

# 联想大数据平台-LEAP

联想大数据

---



# + 目录

**01 联想大数据概述**

**02 联想大数据平台-LEAP**

**03 行业大数据解决方案**

# + 大数据时代企业面临的挑战

随着移动互联网飞速发展，催生“大数据时代”的到来。各种新型智能移动设备的迅速普及、工业互联网概念的兴起及其实践，带来了海量数据的爆炸式增长。



数据存储  
的碎片化



传统烟囱  
式IT架构



管理未充分融  
入企业运营中



缺乏统一  
运维管控



小机+数据库  
技术单一

大数据时代，企业如何管理海量大数据，如何分析应用大数据成为重要的机遇和挑战。

# + 数据资产价值不断增长



Facebook 2012/05/18 日上市，其公司公布的账面资产为66亿美元。

Facebook IPO定价为38美元/股，价值1040亿美元。

Gartner 研究表明：Facebook收集了2.1万亿条“获利信息”，每条信息约为4美分价值，即每个Facebook用户的价值为100美元。



数据的价值并不仅限于特定的用途，它可以为了同一目的而被多次使用，也可以用于其它目的。

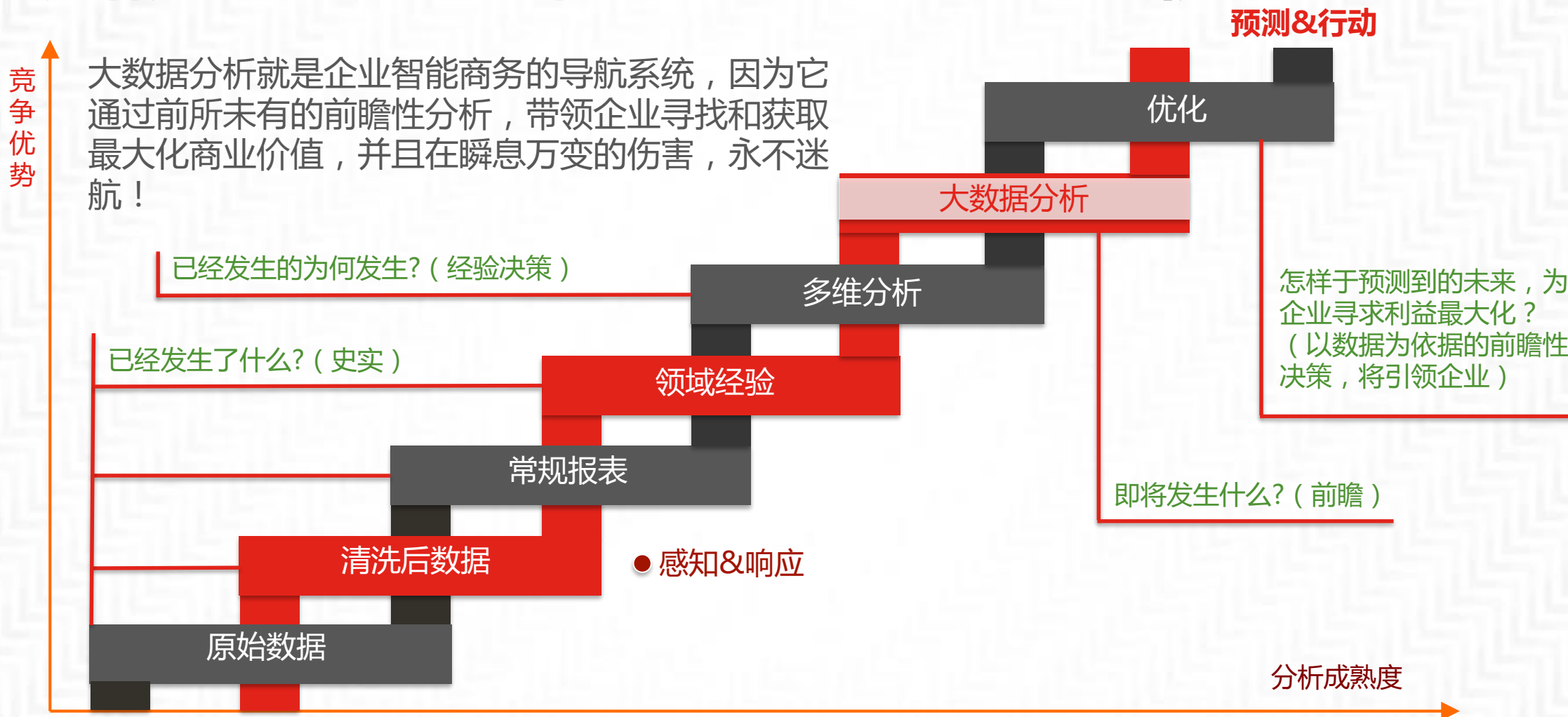
数据的价值并不仅限于特定的用途，它可以为了同一目的而被多次使用，也可以用于其它目的。

