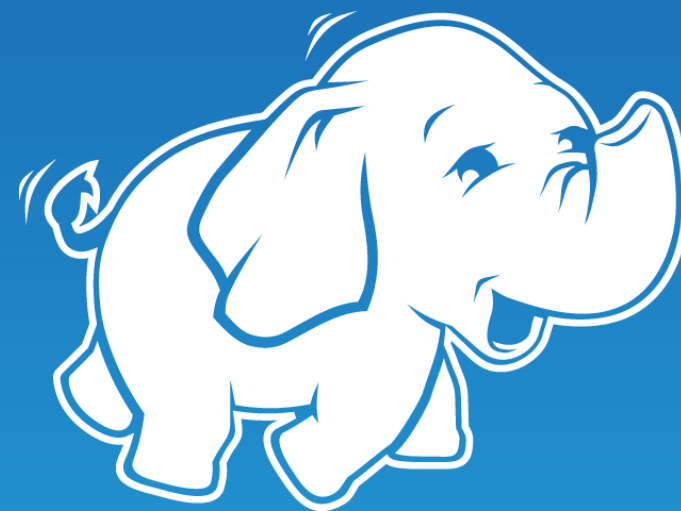


# 大数据和Hadoop介绍



姓名  
职位  
公司

# 议程



什么是大数据?

