

# Big Data Overview

Dec.08,2016 FANG, Jiachao

# Content

- › Story of Big Data
- › Process of Big Data Projects
- › Big Data Infrastructure
- › Case Study: Ford

# Story of Big Data



## Embryonic Stage

1980 - 1996

### Key word: Volume

In this stage, big data only means **massive data**, do not involve any processing technology or storage methods.

- › 1980, the term “**Big Data**” first used by futurist, Alvin Toffler
- › The development of communication industry leads to increase of **information flow**



## Developing Stage

1996 - 2006

### Key word: 3Vs

In this stage, **3Vs(Volume, Velocity, Variety)** became the generally-accepted three defining dimensions of big data.

- › Maturity of **data mining**<sup>1</sup> theories and techniques
- › 2003, the break out of **unstructured data**<sup>2</sup> drives the development of data processing



## Flourishing Stage

2006 - now

### Key word: Big Data Era

In this stage, with the maturity of all kinds of big data techniques, big data is widely applied in almost all fields.

- › The popularity of smartphones leads to rapid growth of **mobile data**
- › Big data application has been grown from theory research to **applied period**

<sup>1</sup> **Data Mining** is a computational process of discovering patterns in large data sets

<sup>2</sup> **Unstructured data** refers to information that either does not have a pre-defined data model or is not organized in a predefined manner

# Process of Big Data Projects in Industry

## 1. Business Understanding

---

- › Understand business objectives
- › Assess the current situation
- › Create data mining goals
- › Establish data mining plans

## 3.Data Preparation

---

- › Data selection
- › Data cleaning
- › Data construction
- › Data formatting

## 5. Evaluation

---

- › Evaluate model results in the context of business objectives
- › New business requirements may be raised
- › “Go or no-go” decision must be made

## 2. Data Understanding

---

- › Initial data collection
- › Data load and integration
- › Examine data properties
- › Examine data quality

## 4. Modeling

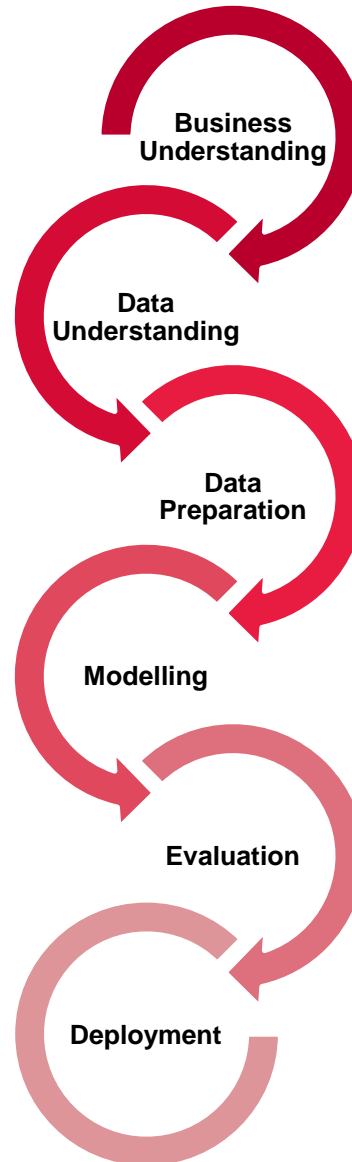
---

- › Select suitable modelling techniques
- › Use test scenario to validate the model
- › Model creation
- › Model assessment