数据的真实价值就像漂浮在海洋中的冰山，第一眼只能看到冰山一角，而绝大部分则隐藏在表面之下。

数据就像一个神奇的钻石矿，是“取之不尽，用之不竭”的

- 通过账面资产来确定企业的价值的方式，已经无法充分反映公司的真实价值；
- 投资者开始注意到数据的潜在价值；
- 拥有数据或者具有数据收集和分析能力的企业，其公司价值会上升；
- 给数据的潜在价值贴上价格标签会给金融部门带来无限的商机。

# + 大数据分析应用的价值越来越受到企业重视



构建面向海量数据的管理与分析能力、实现数据的价值体现正逐渐成为提高企业竞争能力的核心要素之一。联想企业级数据分析平台正是处理企业级大数据场景的高性能一站式分析平台。

# + 早在2011年8月，联想就启动了大数据建设

技术融合  
深度优化



5年！  
300+研发人员持续投入

300名开发工程师  
60名运维工程师  
30名数据科学家



全球部署的  
超大规模集群

全球8个数据中心  
2000台服务器  
3000名操作用户



海量数据分析与  
持续性业务支撑

12Pb数据容量规模  
9Pb数据量  
150亿条记录 / 天  
30TB / 天

持续演进  
工匠情怀



在战斗中成长！

# + 联想大数据提供端到端的企业级大数据产品与服务

## 联想大数据

服务：

业务咨询服务

高级分析服务

IT/DT规划服务

数据管理服务

软件：

工具与应用软件

能力开放平台

大数据计算平台

硬件：

服务器、存储

技术支持、运维

形成端到端的整体解决方案，将处于技术底层的企业数据资产，通过软硬件平台和专业化服务，一步步转化为上层业务价值

### 当客户需要将数据资产转化为业务洞察和商业价值时

我们提供大数据分析及应用业务的咨询服务。我们拥有强大的数据科学家团队，已经支撑了联想全集团、全产品线、全业务流的大数据分析支撑，同时还在制造、零售、能源等行业的大数据分析领域针对客户痛点问题进行数学建模和业务分析，最终实现业务优化。

### 当客户难以管理自己多源、异构、海量的大数据资产时

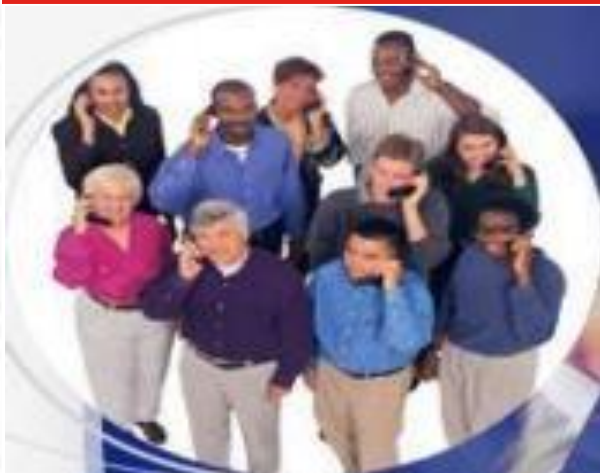
我们提供数据科学管理的咨询与实施。包含了数据质量管理、数据架构咨询、数据安全与隐私、元数据管理，通过一系列业务咨询、软件工具、和技术实施，确保客户数据资产的完整性、正确性、可用性、以及业务连续性。

### 当客户需要一个成熟的、高性能的大数据平台及解决方案时

我们提供联想企业级大数据分析平台（LEAP）产品以及运维服务。LEAP是软件化平台，整合了最先进的大数据开源技术，实现了基于IaaS层联想服务器面向大数据分析的深度优化，同时也包含了面向数据及分析能力开放的各类数据工具组件。

# + 联想大数据已具备多方面的能力

## 专业的分析能力



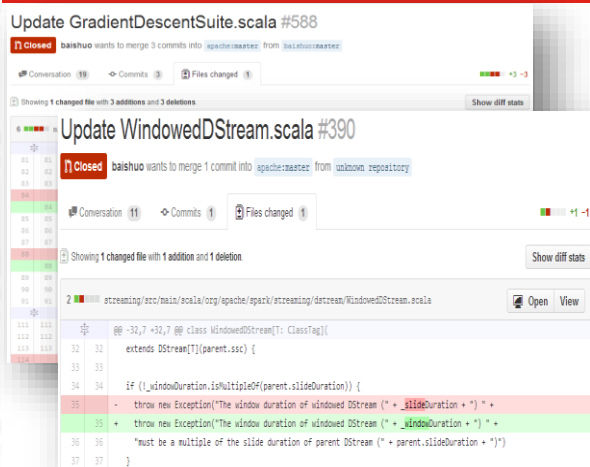
- 联想大数据专家、数据科学家、行业顾问50+余人
- 80%海外留学背景，博士占比近70%
- 在顶级期刊和会议中发表学术论文近百篇，获得国内和美国专利数十项

## 超大的部署规模



- 全球化多中心部署
- 2000台的集群规模
- 3000名操作用户
- 总容量12PB
- 数据总量9PB
- 日新增数据30TB
- 日处理数据4.3PB

## 强大的研发实力



- 300+研发人员的持续投入
- 多年来不断向开源社区做出技术贡献。
- 被北京市发改委于2011年评为“移动互联网系统软件及服务工程 北京市工程实验室”