理解基础

微软和Hadoop

# 什么是大数据?



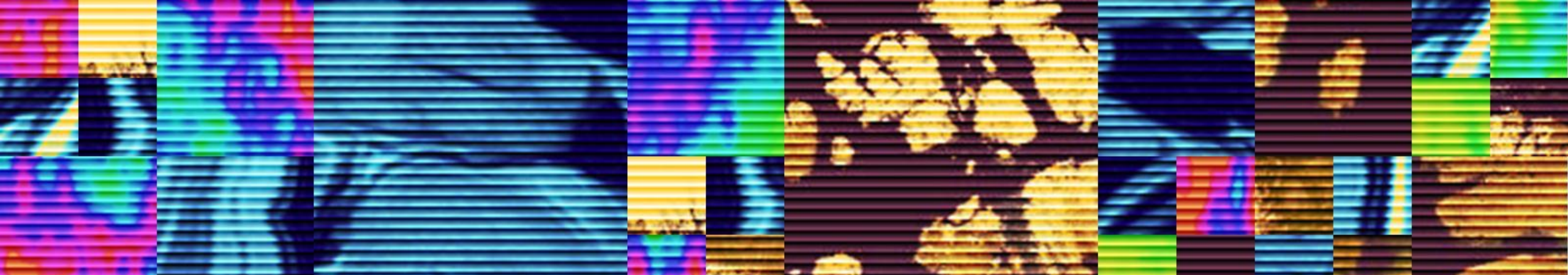


# 1.8 ZETTABYTES

信息将会在2011年创建

来源: CenturyLink 资源中心, 根据readwriteweb报告, 11月17日, 2011

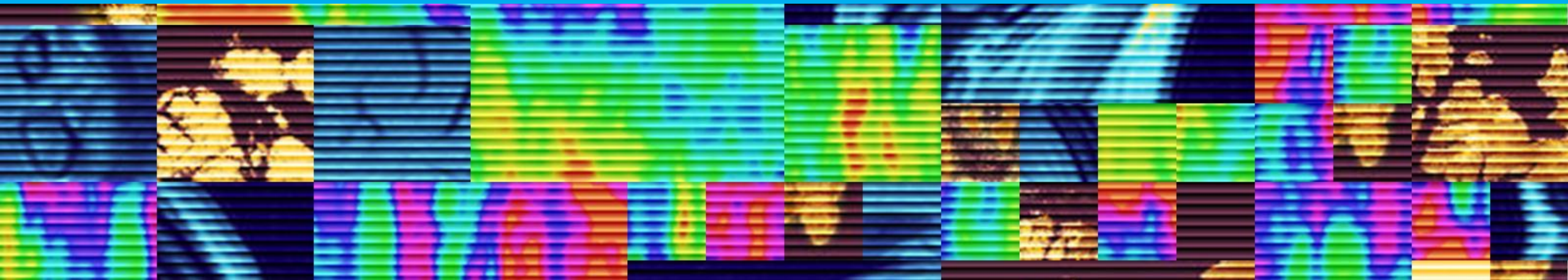




# 7.9 ZETTABYTES

2015前

来源: CenturyLink 资源中心, 根据readwriteweb, 11月17日, 2011报告







Bing 引入的数据 > 7 petabytes  
一个月

推特社区每天生成超过 1  
terabyte 的推特信息

Cisco 预测 2013 年每年互联网流  
量将会达到 667 exabytes

来源: 经济学人, 2010年2月; DBMS2; 微软公司

# The Economist

FEBRUARY 27TH - MARCH 5TH 2010

Economist.com

Obama the warrior

Misgoverning Argentina

The economic shift from West to East

Genetically modified crops blossom

The right to eat cats and dogs

## The data deluge

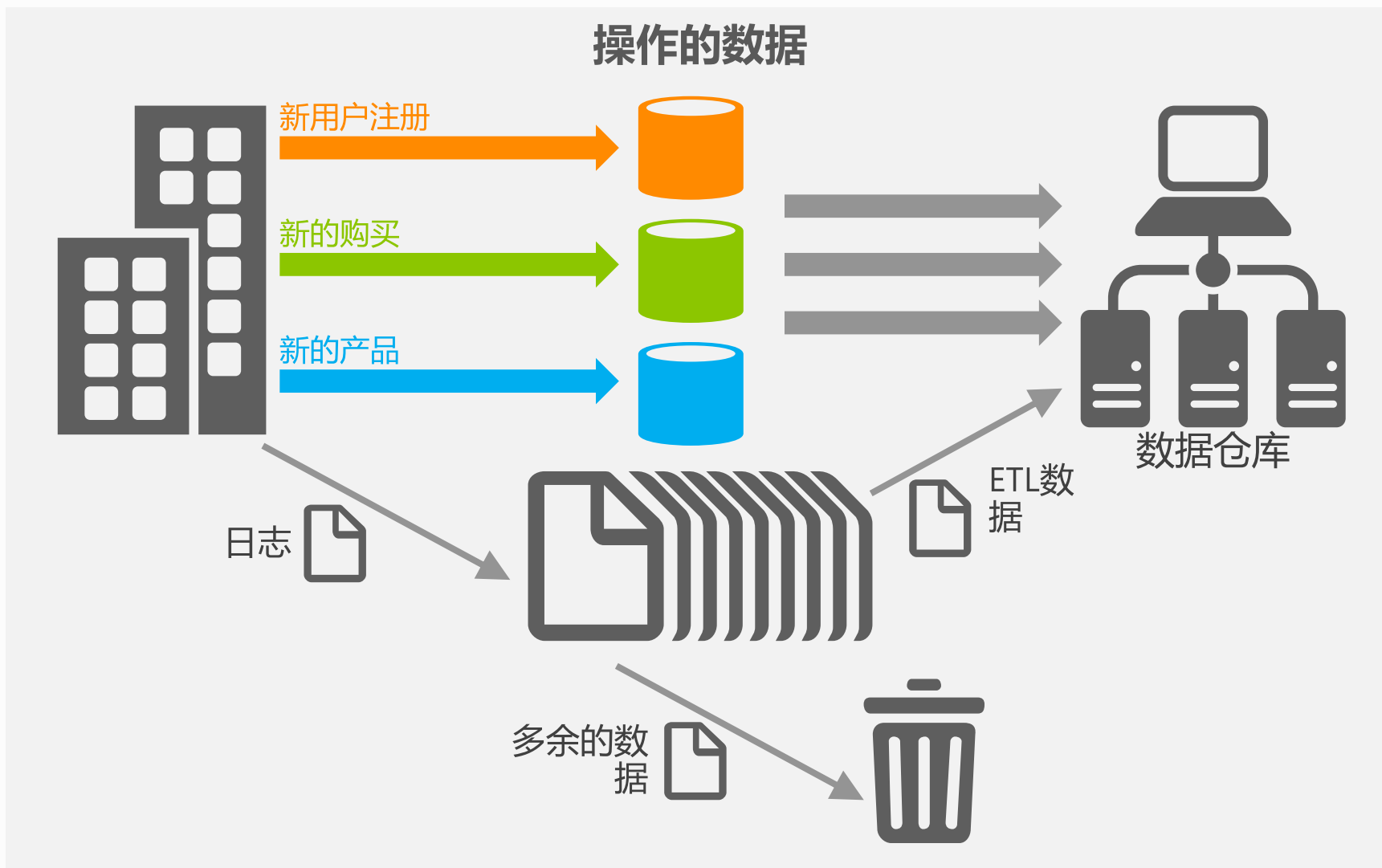
AND HOW TO HANDLE IT: A 14-PAGE SPECIAL REPORT



