

# Deploying Data Science Systems

Dr. Daniel Schien

[Daniel.schien@bristol.ac.uk](mailto:Daniel.schien@bristol.ac.uk)

## Outline

- Recap
- Motivation
- Data Science Process
- Cloud, Docker
- Big Data

# Recap

- Lec 1 Intro
- Lec 2 Data Ingress
- Lec 3 Recommender Systems
- Lec 4 Databases
- Lec 5 Data Wrangling
- Lec 6 Data Fusion
- Lec 7 Data Exploration
- Lec 8 Data Visualisation
- Lec 9 Data sharing, privacy and anonymisation
- **Lec 10 Deploying data science systems**
- Lec 11 The future of data science

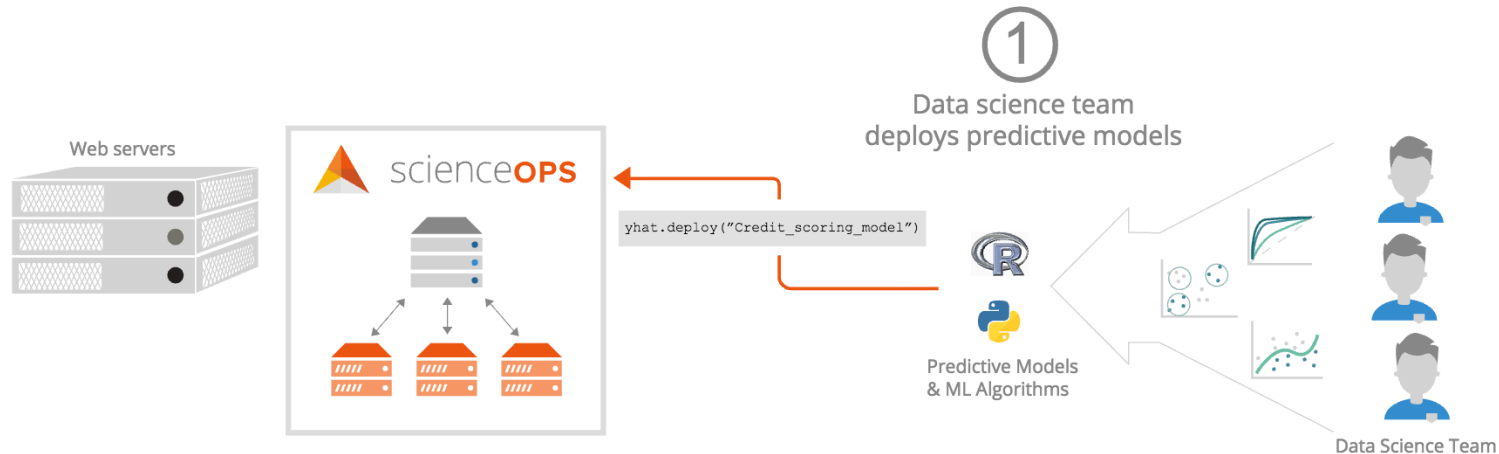
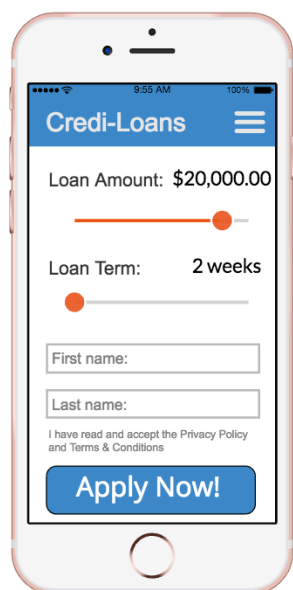
## Learning Outcomes

- deployment and knowledge of **infrastructure** with the role of data scientists
- Have awareness for some key properties in the data science **software process**
- Have awareness for the key concepts in **Big Data** infrastructures
- Have an overview of key **cloud** deployment building blocks
- Be able to **package an ipython notebook** for cloud deployment

- Question: What is deployment?

# Deployment Concerns

- On your machine



- Question: How does deployment overlap with the role of the data scientist?

