

# BigData

## Building Scientific and Consumer Applications



Terrance Boulton



Abhijit Bendale

# Definition of Big Data

No single definition

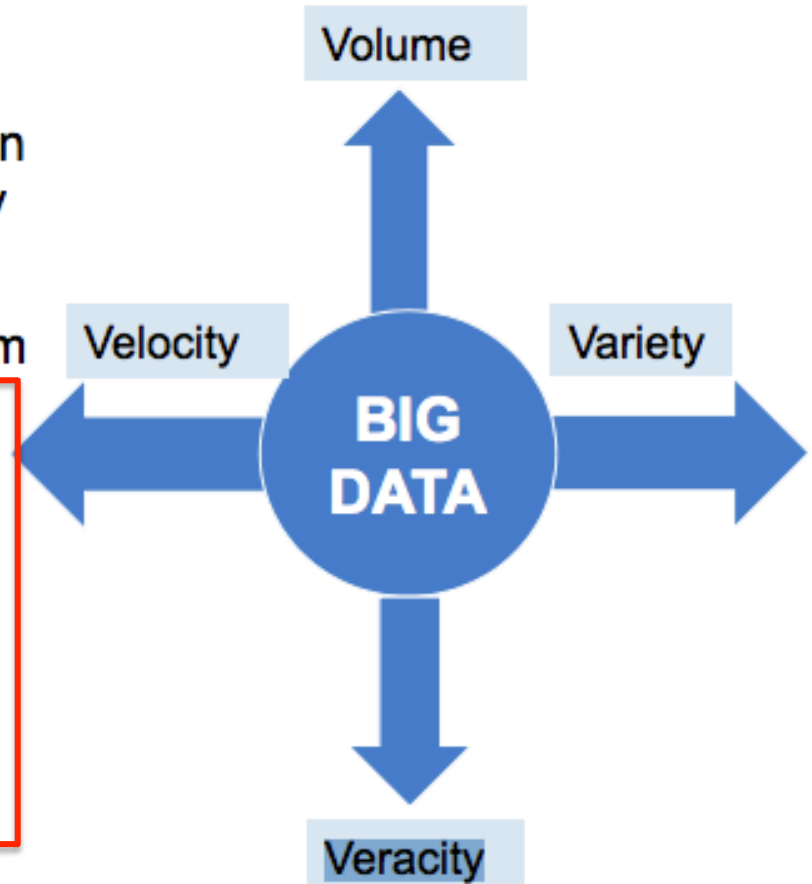
- A collection of large and complex data sets which are difficult to process using common database management tools or traditional data processing applications.
- “Big data refers to the tools, processes and procedures allowing an organization to create, manipulate, and manage very large data sets and storage facilities”

– according to zdnet.com

Big data is not just about size.

- Finds insights from complex, noisy, heterogeneous, longitudinal, and voluminous data.
- It aims to answer questions that were previously unanswered.

The challenges include capturing, storing, searching, sharing & analyzing.



The four dimensions (V's) of Big Data

# Who's Generating Big Data ?

**12+ TBs**  
of tweet data  
every day



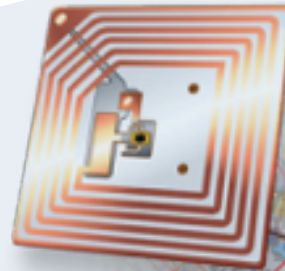
? TBs of  
data every day



**25+ TBs** of  
log data  
every day



**30 billion** RFID  
tags today  
(1.3B in 2005)



**4.6 billion**  
camera  
phones  
world wide



**100s of millions**  
of GPS  
enabled  
devices sold  
annually



**76 million** smart  
meters in 2009...  
200M by 2014

**2+ billion**  
people on  
the Web  
by end  
2011



# Who's Generating Big Data



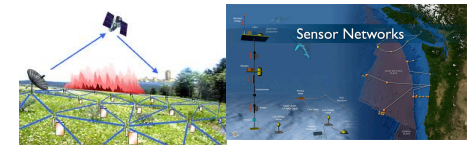
**Social media and networks**  
(all of us are generating data)



**Scientific instruments**  
(collecting all sorts of data)



**Mobile devices**  
(tracking all objects all the time)



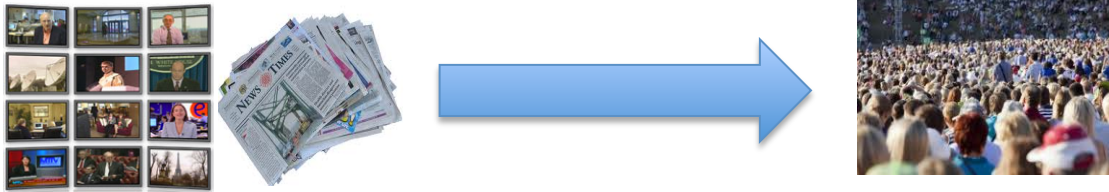
**Sensor technology and networks**  
(measuring all kinds of data)

- The progress and innovation is no longer hindered by the ability to collect data
- But, by the ability to manage, analyze, summarize, visualize, and discover knowledge from the collected data in a timely manner and in a scalable fashion

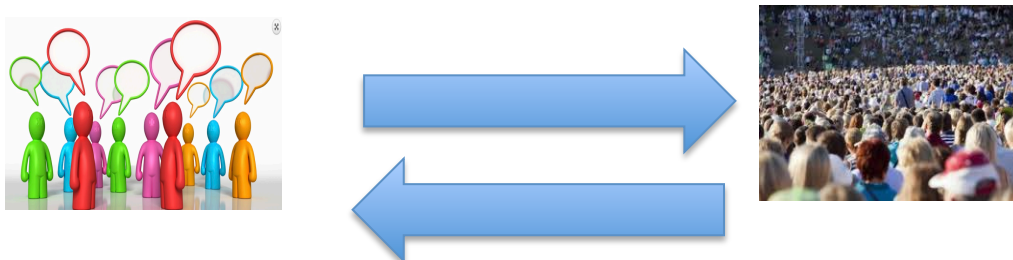
# The Model Has Changed...