## 完善的运维体系



- 全面完整的技术支持保障
- 贯通一致的针对产品功能的高质量技术支持
- 高级服务包含与售后及升级，迁移的咨询服务
- 提供完善的产品文档及自助服务选项



# + 目录

**01 联想大数据概述**

**02 联想大数据平台-LEAP**

**03 行业大数据解决方案**

# + 联想大数据LEAP核心解决的问题

LEAP，联想集团基于开源架构自研的业内领先平台级大数据产品，旨在为企业客户提供端到端的完整大数据解决方案，使企业能够快速构建强大的大数据平台，便捷的部署基于自身业务的大数据分析应用，通过挖掘大数据潜在价值打造企业面向未来的核心竞争力。



**数据**：整合各方数据，沉淀业务知识，LEAP将为客户提供丰富的数据接口与强大的数据资源整合能力。

**平台**：LEAP将为客户提供安全可靠的分式的大数据平台，解决了海量数据的计算，存储实时数据计算等问题。

**管理**：参与到企业运营的各环节，通过对业务数据的分析，发现各种规律趋势，为策略制定提供参考依据。

**运维**：LEAP提供集中的运维管控组件，实现从设备到服务的全方位监控、管理和扩展。

**价值**：联想提供端到端的大数据服务，旨在发现数据潜在价值，帮助客户通过大数据解决商业问题，与客户共同成长。

# + 联想大数据LEAP6大产品线全景图



数据分析应用套件 Nash



数据能力开放平台 Gauss  
*Big Data as a Service*

资源开放  
Riemann



分析武库  
Bayes



数据工厂  
Fourier



大数据计算平台 Descartes  
*大数据技术整合与深度优化*



数据采集转换套件 Euclid



数据资产管理平台

Euler



系统运维监控中心

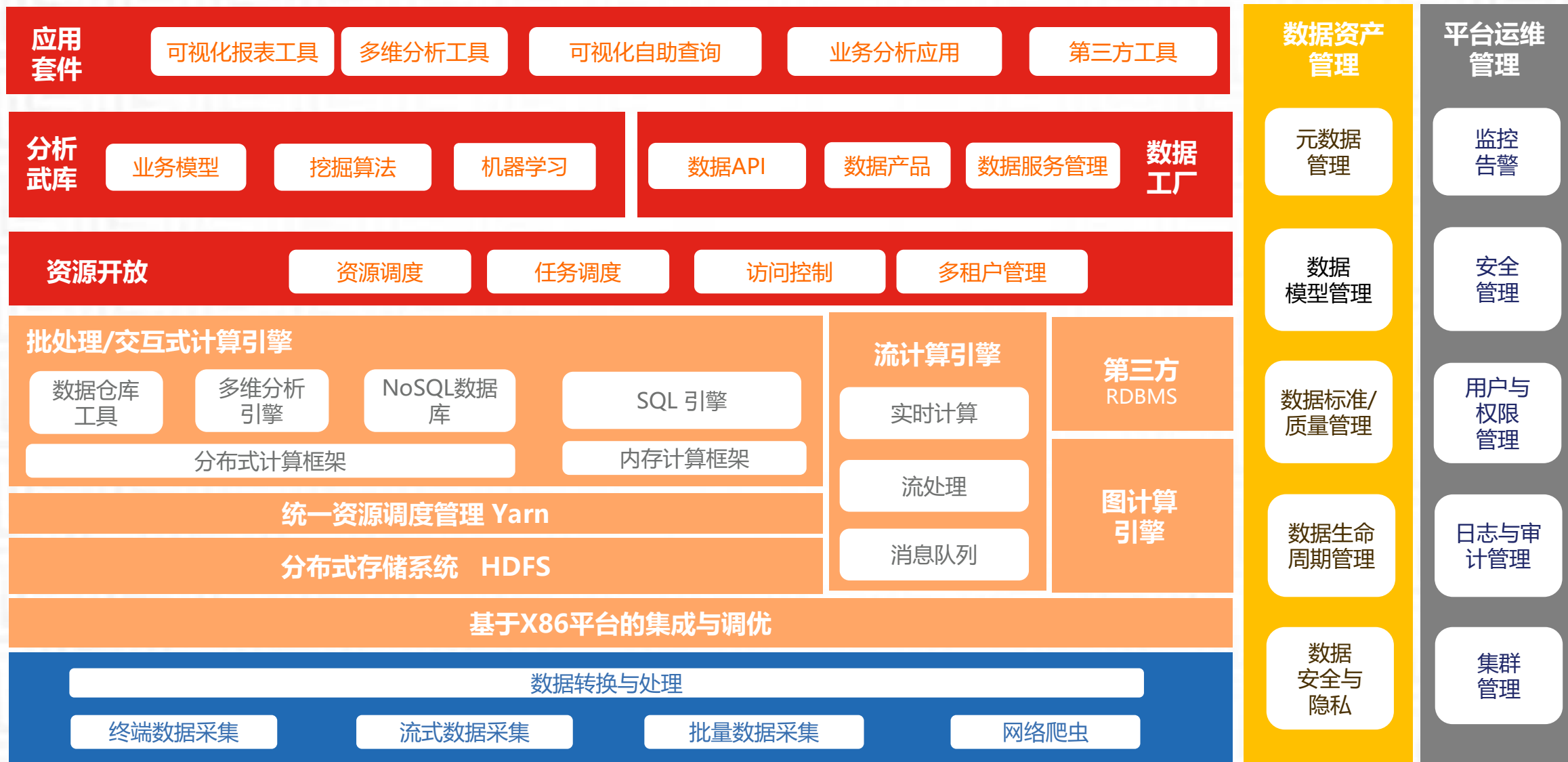
Architon

◆为企业级数据中心提供数据存储和数据处理能力，提供统一的集成平台环境，将硬件和平台软件做有效的集成。

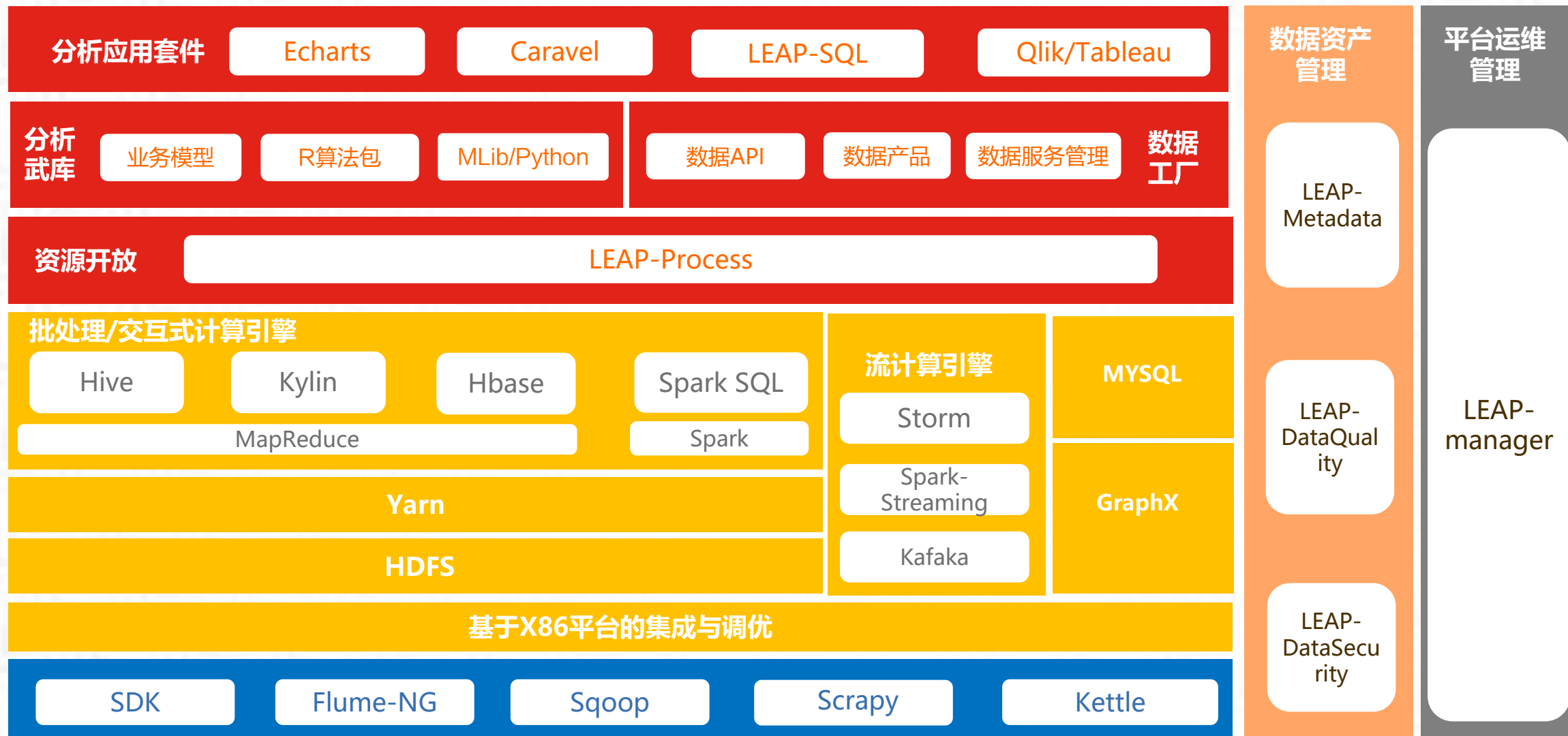
◆为企业级数据中心提供多域的数据模型、标准的元数据、数据处理调度任务、后台处理程序和前台应用程序，以及数据产品；

