# Disco/Hadoop MapReduce

## Benjamin Zaitlen
### *Continuum Analytics*

## PyCon 2013

**CONTINUUM**
**A N A L Y T I C S**

# What is Large Data?

# What is Large Data?

- Can't fit into Excel

# What is Large Data?

- Can't fit into Excel
  - Increase Memory

# What is Large Data?

- Can't fit into Excel

  - Increase Memory

- Can't fit into R

# What is Large Data?

- Can't fit into Excel

  - Increase Memory

- Can't fit into R

  - Increase Memory

# What is Large Data?

- Can't fit into Excel

  - Increase Memory

- Can't fit into R

  - Increase Memory

- Can't fit into Memory

**CONTINUUM**
**ANALYTICS**

# What is Large Data?

- Can't fit into Excel

  - Increase Memory

- Can't fit into R

  - Increase Memory

- Can't fit into Memory

  - Increase Memory

**CONTINUUM**
**ANALYTICS**

# What is Large Data?

- 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

**CONTINUUM**
**A N A L Y T I C S**

# What is Large Data?

- 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.

**CONTINUUM**
**ANALYTICS**

# MapReduce

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

# MapReduce

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

- A mass of data: records

# 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

# MapReduce

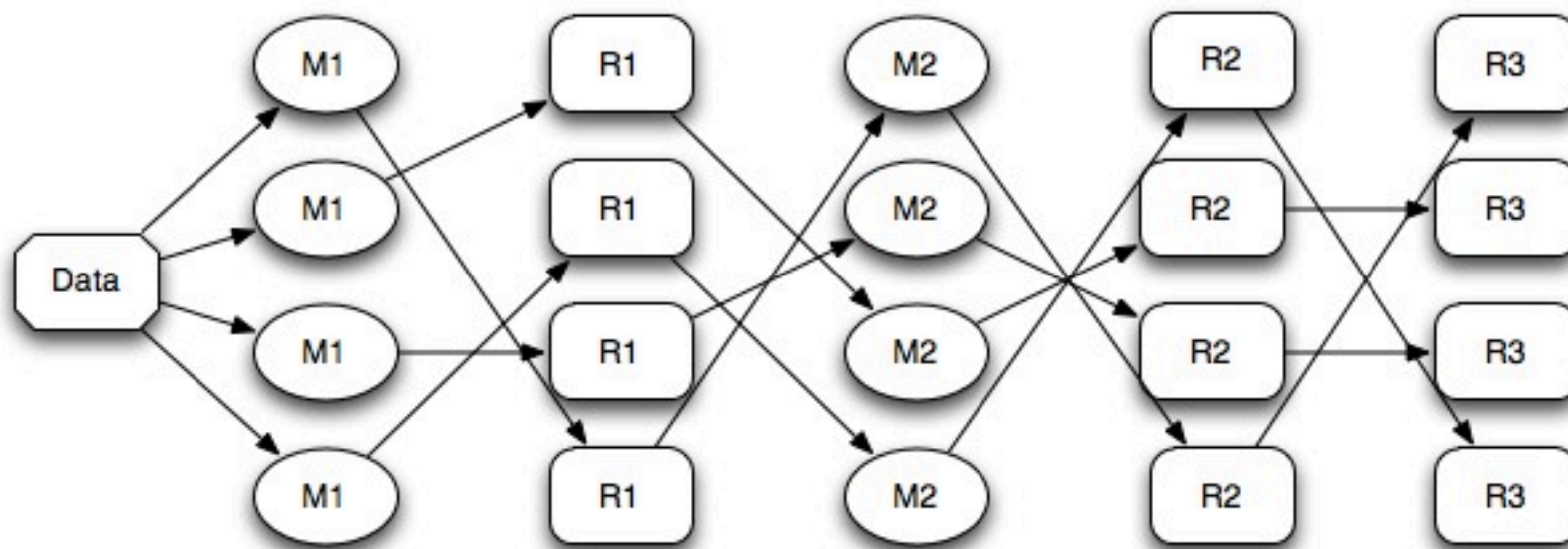Framework to help solve the problem of distributed computation for distributed data