- **The Model of Generating/Consuming Data has Changed**

**Old Model:** Few companies are generating data, all others are consuming data

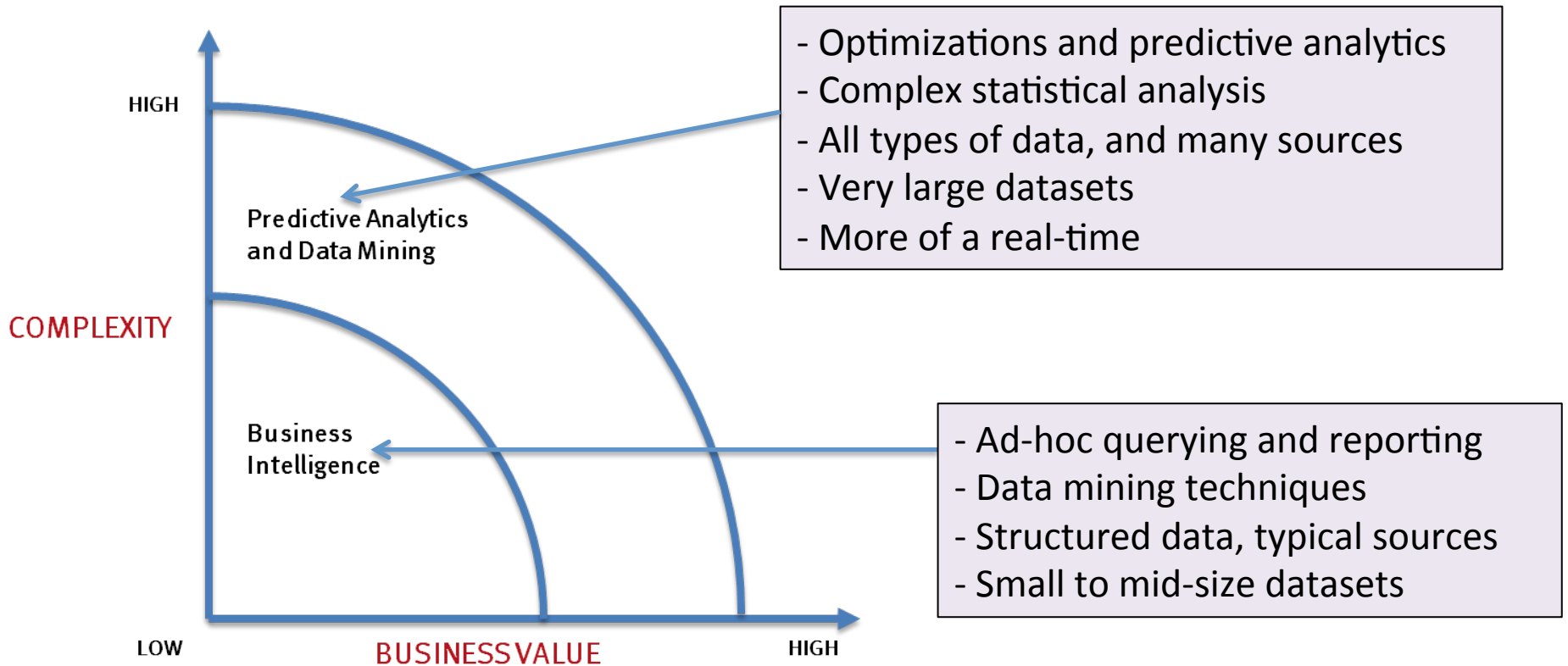


**New Model:** all of us are generating data, and all of us are consuming data



Two way street...!

# What's driving Big Data



# What is driving Bigdata ?

Money to be made..!!





June 2011

# Big data: The next frontier for innovation, competition, and productivity





# Why Big-Data?

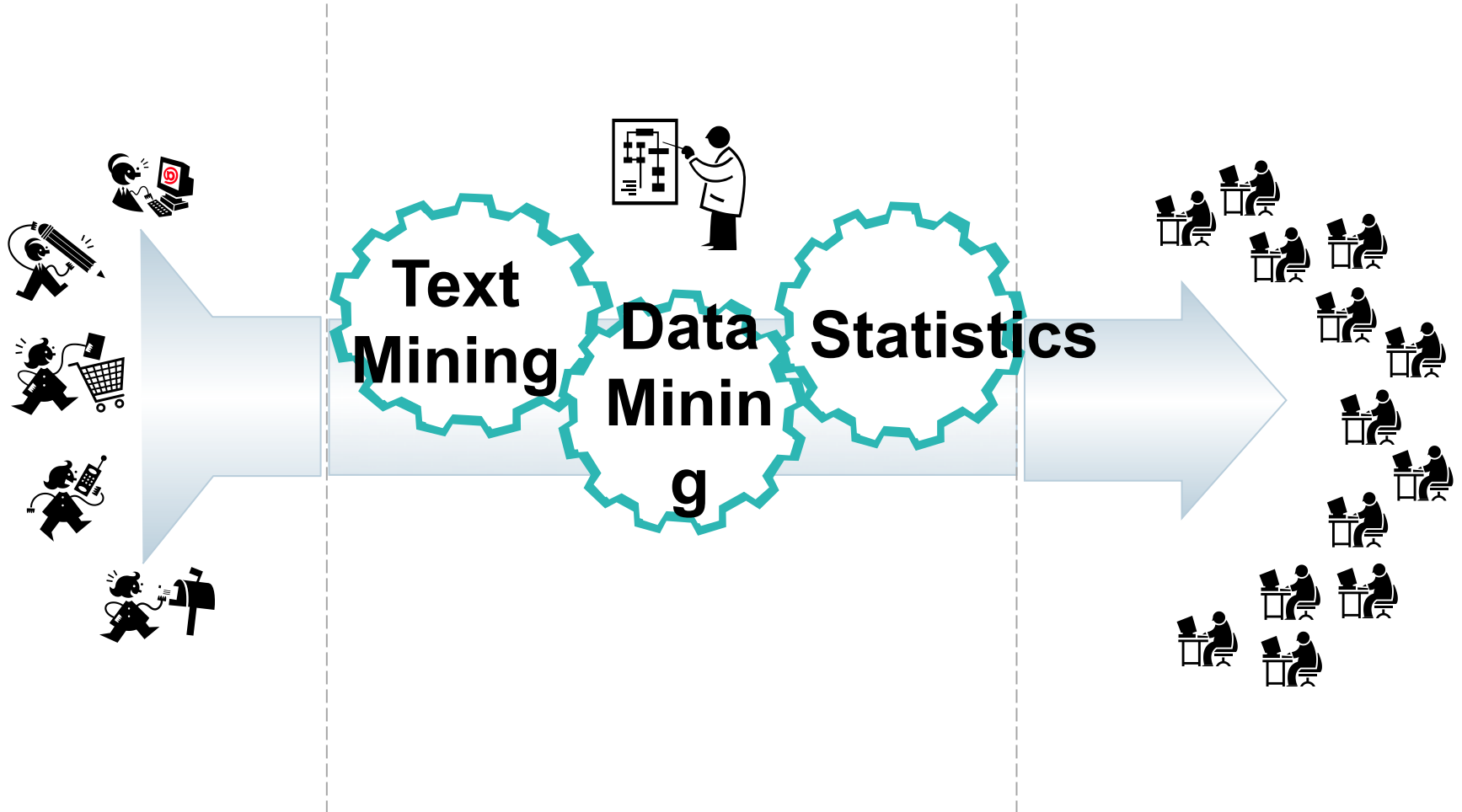
- Key enablers for the appearance and growth of 'Big-Data' are:
  - + Increase in storage capabilities
  - + Increase in processing power
  - + Availability of data

