# Introduction to Map Reduce

Continuum Analytics

# Large Relative to What?

- Can't fit into Excel

# Large Relative to What?

- Can't fit into Excel
  - Increase Memory

# Large Relative to What?

- Can't fit into Excel
  - Increase Memory
- Can't fit into R

# Large Relative to What?

- Can't fit into Excel
  - Increase Memory
- Can't fit into R
  - Increase Memory

# Large Relative to What?

- Can't fit into Excel
    - Increase Memory
- Can't fit into R
    - Increase Memory
- Can't fit into Memory
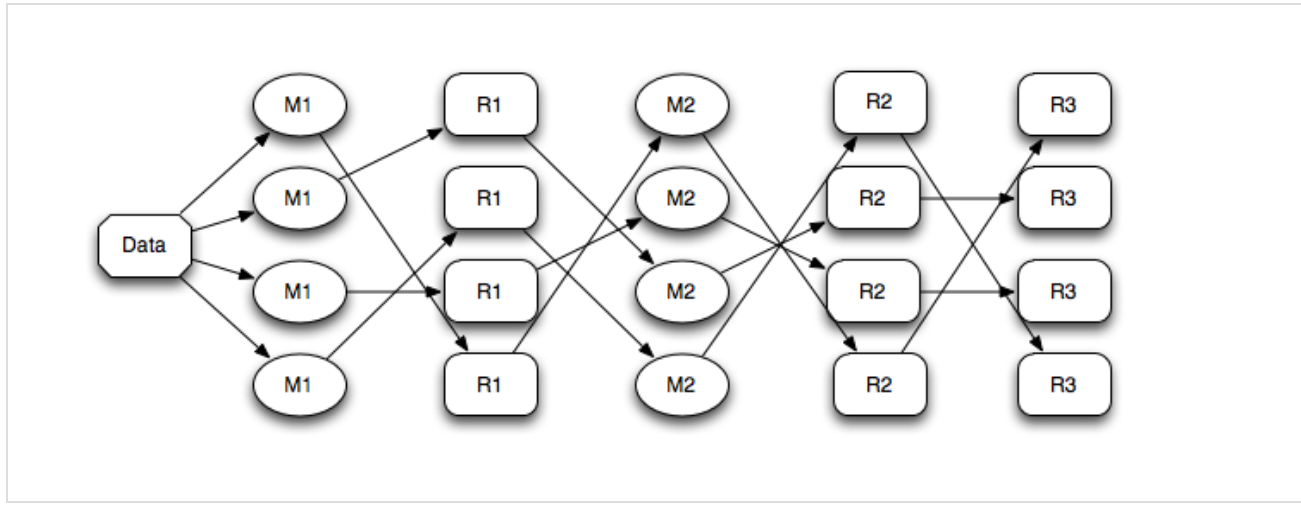    - Increase Memory

# Large Relative to What?

- Can't fit into Excel
  - Increase Memory
- Can't fit into R
  - Increase Memory
- Can't fit into Memory
  - Increase Memory
- Can't fit on a single disk
  - Distributed Filesystem: SAN, HDFS/DDFS, AWS: S3, Redshift, etc.

# MapReduce

Framework to help solve the problem of distributed computation for distributed data

- A mass of data: records
- Split/**Map** records into key-values pairs
- Collect/Partition kv pairs (Optional Sort)
- Buckets are passed to **Reduce** function
- Result is returned