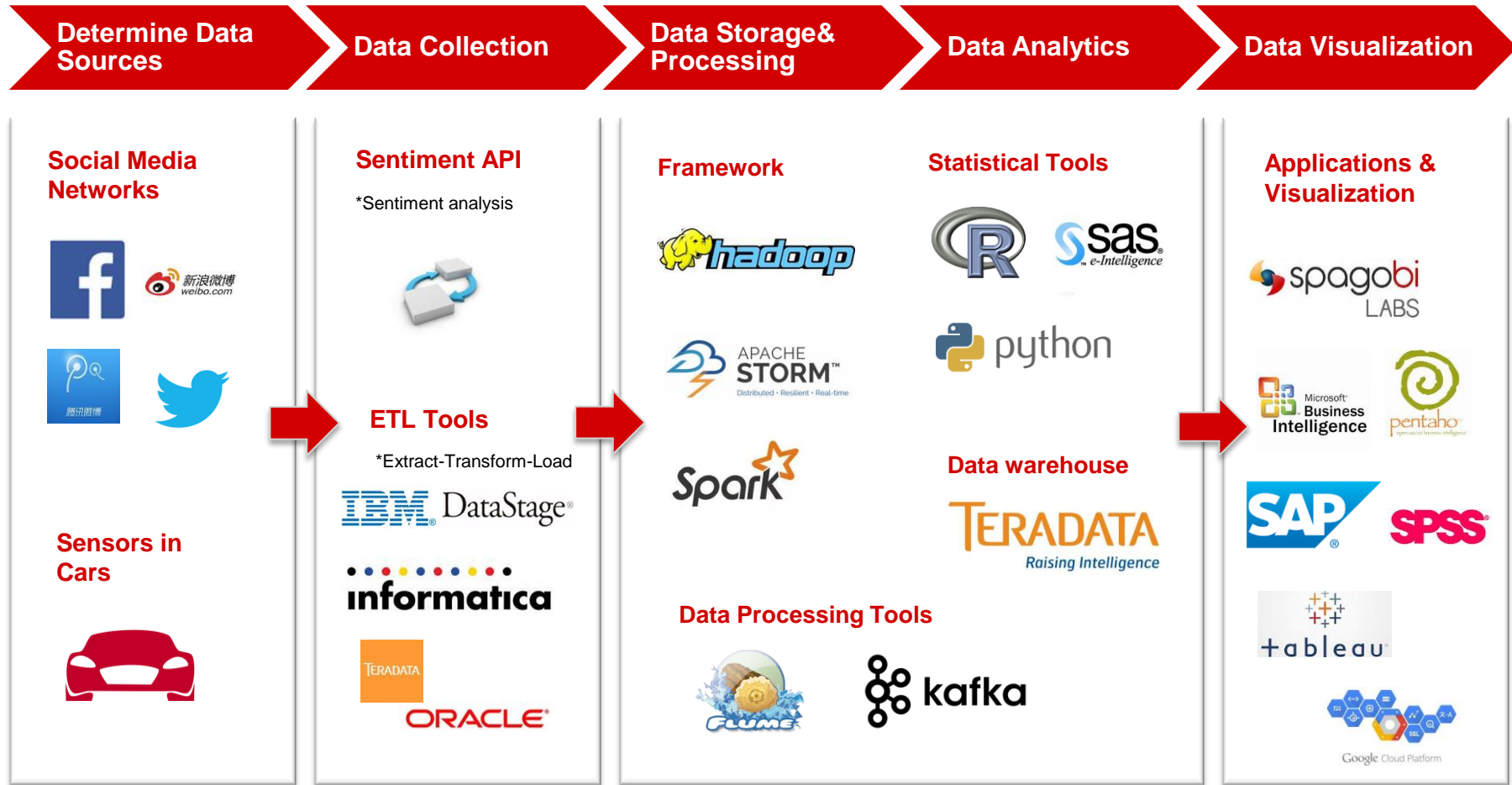
## 6. Deployment

---

- › Create a report or a repeatable data mining process
- › Create deployment ,monitoring and maintenance plan
- › Create a summary project report