# Driving Smarter Business Outcomes

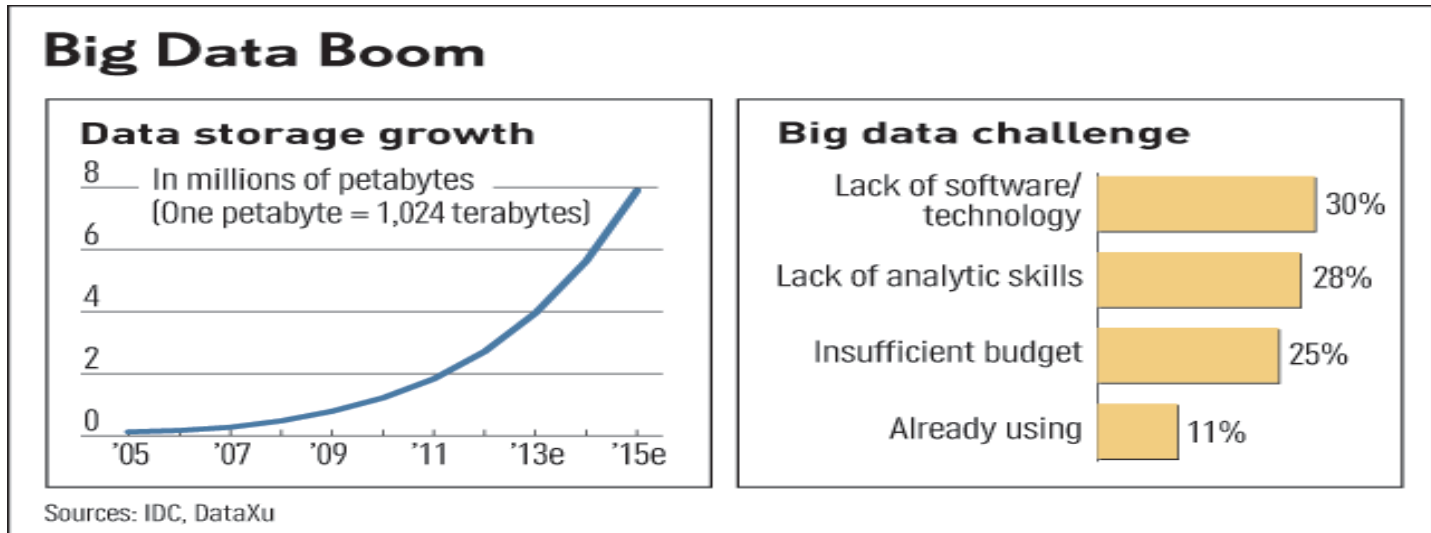
**Capture**

**Predict**

**Act**



# Challenges in Handling Big Data



- **The Bottleneck is in technology**
  - New architecture, algorithms, techniques are needed
- **Also in technical skills**
  - Experts in using the new technology and dealing with big data

Lets join the herd..!

**Big data is like teenage sex:**  
everyone talks about it,  
nobody really knows how to do it,  
everyone thinks everyone else is  
doing it, so everyone claims they  
are doing it...

(Dan Ariely)

# Big Data Landscape

## Vertical Apps



## Ad/Media Apps



## Business Intelligence



## Analytics and Visualization



## Log Data Apps



## Data As A Service



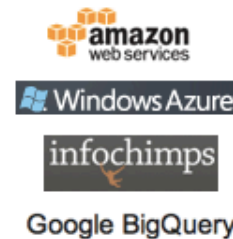
## Analytics Infrastructure



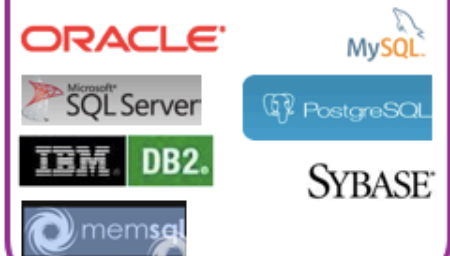
## Operational Infrastructure



## Infrastructure As A Service



## Structured Databases



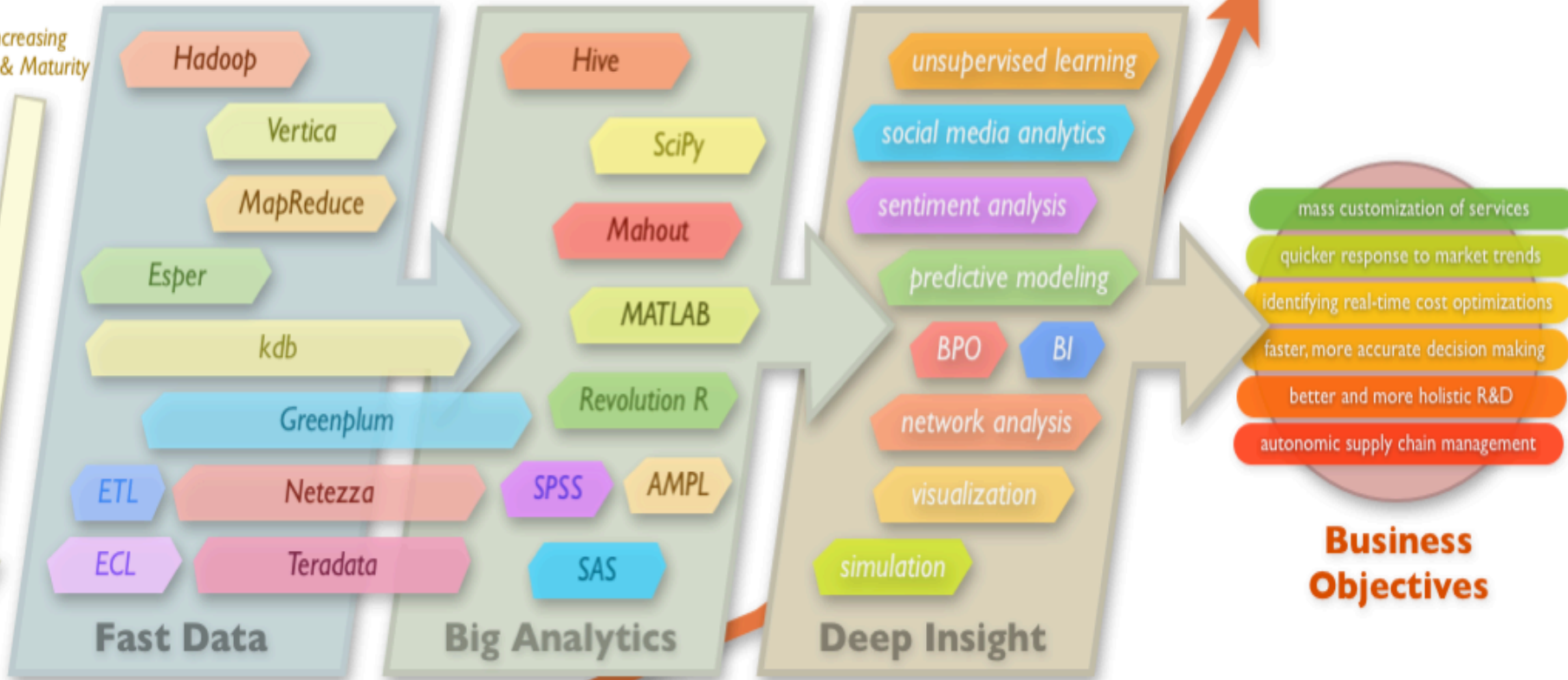
## Technologies





# Big Data: The Moving Parts

Increasing Age & Maturity



From <http://blogs.zdnet.com/Hinchcliffe>

the growth of data will be exponential for the foreseeable future



the amount of data stored by the average company today

# What will you learn ...

- **Problem Analysis and identification**
  - Machine Learning, Statistics analysis techniques to handle big data
  - Theoretical and engineering constraints
  - Parallelizing applications for high throughput
- **Tools and Techniques**
  - Introduction to Machine learning tools like scikits-learn, weka, R etc
  - Parallel Programming using CUDA
  - Introduction to cloud based tools like MapReduce, Hadoop, Pig, Hive etc
- **Case studies from various application domains**
  - Netflix case study
  - Class projects