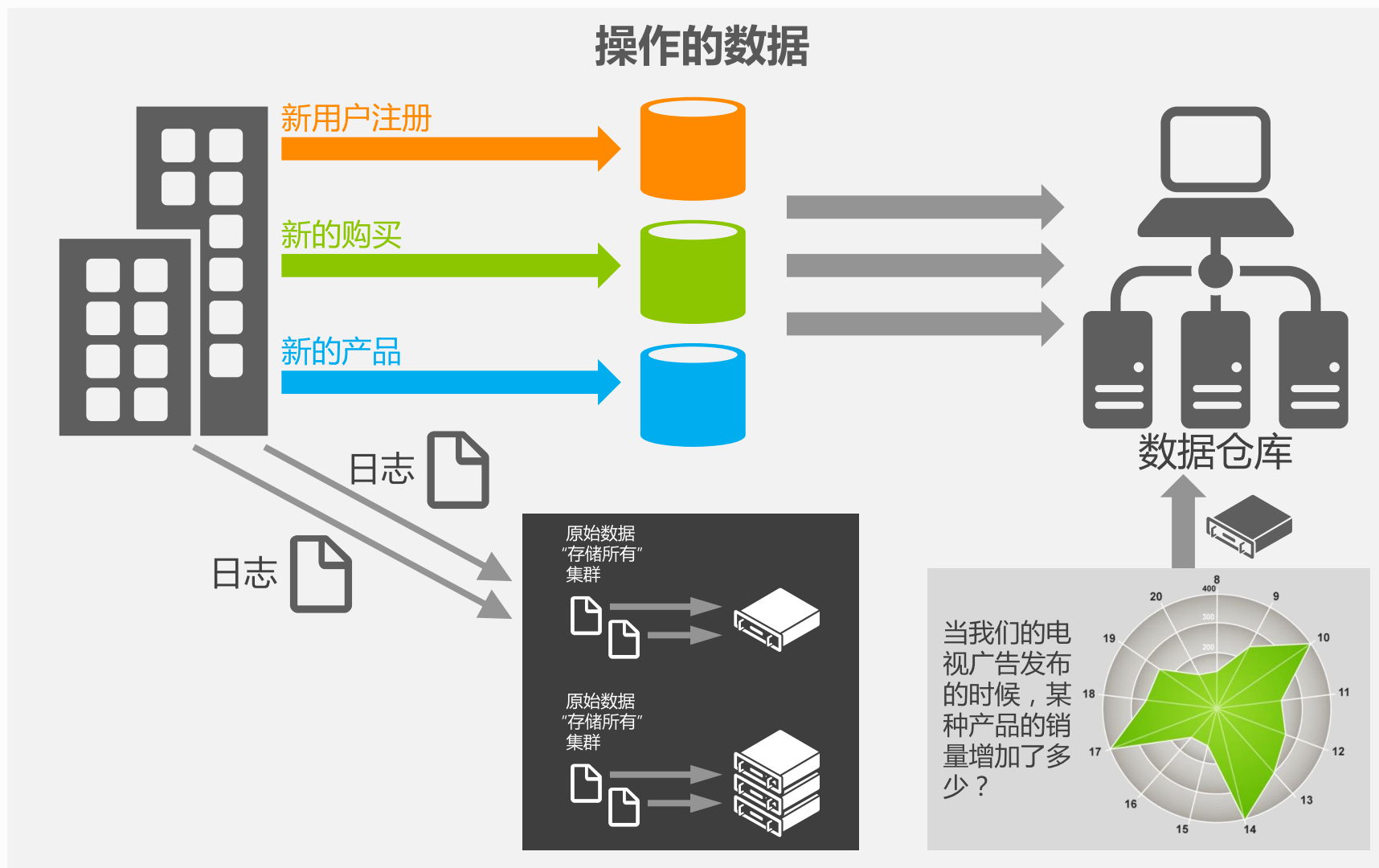
# 举例场景



# 传统电子商务数据流



# 新的电子商务数据流



# 理解基础

## 把计算移向数据





# 大数据的特点



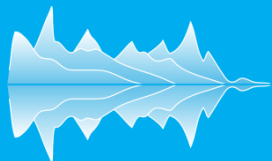
新的数据源



大数据量



新技术



非传统的数据类型



新经济

新的问题和新的视角

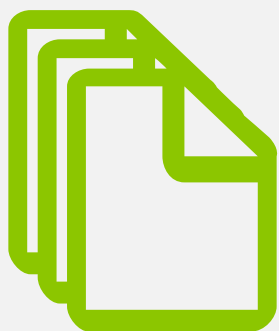


# MapReduce

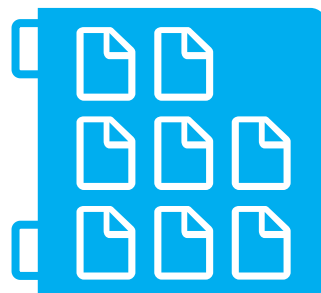
