



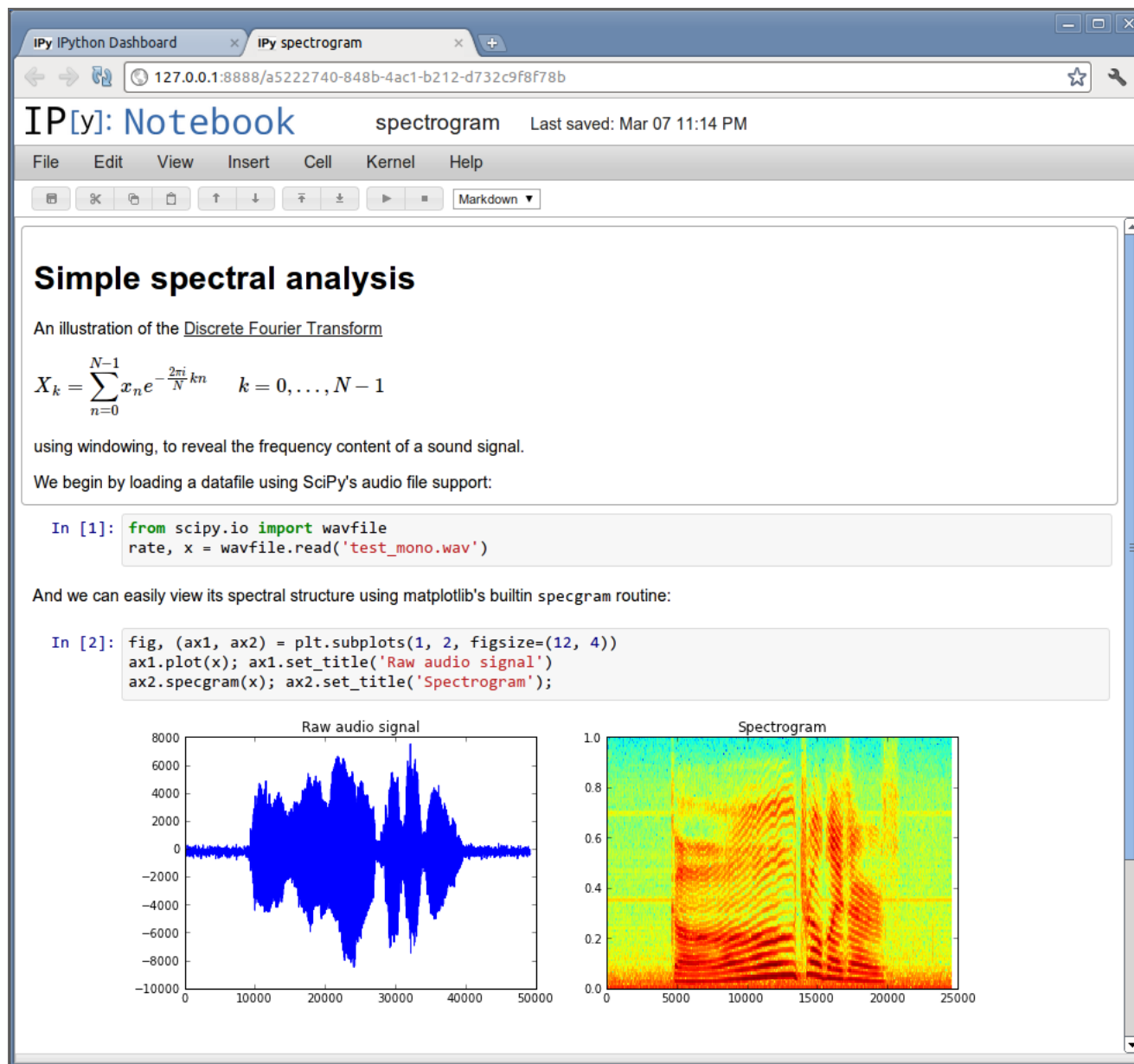
Recap • February 26, 2014





Brian Granger (Cal Poly SLO)

