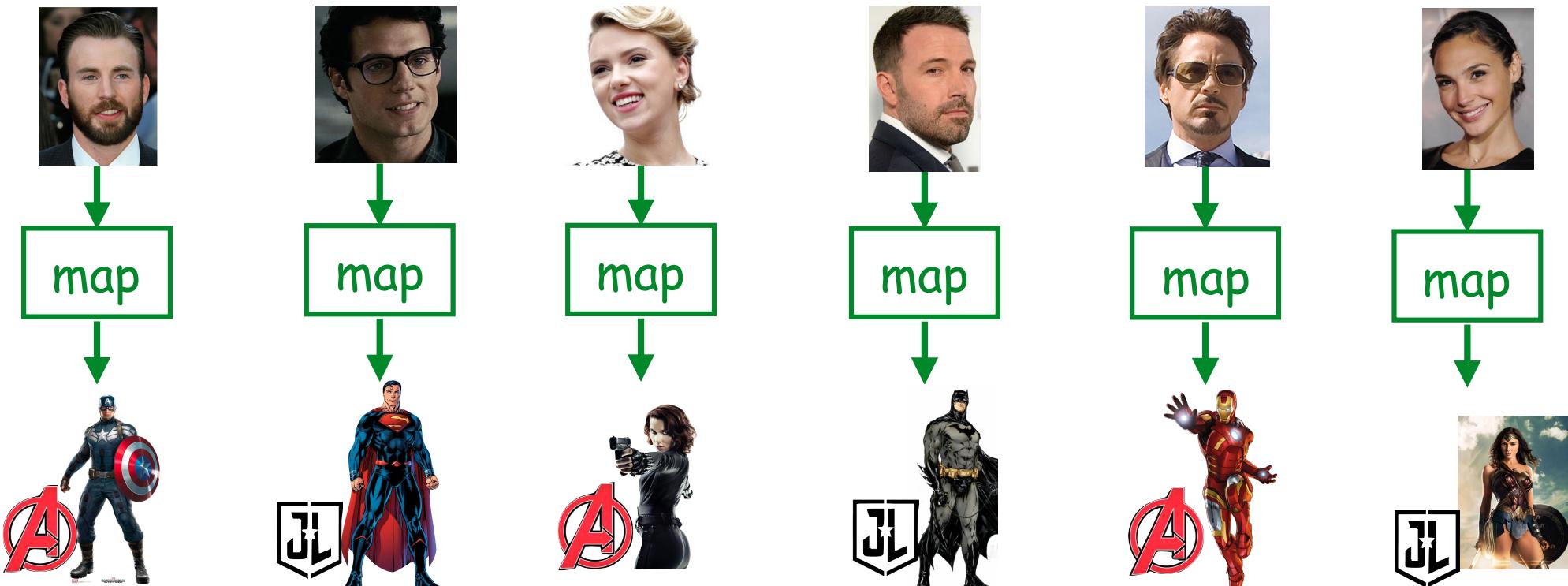


reduce



reduce



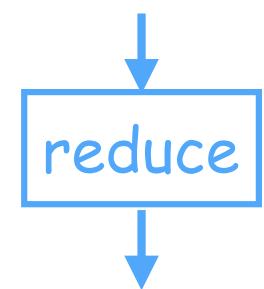
