

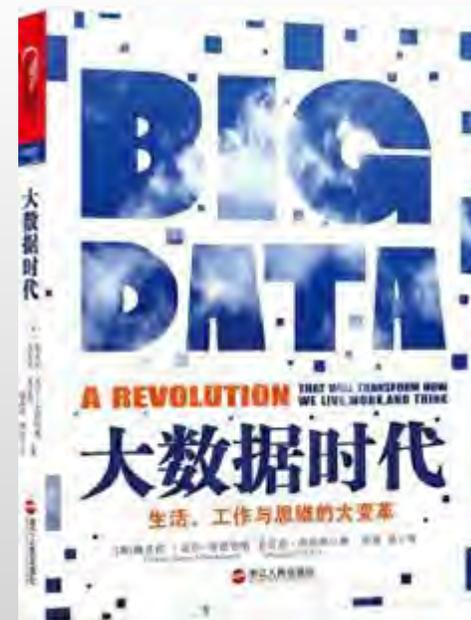
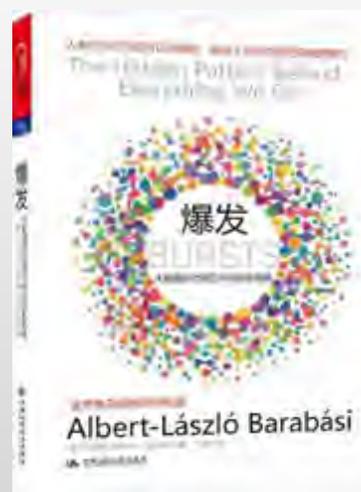
中国大数据发展面临的挑战

主讲人：沈浩

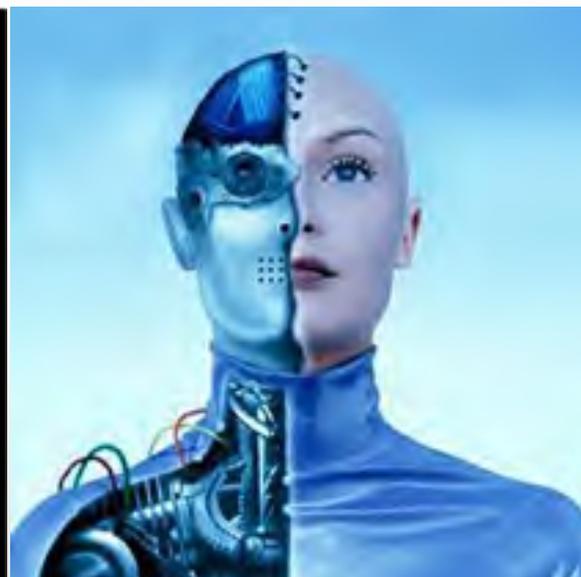
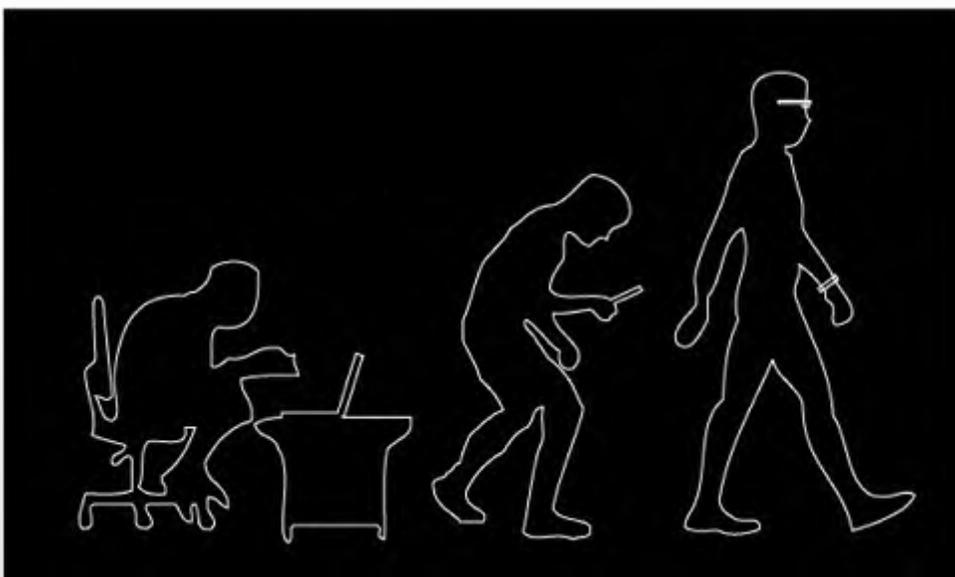
中国传媒大学新闻学院	教授 博导
中国传媒大学调查统计研究所	所长
大数据挖掘与社会计算实验室	主任
中国市场研究行业协会	会长

2016

这是一个令人兴奋的时代，也是一个大数据的时代，数据科学让我们越来越多地从数据中观察到人类社会的复杂行为模式。以数据为基础的技术决定着人类的未来，但并非是数据本身改变了我们的世界，起决定作用的是我们对可用知识的增加。