“Compose and share reproducible stories that involve code and data”



Interactive visualization

coming in IPython 2.0

- Vincent/**Vega**/**d3**
- **mpld3**
- **plotly**

Architecture

- js in browser
- asynchronous, bi-directional JSON messaging (WebSockets/ZeroMQ)
- kernel (Python interpreter)

Multilanguage support

coming soon, a drop down menu for what kernel language you'd like to use

- Julia
- R
- Matlab / Octave
- Scala

- **Brian Granger's slides** are an IPython notebook
- Lots of examples at **nbviewer.ipython.org**

Scikit-learn

Olivier Grisel (Parietal, INRIA)

Machine learning in Python

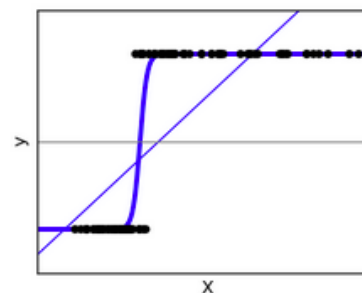
from INRIA

Classification

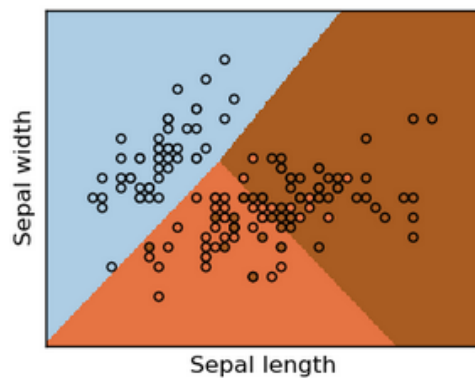
For classification, as in the labeling *iris* task, linear regression is not the right approach as it will give too much weight to data far from the decision frontier. A linear approach is to fit a sigmoid function or **logistic** function:

$$y = \text{sigmoid}(X\beta - \text{offset}) + \epsilon = \frac{1}{1 + \exp(-X\beta + \text{offset})} + \epsilon$$

```
>>> logistic = linear_model.LogisticRegression(C=1e5)
>>> logistic.fit(iris_X_train, iris_y_train)
LogisticRegression(C=100000.0, class_weight=None, dual=False,
                    fit_intercept=True, intercept_scaling=1, penalty='l2',
                    random_state=None, tol=0.0001)
```



This is known as **LogisticRegression**.



dplyr and ggvis

Hadley Wickham (RStudio)

More awesome R packages from Hadley.