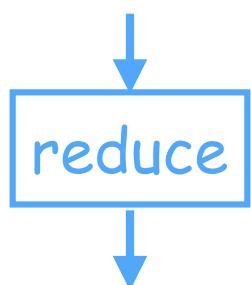
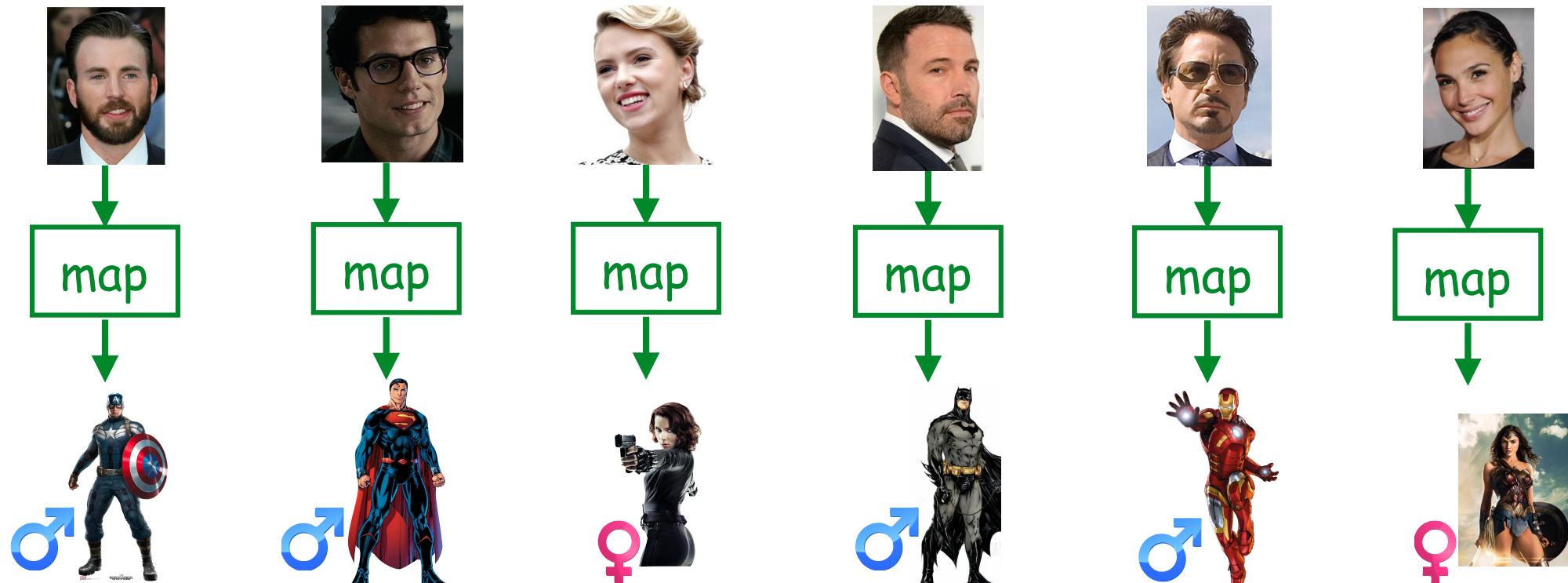