人类行为的93%是可预知的！





BIG DATA

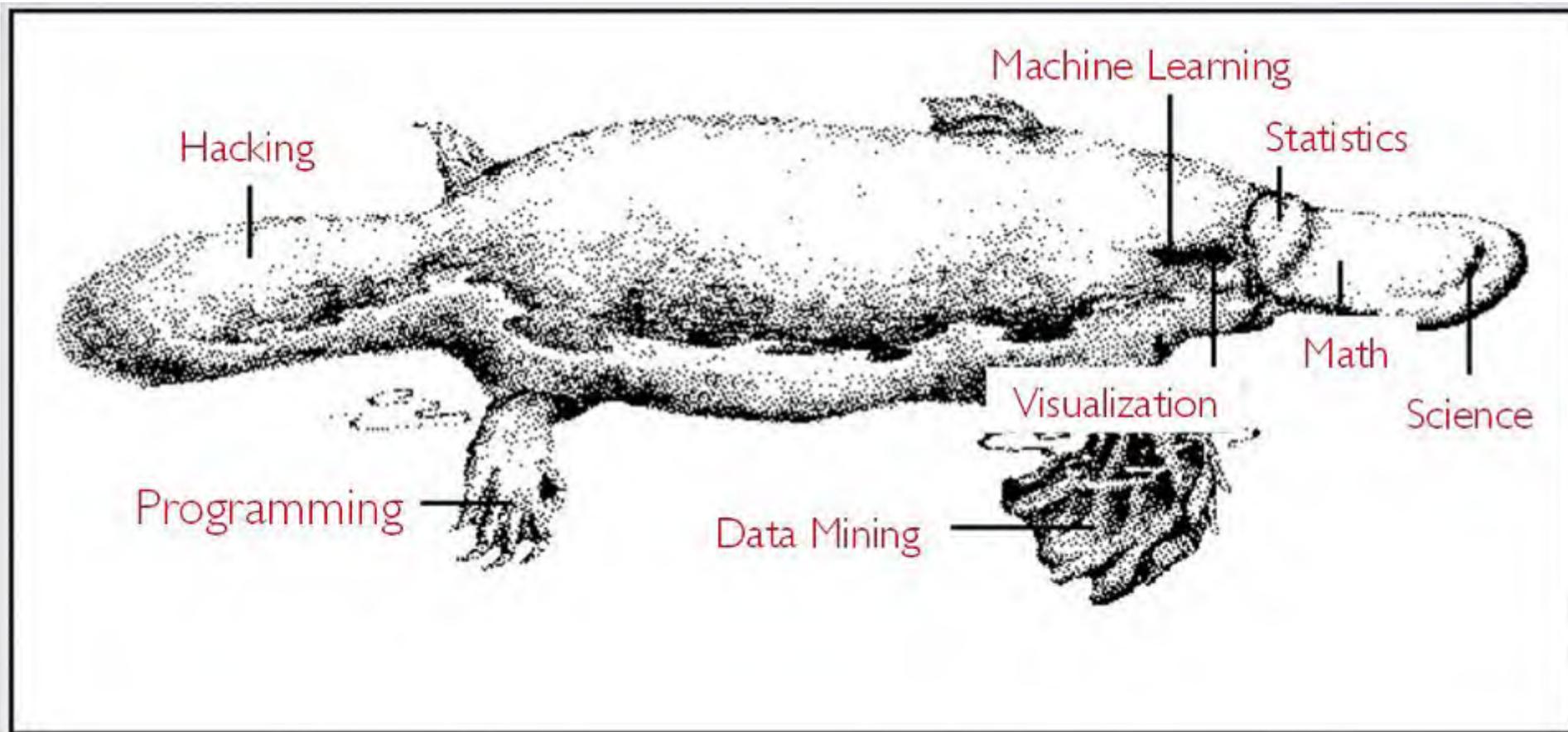
什么是大数据 (Big Data) ? 它将如何改变我们的社会生活? 对政府管理、商业活动、媒介生态、个人生活将产生怎样深刻的影响? 我们该如何拥抱大数据?

其实, 没有多少人在真正的接触大数据, 我们更愿意说这是一个大数据时代, 或许是大数据太热了, 不同学科背景的人都在谈论, 不同行业的都在谈论, 说明大数据时代到了, 全球已经点燃了**大数据时代**。



- Annual data creation in zettabytes (10007 bytes)
- 90% of the world's data created in the last 2 years

大数据冰山一角



大数据时代

大数据的下一个发展阶段？

- 
- 2014，美国白宫：《大数据：抓住机遇、守护价值》
 - 2012，Splunk成为第一家上市的大数据处理公司
 - 2012，瑞士达沃斯：《大数据，大影响》
 - 2011，麦肯锡《大数据：创新、竞争和生产力的下一个新领域》
 - 2010，《经济学人》：“数据，无所不在的数据”
 - 2009，美国政府：Data.gov（开放数据）
 - 2008，计算社区联盟：《大数据计算：在商务、科学和社会领域创建革命性突破》
 - 2005，Hadoop项目诞生

大数据时代



2009年开始,“大数据”作为中国互联网信息技术行业的流行词汇

2011年12月,工信部发布物联网十二五规划,提出的海量数据存储、数据挖掘、图像视频智能分析都是大数据的重要组成部分

2012年7月,阿里巴巴集团在管理层设立“首席数据官”一职,负责全面推进“数据分享平台”战略

2013年,媒体称之为“大数据元年”

2016年3月5日,国务院总理李克强政府工作报告,专门提出“促进大数据、云计算、物联网广泛应用。”这也是自2014年3月5日首次进入政府工作报告以来,大数据连续三年成为我国政府的聚焦点,甚至已被看做“新经济”的高效率引擎。