- data manipulation
- interactive visualization

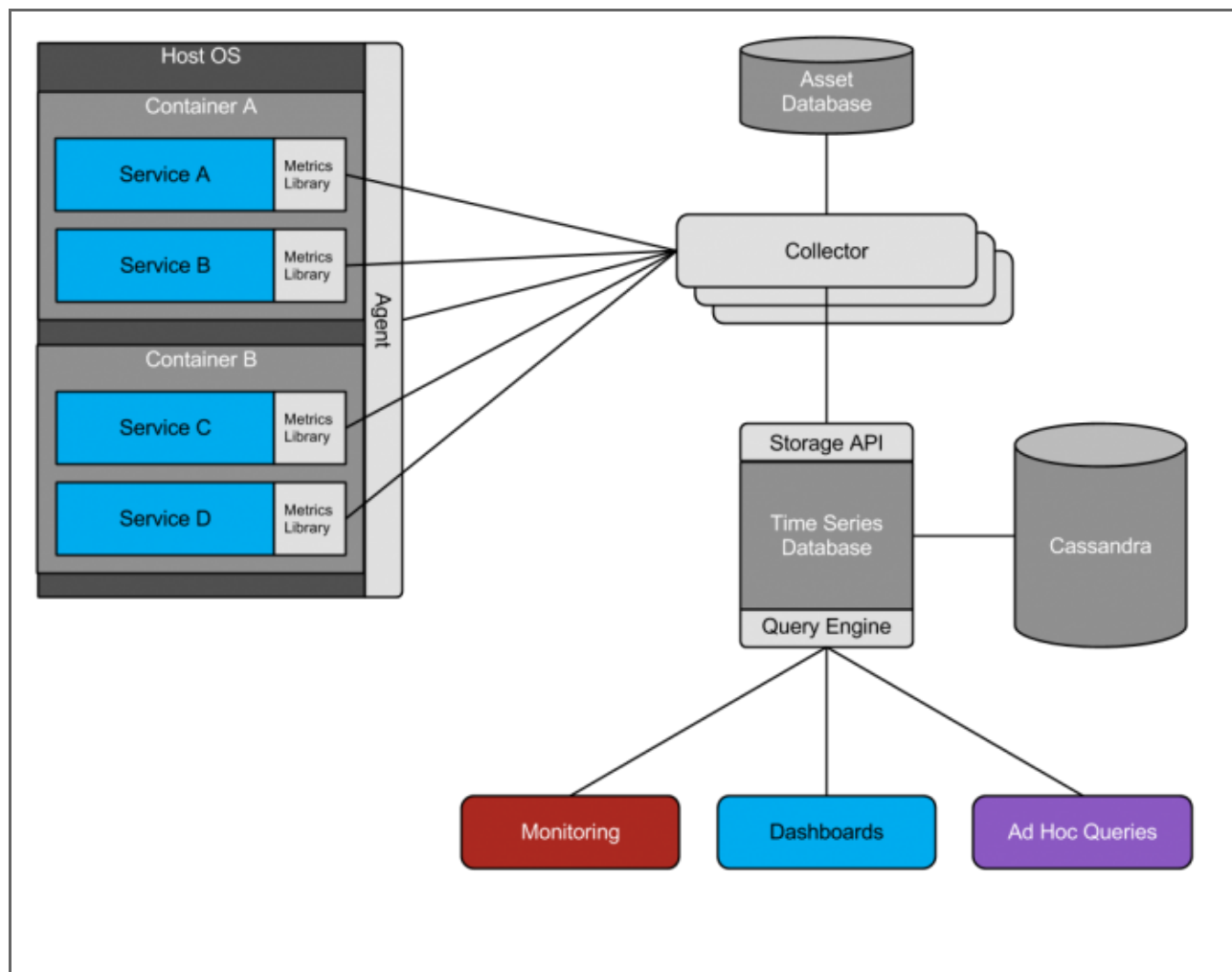
Twitter's Observability stack

Yann Ramin (Twitter)

- “Moving beyond logging” (logging and parsing are expensive)
- Instrument applications to collect runtime statistics

~

Just like you design for testability, design for monitoring



Stack

- instrumentation (**Ostrich, Finagle**)
- storage (Cassandra[†])
- query engine (Cuckoo)
- vizualization
- analytics

[†] cassandra-like key/value store, fronted by an API over caches, queues, storage

Ostrich

Increment a counter

```
Stats.incr("some_important_counter")
```

Ostrich

Store the value of a variable at a particular time

```
Stats.addGauge("current_temperature") {  
  myThermometer.temperature  
}
```