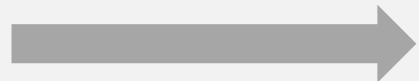


# 如何工作?

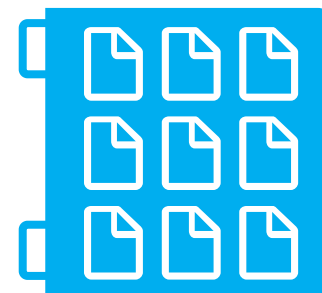
首先，存储数据



文件



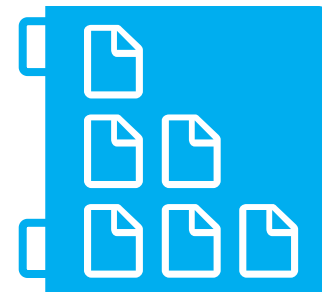
服务器



服务器



服务器



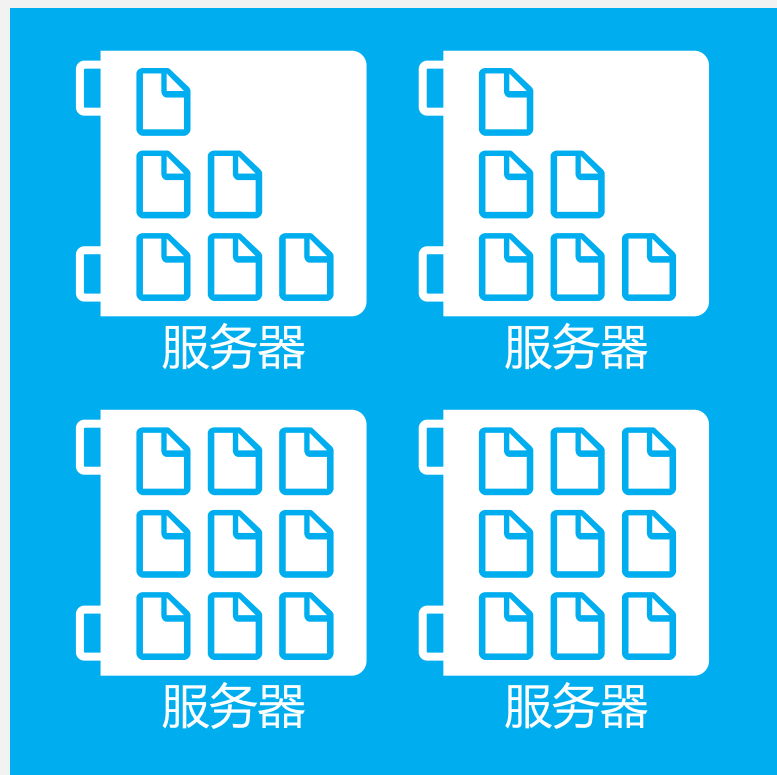
服务器



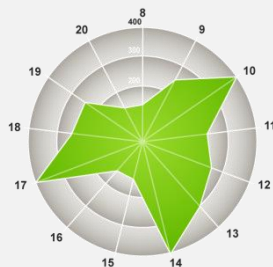
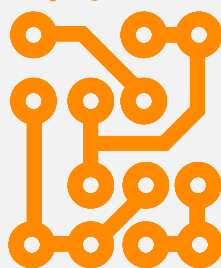
# 如何工作?