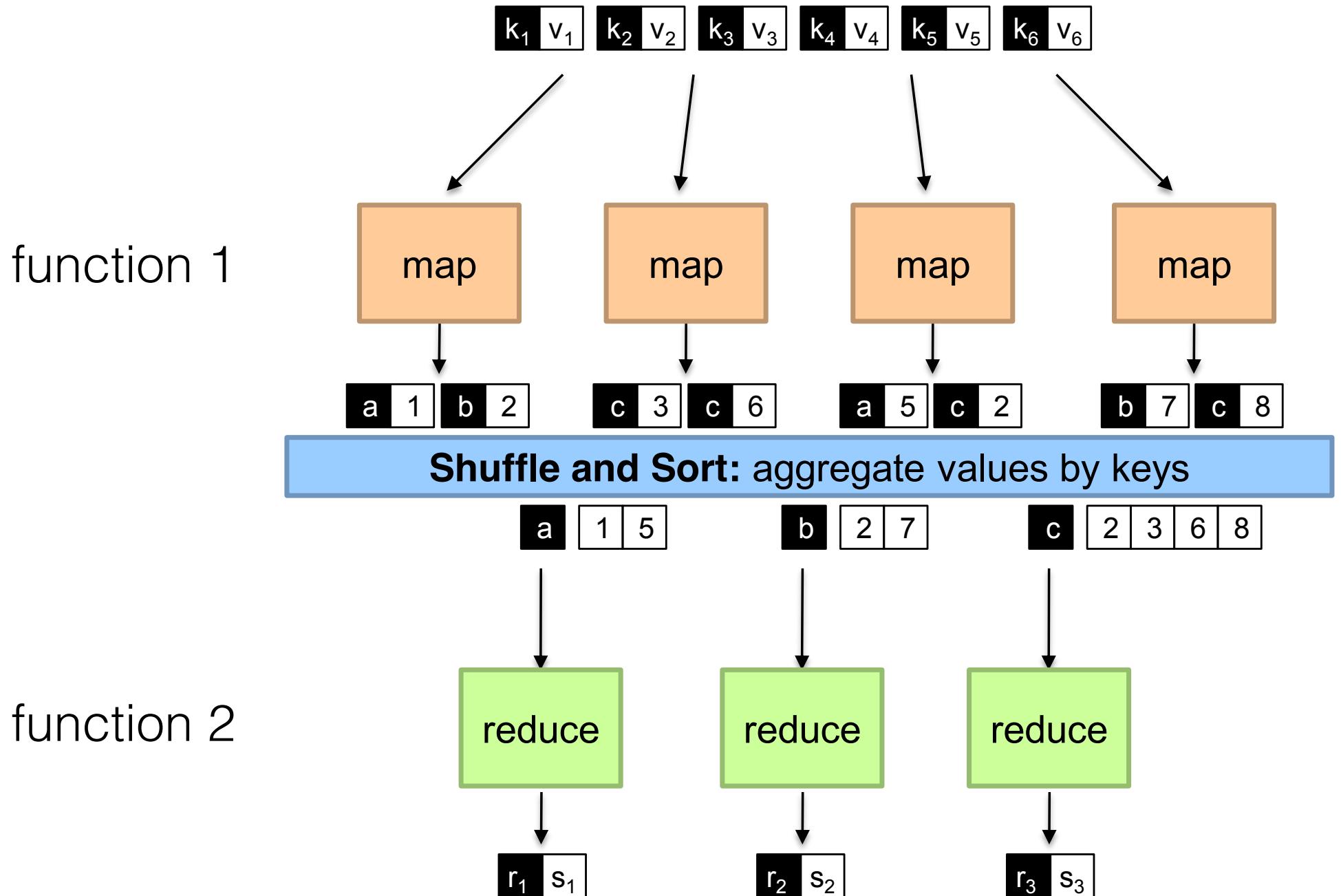
reduce



reduce

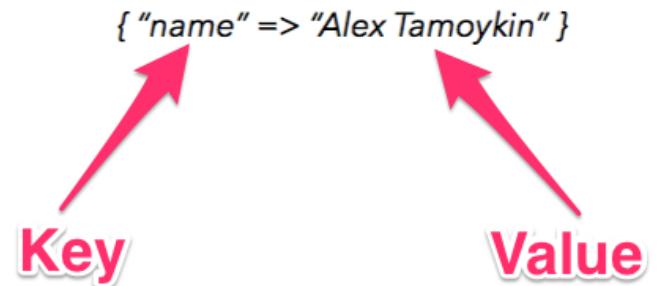




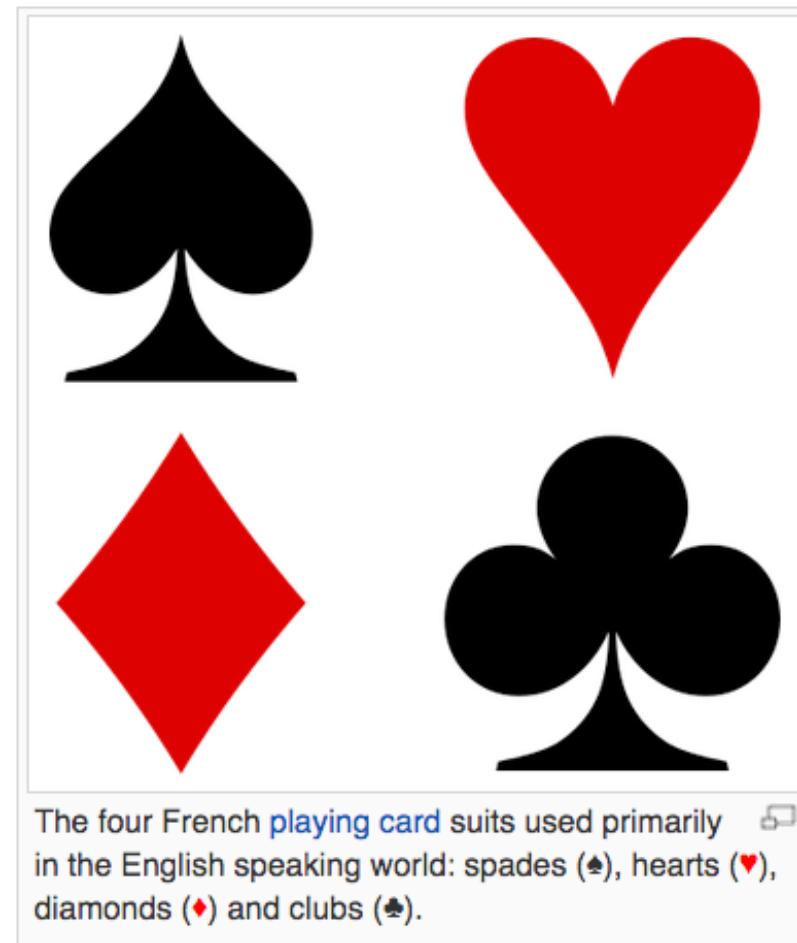


Key- Value Pairs

- Mappers and Reducers are users' code (provided functions)
- Just need to obey the Key-Value pairs interface
- **Mappers:**
 - Consume <key, value> pairs
 - Produce <key, value> pairs
- **Reducers:**
 - Consume <key, <list of values>>