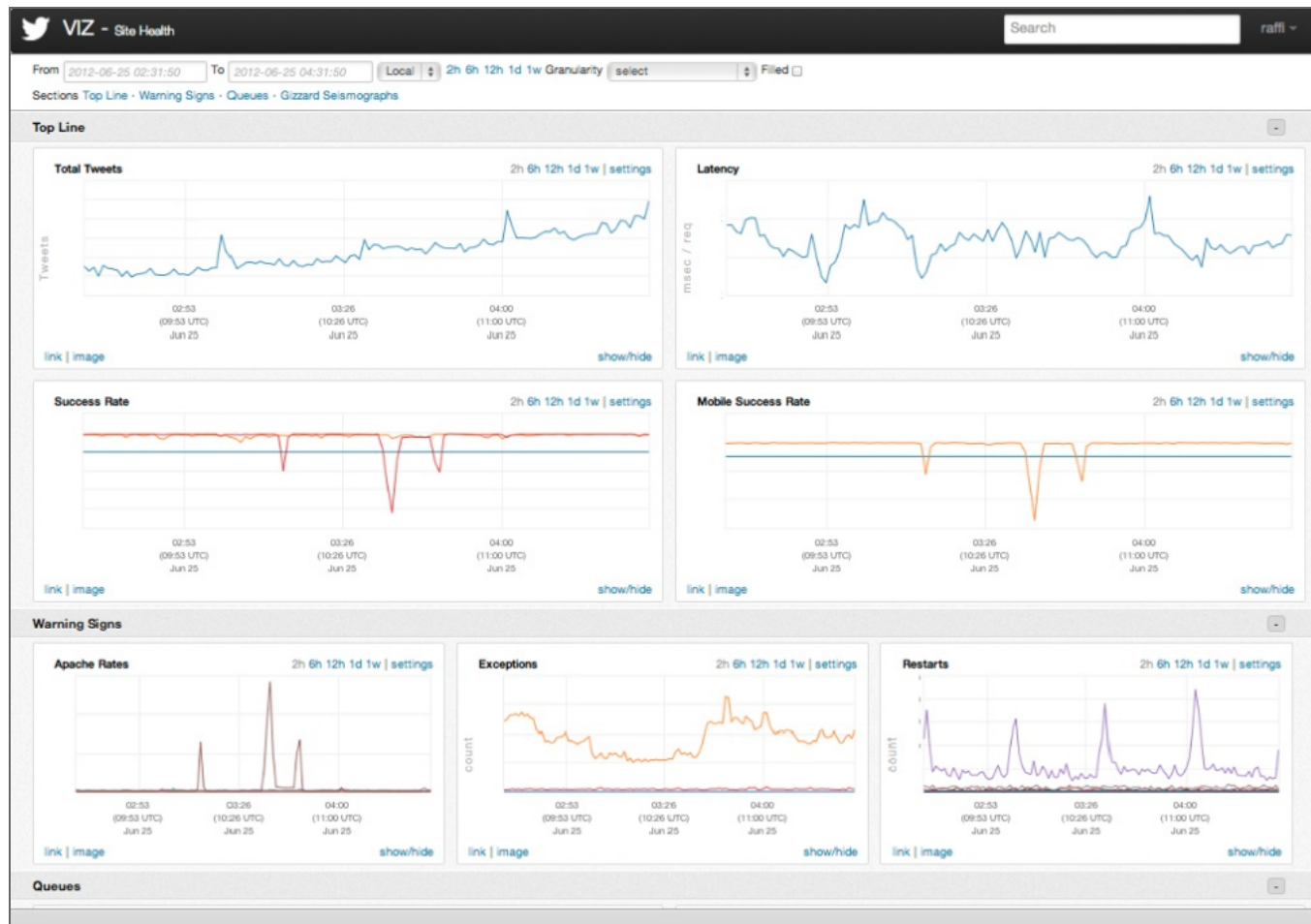

Ostrich

Record running time

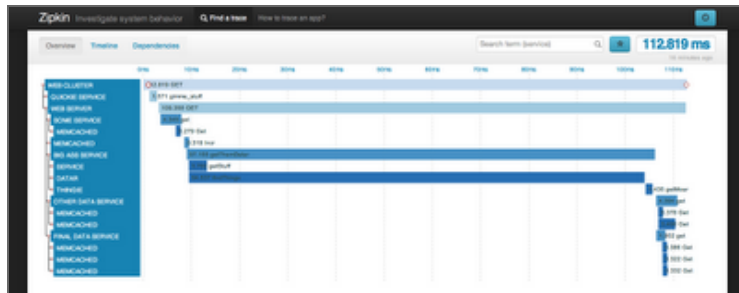
```
Stats.time("translation") {  
  document.translate("de", "en")  
}
```

Keys

(service,source(host,process),metric)



Zipkin tracing for distributed systems (modeled after Google Dapper)