- A mass of data: records

- Split/**Map** records into key-values pairs

- Collect/Partition kv pairs (Optional Sort)
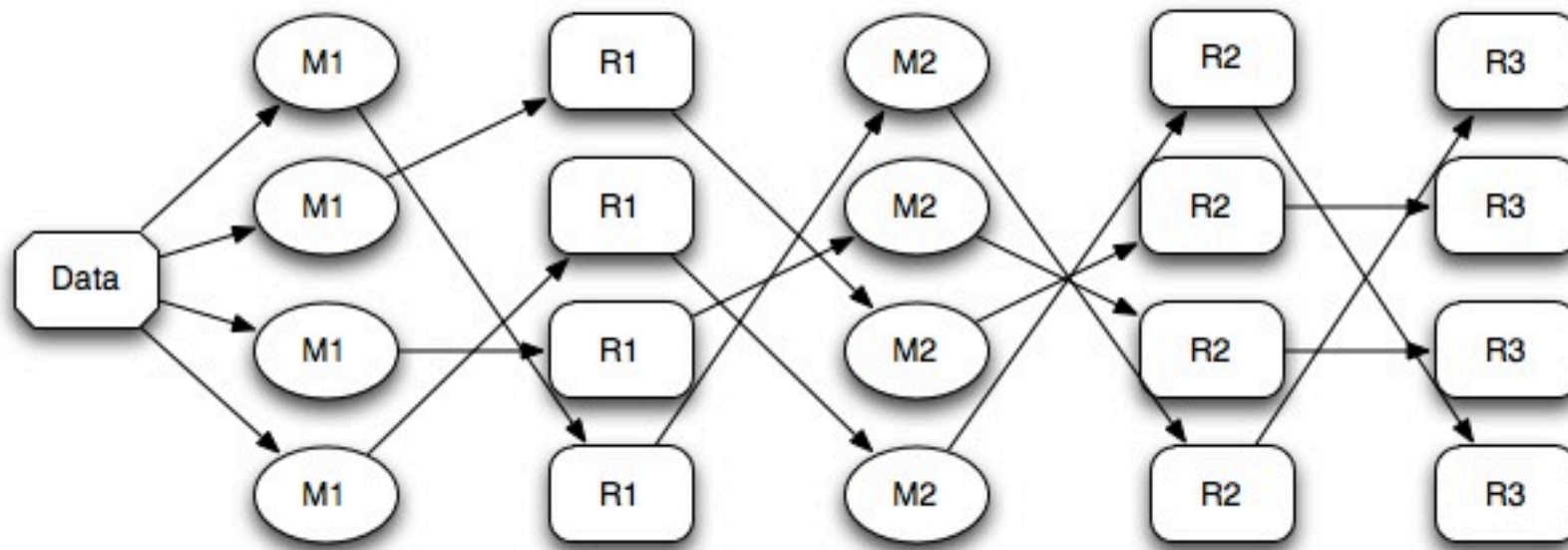
# MapReduce

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

# 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

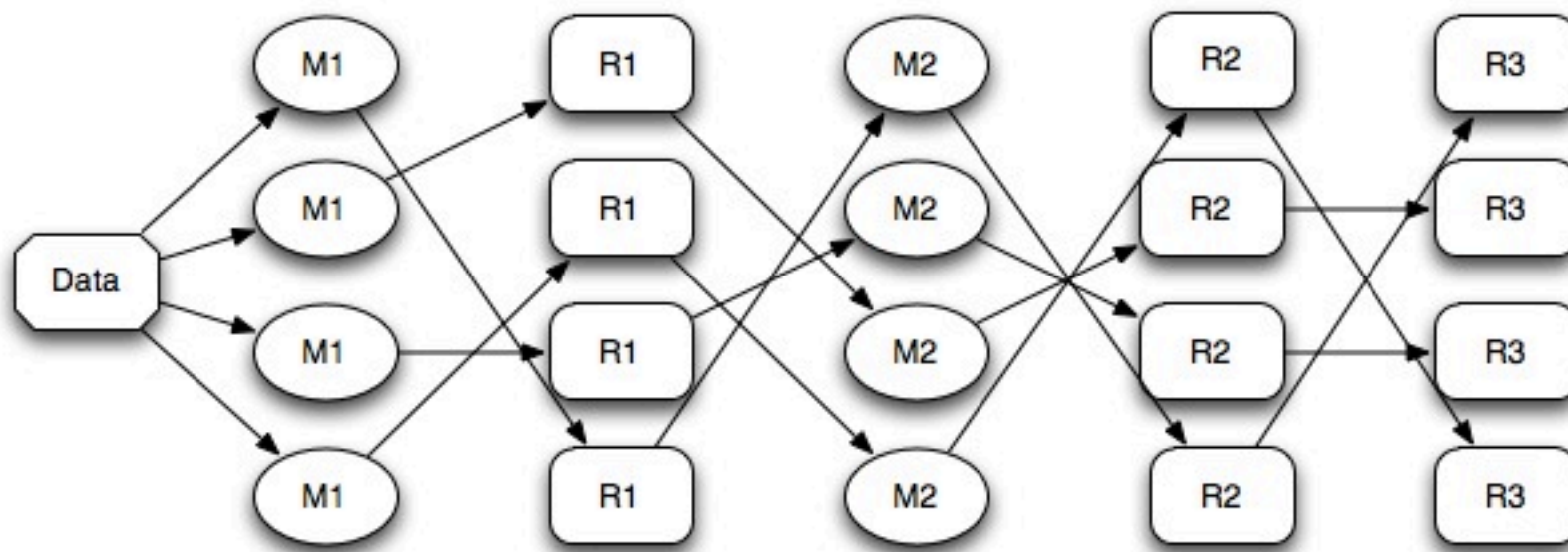- Result is returned

# MapReduce Workflow

# MapReduce Workflow



- Push code to data

# MapReduce Workflow



- Push code to data

- Lots of network traffic

# MapReduce it's a Party

# MapReduce it's a Party

# MapReduce it's a Party

# MapReduce it's a Party

- Disco: Python + Erlang

**CONTINUUM ANALYTICS**

# MapReduce it's a Party

- Disco: Python + Erlang

- Hadoop: Java-Dumbo (Python Streaming)