# Course Structure

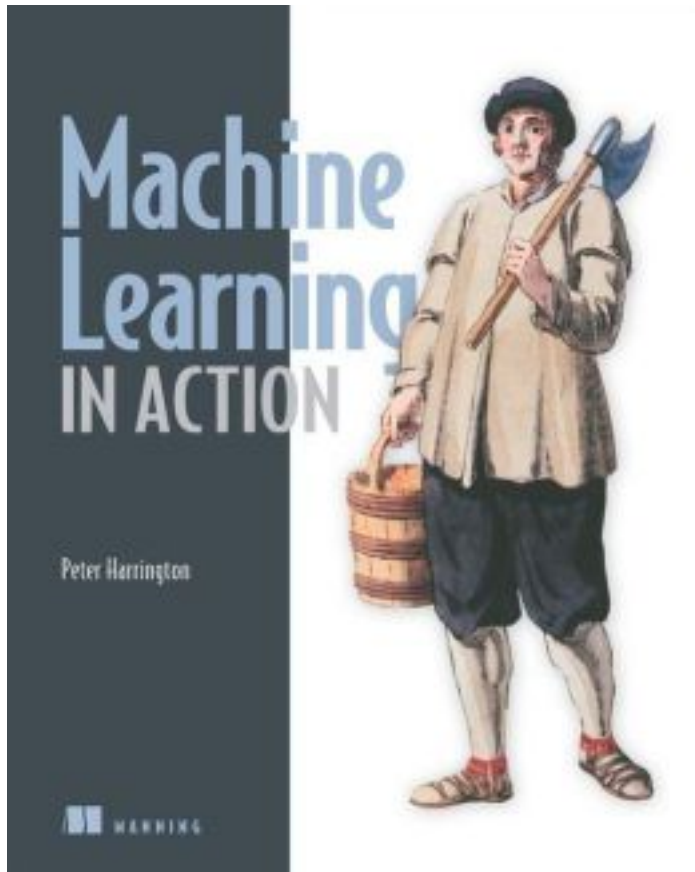


- Pre-requisites: C/C++, Python/Java/R/Matlab, OS
- Lectures (Tue Osborne, 4:45 – 7:30 pm)
- Assignments (2 undergrads, 3 graduate)
- Semester Project
- Grading
  - Assignments (30 %)
  - Project Proposal + Final Project (10 + 40 %)
  - Class participation (20 %)



# Textbooks

amazon.com



“Machine Learning in Action”  
Peter Harrington, Manning Publications  
2012



Programming Massively Parallel Processors:  
A Hands-on Approach  
David Kirk, Wen-mei Hwu  
Morgan Kaufmann 2<sup>nd</sup> Edn 2012

We have  
some copies

# Resources

- We have some GPUs
- Will be installed on machines on campus/  
students will be able to access (more details in  
subsequent lectures)
- Programming
  - C/C++ for CUDA assignments
  - For class projects: C/C++, Java, Python, R etc.

# Application domains and Case Studies

Scientific

Social Media &  
Entertainment

Politics

Finance

Healthcare

Consumer  
Applications

More

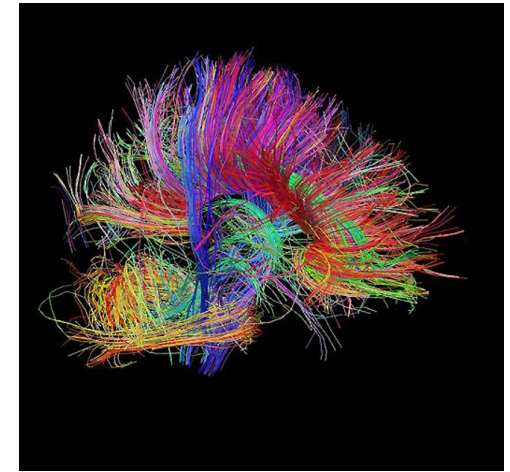
# Scientific Instruments



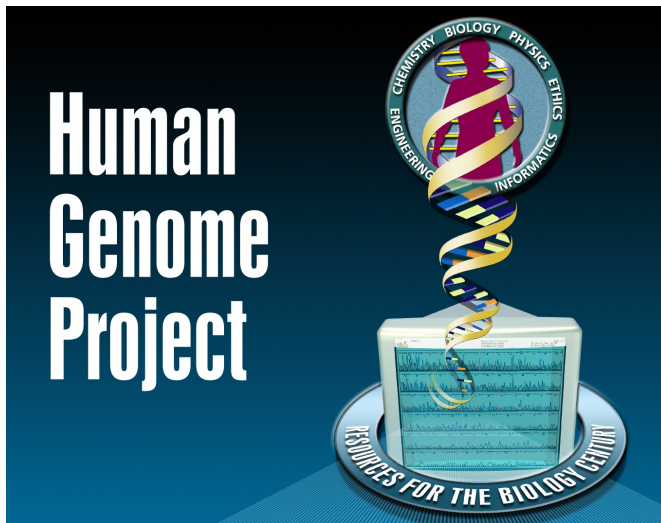
**Murchinson Widefield Array**  
Data produced at the rate of  
16 GB/s



**Scientific instruments**  
(collecting all sorts of data)



**Human Connectome Project**  
12 TB total data, 100M new  
Photographs per day



514 Petabytes, Annotations  
17 TB



**Large Hadron Collider:**  
30 PB annually, stored on 83K physical disks



Processing power of about 110 Petabytes

# Social Media and Images/Video



facebook®

140 billion images  
6 billion added monthly



You Tube

72 hours uploaded  
every minute



flickr  
6 billion images



the simple image sharer  
imgur  
1 billion images  
served daily



3.5 trillion  
photographs

**90%** of net traffic will be visual!



# Retailers..

## Consumer Products Companies

**P&G**



**KRAFT**



## Big Box Stores



**COSTCO**  
WHOLESALE

**STAPLES**

**WAL★MART**  
ALWAYS LOW PRICES.

*Always.*

# Why are they collecting all this data?

## Target Marketing

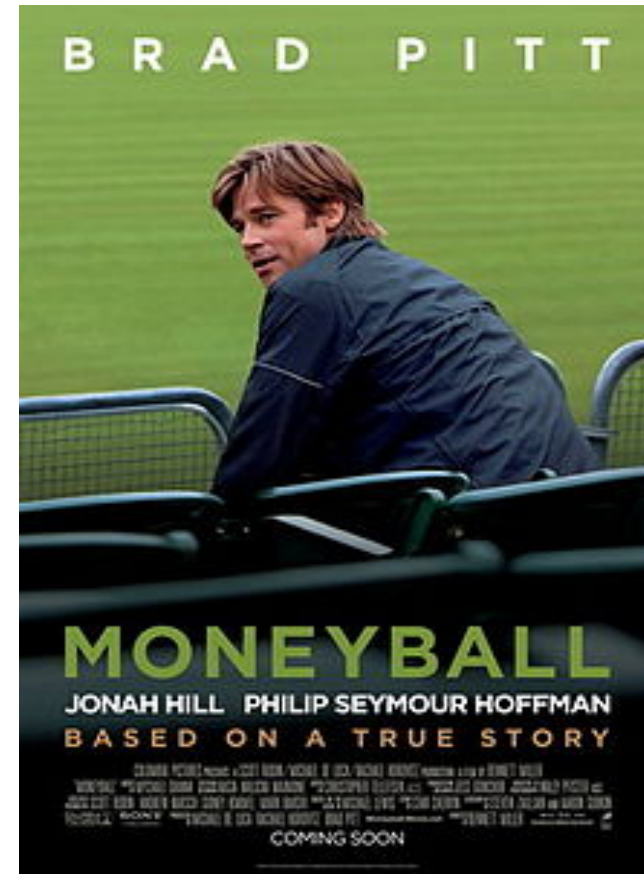
- To send you catalogs for exactly the merchandise you typically purchase.
- To suggest medications that precisely match your medical history.
- To “push” television channels to your set instead of your “pulling” them in.
- To send advertisements on those channels **just for you!**

## Targeted Information

- To know what you need before you even know you need it based on past purchasing habits!
- To notify you of your expiring driver’s license or credit cards or last refill on a Rx, etc.
- To give you turn-by-turn directions to a shelter in case of emergency.