 - Produce <key, value>
- **Shuffling and Sorting:**
 - Hidden phase between mappers and reducers
 - Groups all similar keys from all mappers, sorts and passes them to a certain reducer in the form of <key, <list of values>>



Game

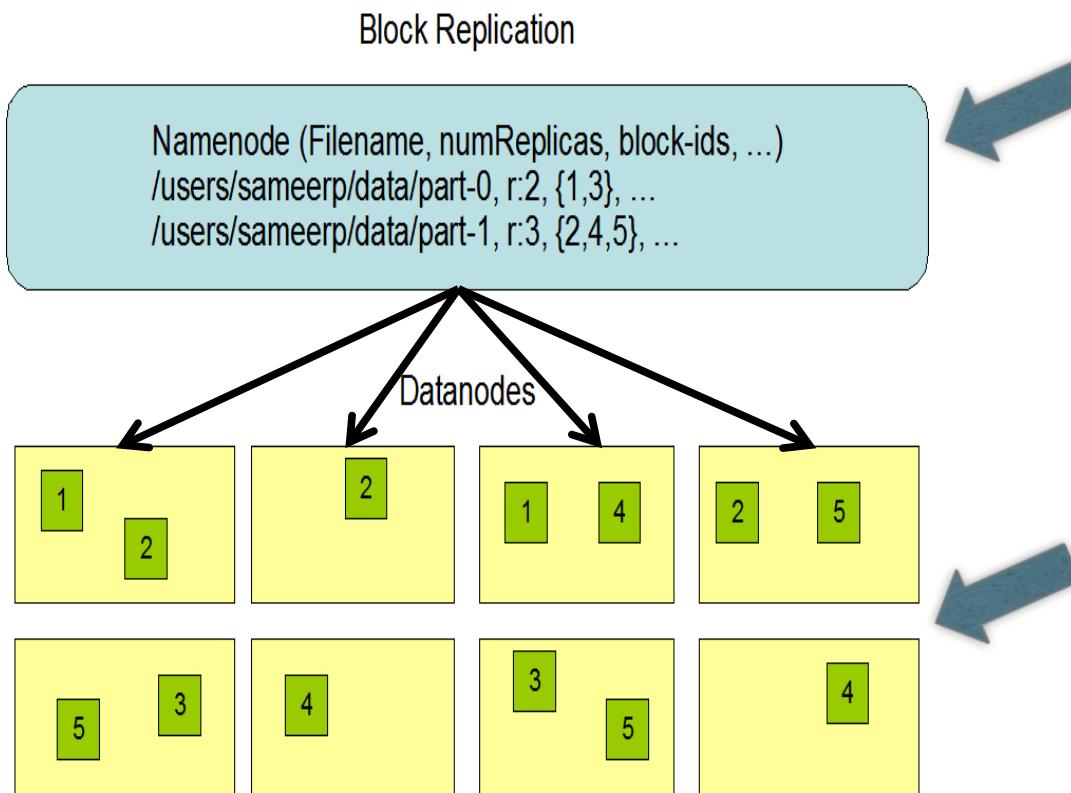


The four French [playing card](#) suits used primarily in the English speaking world: spades (♠), hearts (♥), diamonds (♦) and clubs (♣).