# The data science process

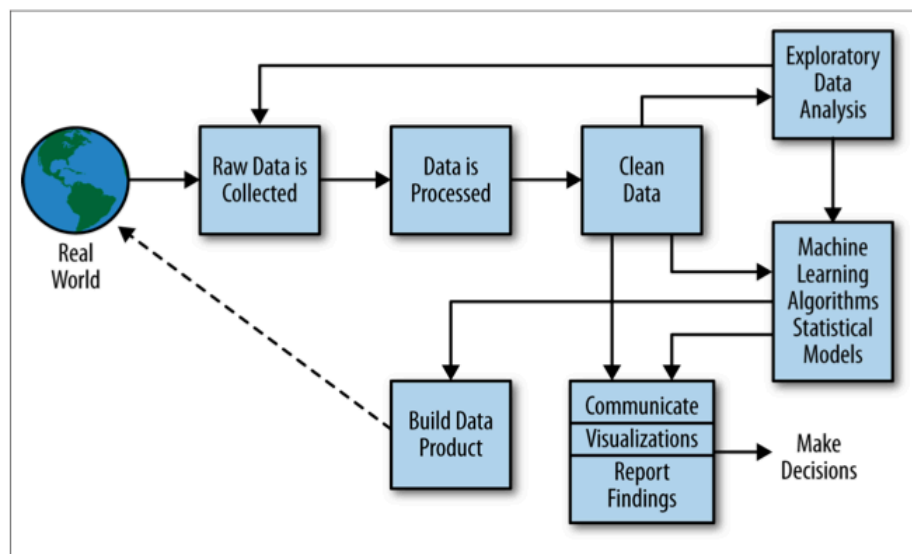


Figure 2-2. The data science process



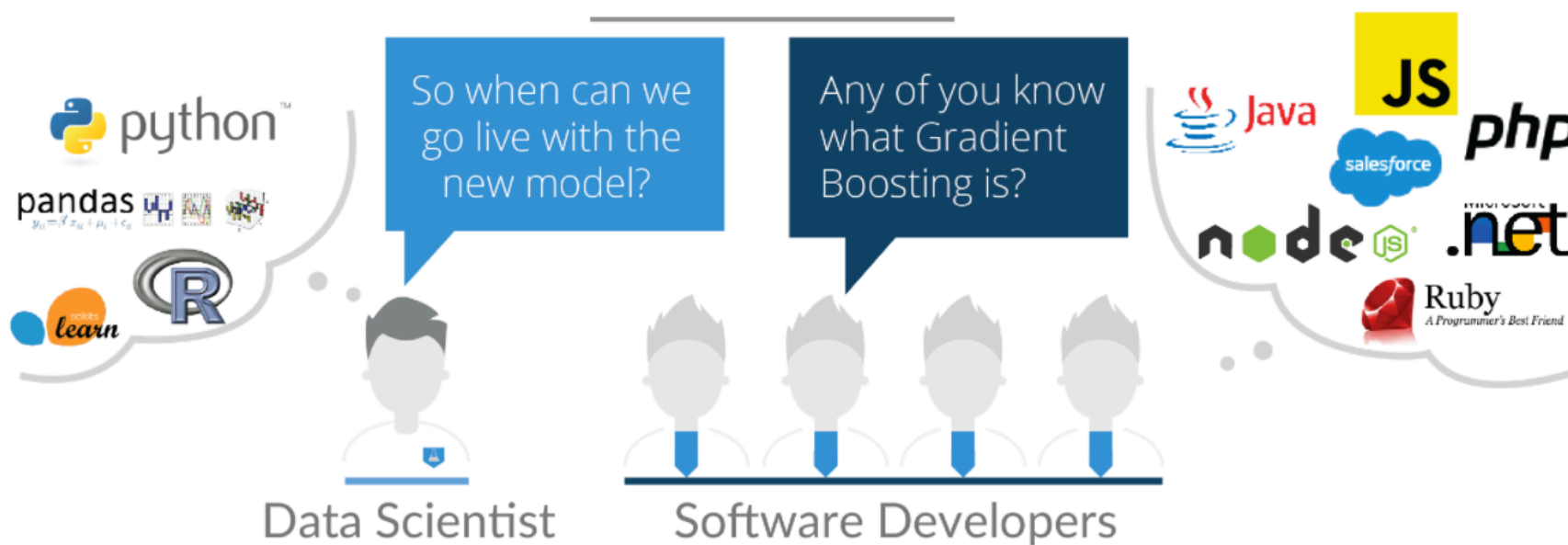
Rachel Schutt & Cathy O'Neil



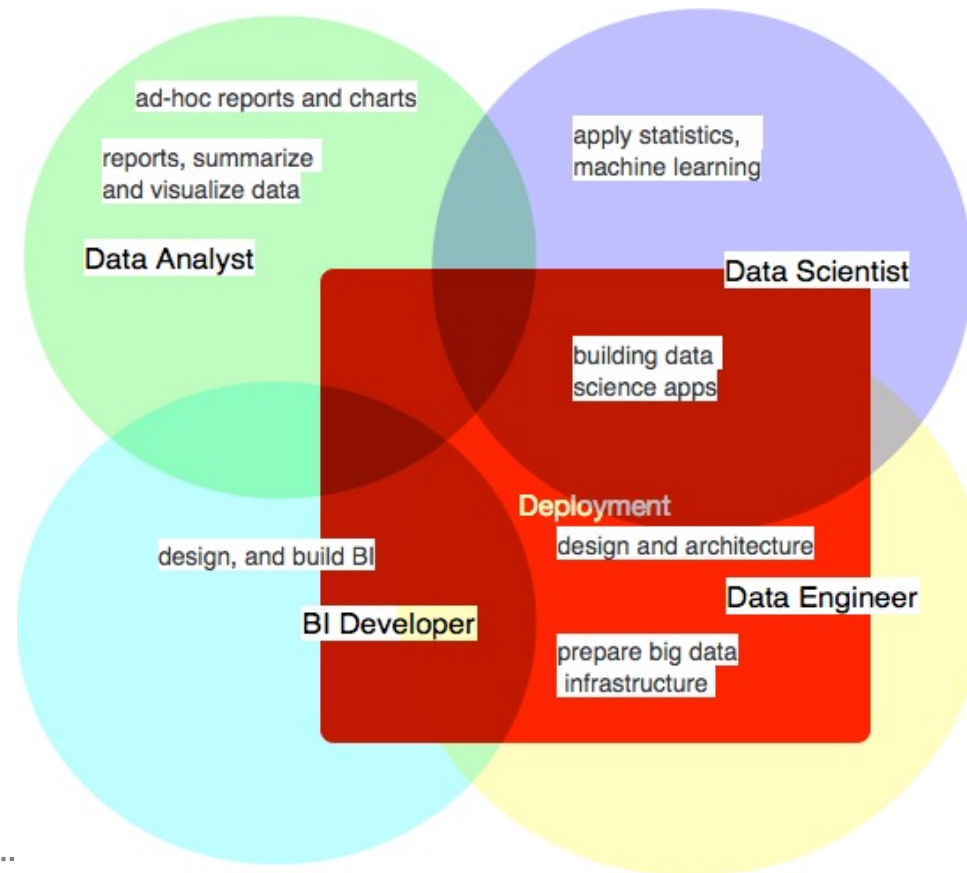
## Data Products

- Business Intelligence and Reports
  - Data Analytics Systems
    - Commercial systems to build data products by analysts
      - Tableau
      - <https://looker.com/product>
      - Knime
      - etc
-