# Big Data and Sports

- **Moneyball: The Art of Winning an Unfair Game**  
Oakland Athletics baseball team and its general manager Billy Beane
- Oakland A's' front office took advantage of more analytical gauges of player performance to field a team that could compete successfully against richer competitors in MLB
- Oakland approximately \$41 million in salary, New York Yankees, \$125 million in payroll that same season. Oakland is forced to find players undervalued by the market,





# Finance

## Credit Card Companies

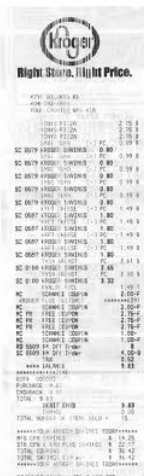


## What data are they getting?

### Airline ticket



### Grocery Bill



### Restaurant check

1	Cab Sauv Corvina	\$15.25
1	Cab Sauv Glass	\$4.00
<b>26 Net Sales Total</b>		<b>\$73.45</b>
1	SUCK MY D... FACE	\$0.00
1	FISH CAKES	\$4.95
2	1/2 Wings Starter	\$7.90
1	MELON AND FARMA HAM	\$3.95
1	Calamari	\$4.95
1	Garlic bro starter	\$2.25
3	Meatball Starter	\$17.85
1	Aub & Feta Starter	\$3.95
1	Cas Barc	\$12.40

### Hotel Bill

Sheraton Gateway Hotel Atlanta Airport  
2999 Sullivan Rd  
Atlanta, GA 30327  
Tel: 770-271-1100 Fax: 770-992-8306

Room No.	Room Description	Room Rate	Tax	Subtotal
03-08C-01	RT313	Room Charge	9.94	9.94
03-08C-01	RT313	Room Tax Occupancy	9.94	19.88
03-08C-02	RT313	Room Charge	9.94	9.94
04-08C-01	RT313	Room Charge	9.94	9.94
04-08C-02	RT313	Room Charge	9.94	9.94
05-08C-01	RT313	Room Charge	9.94	9.94
05-08C-02	RT313	Room Charge	9.94	9.94
06-08C-01	RT313	Room Charge	9.94	9.94
06-08C-02	RT313	Room Charge	9.94	9.94
07-08C-01	RT313	Room Charge	9.94	9.94
07-08C-02	RT313	Room Charge	9.94	9.94
08-08C-01	RT313	Room Charge	9.94	9.94
08-08C-02	RT313	Room Charge	9.94	9.94
09-08C-01	RT313	Room Charge	9.94	9.94
09-08C-02	RT313	Room Charge	9.94	9.94
10-08C-01	RT313	Room Charge	9.94	9.94
10-08C-02	RT313	Room Charge	9.94	9.94
11-08C-01	RT313	Room Charge	9.94	9.94
11-08C-02	RT313	Room Charge	9.94	9.94
12-08C-01	RT313	Room Charge	9.94	9.94
12-08C-02	RT313	Room Charge	9.94	9.94

Total Due: \$119.88

# Usage Example in Big Data

## Data Analysis prediction for US 2012 Election

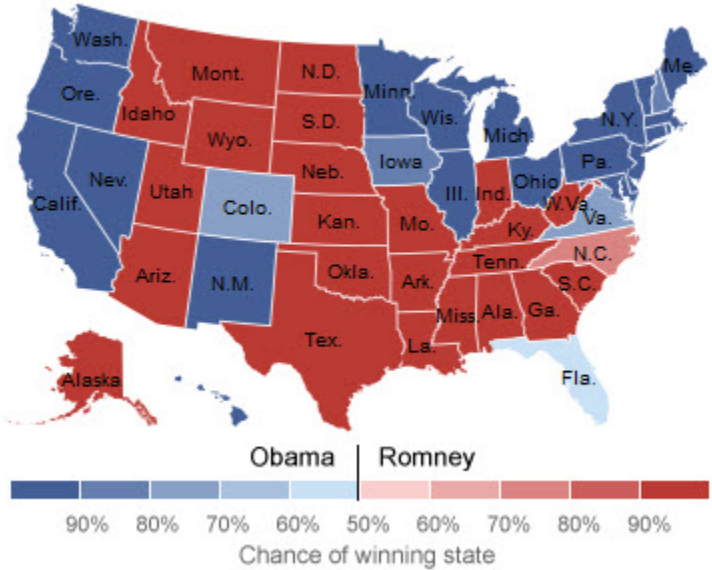
Drew Linzer, June 2012  
332 for Obama,  
206 for Romney

media continue reporting the race as very tight

Nate Silver's, Five thirty Eight blog  
Predict Obama had a 86% chance of winning  
Predicted all 50 state correctly

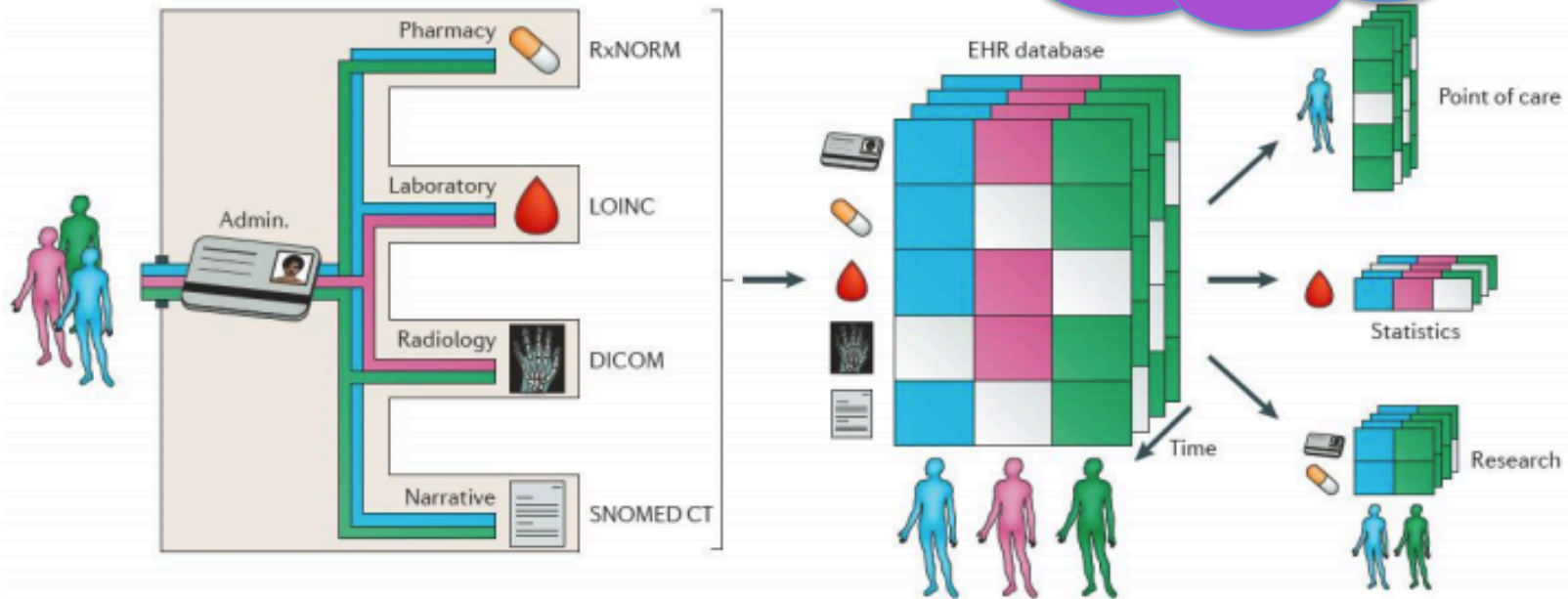
Sam Wang, the Princeton Election Consortium  
The probability of Obama's re-election  
at more than 98%

State-by-State Probabilities



# Data Collection and Analysis

## Healthcare Analytics



Effectively integrating and efficiently analyzing various forms of healthcare data over a period of time can answer many of the impending healthcare problems.

Jensen, Peter B., Lars J. Jensen, and Søren Brunak. "Mining electronic health records: towards better research applications and clinical care." *Nature Reviews Genetics* (2012).