## 第二,处理数据

运行时



代码



```
// Map Reduce function in JavaScript
var map = function (key, value, context) {
  var words = value.split(/[a-zA-Z]/);
  for (var i = 0; i < words.length; i++) {
    if (words[i] !== "")
      context.write(words[i].toLowerCase(),
        1);
  }
};

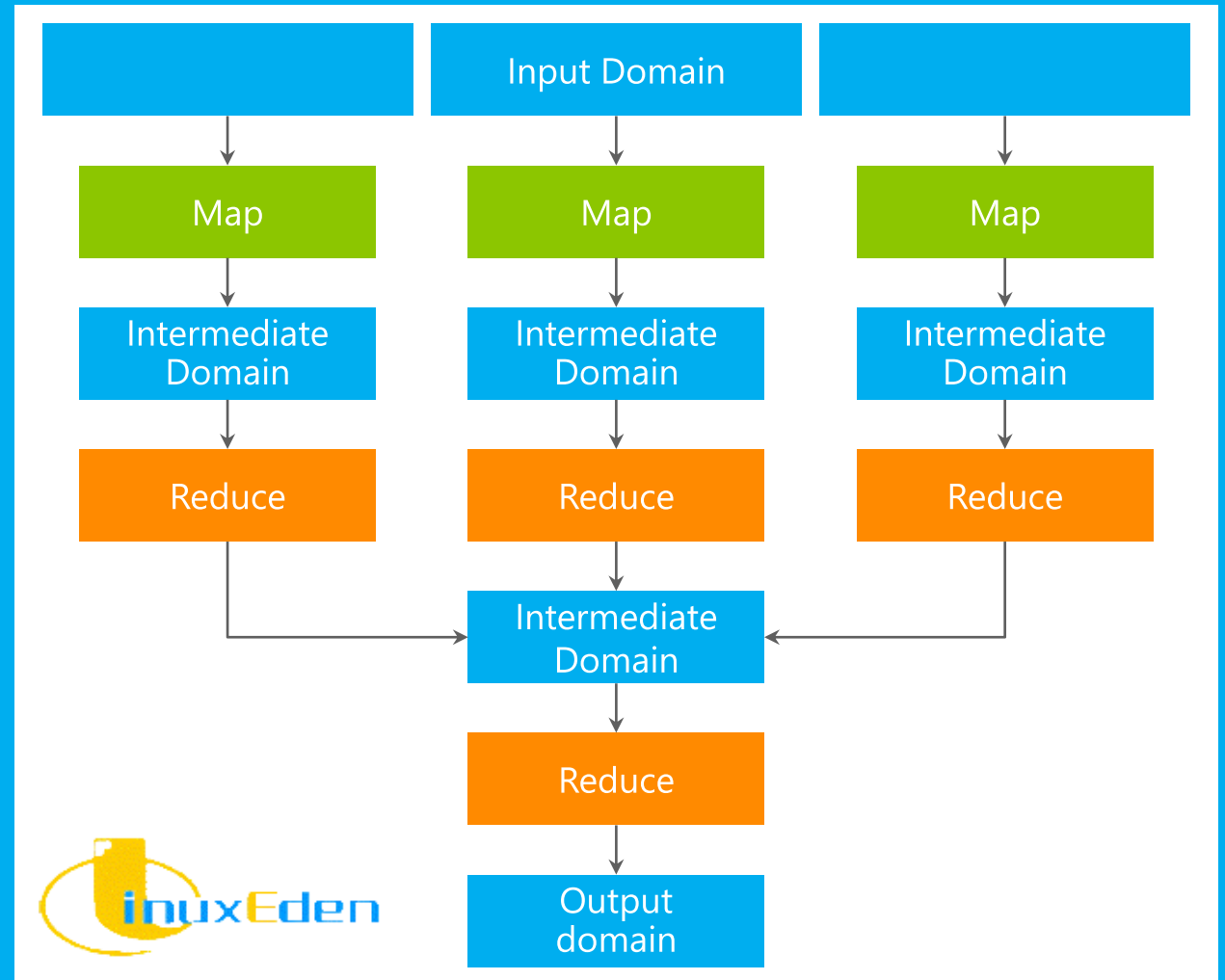
var reduce = function (key, values, context) {
  var sum = 0;
  while (values.hasNext()) {
    sum += parseInt(values.next());
  }
  context.write(key, sum);
};
```