Trained systems

Chris Ré (Stanford)

Examples of trained systems:

- IBM's Watson
- Google's Knowledge Graph
- Netflix recommender engine

- RDBMS + machine learning, stats, optimization
- combining structured and unstructured data

DeepDive

Extracting structured information from unstructured sources such as raw text

DeepDive: Web-scale Knowledge-base Construction using Statistical Learning and Inference

Feature engineering

“Good features allow a simple model to beat a complex model”

Brainwash: A Data System for Feature Engineering

~

† similar to what we observe in challenges

Industry strategy: copy Google

Google	Open
MapReduce	Hadoop
GFS	HDFS
BigTable	HBase
Dremel	Impala
Pregel	Giraph

Giraph

Avery Ching (Facebook)



“iterative graph processing system”

Model of Computation

Two phases

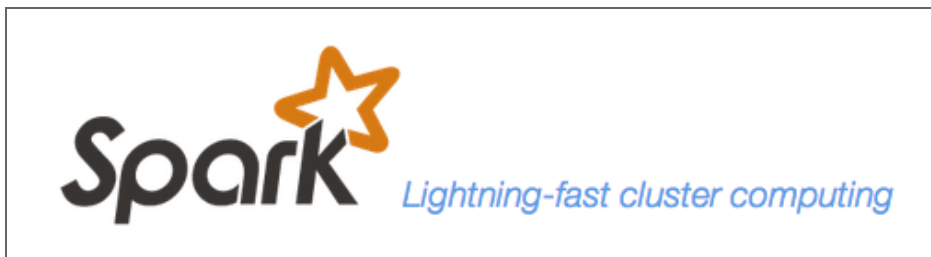
- compute on each vertex
- send messages to other vertices

~

...iterate until halting criteria reached

Examples

- Page rank
- friends of friends



Matei Zaharia (MIT, DataBricks)

Spark is a next-generation Hadoop, in-memory computing

- Written in Scala
- APIs for Scala, Java and Python
- Shark, distributed SQL query engine for Spark (like Hive)
- can read from HDFS, HBase or Cassandra
- Runs on top of Hadoop YARN, Mesos or EC2

