mrjob: part 1

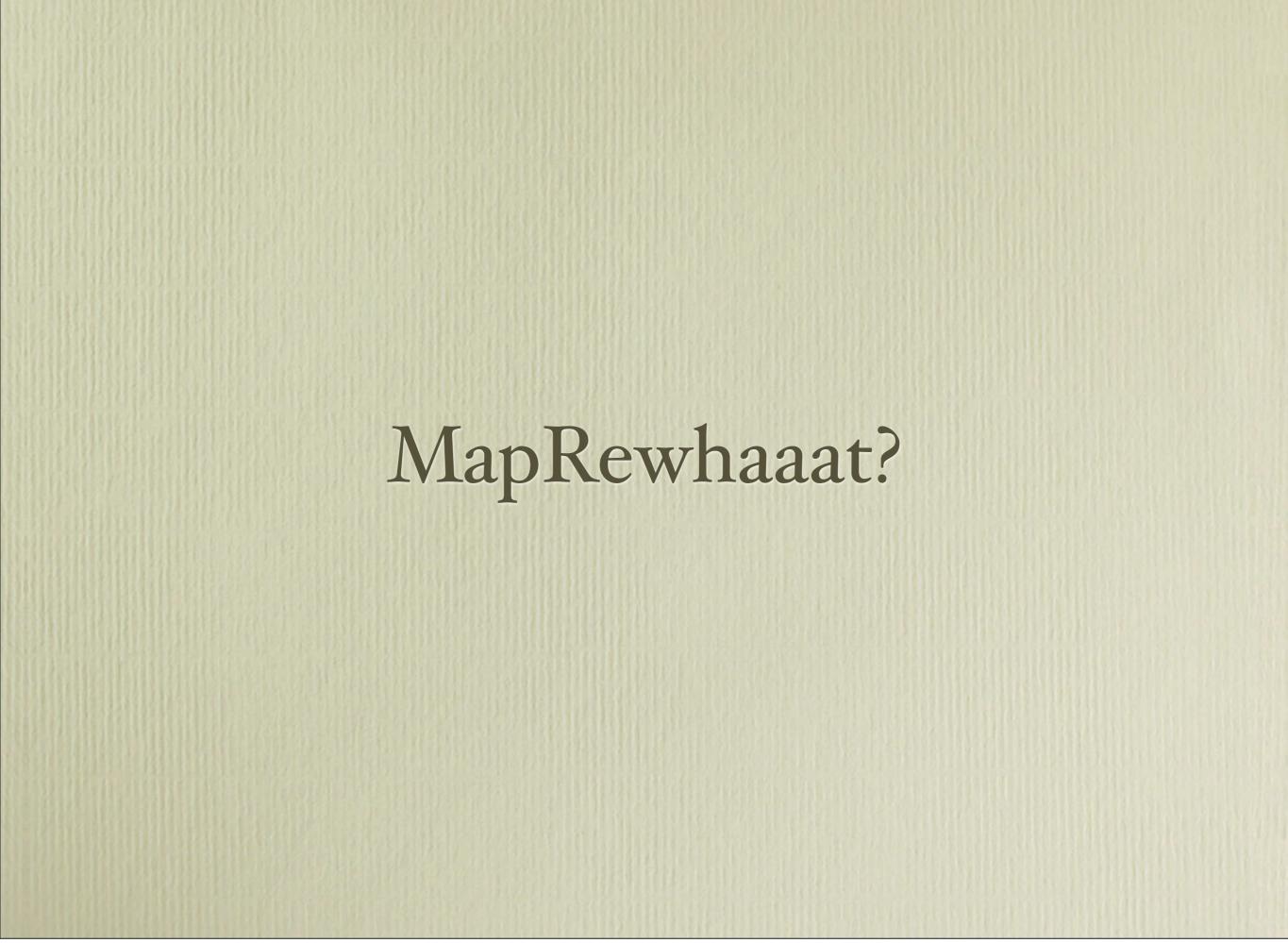
steve johnson - steve@steveasleep.com

### Reference slide comes first

- <a href="http://mrjob.readthedocs.org/">http://mrjob.readthedocs.org/</a> (very latest)
- <a href="http://github.com/yelp/mrjob">http://github.com/yelp/mrjob</a> (read the source if you get bored)
- <a href="http://www.python.org/dev/peps/pep-0008/">http://www.python.org/dev/peps/pep-0008/</a> (if you want to be nice to me when I read your code)

## Setup

- 1. Have a Github account (email your username to steve+bigdive@steveasleep.com)
- 2. clone <a href="https://github.com/irskep/mrjob\_course">https://github.com/irskep/mrjob\_course</a> on Github
- 3.git clone <a href="https://github.com/YOURNAME/">https://github.com/YOURNAME/</a> <a href="mrjob\_course.git">mrjob\_course</a> <a href="mrjob\_course.git">ab. com/YOURNAME/</a>
- 4.pip install -r requirements.txt
- 5. <a href="http://www.yelp.com/dataset\_challenge">http://www.yelp.com/dataset\_challenge</a>: Get The Data
- 6.Put it in mrjob\_course/data/yelp, run make from mrjob\_course



# "map"

• Input looks like this:

<key>\t<value>\n

• In Python, your function looks like this:

```
def mapper(self, key, value):
    # as many times as you want
    yield <new key>, <new value>
```

### "reduce"

• Input looks like this:

```
<key1>\t<value>\n
<key1>\t<value>\n
<key2>\t<value>\n
<key2>\t<value>\n
<key2>\t<value>\n
```

• In Python, your function looks like this:

```
def reducer(self, key, values):
    # values for the same key
    for value in values:
        # do something
    yield <new key>, <new value>
```

### Exercises!



### Protocols

Forward dive\n
Twisting dive\n
Small dive\n

- No key
- Raw lines of text
- Use RawValueProtocol in mrjob
  - first task sees None for key

#### Protocols

dive\_type\tForward dive\n
dive\_type\tTwisting dive\n
dive\_type\tSmall dive\n

- Raw string key separated from value by \t
- Raw string value
- Use RawProtocol in mrjob

### Protocols

```
"dive"\t{"type": "forward"}\n
"dive"\t{"type": "twisting"}\n
"dive"\t{"type": "small"}\n
```

- JSON-encoded keys and values separated by \t
- Use JSONProtocol in mrjob
- JSONValueProtocol exists for reading or writing JSON data without keys

### Exercises!