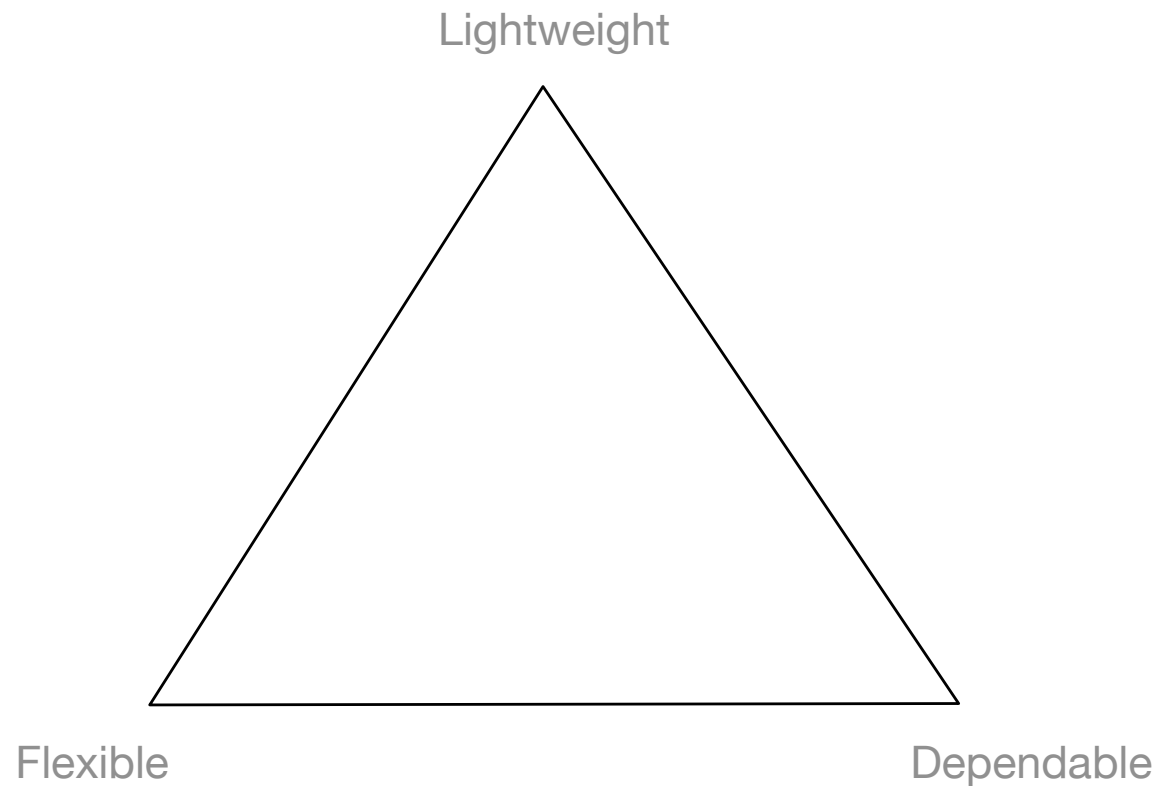
## Now What?



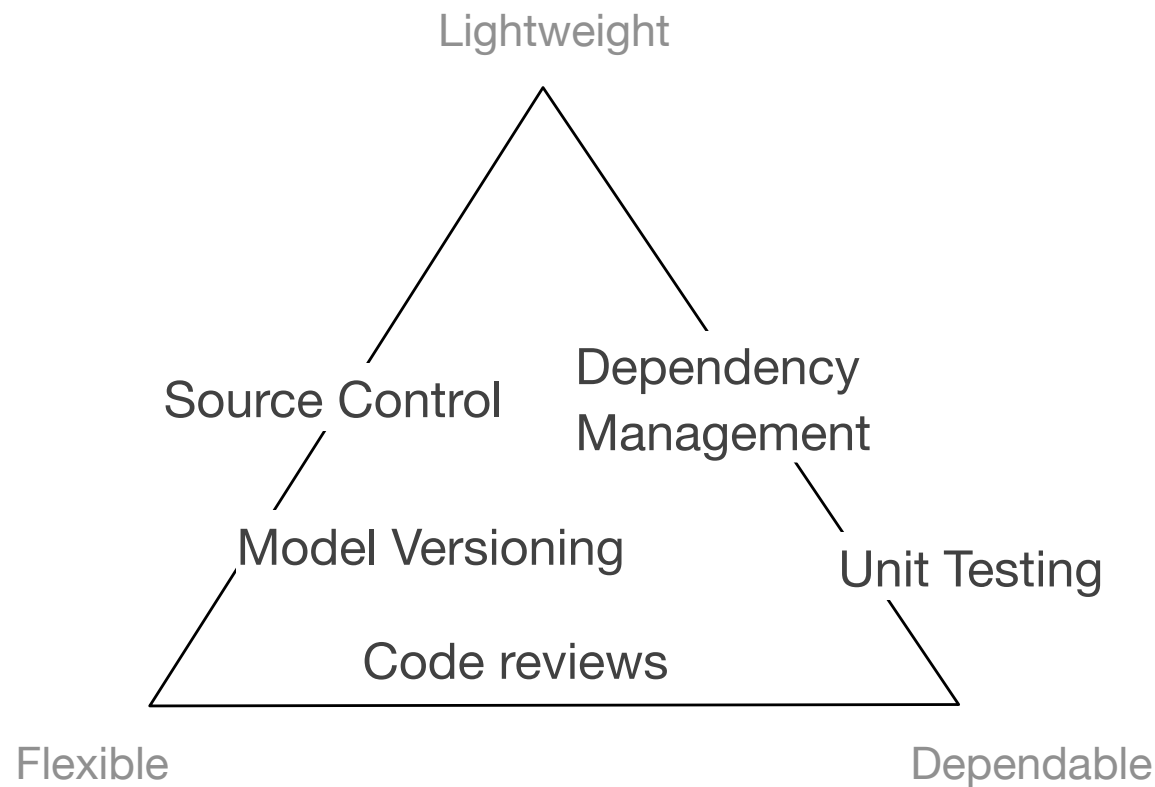
# Data Roles



# Data Science Software Engineering 101



# Data Science Software Engineering 101



# Dependency Management

- Dependency hell
- Reproducibility
- Facilitation to on-board new people

# Model Versioning

## MODEL MANAGEMENT

Test multiple models at once, promote new models into production, and continuously iterate over time.



# Unit Testing

- TDD
- CI
- nosetests, pytest, unittest2
- Source Control – remove output from ipynb's before commit





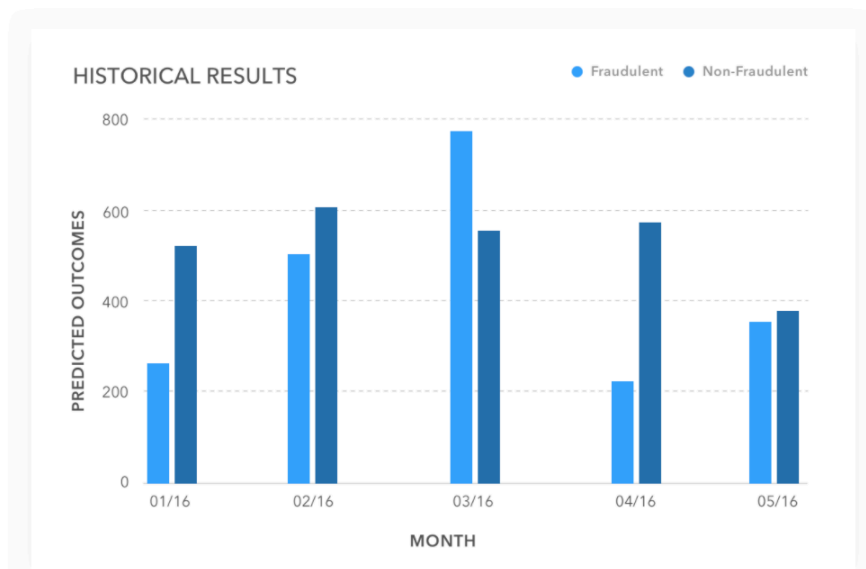
## Quality

- Code reviews
- Pair programming
- Verification
- Validation

n=165	Predicted: NO	Predicted: YES	
	Actual: NO	TN = 50	FP = 10
Actual: YES	FN = 5	TP = 100	105
	55	110	

# Monitoring

- System Health Overview
- Logging



## MODEL MONITORING

Measure model performance and its impact on your business with data from API calls, training, and cross validation.