# Basic Steps of Big Data



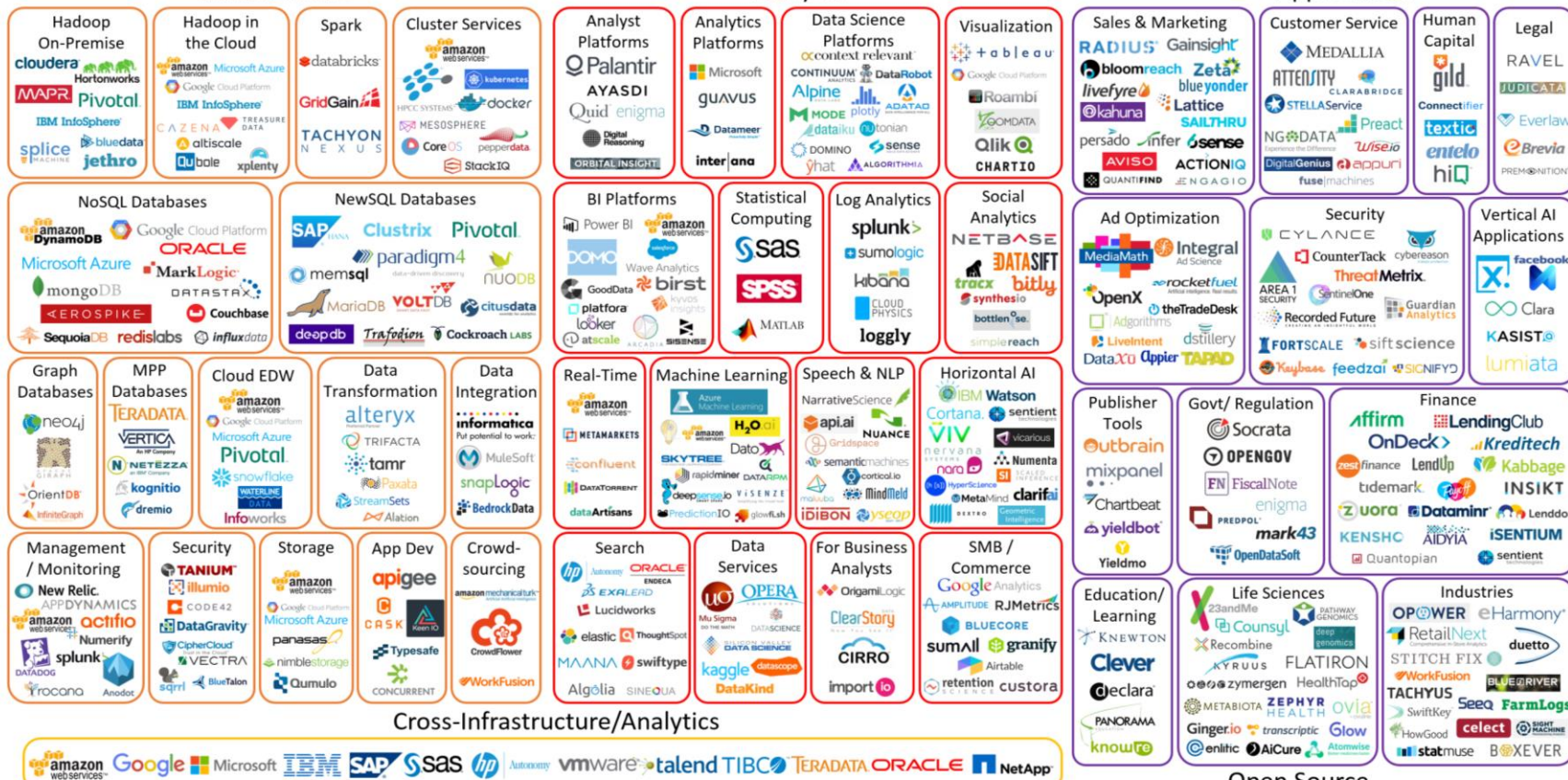


## Big Data Landscape 2016

## Infrastructure

## Analytics

## Applications



## Open Source



## Data Sources & APIs



© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap)




FIRSTMARK

Vorsprung durch Technik 

# Big Data Framework

## What exactly is a big data framework?

It is a combination of technologies and methodologies that is used to **transform big data** in its **raw form** into **refined data** governed by a mature framework that can continuously be used for virtually any big data application — whether it is a batch report, a near-real time data stream analysis or a dashboard.