CONTENT



01

更新技术架构

02

培养专业人才

03

打破数据孤岛

04

权衡开放&隐私

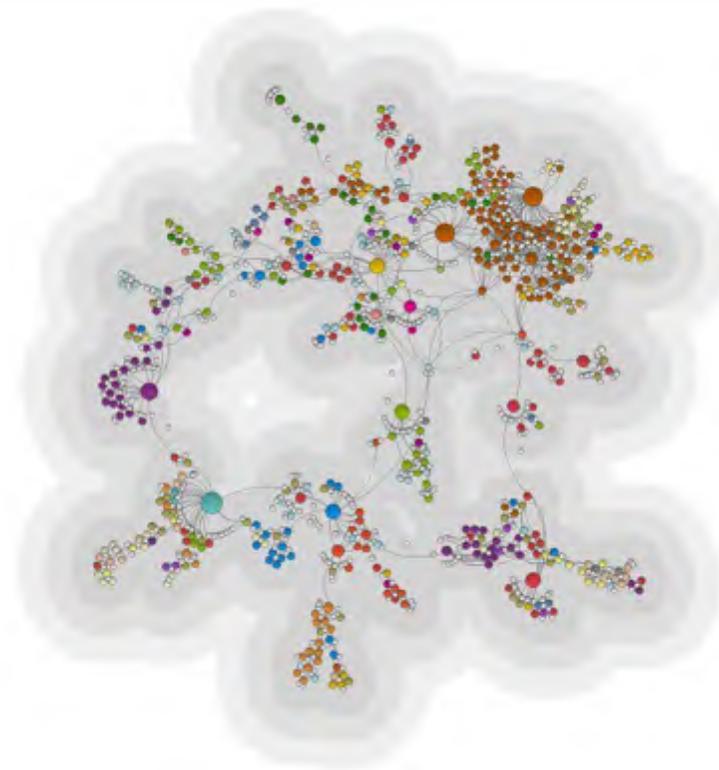
BIG DATA



更新技术架构



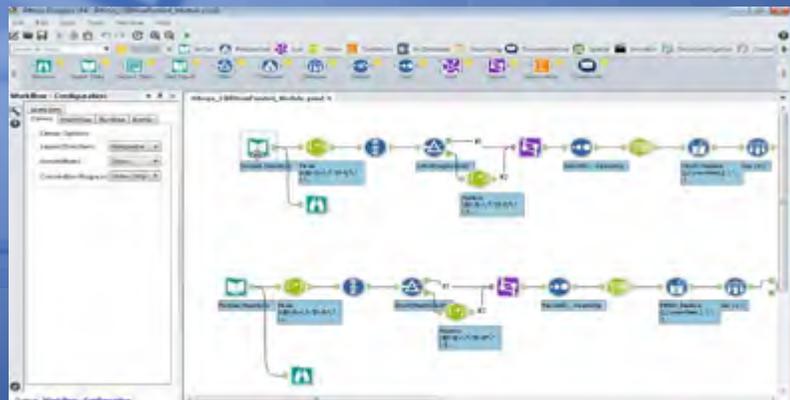
Visualization



让数据模式改变你的心理模式

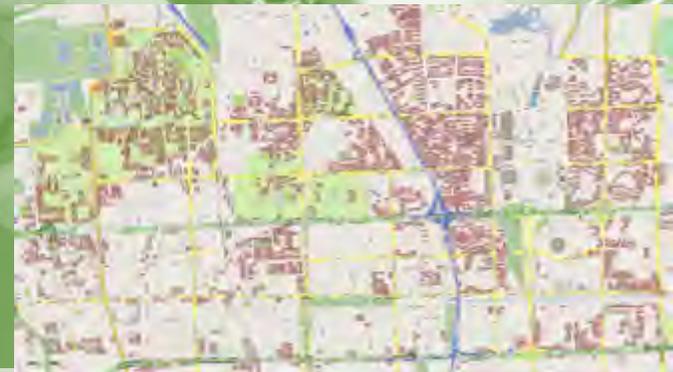
数据科学 Data Science

- 数据挖掘
- 文本挖掘
- 意见挖掘
- 情感挖掘
- 语义挖掘



空间地理科学 Spatial Science

- GIS分析
- 空间地理信息匹配算法
- 城市规划
- 交通与智慧城市



网络科学 Network Science

- 复杂网络
- 社会网络分析
- 传播网络



可视化技术 Data Visualization

- 数据可视化
- 信息可视化
- 交互可视化
- 大屏实时展现



因此，挖掘大数据价值需要技术更新

- 分布式文件系统和并行计算框架
- 复杂算法的高效迭代运算能力
- 数据结构从关系型到非关系型发生改变
- 存储、计算、应用一切尽在云端
- 非结构化和关系数据的分析方法发展
- 机器学习和人工智能大行其道
- 空间地理分析成为大数据分析重点
- 数据可视化重要性上升到极大高度

海量数据与Big Data的图示化理解

海量数据——结构化数据——适合数据挖掘——分类算法

Type	年龄	性别	证券业务	数据业务	客户等级	险种	客户分类
1001	23	女性	79.25	31.26	高端客户	预付保费	VIP客户
1002	47	男性	73.93	56.47	低端客户	预付保费	普通会员
1003	47	男性	89.73	68.94	低端客户	预付保费	VIP客户
1004	28	女性	98.37	72.28	普通客户	预付保费	VIP客户
1005	61	女性	66.93	31.00	低端客户	预付保费	VIP客户
1006	22	女性	67.69	78.65	普通客户	预付保费	普通会员

Big Data——非结构化数据——适合文本挖掘——词云或关联规则——个性化推荐

1001 收入预期对于缺乏一手数据的同学来说是非常困难的，后者是热门的数据挖掘，前者还有文本挖掘、web挖掘、数据挖掘等等

1002 数据是以非结构化和半结构化数据为主，新晋偏向结构化数据和交易数据，现在还有叫社会计算的挖掘应用

传统的结构化数据存储结构无法进行有效分析

语义分析后，转换为Hadoop结构

tweet	Serial	1001	1002
1	1	收入	前者
1	2	预期	是
1	3	对于	以
1	4	缺乏	非
1	5	一手	结构
1	6	数据	化
1	7	的	数据
1	8	同学	方
1	9	来	王
1	10	说	前者
1	11	是	偏向
1	12	非常	说
1	13	困难	结构化
1	14	的	和
1	15	-	交易
2	21	后者	数据
2	22	是	。
2	23	语义	现在
2	24	的	还有
2	25	数据挖掘	叫
2	26	前者	社会
2	27	还有	计算
2	28	文本挖掘	的
2	29	Web挖掘	挖掘
2	30	数据库	应用
2	31	挖掘	等
2	32		

1-可以容易分析各编号(列)的词频，删除高频词(是、的等)后，就可以得到该ID的主要内容词频，完成词云分析和词云图

2-对任何编号ID可以切片存储，这是Hadoop的特点，而且可以分析1001与1002的相关和关联规则，相似性和聚类特性等!