## Infrastructure as Code

- maintainable, versionable, testable, and collaborative.
  - Fabric
  - CloudFormation
  - <https://www.terraform.io/>
  - Chef
  - etc
-

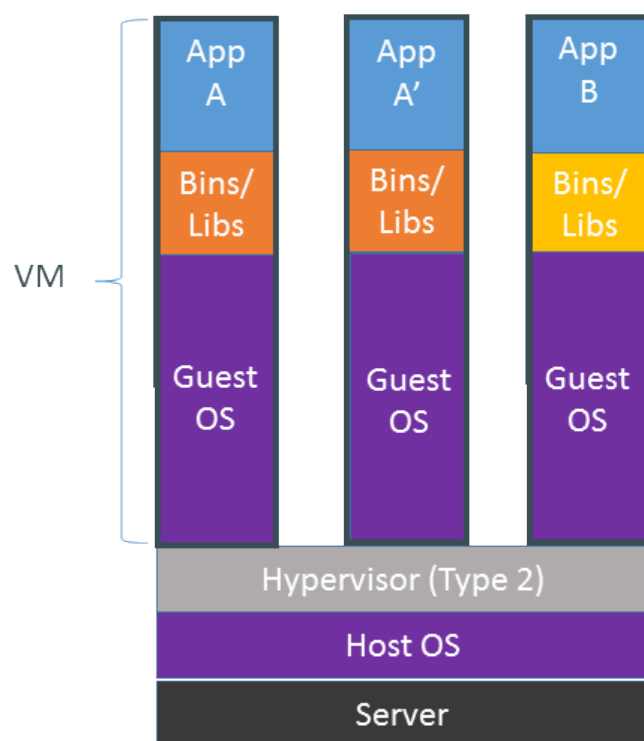
# Important infrastructure

- Your machine
- Cloud
- Big Data

# Docker the-matrix-from-hell

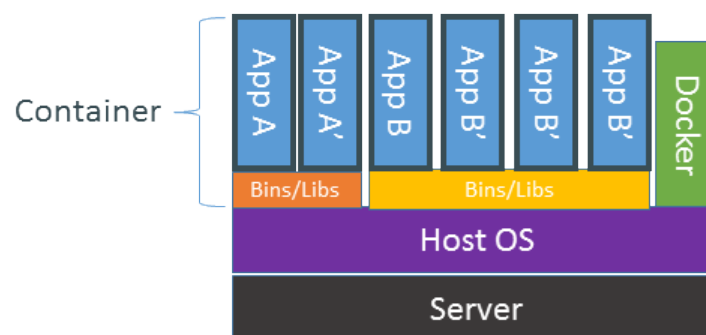
	Static website	?	?	?	?	?	?	?
	Web frontend	?	?	?	?	?	?	?
	Background workers	?	?	?	?	?	?	?
	User DB	?	?	?	?	?	?	?
	Analytics DB	?	?	?	?	?	?	?
	Queue	?	?	?	?	?	?	?
		Development VM	QA Server	Single Prod Server	Onsite Cluster	Public Cloud	Contributor's laptop	Customer Servers
								