CONTINUUM
ANALYTICS

# MapReduce it's a Party

- Disco: Python + Erlang

- Hadoop: Java-Dumbo (Python Streaming)

- Distributed FileSystem: DDFS

# MapReduce it's a Party

- Disco: Python + Erlang

- Hadoop: Java-Dumbo (Python Streaming)

- Distributed FileSystem: DDFS

- Bring your friends

  - NumPy

  - SciPy

  - pandas

  - scikits-learn

  - OpenCV

# Canonical Example

```python
1
2   from disco.job import Job
3   from disco.worker.classic.func import chain_reader
4   from disco.core import result_iterator
5
6   class WordCount(Job):
7       partitions = 2
8       input=["sherlock_complete.txt","poirot_complete.txt"] #collected works
9
10      @staticmethod
11      def map(line, params):
12          import string
13          for word in line.split():
14              yield strippedWord, 1
15
16      @staticmethod
17      def reduce(iter, params):
18          from disco.util import kvgroup
19          for word, counts in kvgroup(sorted(iter)):
20              yield word, sum(counts)
21
22  if __name__ == "__main__":
23      from MapReduce_CountWords_Chain import WordCount
24
25      for (word, counts) in result_iterator(WordCount.wait(show=False)):
26          print word, counts
```
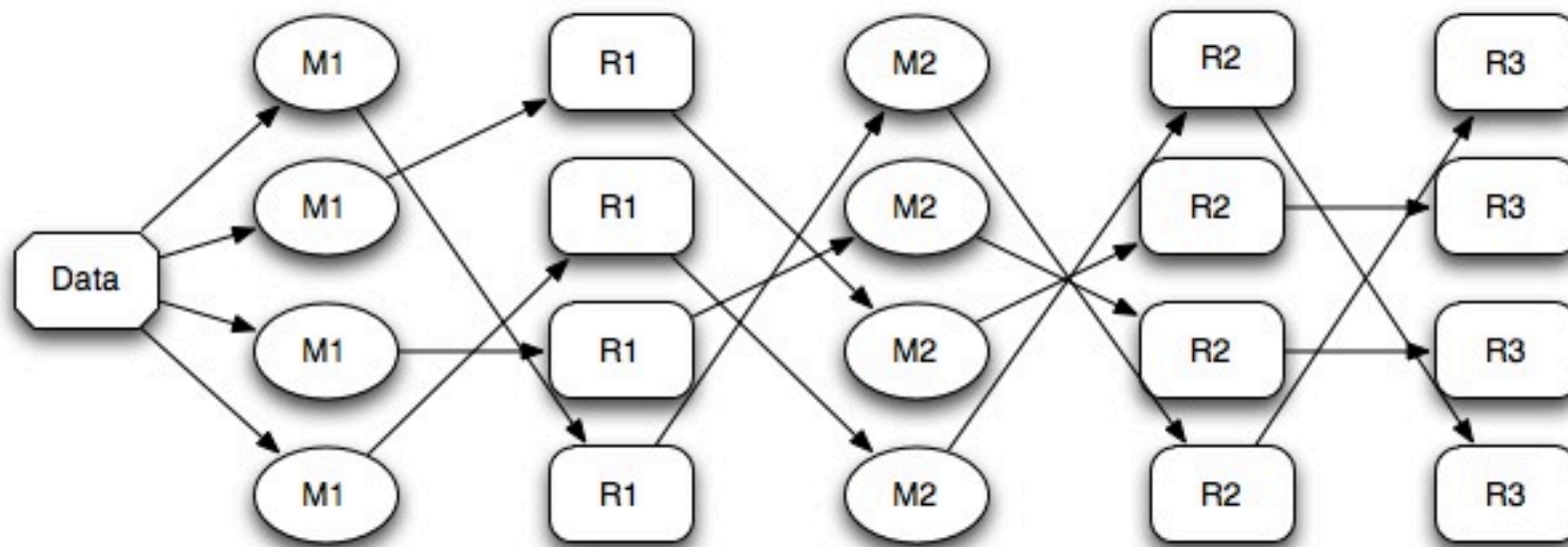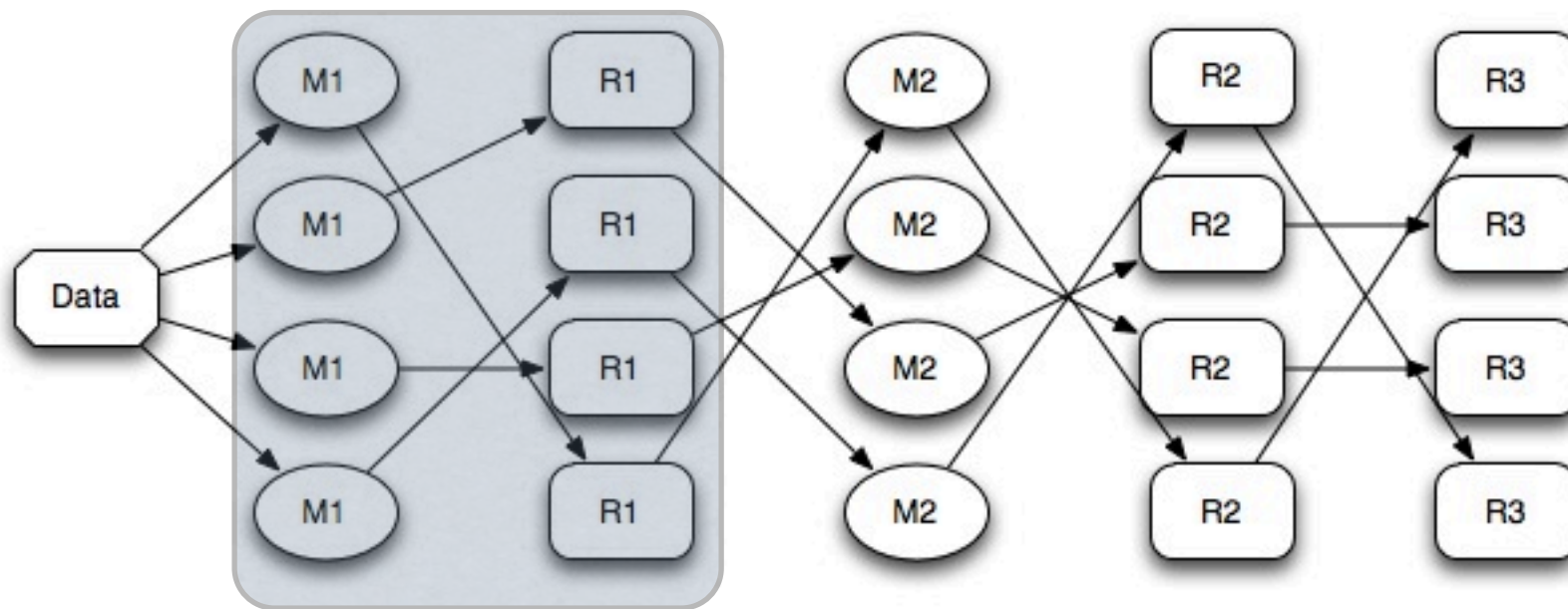
CONTINUUM
ANALYTICS