海量数据：是在社会化数据没出现就有了，商业自动化导致海量数据存储，用于决策的有效信息隐藏在数据中，数据挖掘应运而生!

大数据：是伴随社会化数据出现和大量的在线流量文本(图片、流媒体)数据等，主要为了应对非结构化和半结构化数据;

首先是内容的存储结构发生适应性改变，特征就是Hadoop的出现! 社交挖掘、文本挖掘、Web挖掘等!

海量数据的数据挖掘可特指一种KDD分析方法，有较完善的方法论和数据挖掘软件工具!

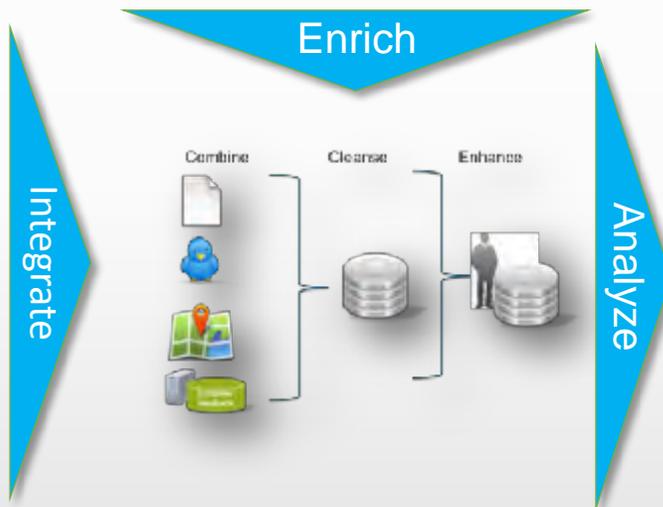
Big Data更多的是宽泛的概念，是一种潮流或IT趋势，内涵和外延比较宽泛，没有特定的方法论限定!