	1	2	3
			
Differences	<b>Data Processing Models</b> <ul style="list-style-type: none"> <li>› best suited for batch processing</li> </ul>	<ul style="list-style-type: none"> <li>› batch processing and real time processing</li> </ul>	<ul style="list-style-type: none"> <li>› supports micro-batching</li> </ul>
	<b>Development</b> <ul style="list-style-type: none"> <li>› written in Java and implemented using Apache pig</li> </ul>	<ul style="list-style-type: none"> <li>› Implemented by Scala tuples, a bit difficult to implement over java</li> </ul>	<ul style="list-style-type: none"> <li>› uses DAG's on every node and data transfer through Storm tuples</li> </ul>
	<b>Scenarios</b> <ul style="list-style-type: none"> <li>› Batch data processing</li> <li>› distributed computing model based on RPC</li> </ul>	<ul style="list-style-type: none"> <li>› Multiple operations of specific data sets</li> </ul>	<ul style="list-style-type: none"> <li>› Offline analysis of massive data</li> <li>› Large scale search of web info</li> </ul>
Similarities	<p>Hadoop, Spark and Storm are preferred choice of frameworks amongst developers for big data applications (based on the requirements) because of their <b>simple implementation methodology</b>.</p> <ul style="list-style-type: none"> <li>› open source processing frameworks</li> <li>› real time BI and big data analytics</li> </ul>		
		<ul style="list-style-type: none"> <li>› Implemented in <b>JVM</b> based programming languages</li> <li>› Provide fault tolerance and scalability</li> </ul>	

# Stat Tools for Big Data

- Statistical tools are used for data processing, statistics and data visualization
- The right choice of statistical tools will largely increase the efficiency of data analysis
- SAS, R and Python are most frequent used among all

1



2



3



## Advantages

- › Market leader in **commercial** analytics space
- › Huge array of statistical **functions**
- › Quickly learnt

- › Open source counterpart of SAS
- › Traditionally been used in **academics and research**
- › Latest techniques get **released quickly**
- › A very **cost-effective** option.

- › An open source **scripting language**
- › Apply to almost any statistical operation / model
- › **Strong performance** in operations on structured data