◆实现对环境中系统资源、软件资源、业务应用、参与人员等各种资源统一管理，综合监控。

# + 联想企业级大数据分析平台（LEAP）功能架构



# + 联想企业级大数据分析平台（LEAP）技术架构



# + 各组件清单列表 ( 1 )

组件	描述	版本
Apache Hadoop	稳定可靠并具有高可扩展性的分布式存储于技术级别架构	hadoop-2.6.0
HDFS	Hadoop分布式文件系统-可扩展、分布式的、高容错性和高吞吐量的数据存储	hdfs-2.6.0
MapReduce2(YARN)	新一代Apache Hadoop分布式计算框架, 引入统一资源管理器YARN	yarn+MR2-2.1.0
Hbase	分布式的、面向列的NoSQL数据库, 在Hadoop之上提供类似于Bigtable的能力	hbase-1.2.0
Hive	基于Hadoop的一个数据仓库工具, 可以将结构化的数据文件映射为一张数据库表, 并提供简单的SQL查询功能, 可以将SQL语句转换为MapReduce任务进行运行	hive-1.1.0
LEAP SQL	可以在浏览器端的Web控制台上与Hadoop集群进行交互来分析处理数据, 例如操作HDFS上的数据, 运行MapReduce Job等等	Priest -SQL-2.0.1
Sqoop	将Hadoop与关系型数据库集成的书籍传输引擎	sqoop-1.4.6
Flume	高可用的、高可靠的、分布式的海量日志采集、聚合和传输的系统, 并具有写到各种数据接受方(可定制)的能力	flume-ng-1.6.0
Kafka	高度可扩展的、容错的发布-订阅消息系统	kafka-2.0.0
Slider	基于Hadoop YARN的应用服务管理框架	slider-0.80.0

组件	描述	版本
spark	快速综合的数据处理引擎, 支持循环数据流的内存计算	spark-2.0.0
pig	处理Hadoop中存储数据的高级数据流语言	pig-0.12.0
Oozie	协调Hadoop活动的流程引擎	oozie-4.1.0
Zookeeper	高度可靠的分布式协同服务	zookeeper-3.4.5
Ambari	提供基于ambari平台的优化方案, 增加自动化监控运维的能力, 另外针对管理层次, 采用松耦合的方式与其他框架进行衔接	Ambari-2.4
Leap Process	仓库计算调度系统: 1、process平台提供了多种组件工具, 技术人员可以根据业务需要, 按流程方式组织计算任务, 设置调度方式, 报警方式。 2、可以新建计算任务流程、查询流程、停用、启用、导出流程。 3、流程定义完成后, 可以查询流程的运行情况, 并可进行补跑等操作。 4、可根据业务需求, 在此系统中管理监控规则, 如果执行情况满足条件, 则根据预警设置的方式, 以短信或邮件的形式进行通知。 5、可以在此系统中查看流程每天运行时间点分布图, 观察同一时间段流程运行密集度。	Leap-Process-3.0.1
storm	实时流数据处理引擎	storm-0.10.0
Solr	自由文本、模糊配对以及多面化搜索引擎	solr-4.10.3

## + 各组件清单列表（2）

组件	描述	版本
监报告警	监视各个模块的cpu,内存,网络等资源使用情况，当资源占用超过设定的阈值时，向系统管理员和业务负责人发出警报	Leap-3.0.1
用户管理	管理平台中的用户角色，配置各个用户的应用访问权限	Leap-Manager-3.0.1
集群管理	配置集群中各个模块的启动，停止，查看模块的状态，同时可以配置模块的具体参数。	Leap-Manager-3.0.1
自动化部署	根据业务的要求，部署计算存储模块，动态生成最优配置	Leap-Manager-3.0.1
访问控制	配置各个用户对数据访问的权限	Leap-Manager-3.0.1
性能管理	从业务的角度分析业务用户的响应时间，优化数据业务过程	Leap-Manager-3.0.1
容灾管理	当机器出现故障的时候，模块系统自动切换到备份的机器	Leap-Manager-3.0.1
智能运维	智能优化模块部署，配置参数，优化系统运行情况	Leap-Manager-3.0.1