# Docker vs VMs

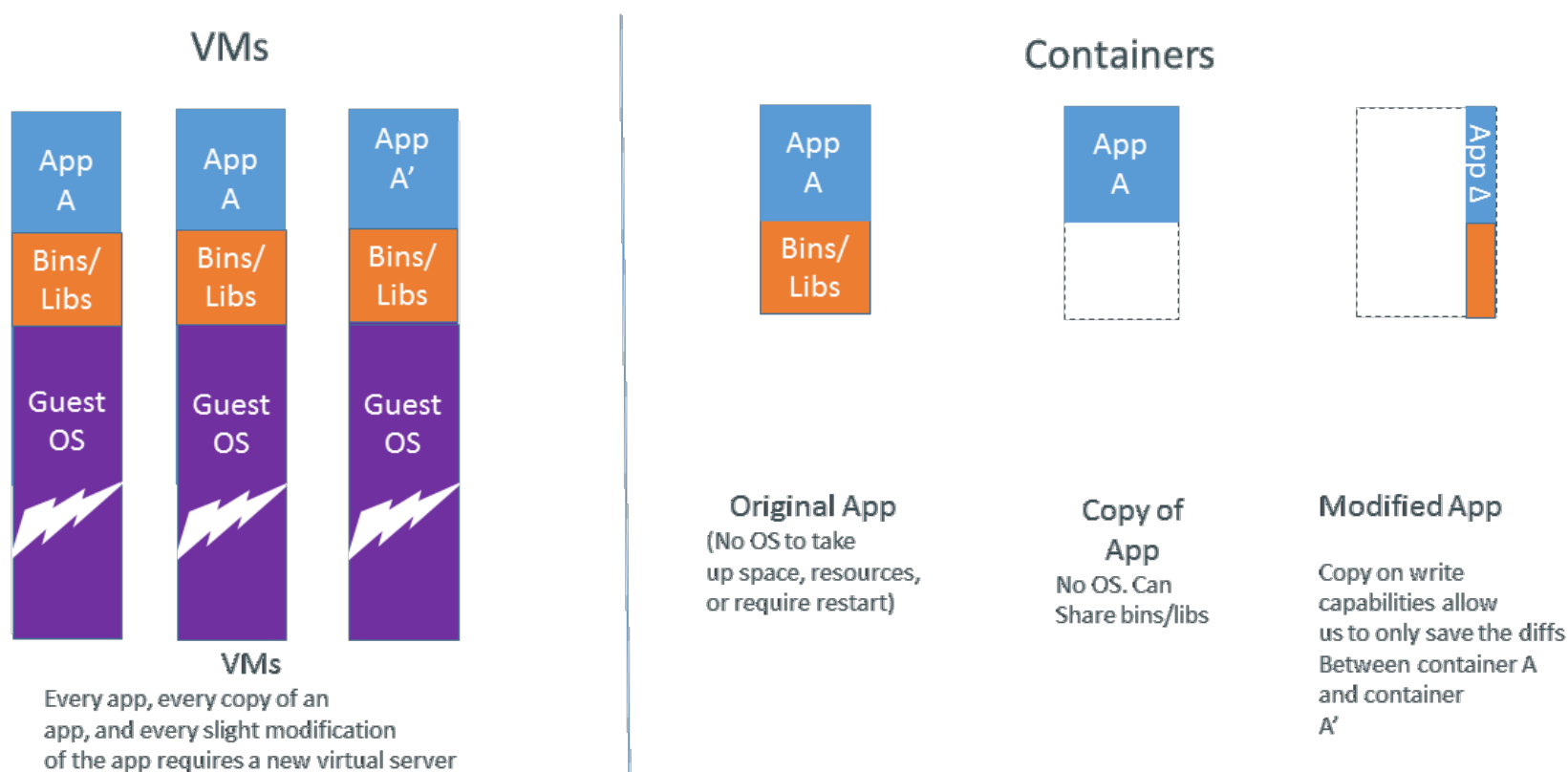


Containers are isolated,  
but share OS and, where  
appropriate, bins/libraries

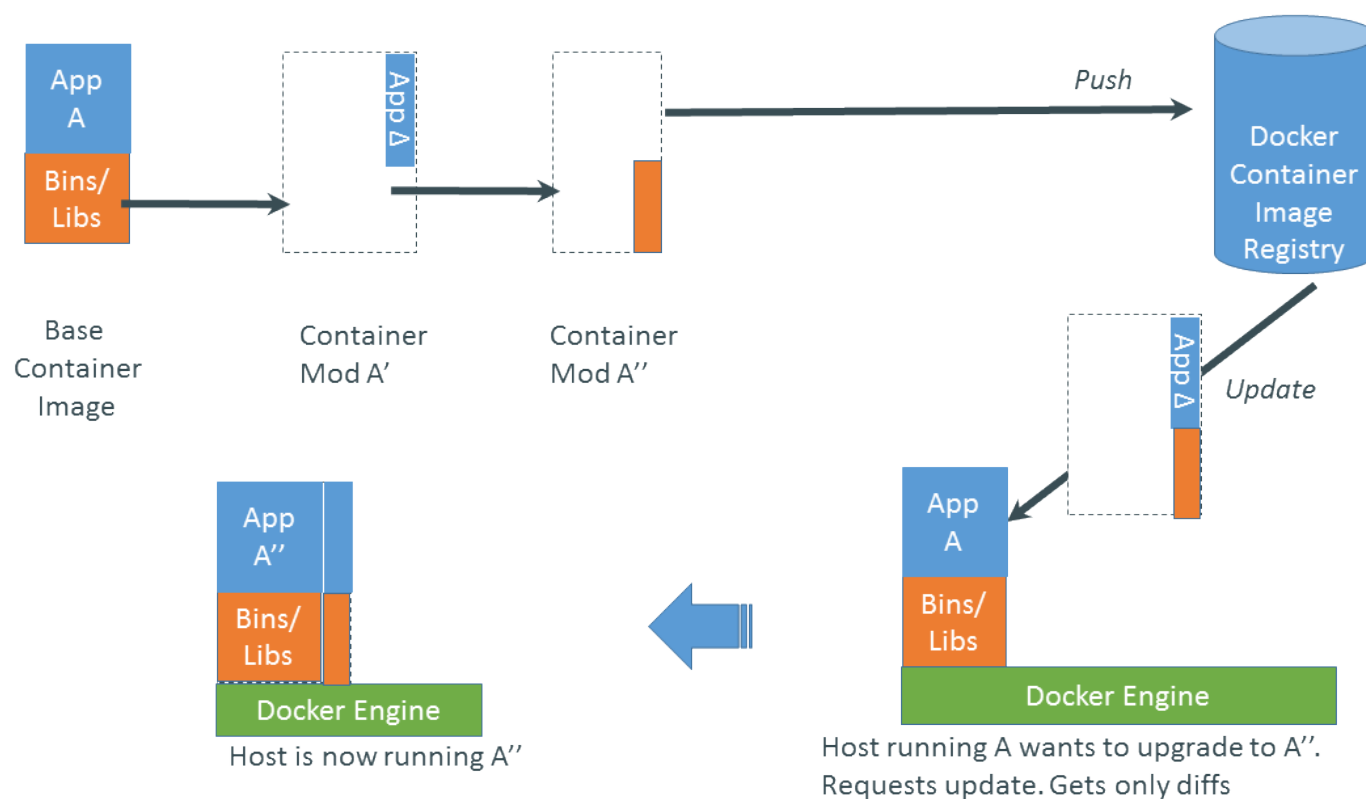
...result is significantly faster deployment,  
much less overhead, easier migration,  
faster restart