# How Can You Avoid *Big Data*?



Always pay in cash



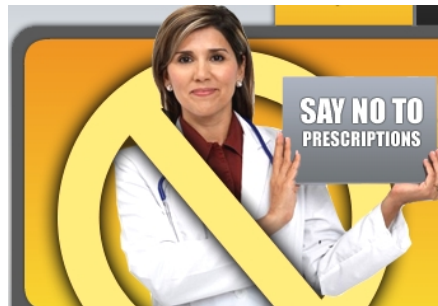
You make no sense



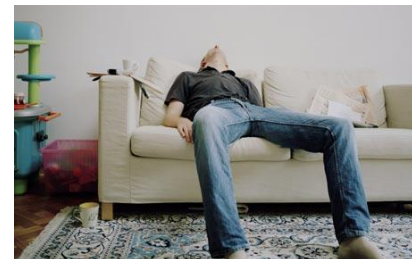
Never go online



Never use club cards



**CELL PHONES!**



Stay home all day

## High Throughput Computing

- Parallel Programming
- CUDA



## Using the Cloud

- Hadoop, MapReduce
- Pig, Hive, Hbase



## Analysis Tools

- Machine learning tools
- Data Mining Tools



Handling  
Big Data



Python



Java



Python



**mipack**  
a scalable c++ machine learning library



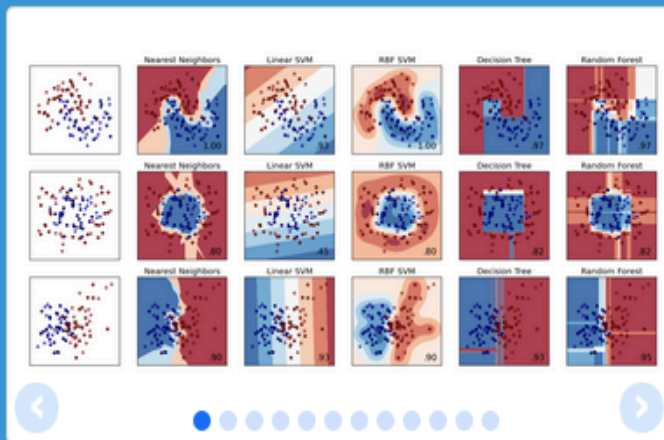


- Library of Machine Learning models
- Simple **fit** / **predict** / **transform** API
- Python / NumPy / SciPy / Cython  
& wrappers for libsvm / liblinear
- Model Assessment, Selection & Ensembles
- Some support for multi-core

# scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license



## Classification

Identifying to which set of categories a new observation belong to.

**Applications:** Spam detection, Image recognition.

**Algorithms:** *SVM, nearest neighbors, random forest, ...* — Examples

## Regression

Predicting a continuous value for a new example.

**Applications:** Drug response, Stock prices.

**Algorithms:** *SVR, ridge regression, Lasso, ...* — Examples

## Clustering

Automatic grouping of similar objects into sets.

**Applications:** Customer segmentation, Grouping experiment outcomes

**Algorithms:** *k-Means, spectral clustering, mean-shift, ...* — Examples

## Dimensionality reduction

Reducing the number of random variables to consider.

**Applications:** Visualization, Increased efficiency

**Algorithms:** *PCA, Isomap, non-negative matrix factorization.* — Examples

## Model selection

Comparing, validating and choosing parameters and models.

**Goal:** Improved accuracy via parameter tuning

**Modules:** *grid search, cross validation, metrics.* — Examples

## Preprocessing

Feature extraction and normalization.

**Application:** Transforming input data such as text for use with machine learning algorithms.

**Modules:** *preprocessing, feature extraction.* — Examples



# Use off-the-shelf tools: Consistent API

```
from sklearn.svm import SVC

clf = SVC()
clf.fit(X_train, y_train)

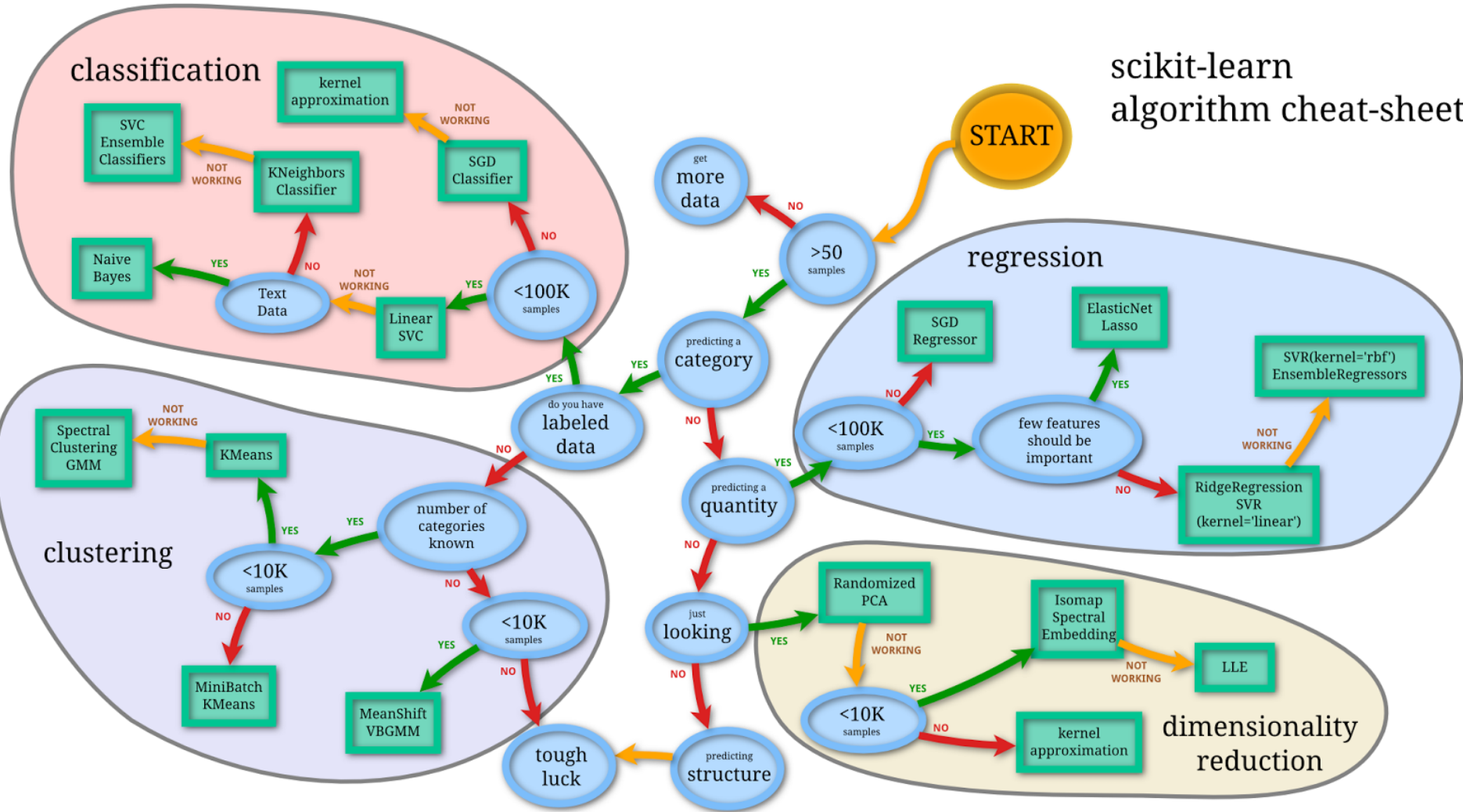
y_pred = clf.predict(X_test)
```

```
from sklearn.ensemble import RandomForestClassifier

clf = RandomForestClassifier()
clf.fit(X_train, y_train)

y_pred = clf.predict(X_test)
```