# + 特性说明

01



○ 全源数据整合能力，快速汇聚各类数据

02



○ 高性能数据存储和计算平台，快速处理与分析

03



○ 一站式大数据开发应用环境，快速构建分析应用

04



○ 数据深度分析引擎，挖掘数据价值

05



○ 全球超大规模实践，验证系统高可靠性

06

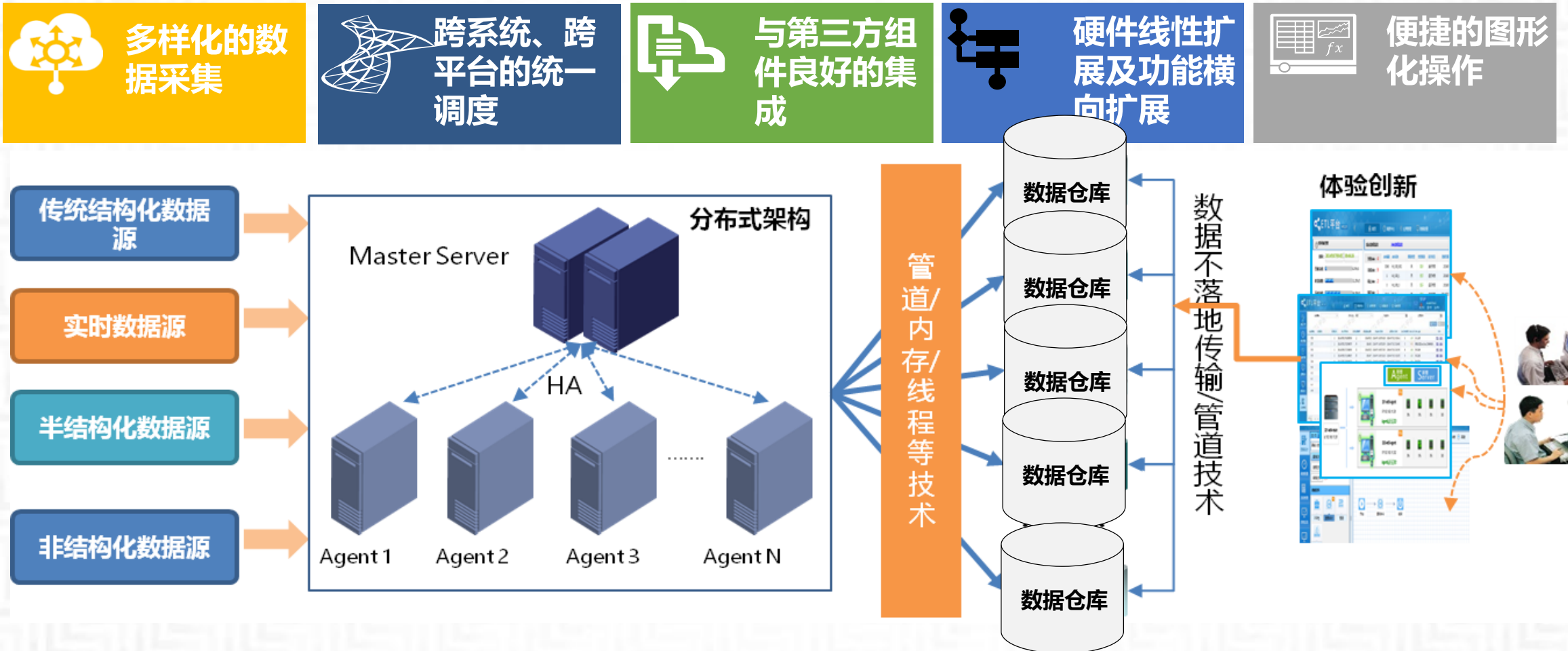


○ 自动化智能运维，易用易管理

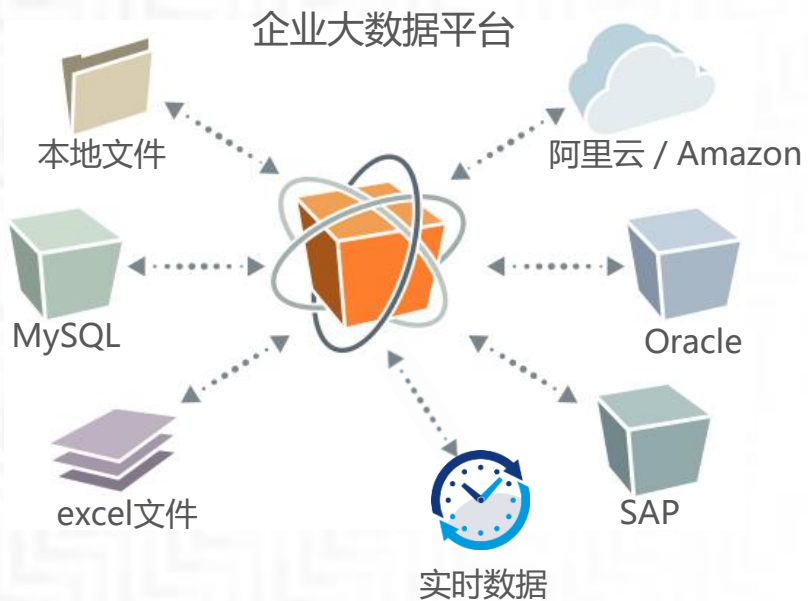


# + 特性1：全源数据整合能力，快速汇聚各类数据

通过统一的数据采集转换套件，可降低多套采集系统运维的复杂度，实现数据采集的统一调度。通过图形化的配置监控界面，快速完成各类采集规则的配置和发布。同时结合多样化的数据采集能力，对各类表、文件、消息等多种数据的实时采集和批量采集。