# Canonical Example

```python
from disco.job import Job
from disco.worker.classic.func import chain_reader
from disco.core import result_iterator

class WordCount(Job):
    partitions = 2
    input=["sherlock_complete.txt","poirot_complete.txt"] #collected works

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

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

CONTINUUM
A N A L Y T I C S

# Canonical Example

```python
1
2    from disco.job import Job
3    from disco.worker.classic.func import chain_reader
4    from disco.core import result_iterator
5
6    class WordCount(Job):
7        partitions = 2
8        input=["sherlock_complete.txt","poirot_complete.txt"] #collected works
9
10       @staticmethod
11       def map(line, params):
12           import string
13           for word in line.split():
14               yield strippedWord, 1
15
16       @staticmethod
17       def reduce(iter, params):
18           from disco.util import kvgroup
19           for word, counts in kvgroup(sorted(iter)):
20               yield word, sum(counts)
21
22   if __name__ == "__main__":
23       from MapReduce_CountWords_Chain import WordCount
24
25       for (word, counts) in result_iterator(WordCount.wait(show=False)):
26           print word, counts
```
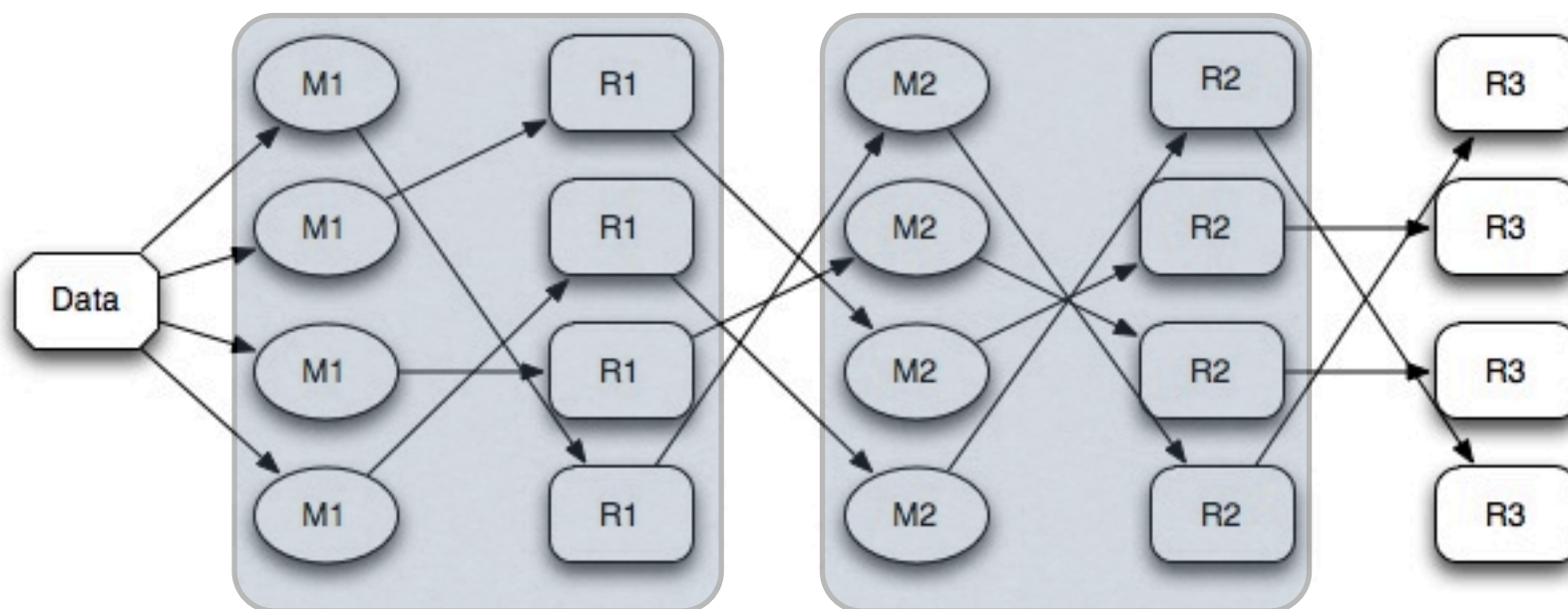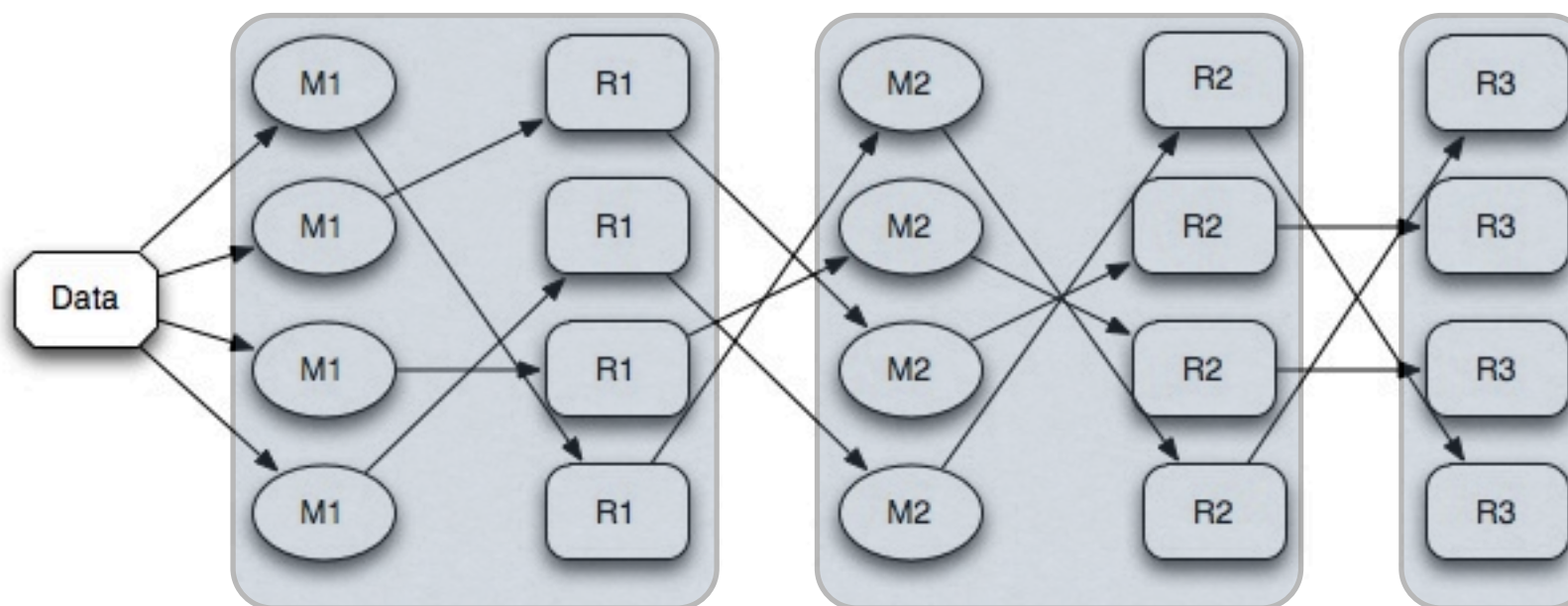
**CONTINUUM**
**A N A L Y T I C S**

# Chain Jobs

# Chain Jobs

# Chain Jobs