# MapReduce – workflow

一个MapReduce 工作通常把输入数据集分割成独立的块，这些独立的块用map任务以完全并行方式进行处理

框架把“map”的输出进行排序，然后输入到“reduce”任务中

框架关心排程的任务，管理和重新执行失败的任务



# MapReduce – workflow

数据获取  
和建模

合作  
和可视化

分析  
和数据挖掘

传播，分享  
和保持

“It takes more time to hand a project from the seismic guys to me to the engineers in production than it does to figure out the oil field plays.”

**Geologist,**  
Major oil and gas company

“Our weather model and resulting data sets should be accessible to universities and other institutions.”

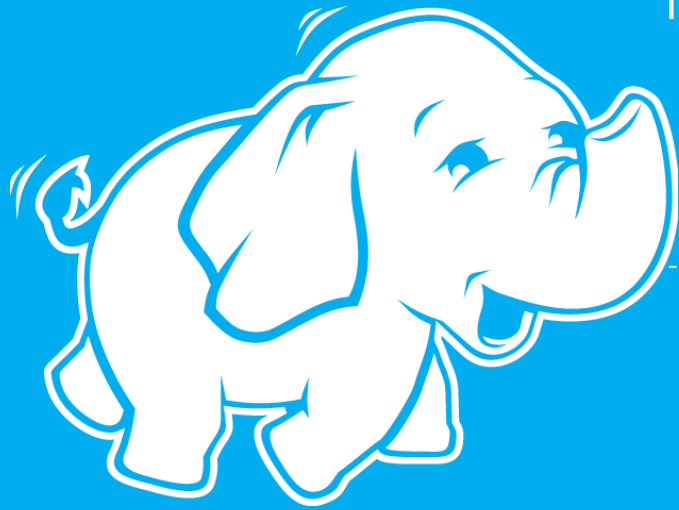
**Aerospace Development Manager,**  
U.S. Federal Government