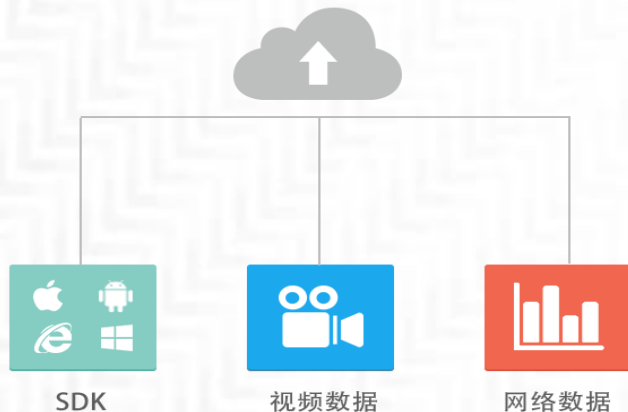


# + 多样化的数据接入采集能力



## 支持的数据源

IBM DB2	Oracle	Informix	MySQL	SAP ERP System	MS Access
Sybase	MS SQL Server	PostgreSQL	Teradata	SQLite	Hadoop Hive 2
AS/400	MaxDB	ExtenDB	OpenERP Server	H2	dBase 5
Apache Derby	MonetDB	Firebird SQL	Oracle RDB	Hadoop Hive	Ingres
Borlan Interbase	Native Mondrian	Generic database	Palo MOLAP Server	Hypersonic	Intersystems Cache
Calpont InfiniDB	Neoview	Greenplum	UniVerse database	Impala	KingbaseES
Exasol 4	Netezza	Gupta SQL Base	Vertica	Infobright	LucidDB
Aliyun	Amazon AWS	Google Cloud	Adobe Cloud	FTP	.....



- SDK

移动端：Android SDK, iOS SDK

Web端：Web SDK

PC端：Windows SDK

- 网络数据

通过网络爬虫来获取企业外部的相关数据

- 视频数据

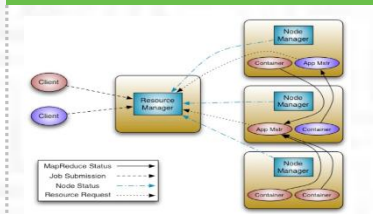
视频采集分析模块，为企业提供基于视频的数据采集



# 特性2：高性能数据存储和计算平台，快速处理与分析

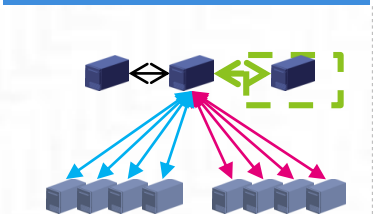
LEAP 平台集成了业界最先进的批量、流式、实时计算技术，采用灵活、高扩展性的数据处理架构，支持通用X86平台，同时面向联想x86 服务器进行了深度集成与优化，实现了超高性能的大数据分析技术平台。

## 设备资源统一管理



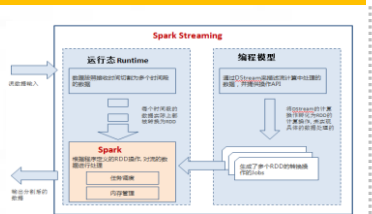
通过资源调度技术，将共享的各类计算资源按需动态分配给不同负载的应用。实现一个集群支撑多套同的应用运行，从而提升集群设备资源利用率。

## 无限的水平扩展



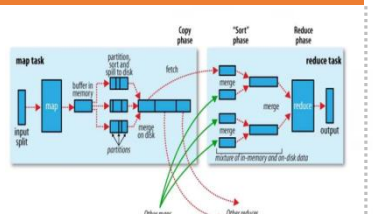
系统可线性扩充存储容量或提高处理性能，只需向向集群中增加机器，无需停机。有效解决企业由于数据增长导致的处理性能缓慢的问题。

## 实时数据处理



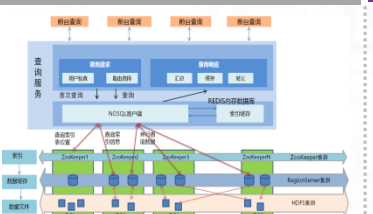
实时流处理引擎提供强大的流计算表达能力，可支持复杂的实时处理逻辑，满足企业实时告警、风险控制、在线统计和挖掘等应用需求。

## 批量数据处理



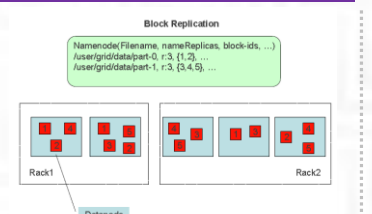
批量处理将海量结构化、半结构化、非结构化等多种类型数据进行批量计算和存储，实现对OLTP事务、全文搜索等统计业务的支撑。