# Docker - Diffs and Layers

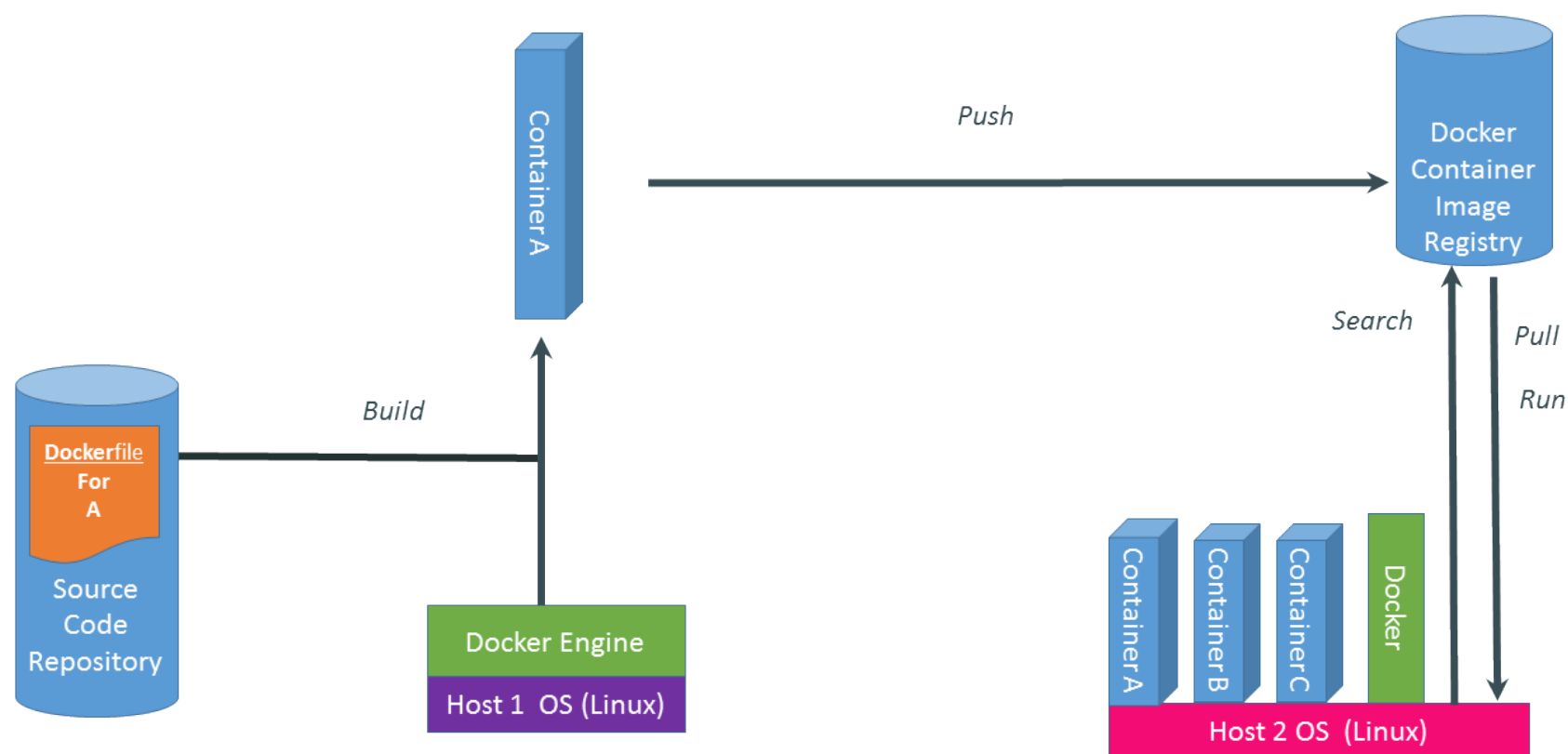


# Docker - Change Cycle





# Docker - Deployment



## Basic Deployment

- Ipython notebook in a container
- Deploy your jupyter notebook online
- simple notebook with widgets
- Dockerfile
- build container
- run with mount and password
- deploy remote with fabric - [fabfile.org](http://fabfile.org)



# Demo

[https://github.com/dschien/ads\\_notebook](https://github.com/dschien/ads_notebook)

```
→ ads git:(master) ✗ docker build -t dschien/ads .
Sending build context to Docker daemon 179.2 kB
Step 1/9 : FROM jupyter/minimal-notebook
----> d9b21ff9ceac
Step 2/9 : MAINTAINER Dan Schien <dschien@gmail.com>
----> Using cache
----> 23087455a6fb
Step 3/9 : COPY requirements.txt /opt/app/requirements.txt
----> Using cache
----> 7db351f9a3b1
Step 4/9 : WORKDIR /opt/app
----> Using cache
----> 505a4a3adc51
Step 5/9 : RUN pip install -r requirements.txt
----> Using cache
----> 2c01e4710577
Step 6/9 : RUN pip install jupyter_dashboards
----> Using cache
----> 112ed28b3bc3
Step 7/9 : RUN jupyter dashboards quick-setup --sys-prefix
----> Using cache
----> 162d98b873ed
Step 8/9 : RUN jupyter nbextension enable --py --sys-prefix widgetsnbextension
----> Using cache
----> 79c769041a1d
Step 9/9 : USER jovyan
----> Using cache
----> 31cc2868e89e
*** Successfully built 31cc2868e89e
```

```
Dockerfile x
FROM jupyter/minimal-notebook

MAINTAINER Dan Schien <dschien@gmail.com>

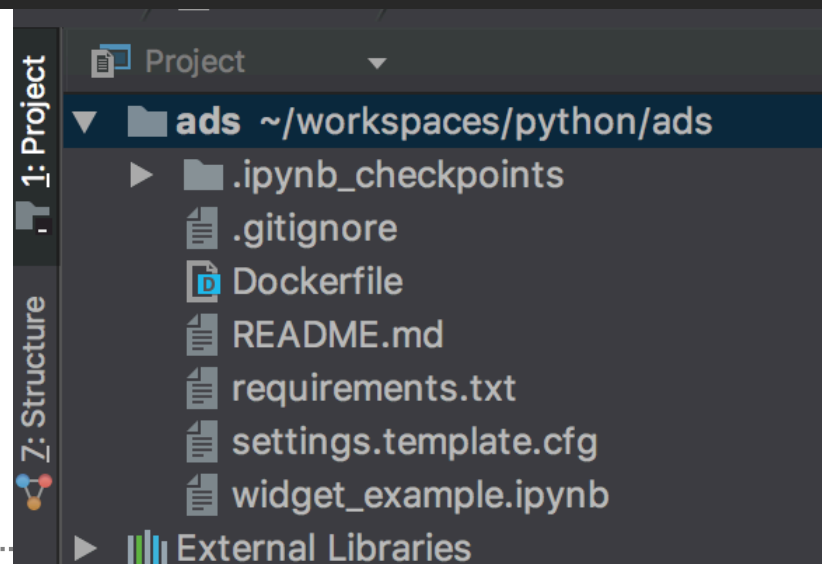
COPY requirements.txt /opt/app/requirements.txt
WORKDIR /opt/app

RUN pip install -r requirements.txt

RUN pip install jupyter_dashboards
RUN jupyter dashboards quick-setup --sys-prefix

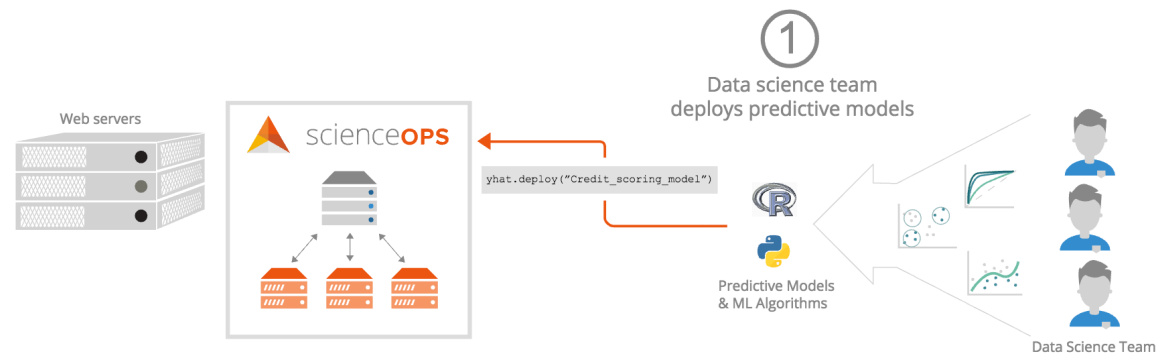
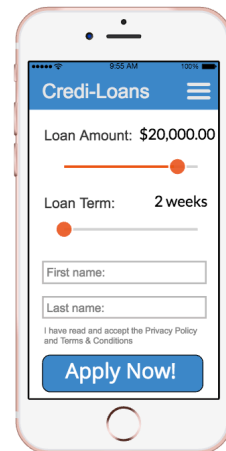
RUN jupyter nbextension enable --py --sys-prefix widgetsnbextension

USER jovyan
```