# Hadoop

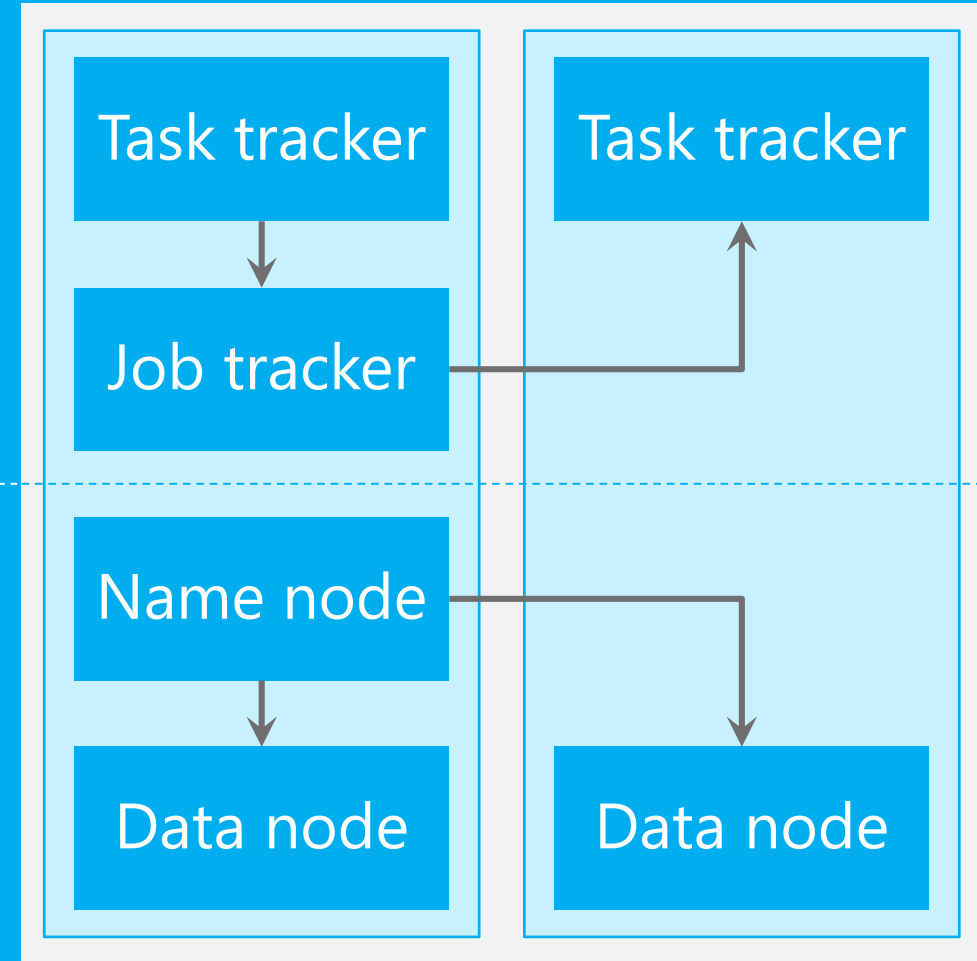


# Hadoop 架构



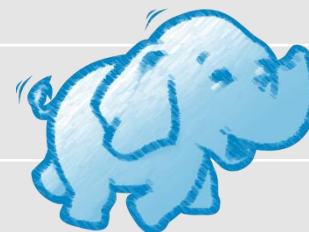
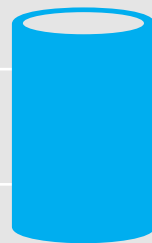
MapReduce  
层

HDFS  
层



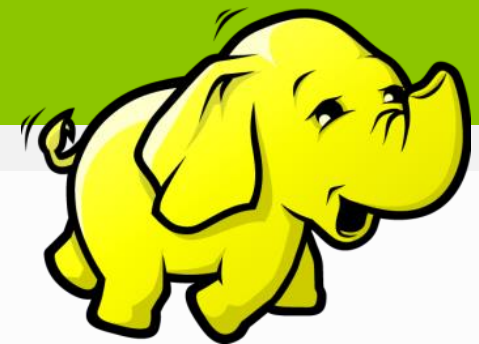
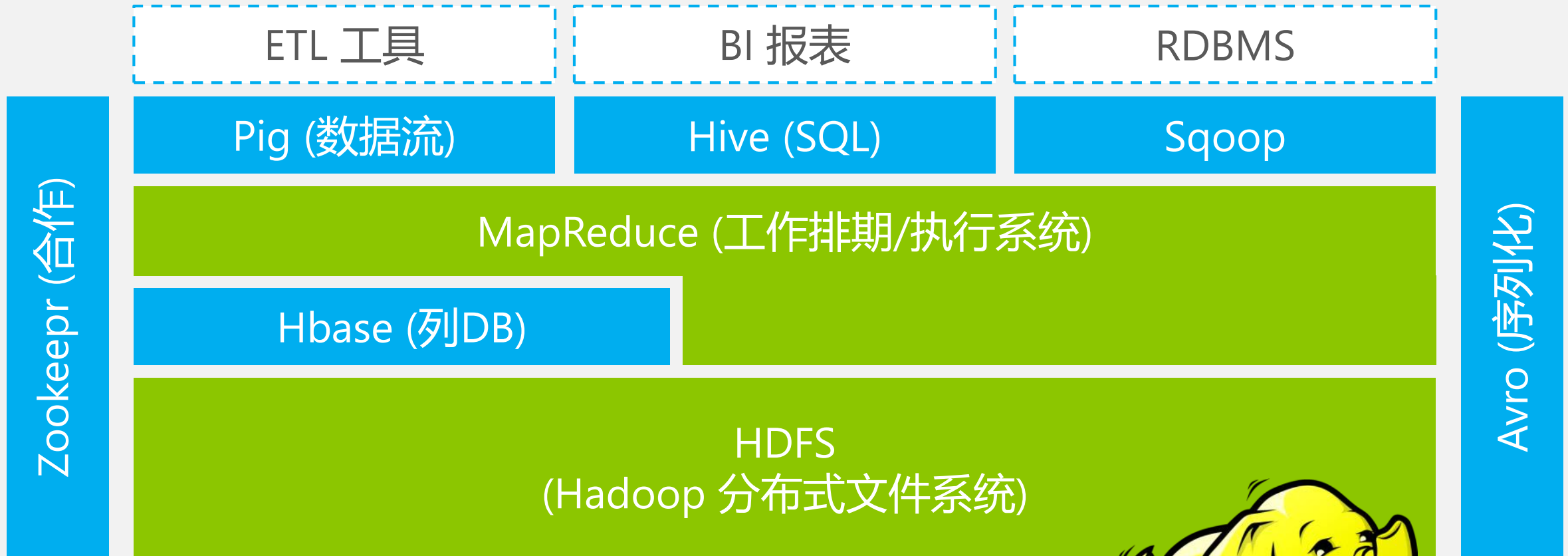
# 传统的RDBMS vs. MapReduce