# MapReduce Workflow



- Push Code to Data

- Lots of Network Traffic

# MR Implementations

- Disco: Python + Erlang
  - Distributed FileSystem: DDFS

- Hadoop: Java
  - Streaming with Python
  - Dumbo
  - MRJob
  - Hadoopy

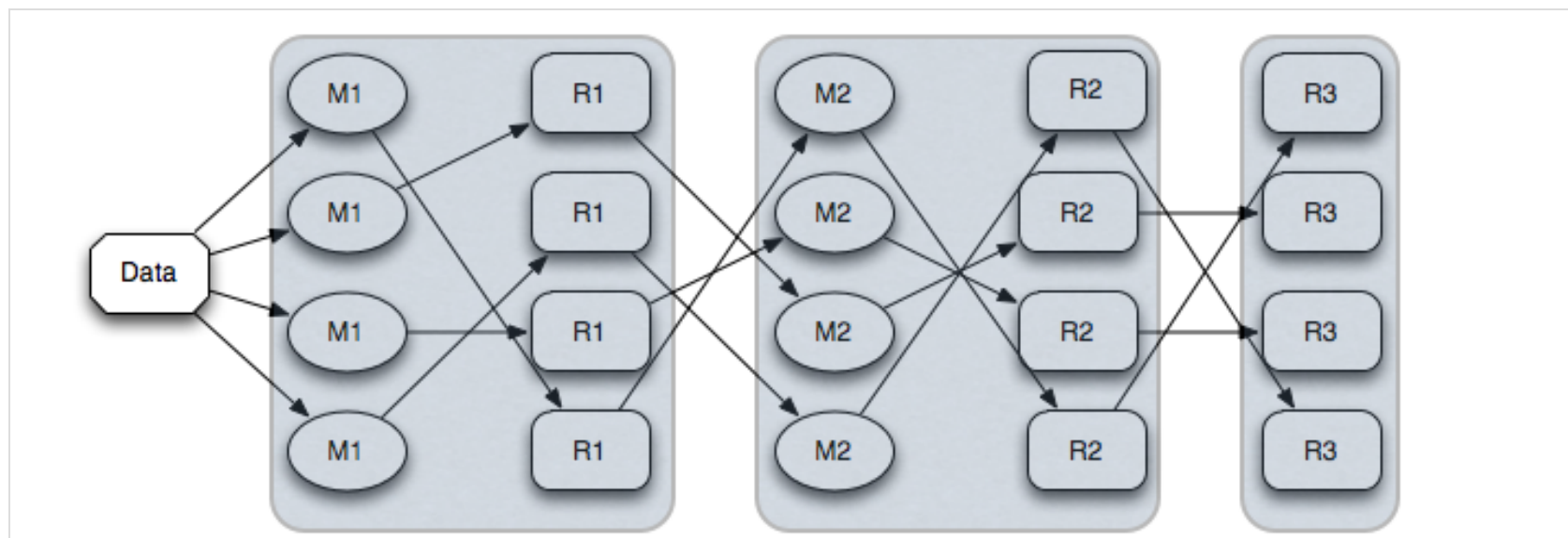# MapReduce: It's a Party

# Buddies Included

- NumPy
- SciPy
- pandas
- scikits-learn
- OpenCV
- ...