## Further Packaging your APP

- Provide the results but not the model
- Custom App with Flask
- REST Api



# Compute Developments

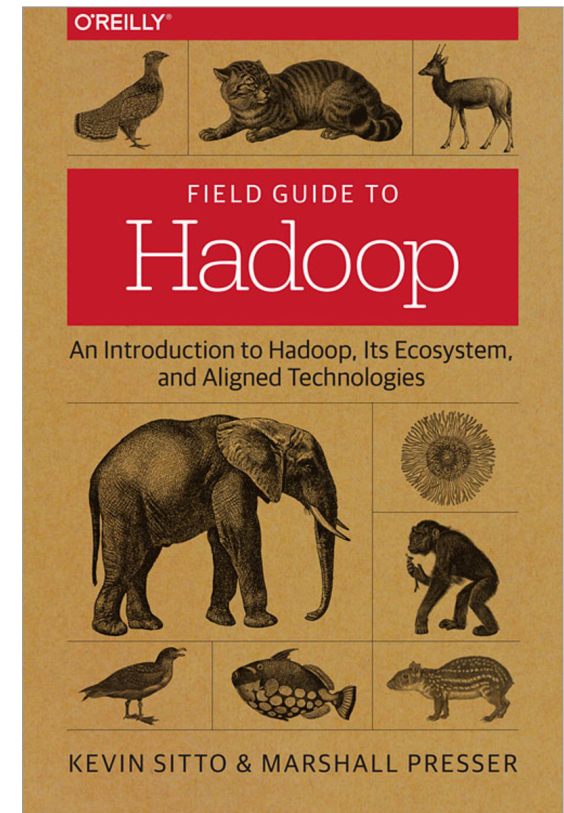
- CPUs don't get faster -> since about 2005 they have stabilised around 3Gh, since then we need to parallelise
- if your datasets get very large, they no longer fit into your (or any commodity hardware) machine -> parallelise
- how can you efficiently **parallelise**
  - Many machines
  - Many cores
- at scale, machines fail -> you need to know that because your algorithms will be affected by it (or rather, you can design algorithms accordingly)

# Model Optimisation

- Bottlenecks,
- stragglers
- hot spots
- Simple timing approaches, benchmarking

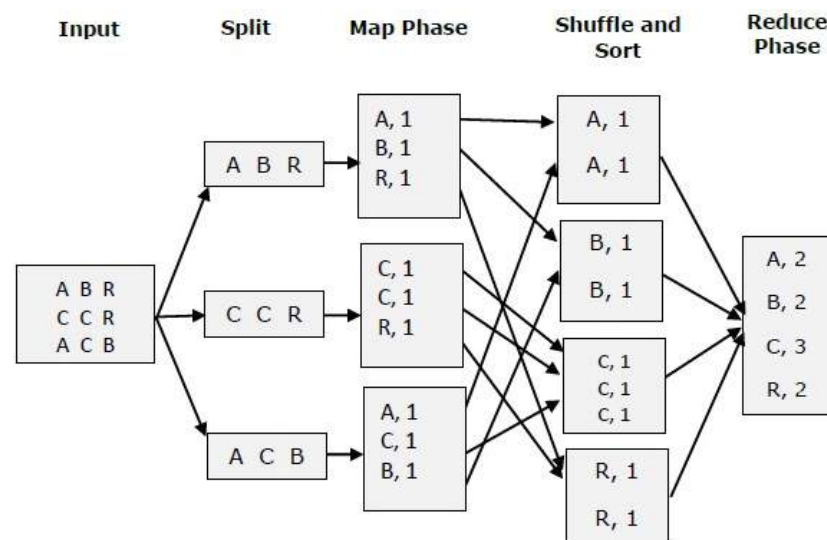
# Hadoop

- **HDFS**- The default storage layer
  - splits data up
  - provides failure recovery through duplication of contents among separate nodes
- **MapReduce**
  - Googles key innovation to process HTML file indexes at web scale
- **YARN**- Responsible for cluster management and scheduling user applications
- Hosted MapReduce available
  - [Cloudera](#), [Hortonworks](#) and [Splice](#)