## Disadvantages

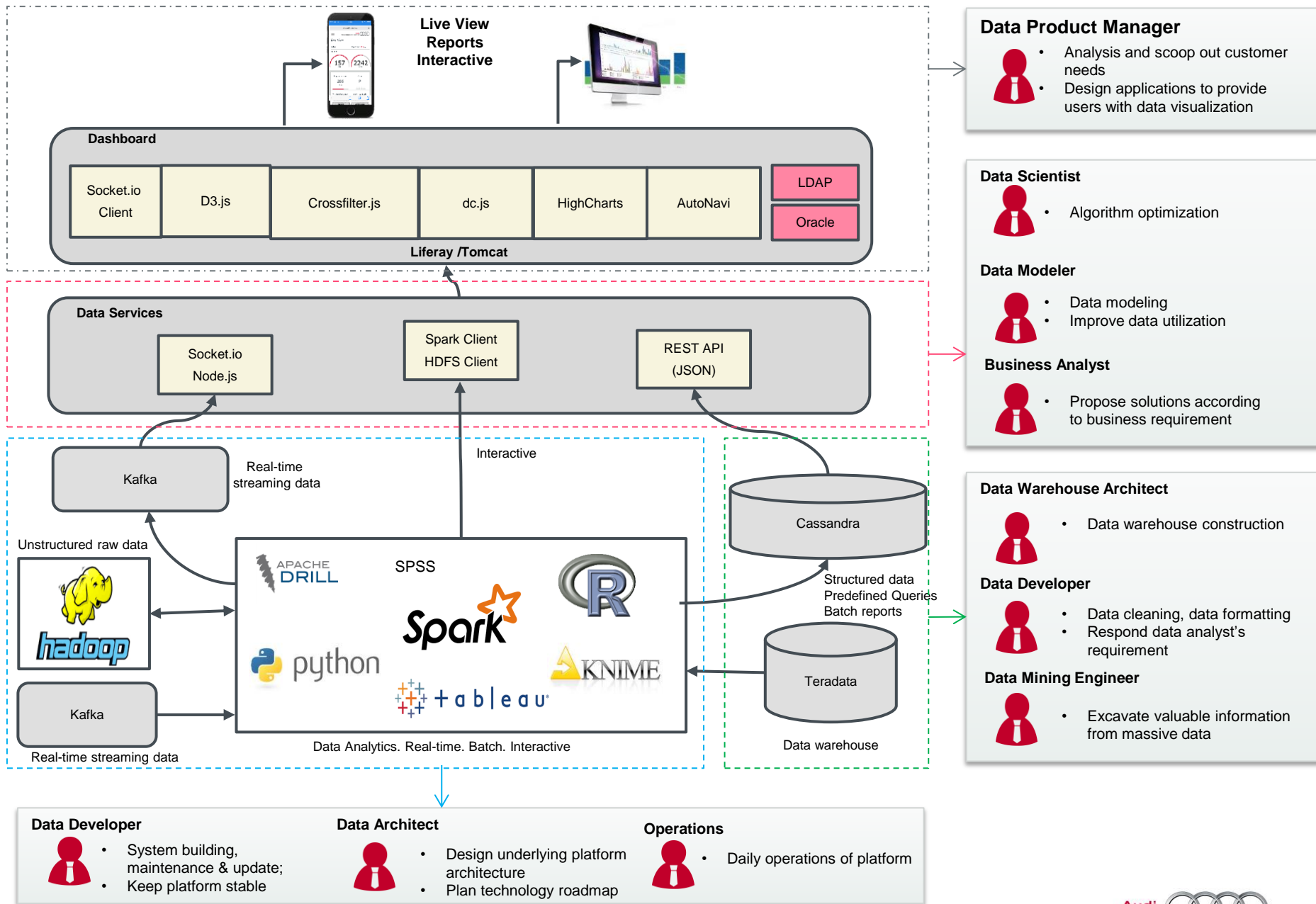
- › **Expensive**
- › not always enriched with latest statistical functions
- › Low graphical capacity

- › **Difficult to learn**
- \* Documentation is sometimes patchy and terse, and impenetrable to the non-statistician

- › **Smaller pool of Python developers** compared to other languages, such as Java
- › **Lack of true multiprocessing** support