大数据技术更新

All Relevant Data



Packaged Market & Customer Data & API & Census



Integrate any data source
Spark Streaming



Rapid design of predictive analytics with unique spatial understanding



大数据技术更新

Visualization

Data, business intelligence, cloud

Big Data

Machine Learning, blending, extract-transform-load

Predictive Analytics

Models, one-click data mining, decision management

Data Warehousing

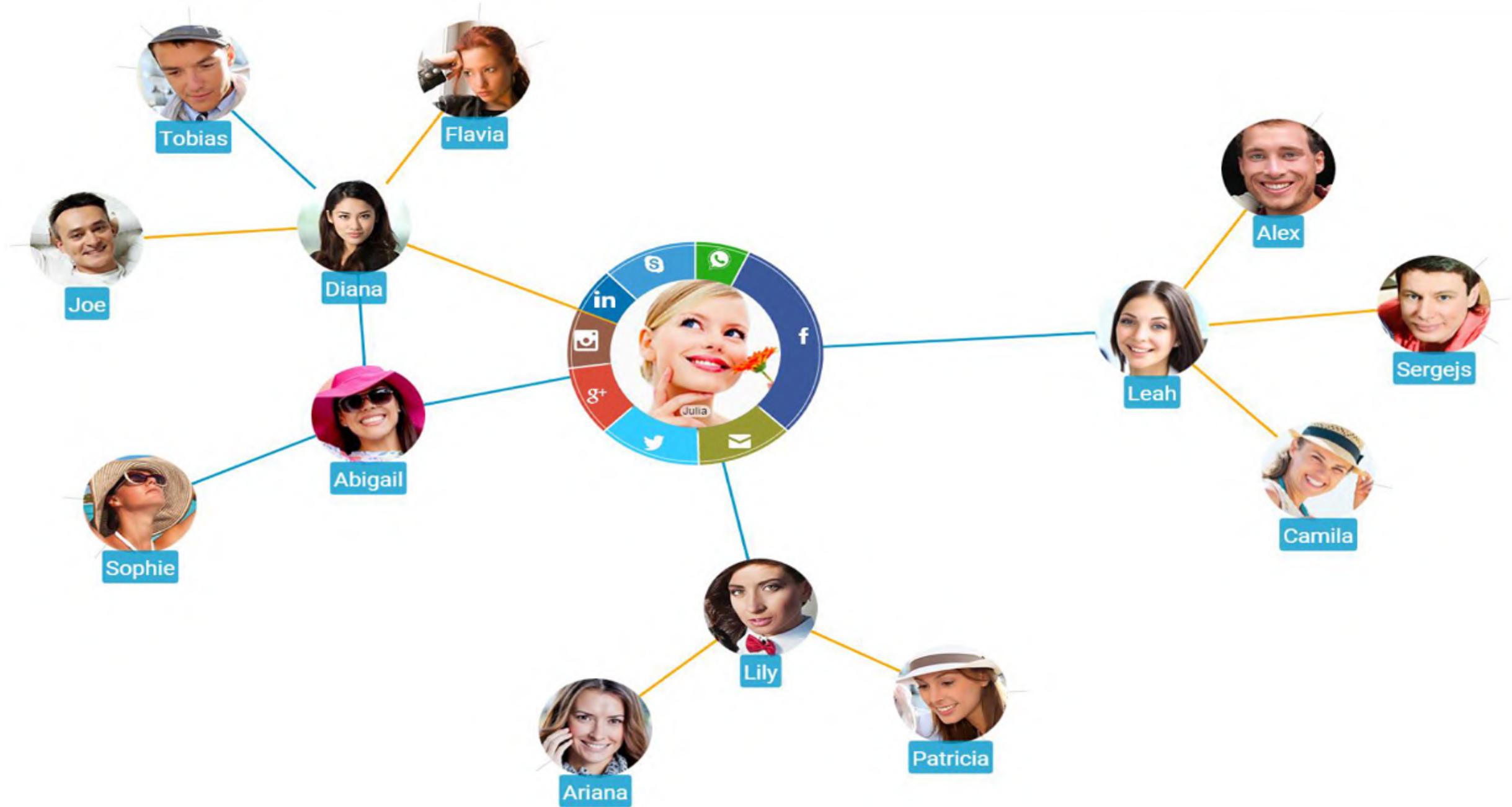
Data storage, data modeling



大数据技术更新

The screenshot displays the Microsoft Azure Machine Learning environment. The main workspace shows a workflow titled "Income Prediction" with the following steps: "Adult Census Income Binary..." (dataset), "Split Data", "Two-Class Boosted Decision..." (model), "Train Model", "Score Model", and "Evaluate Model". A sidebar on the left lists various tools and services like "Saved Datasets", "Trained Models", and "Machine Learning". A right-hand pane shows "Properties" and "Project" details, including "Experiment Properties" (Start Time: 4/18/2016, End Time: 4/18/2016, Status Code: Finished) and a "Description". Below the main interface, a Windows File Explorer window shows the local file system, and a "Workflow - Configuration" window displays a detailed flowchart of the training process.

The screenshot shows the IBM Watson Analytics dashboard. The top navigation bar includes "Explore", "Predict", "Assemble", "Social Media", and "Refine". The main area displays three data quality cards: "23 LOW QUALITY" for "Social Media data set", "70 MEDIUM QUALITY" for "mysocial", and "Employee Performance". Below these, a workflow diagram is shown with the following stages: "File Reader" (Read IMDb reviews from CSV file), "Document Creation" (Transformation of strings to documents), "Preprocessing" (Preprocessing of documents), and "Feature Creation" (Creation of document vectors of frequent 1grams and 2grams). The workflow is split into two paths: "1-gram features" and "1- and 2-gram features". Both paths use a "Decision Tree Learner" (Node 308 and Node 309), a "Decision Tree Predictor", and a "Scorer". The final steps are "Joiner" (Join class probabilities) and "ROC Curve" (Score decision tree models).



掌握的主要技术和方法



掌握的业务流程



BIG DATA



培养专业人才

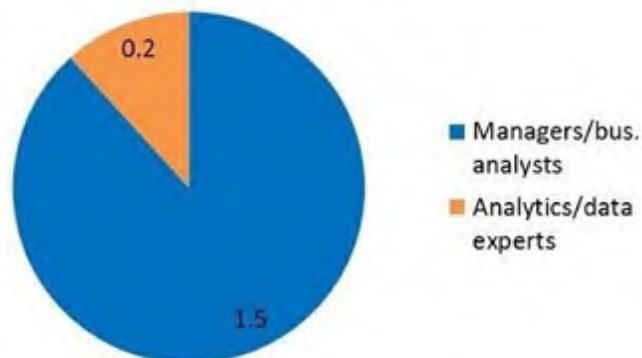
大数据人才需求

Gartner®

2015年，全球将新增440万个与大数据相关的工作岗位，且会有25%的组织设立首席数据官职位。未来，大数据将会出现约100万的人才缺口。



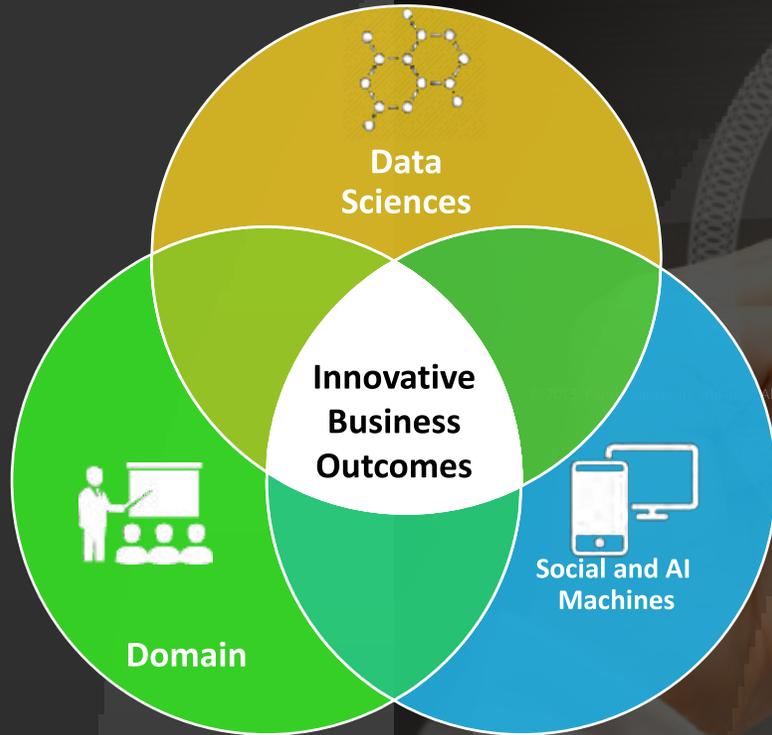
Big Data Skills Gap: 1.7 Million Workers By 2018



Source: McKinsey Global Institute report: "Big Data: The Next Frontier for Innovation, Competition and Productivity"

麦肯锡全球研究院的研究预测在未来6年，仅在美国本土就可能面临缺乏14万至19万具备深入分析数据能力人才的情况，同时具备通过分析大数据并为企业做出有效决策的数据的管理人员和分析师也有150万人的缺口。