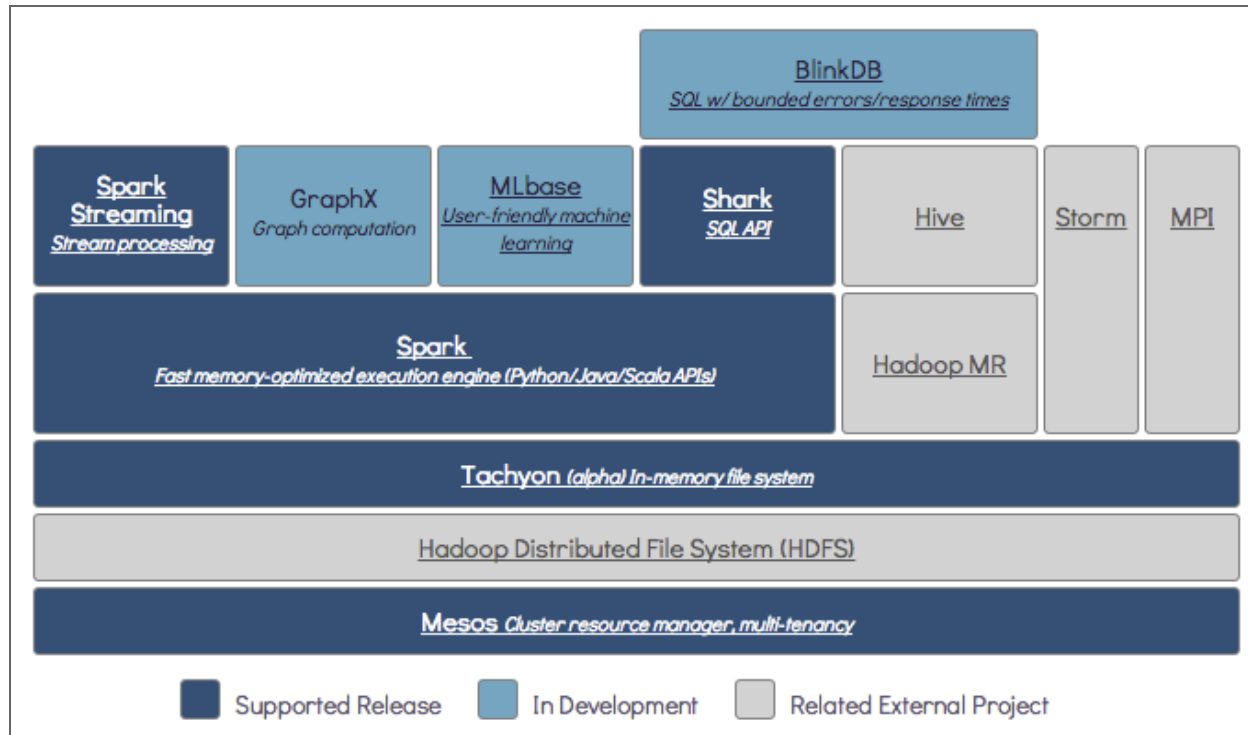
Berkeley Data Analytics Stack

Ameet Talwalkar (UC Berkeley)

BDAS is a suite of software created by the **AMPLab**

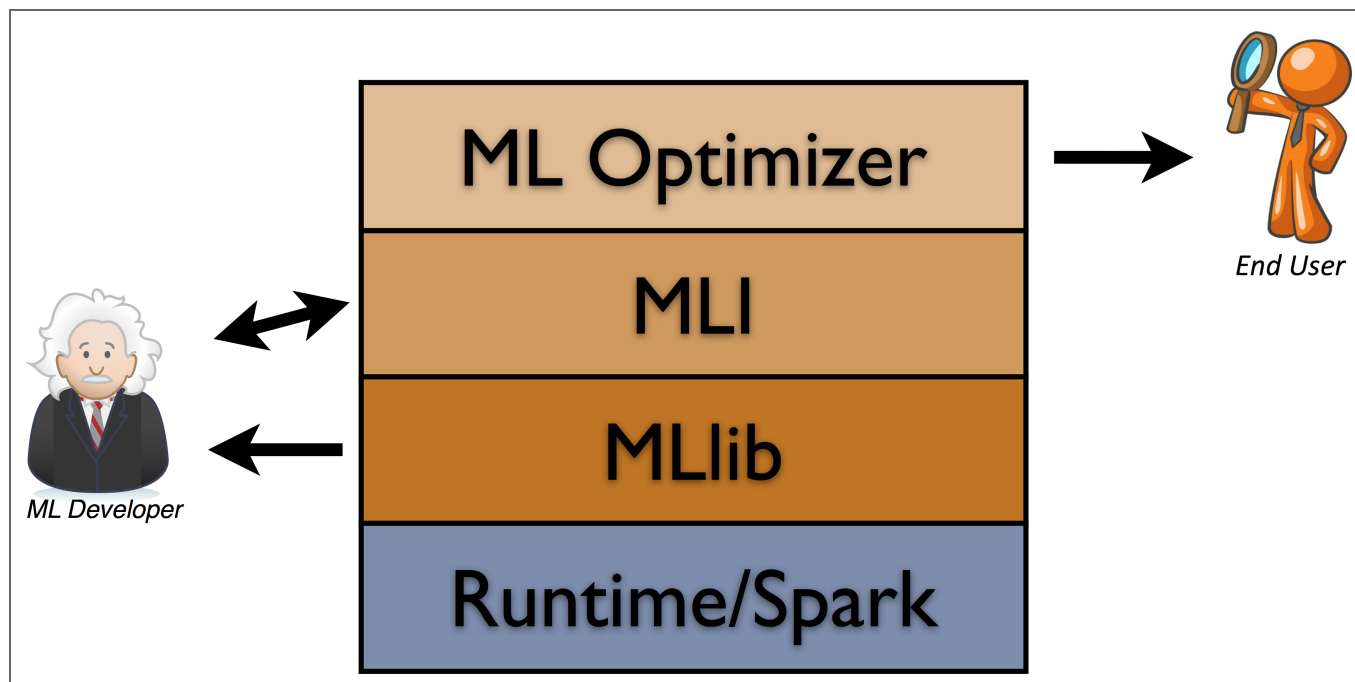
Data-intensive distributed computing stack