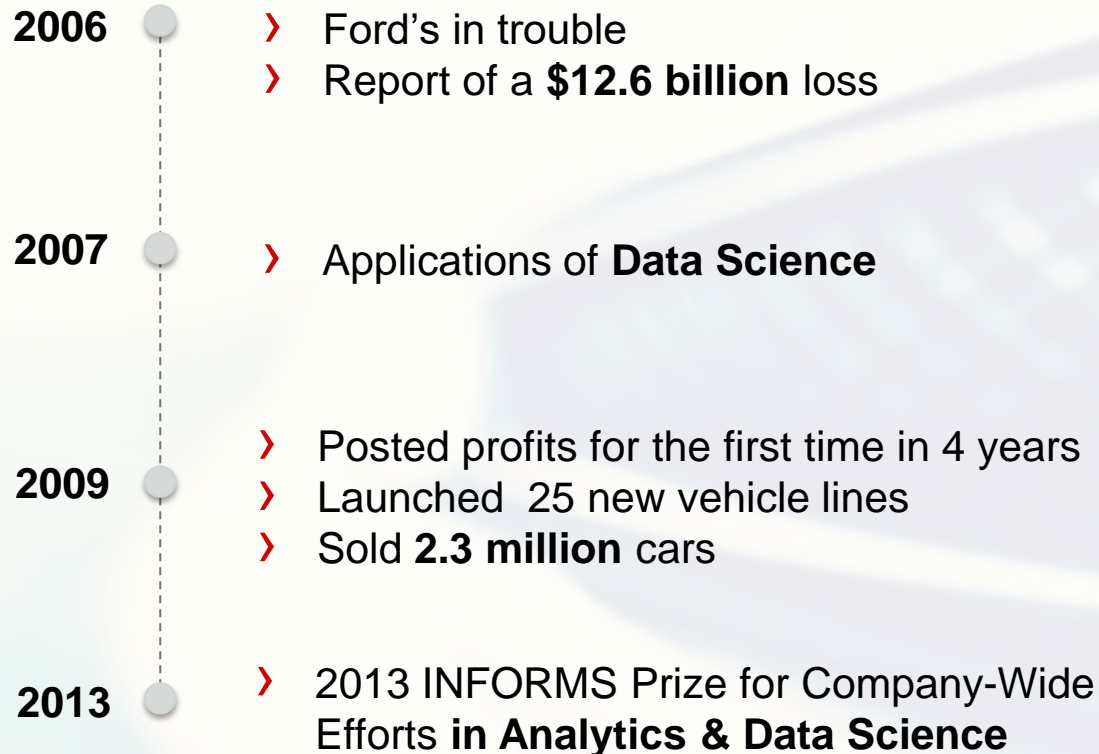
# Big Data Infrastructure & Team



# Case study: Ford

## How Big Data Brought Ford Back from the Brink

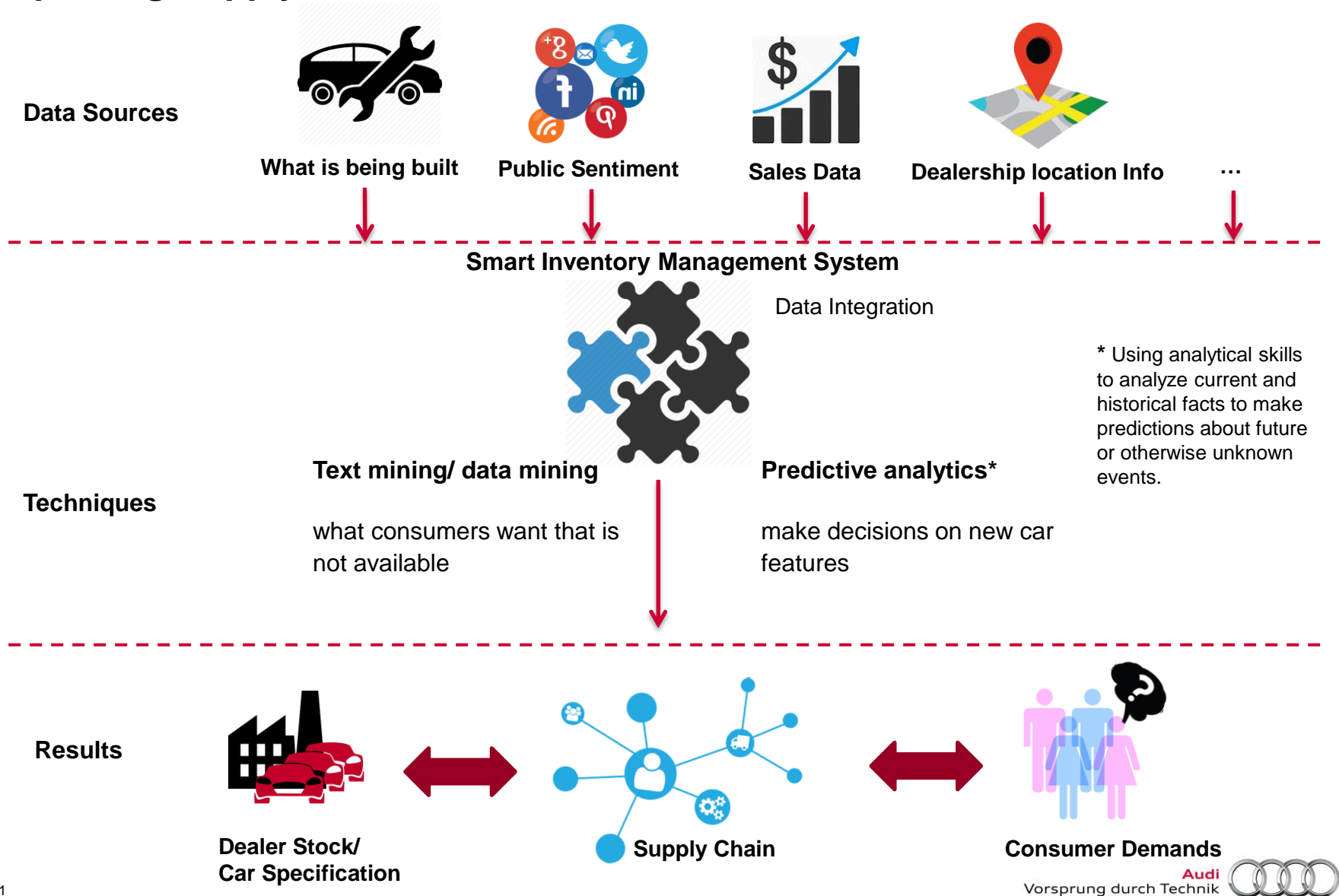
### Background

- 
- A vertical timeline on the left side of the slide, marked with years 2006, 2007, 2009, and 2013. Each year is accompanied by a grey circular marker and a dashed vertical line extending downwards. To the right of each marker, red chevron symbols (›) introduce a list of events.
- 2006**
    - › Ford's in trouble
    - › Report of a **\$12.6 billion** loss
  - 2007**
    - › Applications of **Data Science**
  - 2009**
    - › Posted profits for the first time in 4 years
    - › Launched 25 new vehicle lines
    - › Sold **2.3 million** cars
  - 2013**
    - › 2013 INFORMS Prize for Company-Wide Efforts in **Analytics & Data Science**



# Case study: Ford

## Improving Supply & Demand



# Case study: Ford

## Improving Efficiency on the Factory Floor

### Before



high cost of producing prototypes

### Solution



#### Prototype Optimization Model

- › Computing the **minimum number of vehicles** required to perform the **maximum amount of tests**

### After



- › An estimated **\$12 million** in savings first used on the European Transit vehicle
- › An estimated **\$250 million** in savings if rolled out the whole company



# Case study: Ford

## Anticipating Future Energy Supply and Demand

### Data Sources

- › performance data from sensors in electric cars



Where customers plug in



How many gas miles they drove



How many electric mile



How often they take trips

...

### Solution



### Results



- › Detailed **reports for drivers** about car performance and driving conditions