Data Artisans 数据工匠



Data
Sciences



Domain

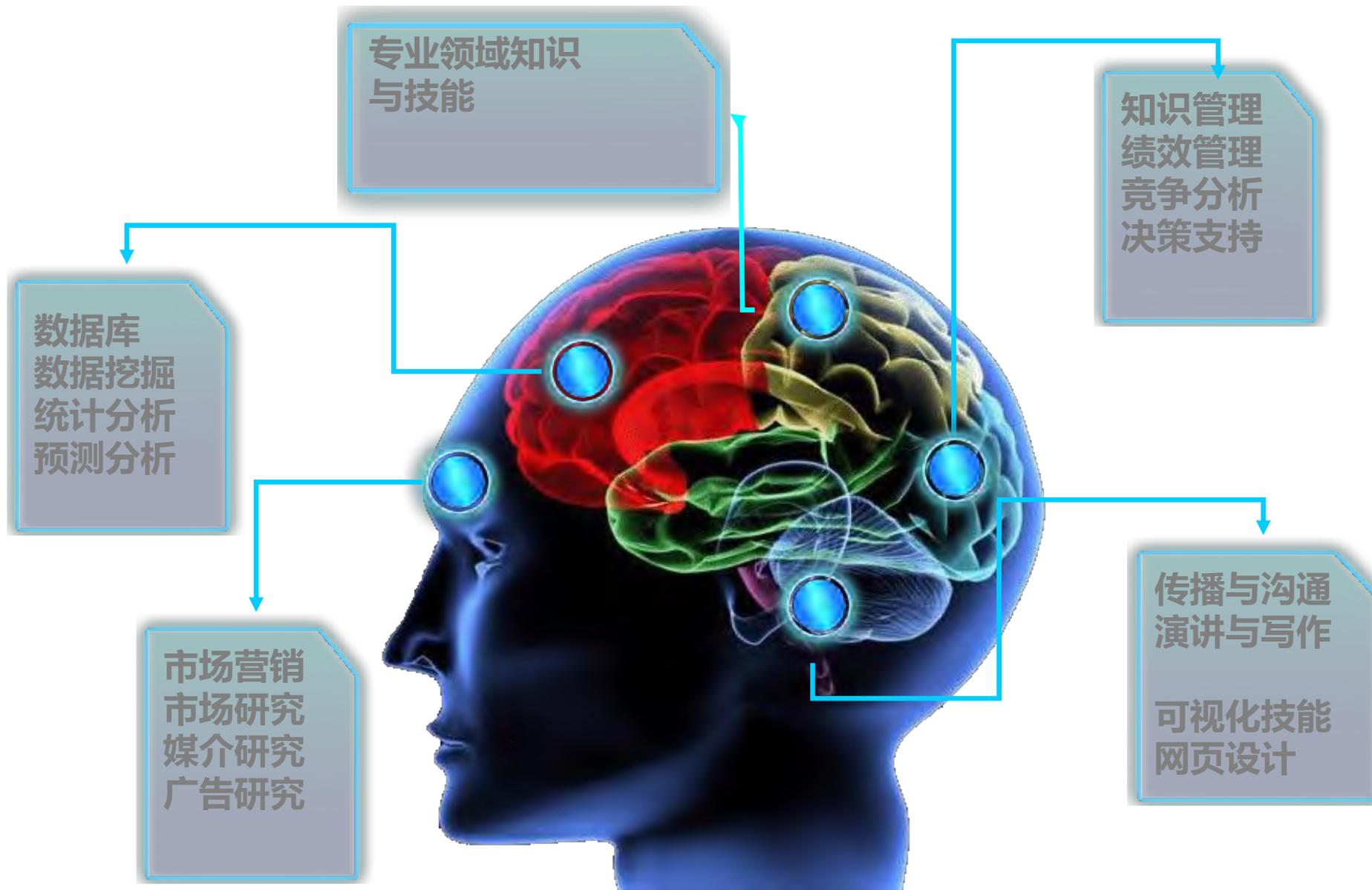


Social and AI
Machines

Innovative
Business
Outcomes

All rights reserved. Pluto7 Confidential

培养专业人才



大数据人才培养

组建师资队伍
开设专业课程
寻找实训基地

高校



知识学习



&



企业



技术实践

建立合作关系
开发合作项目
实现人才交流

保持紧密高效联系，实现联手双赢局面

BIG DATA



打破数据孤岛

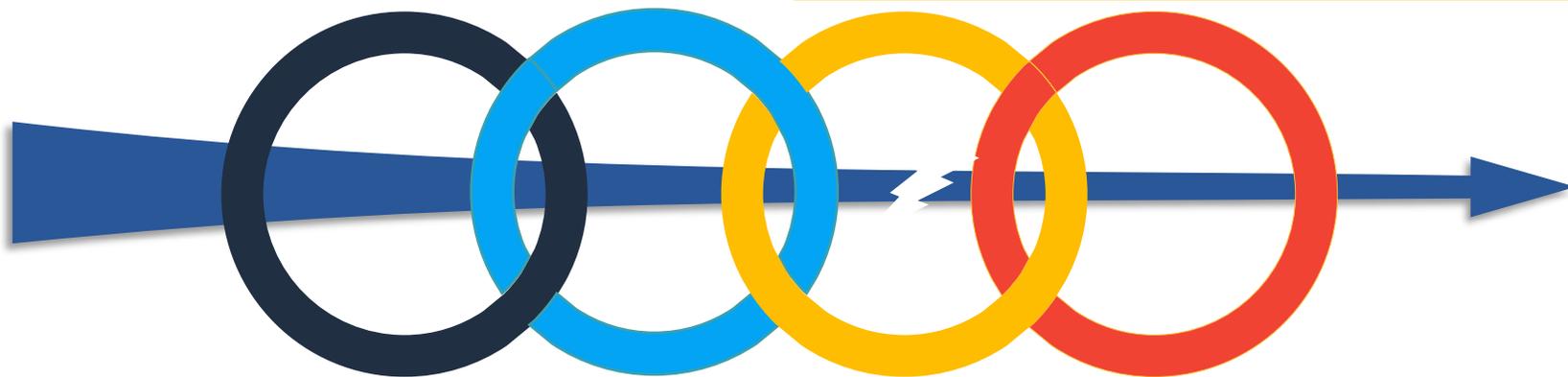


风景虽然美丽
现实问题严重

数据孤岛：指的是一个个相对独立的不同类型的数据资源系统。

1. 企业知道他们能够从信息和数据中获得更有价值的洞察，但不知道怎么做。