# Chain Jobs

# MapReduce Thoughts

# MapReduce Thoughts

- Data Cleansing

  - Everyone's pain point
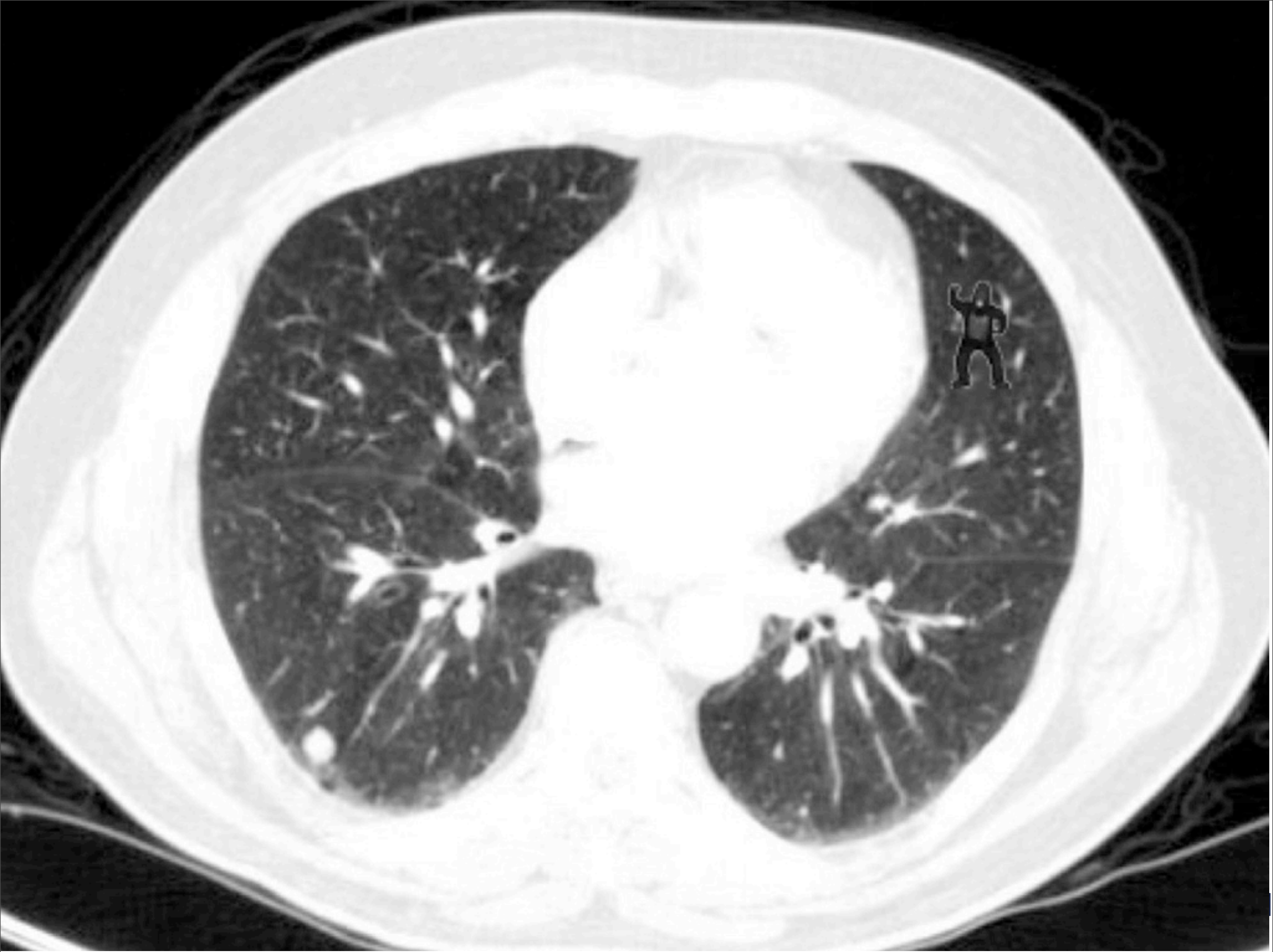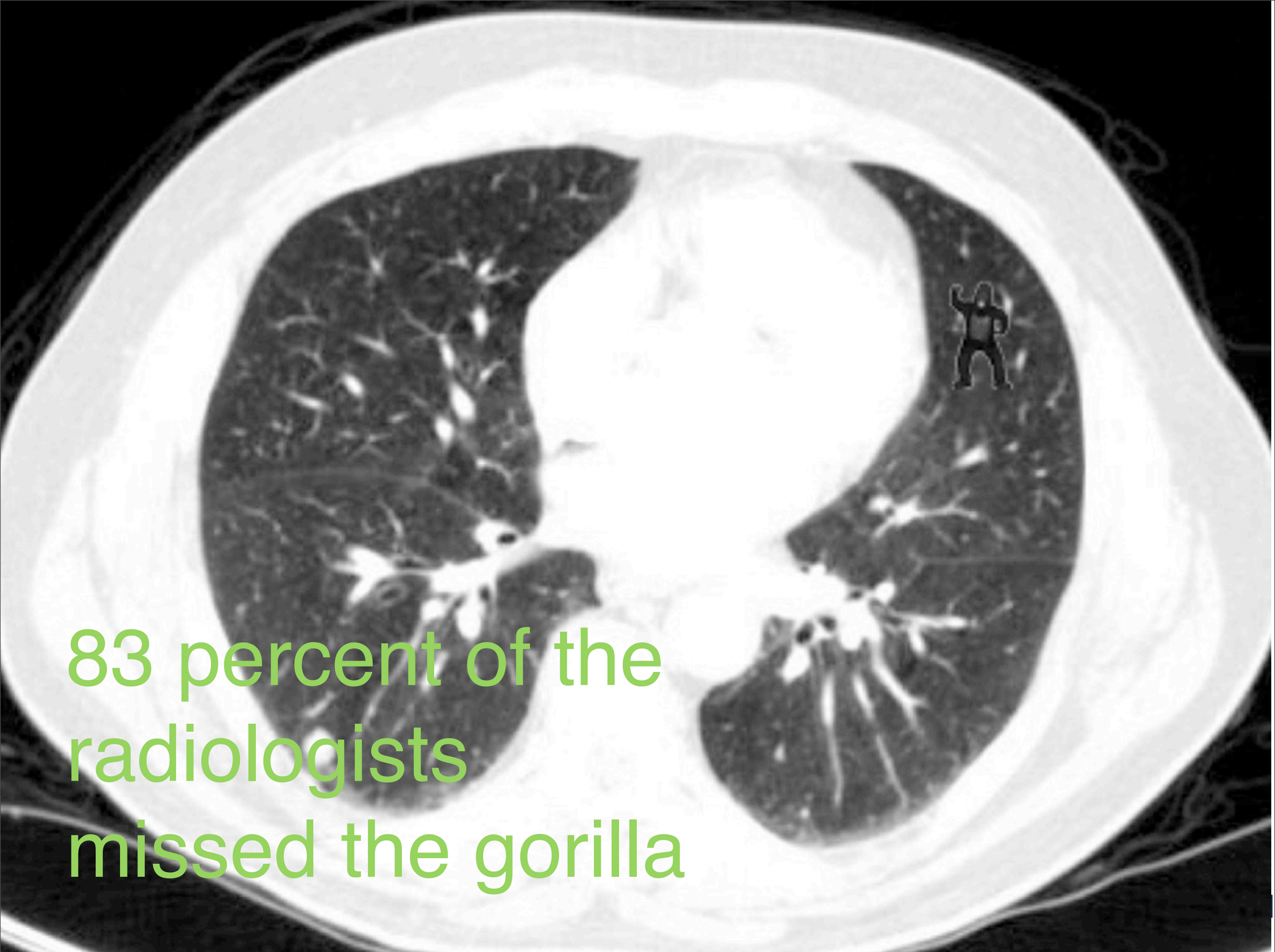
# MapReduce Thoughts

- Data Cleansing

  - Everyone's pain point

- Task Deconstruction

  - Good for code management

  - Hides -- in a good way -- data management

# 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
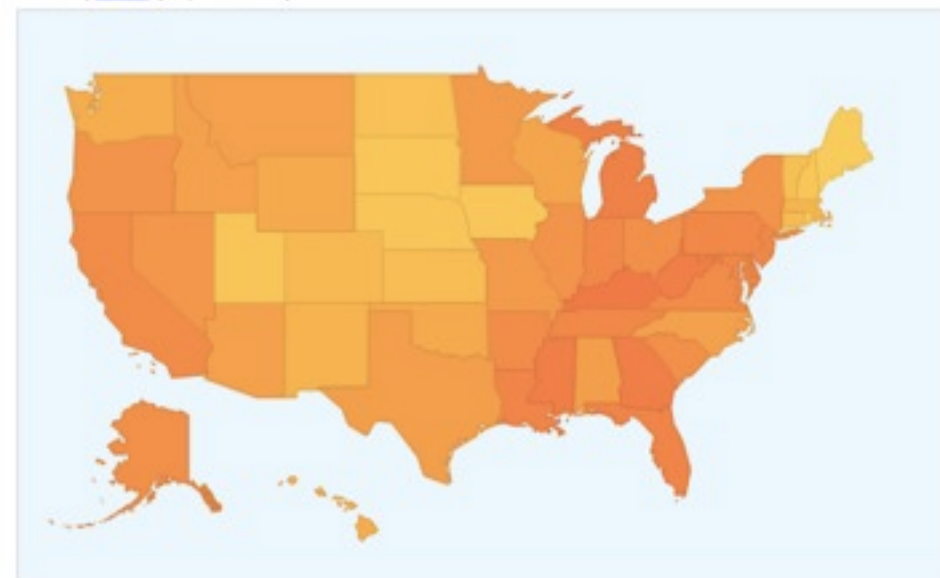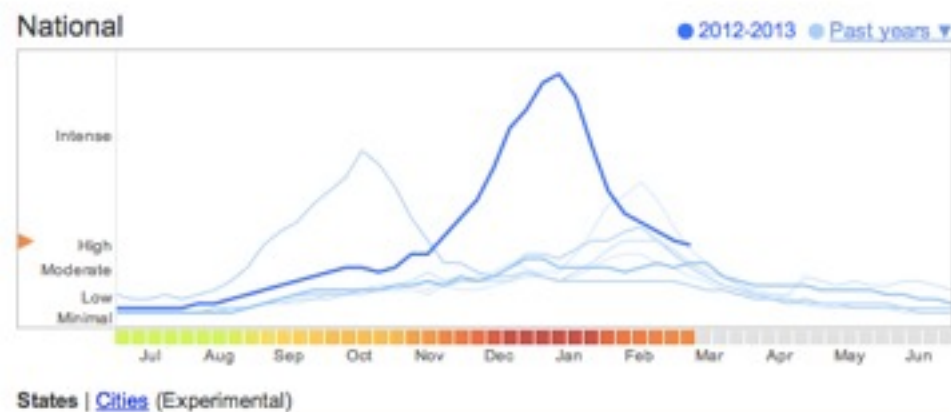