- › Developing vehicles using a range of **new fuel technologies**



- › Combined with **internal process** and knowing where to improve

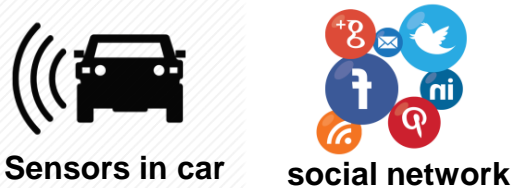
Audi  
Vorsprung durch Technik



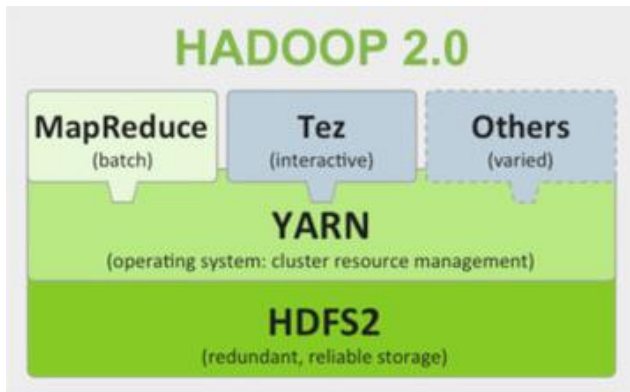
# Case study: Ford

## Ford's Big Data Strategy

### Data Sources



### Analytics Platform



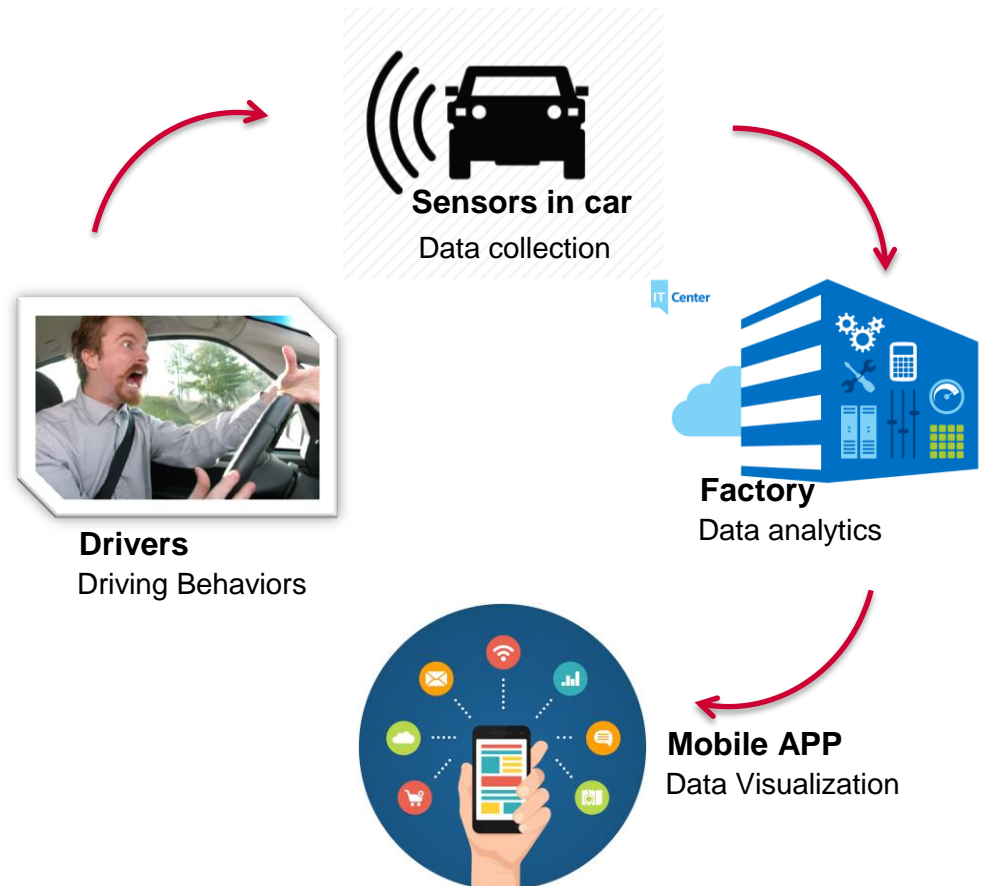
### Tools

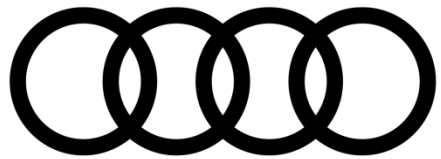
- › R :statistical analysis

### Methods

- › sentiment analysis; real-time analysis; data mining; text mining

### Instance: Behavior Monitoring





► Thank you!