Main Properties of HDFS

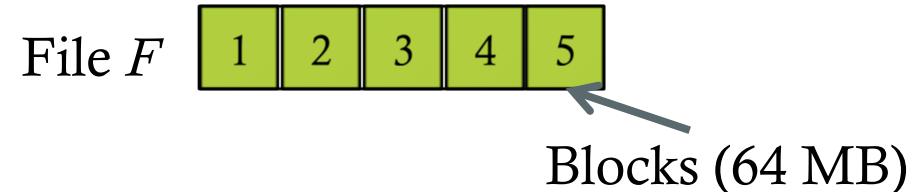
- **Large:** A HDFS instance may consist of thousands of server machines, each storing part of the file system's data
- **Replication:** Each data block is replicated many times (default is 3)
- **Failure:** Failure is the norm rather than exception
- **Fault Tolerance:** Detection of faults and quick, automatic recovery from them is a core architectural goal of HDFS
 - Namenode is consistently checking Datanodes

Hadoop Distributed File System (HDFS)



Centralized namenode

- Maintains metadata info about files



Many datanode (1000s)

- Store the actual data
- Files are divided into blocks
- Each block is replicated N times
(Default = 3)

MRjob package

 pythonhosted.org/mrjob/guides/quickstart.html

mrjob v0.4.2 documentation

[Home](#) » [Guides](#)

← [Why mrjob?](#) | [Concepts](#) →

Table Of Contents

Fundamentals

- Installation
- Writing your first job
 - What's happening
- Running your job different ways
- Writing your second job
- Configuration

Need help?

Join the mailing list by visiting the [Google group page](#) or sending an email to mrjob+subscribe@googlegroups.com.

Fundamentals

Installation

Install with pip:

```
pip install mrjob
```

or from a [git](#) clone of the [source code](#):

```
python setup.py test && python setup.py install
```

Writing your first job

Open a file called `word_count.py` and type this into it:

```
from mrjob.job import MRJob

class MRWordFrequencyCount(MRJob):

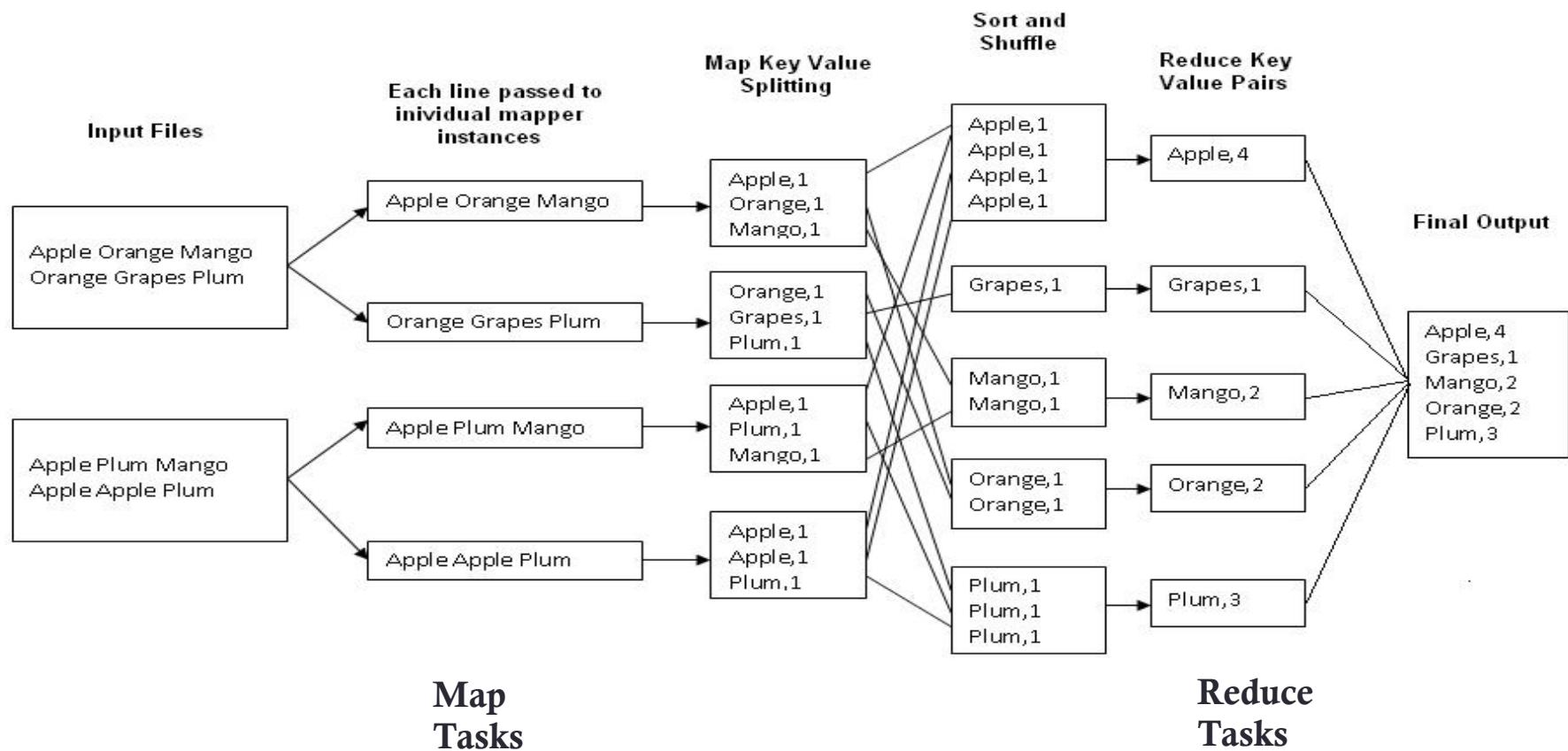
    def mapper(self, _, line):
        yield "chars", len(line)
        yield "words", len(line.split())
        yield "lines", 1

    def reducer(self, key, values):
        yield key, sum(values)
```