  - Job organization

# Data Thoughts

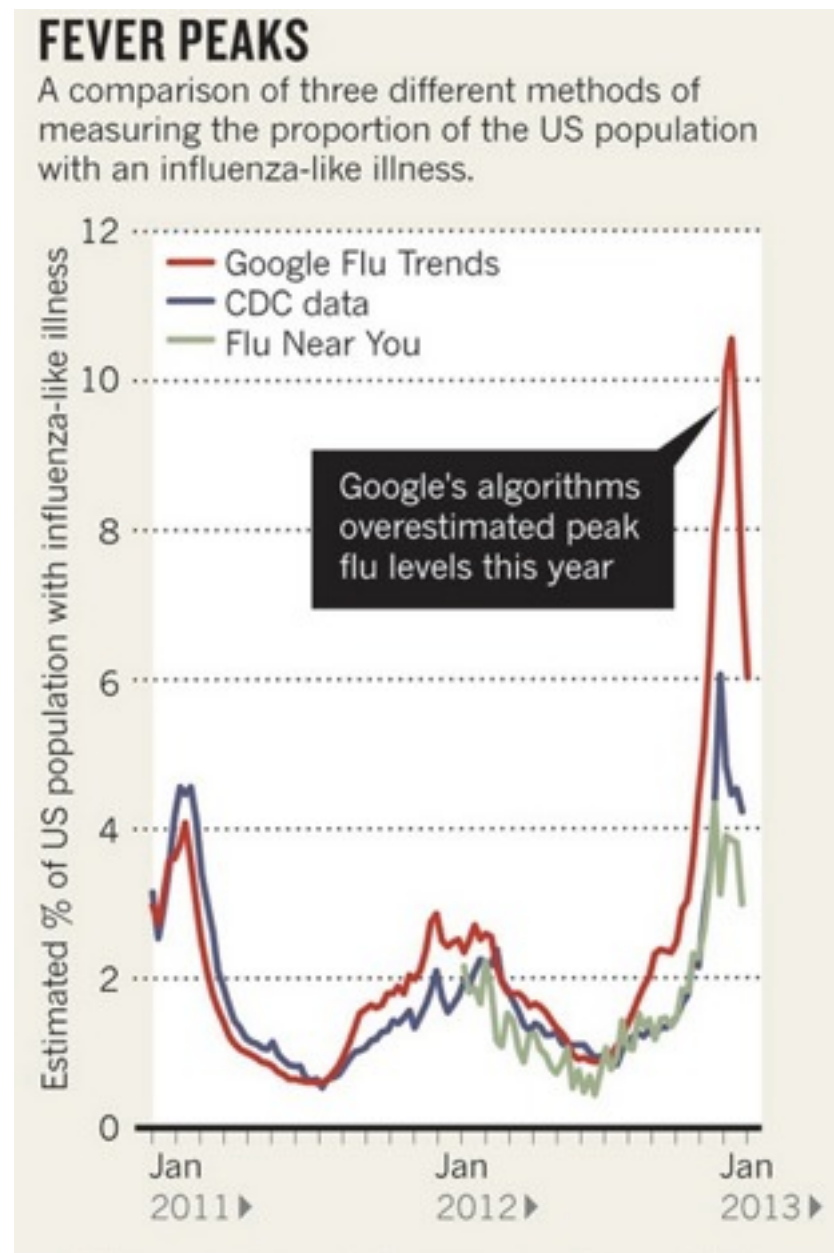83 percent of the radiologists missed the gorilla

# Google Flu



- Data Mining

- Faster than CDC

# Google Get's 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

# Wrapping Up

- Invisible Gorillas will stay invisible

- Turnkey analytics is dangerous

- Good Analysis

  - Requires iterative exploration

  - Peer review and collaboration

**CONTINUUM**
**A N A L Y T I C S**