3. 数据散落在不同部门，存在不同的数据仓库中，不同部门的数据技术也有可能不一样，这导致企业内部自己的数据都没法打通。



2. 虽然没有明确的大数据业务需求，但希望可以整合企业数据，保护数据资产。

4. 大数据需要不同数据的关联和整合才能更好的发挥理解客户和理解业务的优势。将不同部门的数据打通，并且实现技术和工具共享，才能更好的发挥企业大数据的价值。

数据湖：表面上看，数据都是承载在基于可向外扩展的HDFS廉价存储硬件之上的。但数据量越大，越需要各种不同种类的存储，并且，不是所有的企业数据都是适合存放在廉价的HDFS集群之上的。



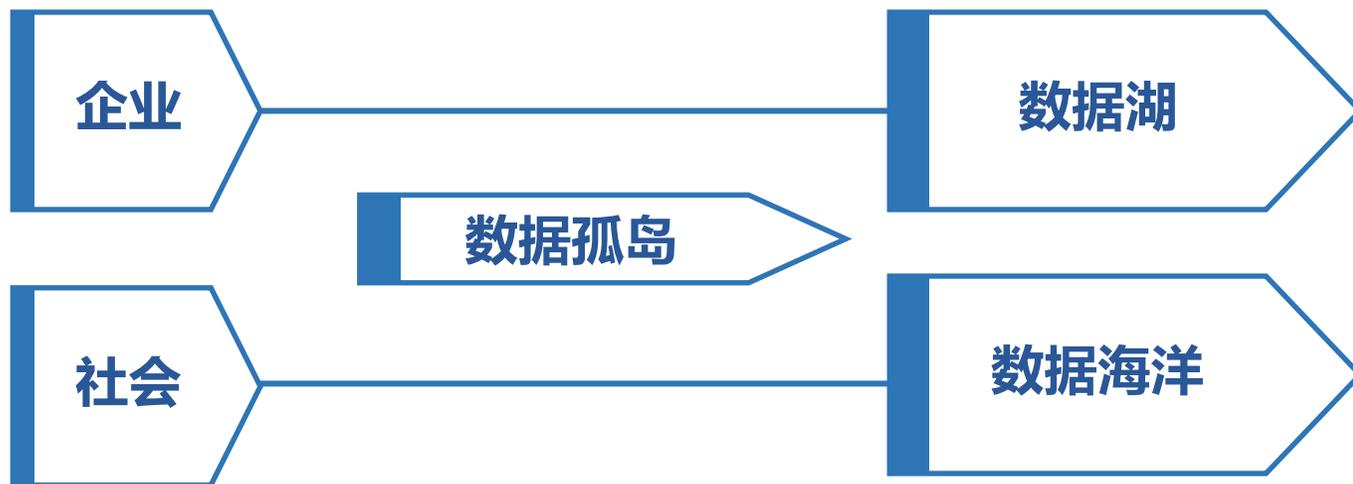
数据湖近在眼前，有人持怀疑的态度，有人热情拥抱。怀疑者认为，数据湖无非另一种将全部数据整合至单一位置的存储形式，支持者认为，数据湖不仅预示着前所未见的存储效率，还让分析成为可能，让每个组织都可用。