# scikit-learn algorithm cheat-sheet



# Face Recognition in 5 mins..!

Total dataset size:

n\_samples: 1288, n\_features: 1850, n\_classes: 7

Extracting the top 150 eigenfaces from 966 faces  
done in 0.466s

Projecting the input data on the eigenfaces orthonormal basis  
done in 0.056s

Fitting the SVM classifier to the training set  
done in 18.549s

Predicting people's names on the test set  
done in 0.062s

	precision	recall	f1-score	support
Ariel Sharon	0.90	0.75	0.82	12
Colin Powell	0.78	0.94	0.85	62
Donald Rumsfeld	0.86	0.72	0.78	25
George W Bush	0.89	0.96	0.92	141
Gerhard Schroeder	0.92	0.74	0.82	31
Hugo Chavez	0.90	0.53	0.67	17
Tony Blair	0.81	0.74	0.77	34
avg / total	0.86	0.86	0.86	322

# Learned Eigen Faces

eigenface 0



eigenface 1



eigenface 2



eigenface 3



eigenface 4



eigenface 5



eigenface 6



eigenface 7



eigenface 8



eigenface 9



eigenface 10



eigenface 11



# Predicting...

predicted: Powell  
true: Powell



predicted: Rumsfeld  
true: Rumsfeld



predicted: Bush  
true: Bush



predicted: Chavez  
true: Chavez



predicted: Bush  
true: Bush



predicted: Bush  
true: Bush



predicted: Bush  
true: Bush



predicted: Schroeder  
true: Schroeder



predicted: Powell  
true: Powell



predicted: Sharon  
true: Sharon



predicted: Blair  
true: Schroeder



predicted: Rumsfeld  
true: Rumsfeld



# High Throughput Computing



  
Share

LEARNING BY DOING

# Moore's Law

Motivation



Share

- The most economic number of components in an IC will double every year
- Historically – CPUs get faster
  - Hardware reaching frequency limitations
- Now – CPUs get wider

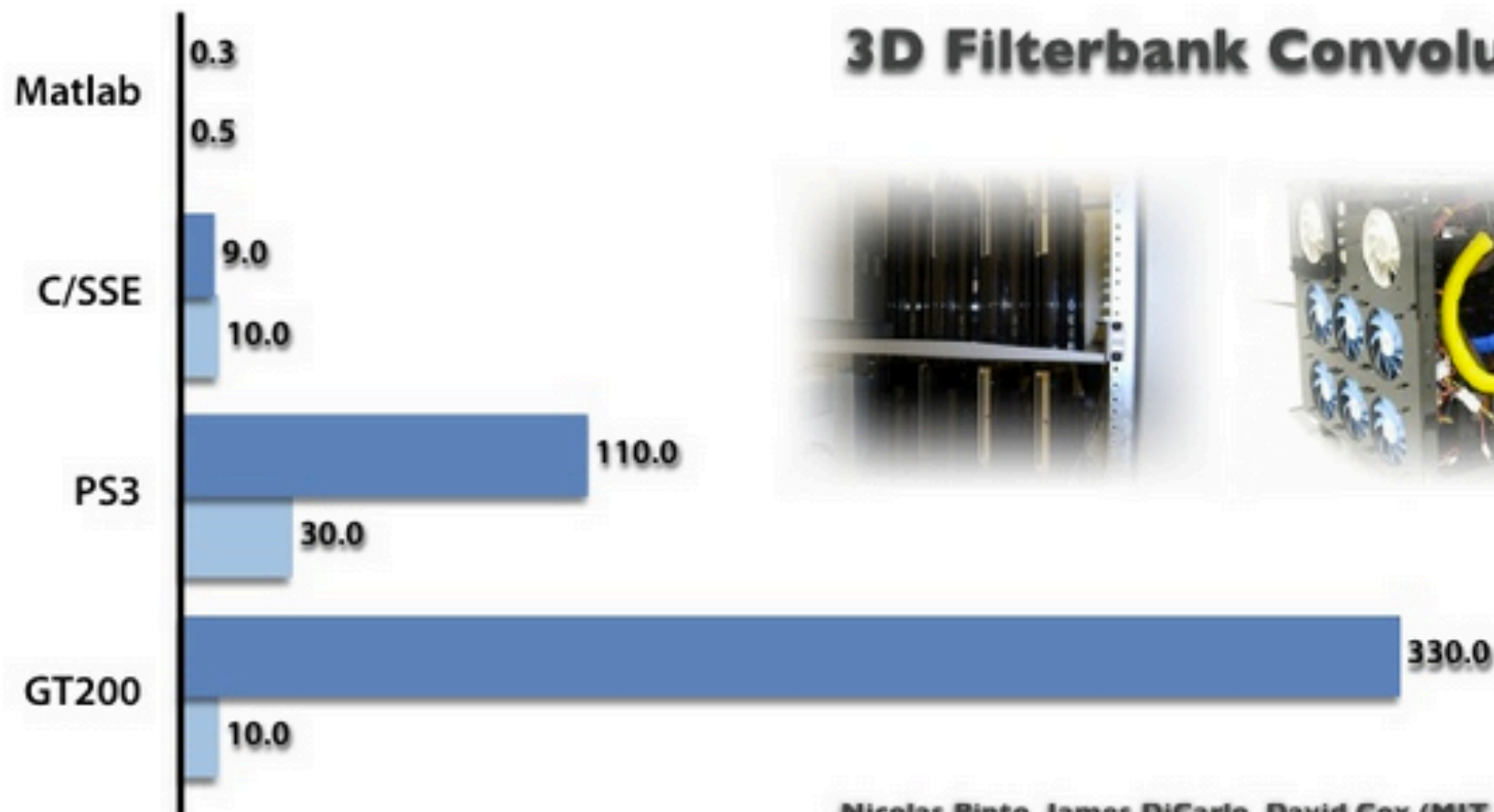
# GPUs are REALLY fast

# GPU?

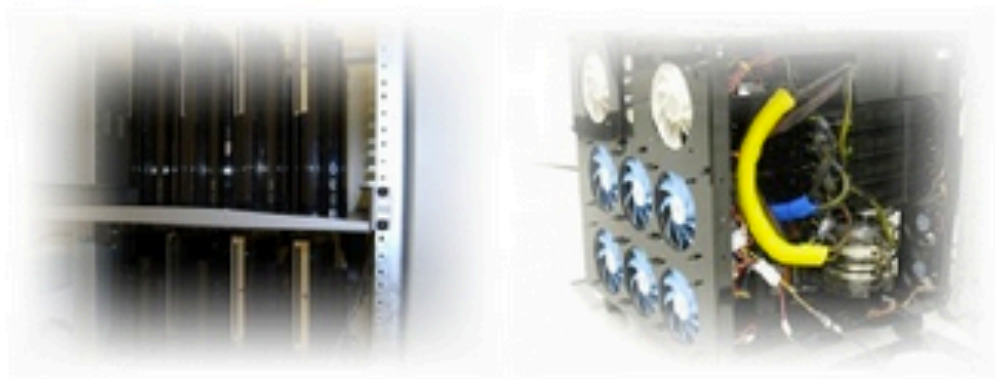


Share

■ Performance (gflops)    ■ Development Time (hours)



## 3D Filterbank Convolution



Nicolas Pinto, James DiCarlo, David Cox (MIT, Harvard)



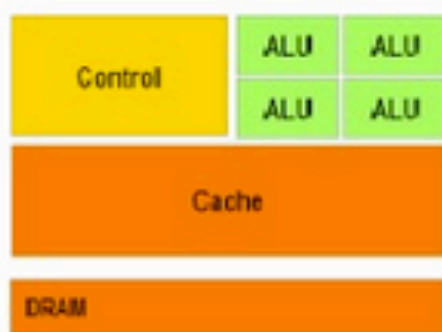
# Why so fast?