## Roll your own MapReduce

- 2 Python scripts, mapper, reducer
- Word count example





# Command Line MapReduce

## Runtime ~8s (single core)

```
1 #!/usr/bin/env python
2 import sys
3
4 # --- get all lines from stdin ---
5 for line in sys.stdin:
6     # --- remove leading and trailing whitespace---
7     line = line.strip()
8
9     # --- split the line into words ---
10    words = line.split()
11
12    # --- output tuples [word, 1] in tab-delimited format---
13    for word in words:
14        print '%s\t%s' % (word, "1")
15
```

```
1 #!/usr/bin/env python
2 import sys
3
4 # maps words to their counts
5 word2count = {}
6
7 # input comes from STDIN
8 for line in sys.stdin:
9     # remove leading and trailing whitespace
10    line = line.strip()
11
12    # parse the input we got from mapper.py
13    word, count = line.split('\t', 1)
14    # convert count (currently a string) to int
15    try:
16        count = int(count)
17    except ValueError:
18        continue
19
20    try:
21        word2count[word] = word2count[word] + count
22    except:
23        word2count[word] = count
24
25 # write the tuples to stdout
26 # Note: they are unsorted
27 for word in word2count.keys():
28     print '%s\t%s' % (word, word2count[word])
29
```

```
➔ mapreduce_wc time cat miserables.txt | ./mapper.py | sort | ./reducer.py | sort -k 2 -r -n > result.txt
cat miserables.txt  0.00s user 0.00s system 1% cpu 0.358 total
./mapper.py  0.34s user 0.02s system 96% cpu 0.370 total
sort  6.74s user 0.03s system 86% cpu 7.791 total
./reducer.py  0.77s user 0.02s system 10% cpu 7.844 total
sort -k 2 -r -n > result.txt  0.47s user 0.01s system 5% cpu 8.324 total
```

## Spark

- Holds all data in memory
  - creates query plans that are optimised
  - Does Repeat analysis without going to disk -> speed up
  - API to manipulate data frames
  - Scala /Python REPL
  - makes interactive queries possible
-

# Spark Demo

- install brew scala, brew spark
- run `/usr/local/Cellar/apache-spark/2.1.0/libexec/sbin/start-master.sh`
- add to env
  - if `which pyspark > /dev/null`; then
  - `export SPARK_HOME="/usr/local/Cellar/apache-spark/2.1.0/libexec/"`
  - `export PYTHONPATH=$SPARK_HOME/python:$SPARK_HOME/python/build:$PYTHONPATH`
  - `export PYTHONPATH=$SPARK_HOME/python/lib/py4j-0.10.4-src.zip:$PYTHONPATH`
  - fi
- start master
  - `/usr/local/Cellar/apache-spark/2.1.0/libexec/sbin/start-master.sh`
  - check
    - <http://localhost:8080/>
- start slave
  - `/usr/local/Cellar/apache-spark/2.1.0/libexec/sbin/start-slave.sh spark://it033887.fen.bris.ac.uk:7077`
- start ipython notebook
  - `pyspark --master spark://it033887.fen.bris.ac.uk:7077`

# Time down to ~1s – uses all 8 cores of my MBP



## 3 Spark Word Count

```
In [1]: from operator import add
```

```
In [3]: def spark_wc(file_name):
        distFile = sc.textFile(file_name)
        counts = distFile.flatMap(lambda line: line.split(" ")) \
            .map(lambda word: (word, 1)) \
            .reduceByKey(add)

        sorted_wcs = counts.sortBy(lambda word_count_pair: word_count_pair[1], ascending=False)
        sorted_wcs.saveAsTextFile("./spark_wc_result")
```

```
1.0379528559860773
```

```
In [ ]: from timeit import default_timer as timer
```

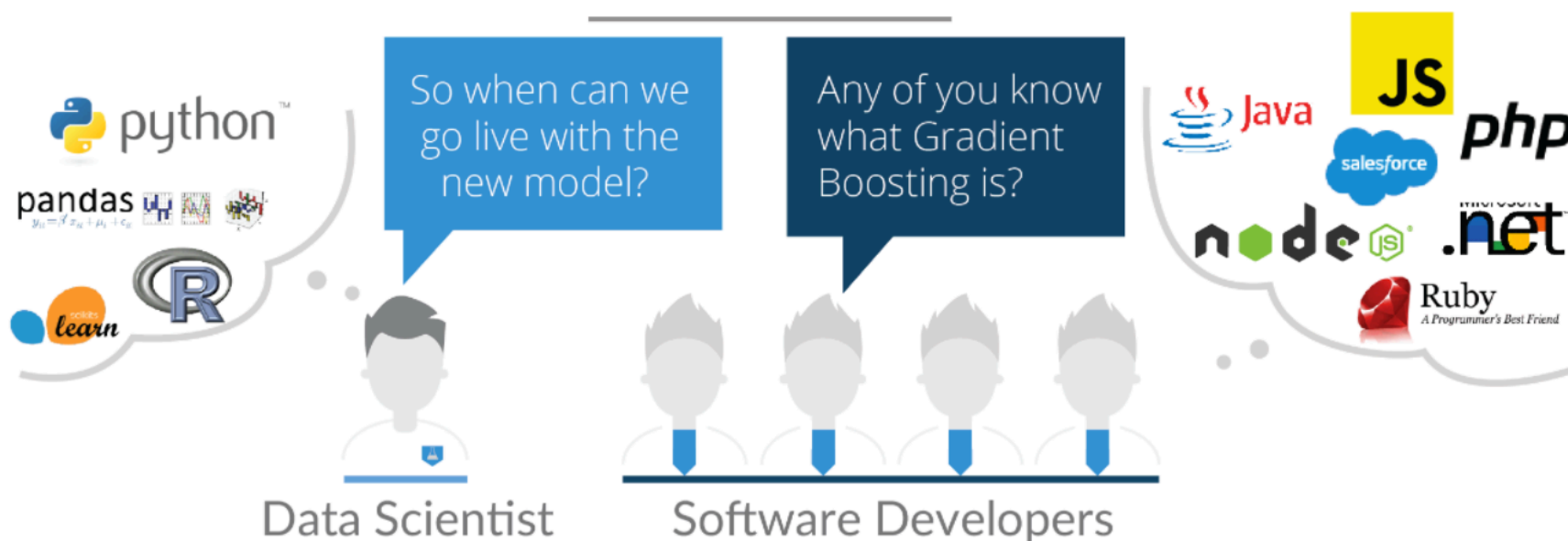
```
start = timer()
spark_wc("mapreduce_wc/miserables.txt")
end = timer()
print(end - start)
```

## Stream Processing

- Spark Streaming with micro batches - convenience of familiar paradigm, sufficient for 98% of use cases
- Storm - for real single event processing

# Horizon Scanning: Predictive Model Markup Language (PMML)

Now What?



## Summary

- In startup environments, data scientists can contribute to **deployment** on several fronts
- knowledge of distributed computing paradigms is important for efficient **implementations**, better models
- for sharing in **trusted** environments, docker packaged notebooks work great

Thank you for your attention