数据湖概念的提出，让我看到打破数据孤岛的重要性。

从宏观层面上看，数据孤岛：**开放数据**

指的是政府、企业和行业信息化系统建设往往缺少统一规划和科学论证，系统之间缺乏统一的标准，形成了众多“数据孤岛”。

因此，它的意义不仅仅是数据开放，更重要的是能够制定一个可以遵循的**数据存储、读取**的管理规则，为所有人使用。

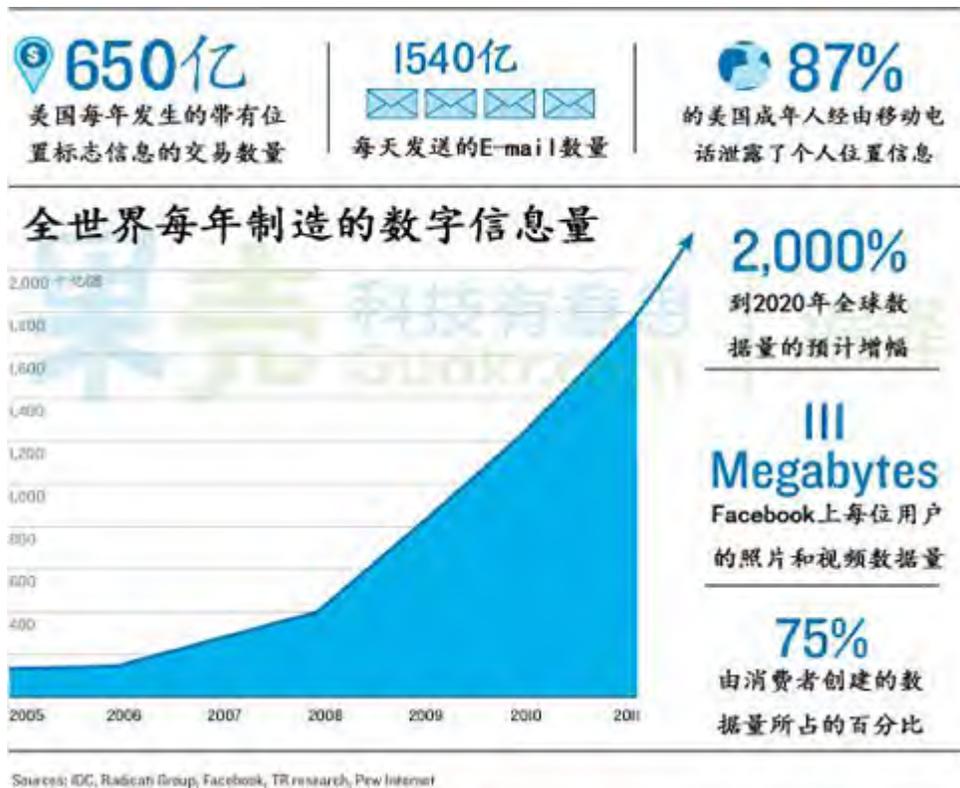


BIG DATA



权衡开放&隐私

数据开放与隐私

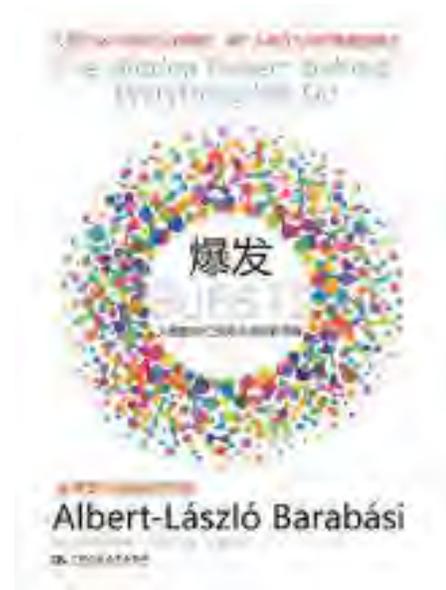


图片：MIT TechnologyReview (2013)



尽管这其中的大部分数据都是不可见的，似乎也并不携带任何个人信息，但事实并非如此。现代数据科学已经发现几乎任何类型的数据都可能用来识别创造它的人，实际上，数据越多，其中可以称得上隐私的就越少。

“我点击了自己的名字，页面上出现了一张熟悉的照片——是我穿着一件蓝色衬衫的照片，旁边配有我的基本履历资料……我点开了一个最近更新的链接，地点是波士顿的马萨诸塞大街……两秒钟后，我在视频中看到了自己推开了地铁站那厚重的大门……每次看到自己出现在视频中，我都会浑身不自在。但现在可好，我的一举一动已经被LifeLinear网的系统给记录了下来……”

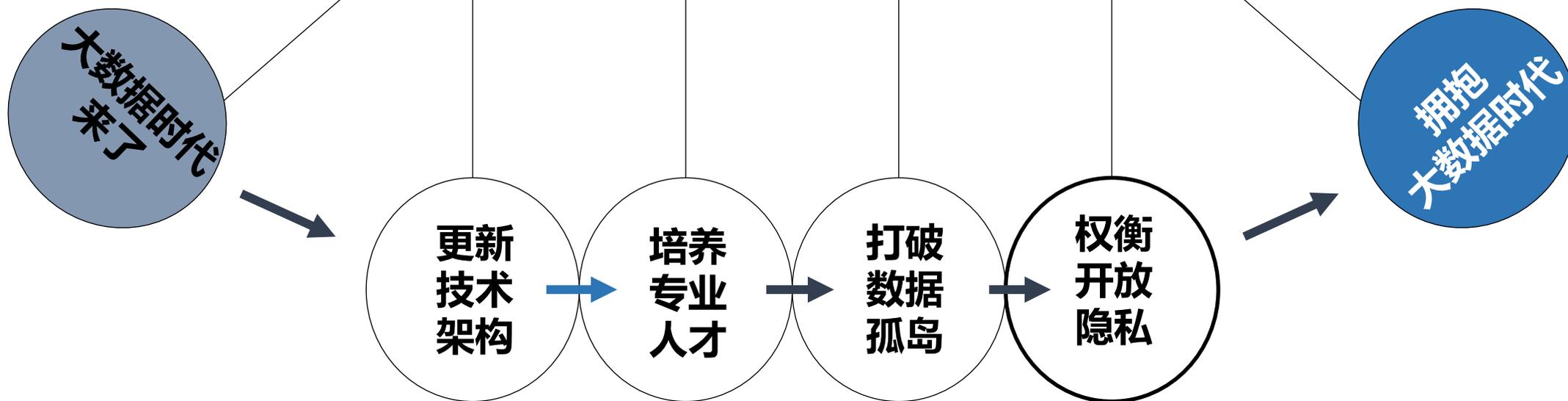


开放是基石，隐私是关键

加大立法制度

加强行业管控

关于今天的主题



大数据时代的到来，让我们认识到现有技术的缺陷，而新技术、新产品、新公司如雨后春笋般涌现；从不缺少技术的更新，当然也需要人才的跟进；当然，对数据资源的需求也在逐步增加，需要数据的开放与相互之间的连接；此时，我们应该回归本体，注重对数据隐私的保护。

我们生活在中，就不得不同数据打交道。我们也是数据的一部分，不论我们想不想与大数据牵扯到一起，数据都会找到我们，覆盖我们。

大数据时代已经来临，如何从海量数据中发现知识，寻找隐藏在大数据中的模式、趋势和相关性，揭示社会现象与社会发展规律，以及可能的商业应用前景，都需要我们拥有更好的数据洞察力。



中国大数据发展面临的挑战

THANKS.

2016