	传统RDBMS	MAPREDUCE
数据量	Gigabytes ( <i>Terabytes</i> )	Petabytes ( <i>Hexabytes</i> )
访问	交互和批量	批量
更新	多次读写	一次写，多次读
结构	静态结构	动态结构
数据完整性	高(ACID)	低
规模扩展	非线性	线性
DBA 比例	1:40	1:3000



Reference: Tom White's Hadoop: The Definitive Guide

# The Hadoop 子系统



Reference: Tom White's Hadoop: The Definitive Guide



# 微软和Hadoop

# 具体提供



深刻见解

Hive ODBC Driver 和 Hive Add-in for Excel  
与Microsoft PowerPivot集成



企业就绪

Windows Server & Azure上基于Hadoop的分发  
与Hortonworks战略合作



更广泛的访问

支持Hadoop的JavaScript框架  
RTM of Hadoop connectors for SQL Server and PDW

# 微软大数据解决方案

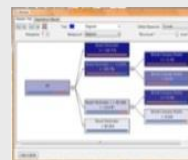
熟悉的最终用户工具



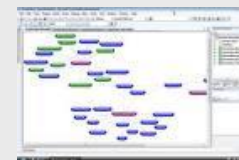
Power View



Excel和PowerPivot



预测分析



集成的BI

BI 平台



SSAS



SSRS



Hadoop On  
Windows Azure



Hadoop On  
Windows Server

连接器



Microsoft EDW

非结构化  
和结构化的数据



Sensors



Devices



Bots



Crawlers



ERP



CRM



LOB



APPs

# 在Azure上部署和集成 Hadoop集群

## 演示





# Windows上的Hadoop

## 所有用户都能对新的数据类型有深刻见解



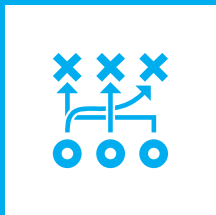
# 微软大数据路线图



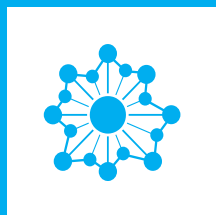
微软正在扩展他在商业智能和数据仓库的领导地位，把任何大小的数据的新的类型提供给所有用户



为了加速微软的基于Windows Server和Windows Azure的Hadoop解决方案，微软和Hortonworks成为合作伙伴



微软发布了大数据端对端的路线图，包含了在Apache Hadoop™企业级的，在Windows Server和Windows Azure上的基于Hadoop的解决方案



微软承诺让所有规模组织的最终用户，开发人员和IT人员访问和使用Hadoop。

# 资源

<http://www.hadooponazure.com/>

<http://hadoop.apache.org/>



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries.

The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft cannot guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

Translated to Chinese Simplified Version by Shanghai Yungoal Info Tech Co., Ltd. [YunGoal](#)