<http://pythonhosted.org/mrjob/guides/quickstart.html>

Example 1: Word Count

- Job: Count the occurrences of each word in a data set



The Famous Word Count Example

```
from mrjob.job import MRJob

class mrWordCount(MRJob):

    def mapper(self, key, line):
        for word in line.split(' '):
            yield word.lower(), 1

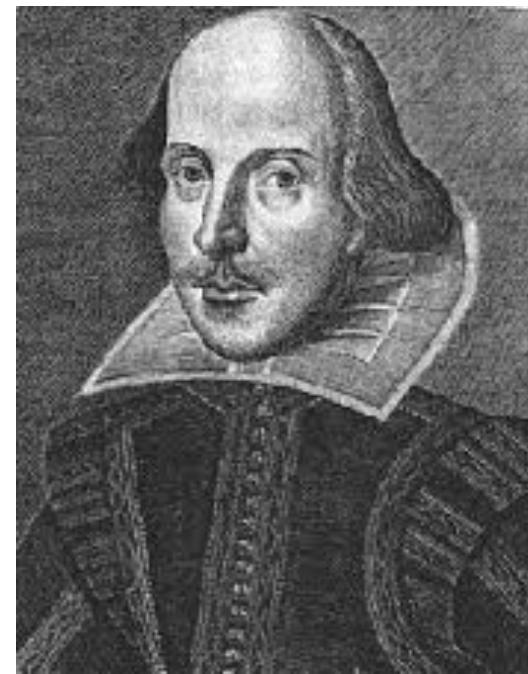
    def reducer(self, word, occurrences):
        yield word, sum(occurrences)

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

Example Input Data



Hamlet



Shakespeare

Launching the Job

Interpreter

Input
Redirection

Output
Redirection

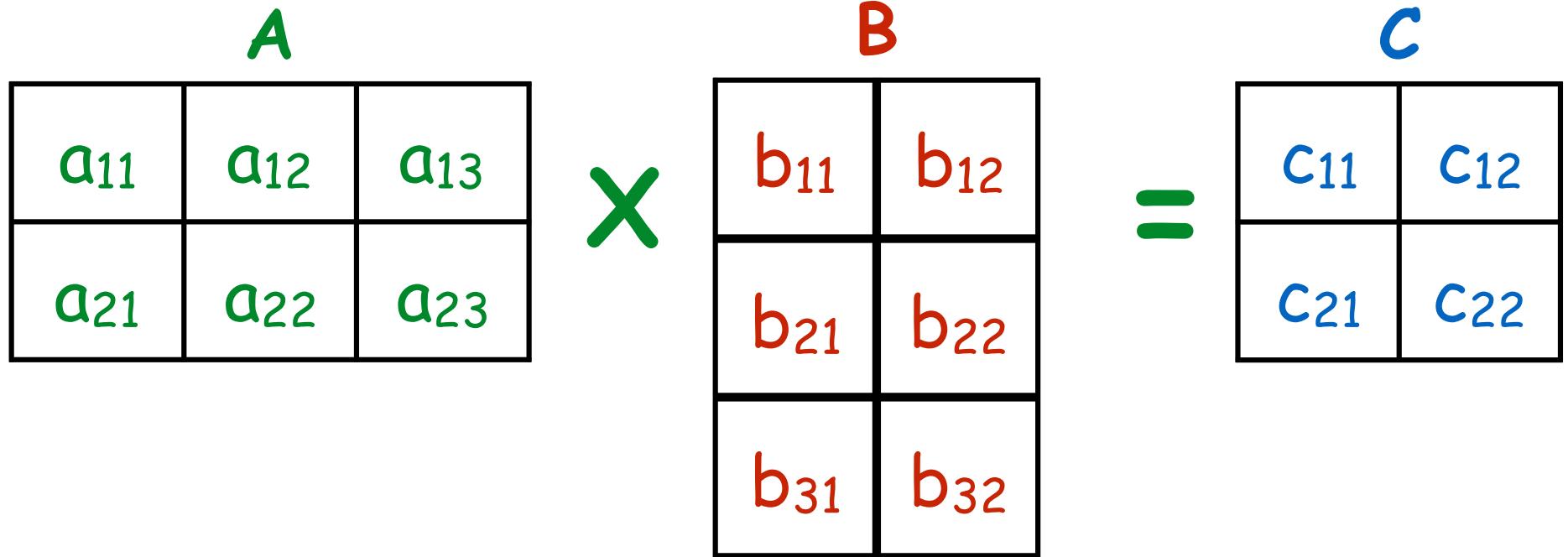
```
python myscript.py < inputfile.txt > outputfile.txt
```

Your