# GPU?

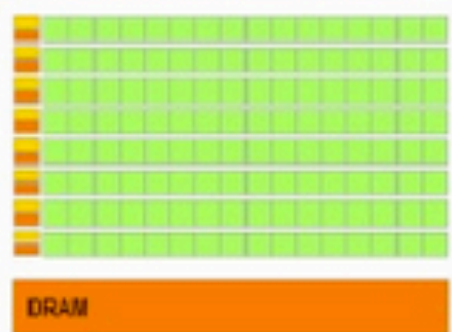


Share

- Designed for math-intensive, parallel problems:



CPU



GPU

- More transistors dedicated to ALU than flow control and data cache

# Is it free?

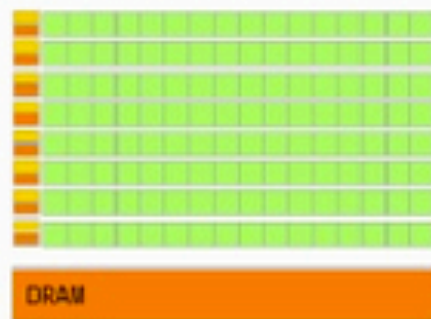
# GPU?



- What are the consequences?
- Program must be more predictable:
  - Data access coherency
  - Program flow



CPU



GPU

# Parallel Computing

**GPU?**



Share

- Rather than expecting CPUs to get twice as fast, expect to have twice as many!
  - Parallel processing for the masses
  - Unfortunately: Parallel programming is hard.
- Algorithms and Data Structures must be fundamentally redesigned

# CPU vs. GPU

---

**GPU?**



Share

- CPU
  - Really fast caches (great for data reuse)
  - Fine branching granularity
  - Lots of different processes/threads
  - High performance on a single thread of execution
- GPU
  - Lots of math units
  - Fast access to onboard memory
  - Run a program on each fragment/vertex
  - High throughput on parallel tasks
- CPUs are great for *task* parallelism
- GPUs are great for *data* parallelism

# CUDA Software Development

CUDA Optimized Libraries:  
math.h, FFT, BLAS, ...

Integrated CPU + GPU  
C Source Code

NVIDIA C Compiler

NVIDIA Assembly  
for Computing (PTX)

CPU Host Code

CUDA  
Driver

Profiler

Standard C Compiler

GPU

CPU

A signpost with two blue signs. The top sign is tilted and contains the text 'BIG DATA' in white, uppercase letters. The bottom sign is also tilted and contains the text 'CLOUD' in white, uppercase letters. The signpost is a silver metal pole with a decorative top. The background is a cloudy sky.

**BIG DATA**

**CLOUD**

What exactly is  **hadoop** ?

---

- Actually a growing collection of subprojects; focus on two right now



# An overview of Hadoop Map-Reduce

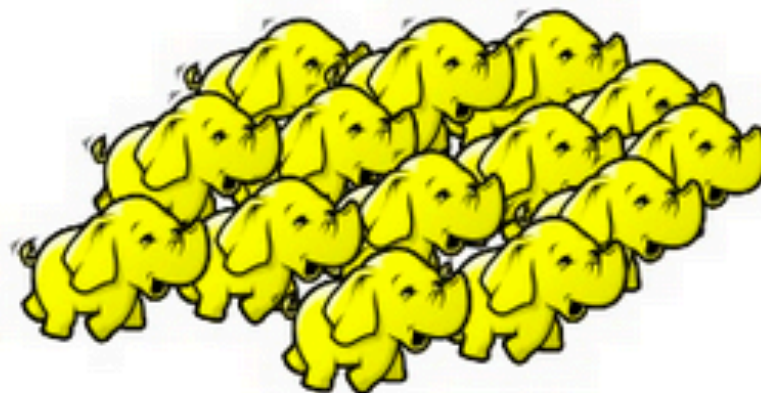
---

Traditional  
Computing



*(one computer)*

Hadoop



*(many computers)*



# An overview of Hadoop Map-Reduce

---

(Actually more like this)



*(many computers, little communication,  
stragglers and failures)*

## Map-Reduce: Three phases

---

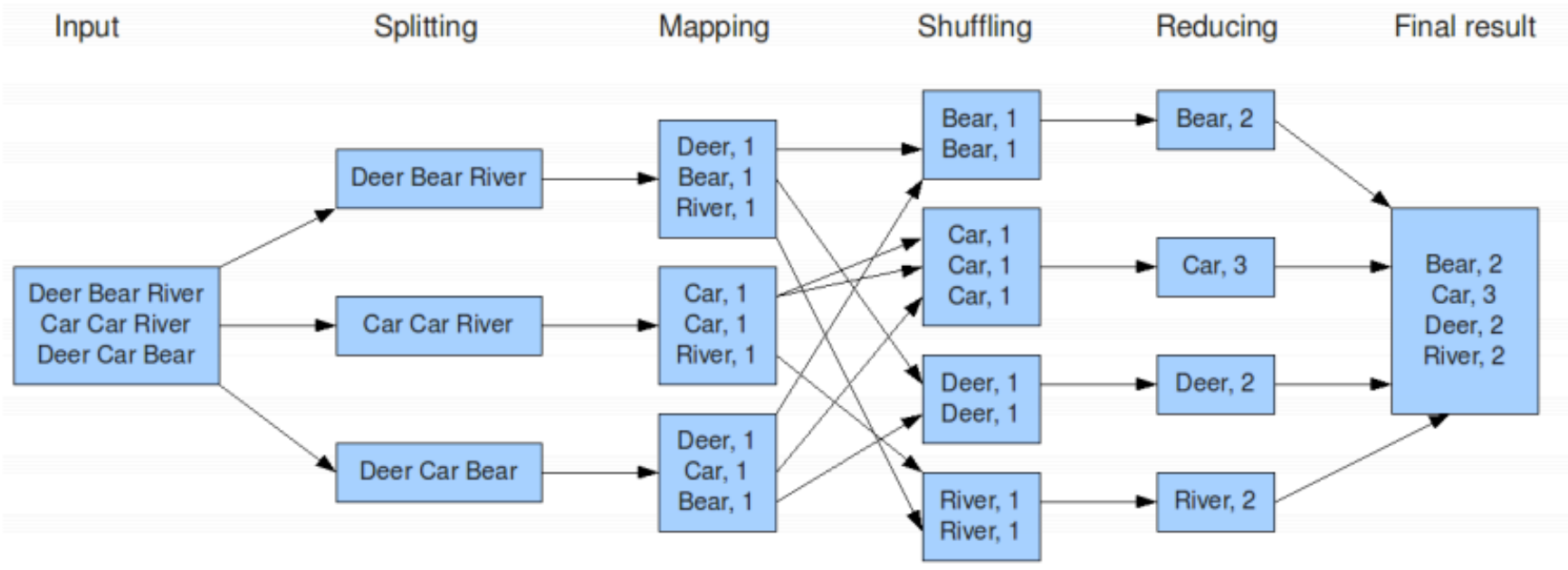


1. Map

2. Sort

3. Reduce

The overall MapReduce word count process





Share

# Map-Reduce: Imagine word-count on the Web

---

## MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

*Google, Inc.*

### Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with these issues.

As a reaction to this complexity, we designed a new

## High Throughput Computing

- Parallel Programming
- CUDA



## Using the Cloud

- Hadoop, MapReduce
- Pig, Hive, Hbase



## Analysis Tools

- Machine learning tools
- Data Mining Tools



Handling  
Big Data

# Always use the right tool !



Share





Share

It would be a win-win-win situation!

WORLD'S  
BEST  
BOSS

Office Season 2, Episode 27: Conflict Resolution)



Stay Tuned.  
Coming Soon!