- Queries / Machine Learning / Graph processing
- Task scheduler
- Distributed storage
- Resource management



MLBase

MLBase is a machine learning libraries that run over Spark (like Mahout/Hadoop)



Jawbone UP

Monica Rogati (Jawbone)

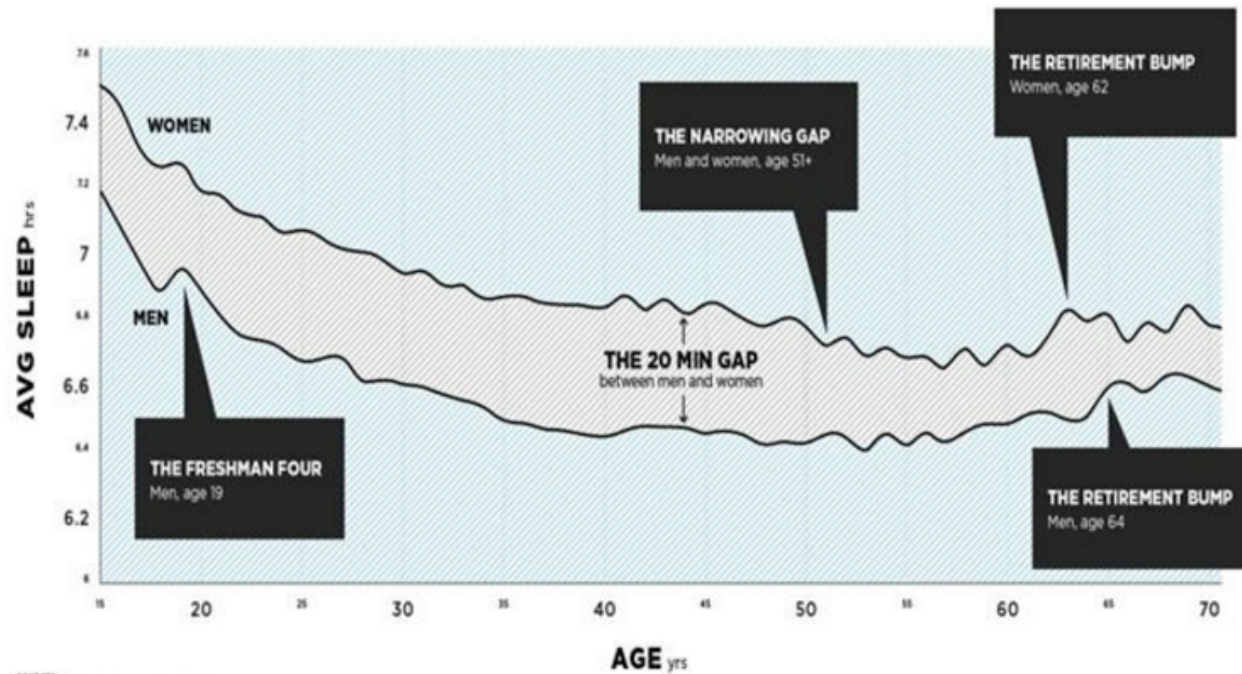


Quantified Self

- 50 million nights of sleep
- 500 billion steps

AGE, GENDER & SLEEP

ON AVERAGE, WOMEN SLEEP 20 MIN MORE THAN MEN



SOURCES:
 Jawahar D. Arora et al. Sleep data from thousands of users
 Center for Retirement Research at Boston College
 The National Institute on Aging

THE END