Script

Input
File

Output
File

Matrix Multiplication



$$c_{11} = a_{11} \times b_{11} + a_{12} \times b_{21} + a_{13} \times b_{31}$$

$$c_{ij} = a_{i1} \times b_{1j} + a_{i2} \times b_{2j} + a_{i3} \times b_{3j}$$

A		
1	2	3
4	5	6

X

B	
7	10
8	11
9	12

C	
50	68
122	167

Input File

A, 1, 1, 1
A, 2, 1, 4
...
B, 1, 1, 7
B, 2, 1, 8
...



Output File

C, 1, 1, 50
C, 2, 1, 122
...

Mapper

Input:

one element
of an input
matrix

$A, 1, 1, 1$

$A[1,1]$

1	2	3
4	5	6

$A, 2, 1, 4$

$A[2,1]$

1	2	3
4	5	6

$B, 1, 1, 7$

$B[1,1]$

7	10
8	11
9	12

$B, 2, 2, 11$

$B[2,2]$

7	10
8	11
9	12

Output:

$C[1,1] \ A[1,1]$

11	12
21	22

1	2	3
4	5	6

key value

$C[2,1] \ A[2,1]$

11	12
21	22

1	2	3
4	5	6

key value

$C[1,1] \ B[1,1]$

11	12
21	22

7	10
8	11
9	12

key value

$C[1,2] \ B[2,2]$

11	12
21	22

7	10
8	11

key value

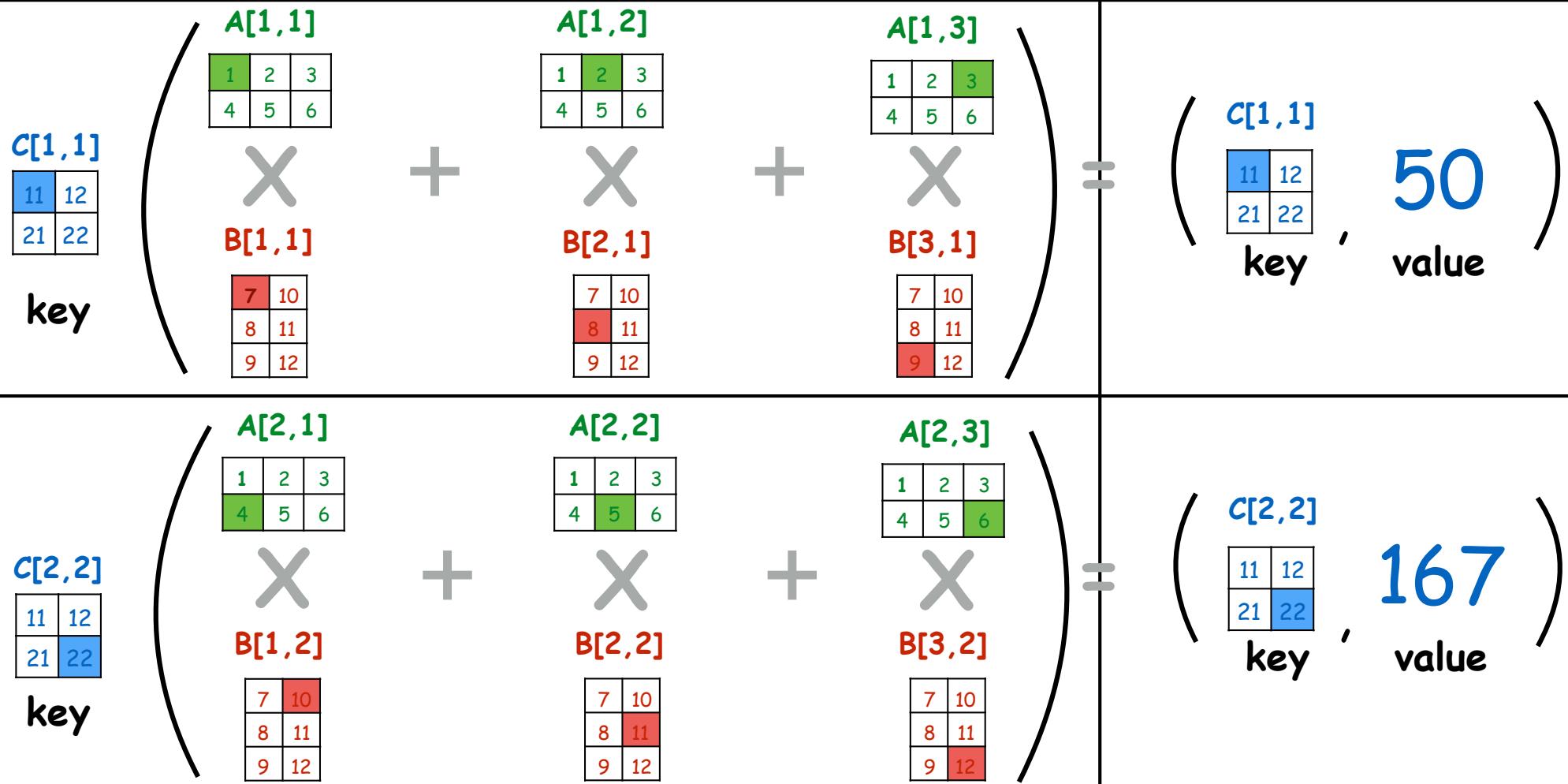
Reducer

Input:

a collection of data records

Output:

one element of
matrix C



Map

Input

A, 1, 1, 1

B, 1, 1, 7

key value

_, (A, 1, 1, 1)

_, (B, 1, 1, 7)

key value

(C, 1, 1), (A, 1, 1, 1)

(C, 1, 1), (B, 1, 1, 7)

key value ...

(C, 1, 2), (A, 1, 1, 1) ...

(C, 2, 1), (B, 1, 1, 7) ...

Reduce

key

value list

(C, 1, 1), [(A, 1, 1, 3), (B, 1, 1, 7), ..., (B, 3, 1, 9)]

(C, i, j), [(A, i, 1, 4), ..., ..., (B, 1, j, 7), ..., ...]

key value

(C, 1, 1), 50

(C, 1, 2), 68

Output

C, 1, 1, 50

C, 1, 2, 68

$$c_{ij} = a_{i1} \times b_{1j} + a_{i2} \times b_{2j} + a_{i3} \times b_{3j}$$