# Canonical Example

```python
1   from disco.job import Job
2   from disco.core import result_iterator
3
4   class WordCount(Job):
5
6       partitions = 3
7       input=["sherlock.txt","poirot.txt","clouseau.txt"]
8
9       @staticmethod
10      def map(line, params):
11          import string
12          for word in line.split():
13              yield word, 1
14
15      @staticmethod
16      def reduce(iter, params):
17          from disco.util import kvgroup
18          for word, counts in kvgroup(sorted(iter)):
19              yield word, sum(counts)
20
21  if __name__ == "__main__":
22      from disco_words import WordCount
23
24      wordcount = WordCount().run()
25
26      for (word, counts) in result_iterator(wordcount.wait(show=True)):
27          print word, counts
28
```
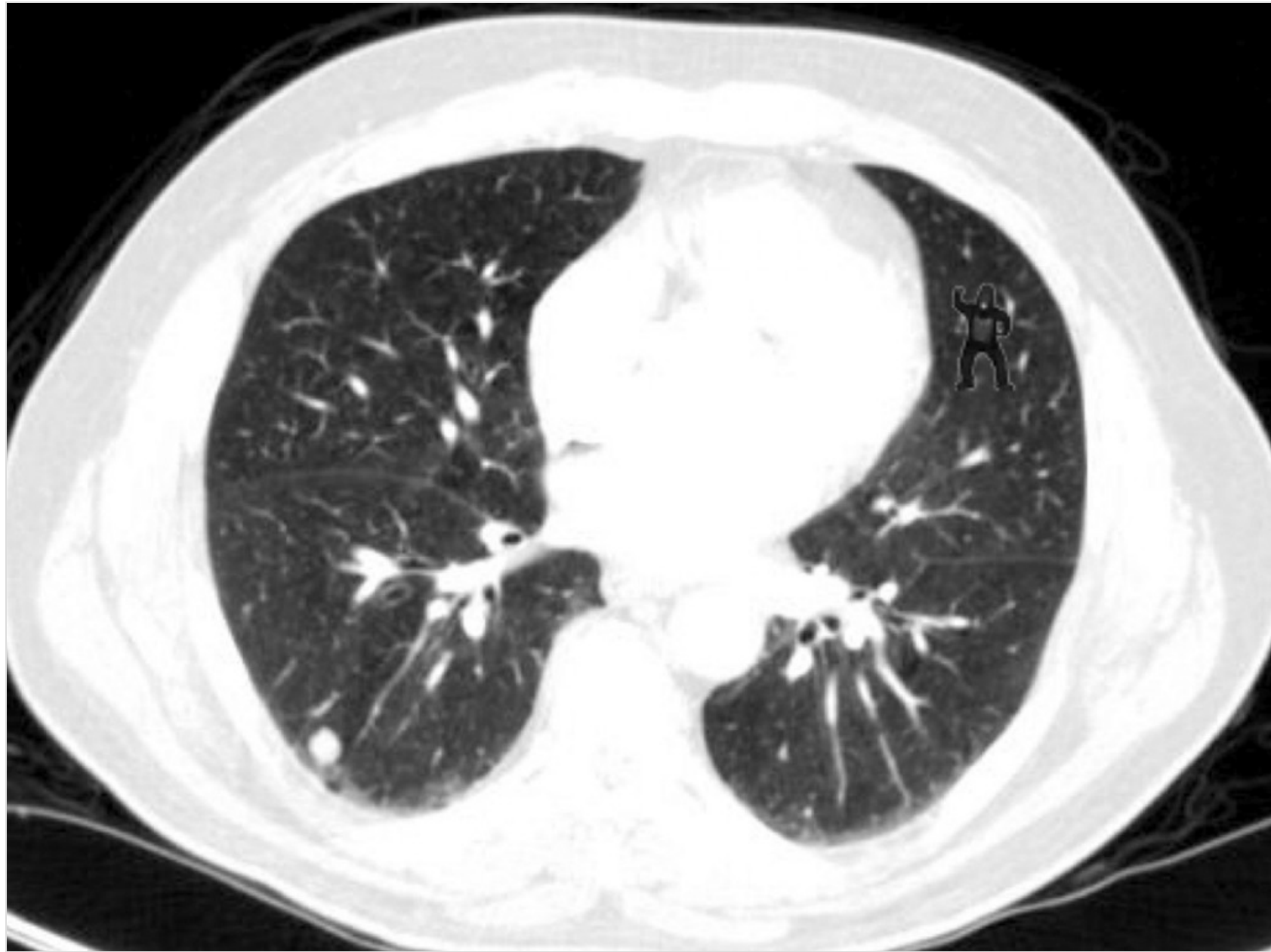
CONTINUUM
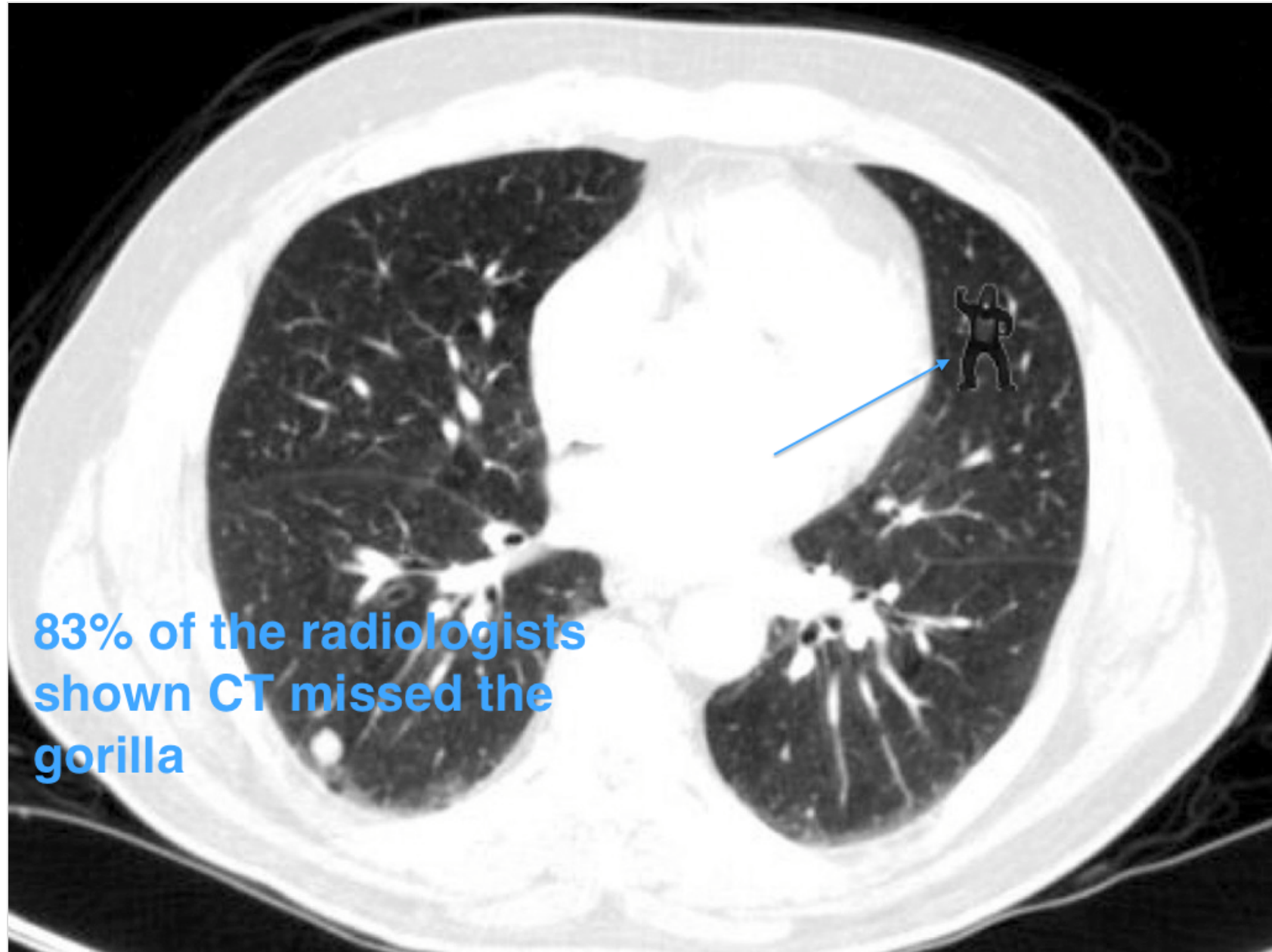ANALYTICS

# Demo 1

# Chaining Jobs

# MapReduce Thoughts

- Data Cleansing
  - Everyone's pain point
- Task Deconstruction
  - Good for code management
  - Hides -- in a good way -- data management
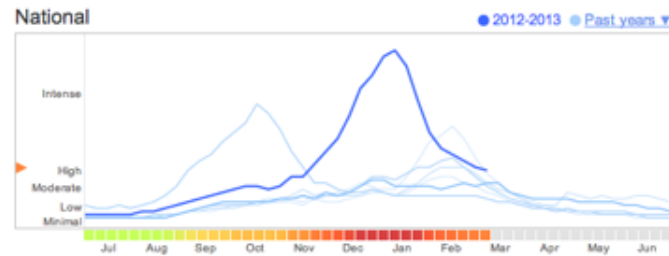- Can Be Inefficient
  - Network traffic
  - Job organization

CONTINUUM
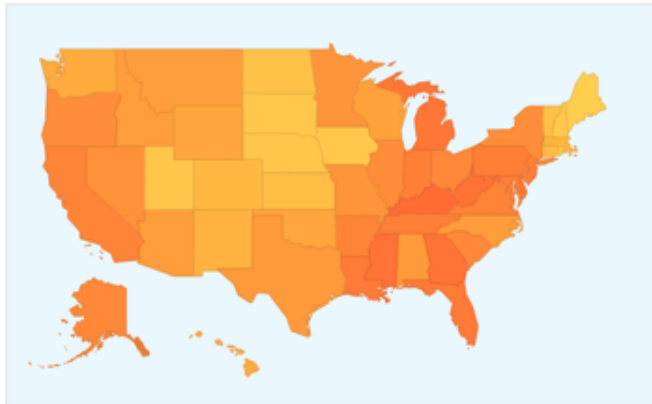ANALYTICS

83% of the radiologists shown CT missed the gorilla
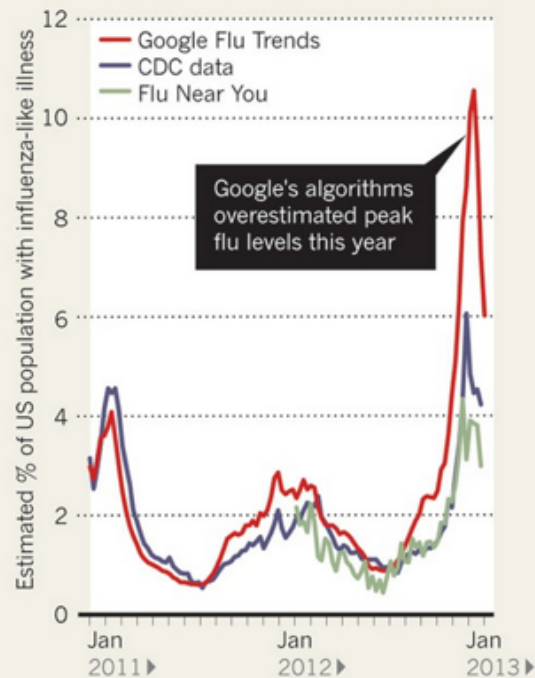
# Google Flu



- Data Mining

- Faster than CDC

# Google Get's It Wrong



**FEVER PEAKS**
A comparison of three different methods of measuring the proportion of the US population with an influenza-like illness.

- Google Flu Trends
- CDC data
- Flu Near You

Google's algorithms overestimated peak flu levels this year

- Typically, prediction is great!

- This year not so much

- Google: No comment!

- Feedback mechanism from hype-up media

CONTINUUM
ANALYTICS

# Data Philosophy

- Invisible Gorillas will stay Invisible
  - Inattentional Blindess
- Machine Learning without Oversight
  - Turnkey analytics is dangerous
- Good Analysis
  - Requires iterative exploration
  - Peer review and collaboration

# Canonical Example

```python
class WordCount(Job):
    partitions = 3
    input=["sherlock.txt","poirot.txt","clouseau.txt"]

    @staticmethod
    def map(line, params):
        import string
        for word in line.split():
            yield word, 1

    @staticmethod
    def reduce(iter, params):
        from disco.util import kvgroup
        for word, counts in kvgroup(sorted(iter)):
            yield word, sum(counts)

if __name__ == "__main__":
    from count_words import WordCount

    for (word, counts) in result_iterator(WordCount.wait(show=True)):
        print word, counts
```