## 海量数据快速查询



通过低成本的硬件提供高性能的数据加载、索引和查询能力，具备对海量数据（PB级）的存储，提供毫秒级的查询响应能力，从而提升客户体验。

## 海量数据存储



通过分布式文件系统，可将海量各类原始数据、结果数据进行快速存储。并通过自带副本机制，完成对数据的多份备份。

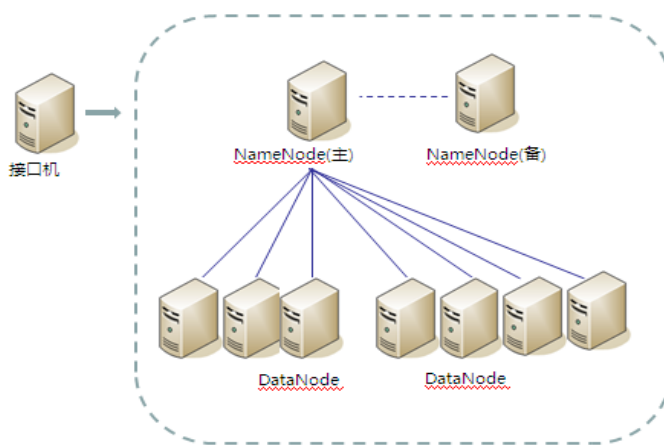
# + TPC-DS测试结果

通过上述的TPC-DS测试结果可以看到，联想的交互式分析引擎即便对于大数据集的查询也可以在**分钟级完成**。

TPC-DS SQL ID <sup>o</sup>	10GB 数据规模执行时间(秒) <sup>o</sup>	100GB 数据规模执行时间(秒) <sup>o</sup>	1TB 数据规模执行时间(秒) <sup>o</sup>
1 <sup>o</sup>	19.729 <sup>o</sup>	29.091 <sup>o</sup>	20.625 <sup>o</sup>
2 <sup>o</sup>	3.46 <sup>o</sup>	3.738 <sup>o</sup>	23.407 <sup>o</sup>
3 <sup>o</sup>	1.491 <sup>o</sup>	1.678 <sup>o</sup>	5.844 <sup>o</sup>
4 <sup>o</sup>	15.648 <sup>o</sup>	31.756 <sup>o</sup>	221.313 <sup>o</sup>
5 <sup>o</sup>	5.743 <sup>o</sup>	14.776 <sup>o</sup>	71.749 <sup>o</sup>
6 <sup>o</sup>	2.364 <sup>o</sup>	4.325 <sup>o</sup>	12.428 <sup>o</sup>
7 <sup>o</sup>	4.427 <sup>o</sup>	4.961 <sup>o</sup>	77.647 <sup>o</sup>
8 <sup>o</sup>	3.394 <sup>o</sup>	5.139 <sup>o</sup>	24.835 <sup>o</sup>
9 <sup>o</sup>	9.846 <sup>o</sup>	22.183 <sup>o</sup>	113.896 <sup>o</sup>
10 <sup>o</sup>	10.153 <sup>o</sup>	19.854 <sup>o</sup>	160.013 <sup>o</sup>
11 <sup>o</sup>	4.885 <sup>o</sup>	5.802 <sup>o</sup>	19.2 <sup>o</sup>
12 <sup>o</sup>	2.691 <sup>o</sup>	6.199 <sup>o</sup>	15.235 <sup>o</sup>
13 <sup>o</sup>	6.242 <sup>o</sup>	19.472 <sup>o</sup>	250.68 <sup>o</sup>
14 <sup>o</sup>	2.712 <sup>o</sup>	5.122 <sup>o</sup>	13.626 <sup>o</sup>
15 <sup>o</sup>	13.465 <sup>o</sup>	16.445 <sup>o</sup>	18.646 <sup>o</sup>

SQL特征	查询数量
子表达式	31
关联的子查询	15
不相互关联的子查询	76
Group by	78
Order by	64
Rollup	9
Partition	11
Exists	5
Union	17
Intersect	2
Minus	1
Case	24
Having	5

测试结果表明：SQL兼容性方面，联想大数据平台已兼容TPC-DS 99中的**89条**，处较高的水平(Spark 1.6.1支持40条，Hive 1.2.1支持45条)；



测试环境名称	SQL兼容性测试
操作系统	CentOS 6.5
Linux Kernel版本号	2.6.32-431
HADOOP软件厂商	联想LEAP
硬件设备配置	2CPU/128GRAM/17TBHD /万兆网卡
设备数量	10台

# + 特性3：一站式的大数据开发应用环境，快速构建分析应用

通过DataHub, 简化原有数据库系统迁移到大数据平台

通过Matrix, 进行元数据和数据字典管理

通过LEAP- SQL, 实现几秒钟从巨量的数据仓库中找到需要的分析结果

通过LEAP- Process, 执行复杂的数据计算流程, 生成定制化的报表

通过可视化工具, 构建快速的商业智能分析平台, 构建酷炫的可视化能力

通过Sense, 无需了解复杂的分析算法, 即可进行深度的数据分析和挖掘

## 一站式的大数据开发环境

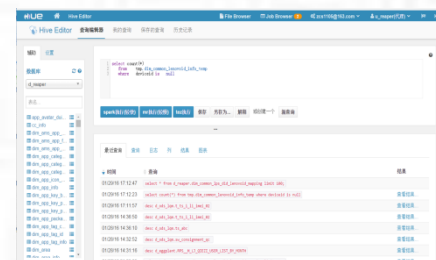
- Matrix**: The metadata system, holds all table profiles.
- Priest SQL**: The Big feed in Lenovo.
- Priest Process**: The graphic data process flows. Defined once and executed periodically.
- Arthas**: A fast query system, run the typed sql and get the results in A short time.
- Athena**: A visualization tool for query data.
- R**: A free software environment for statistical computing and graphics.
- RealtimeLog**: Dump realtime logs collected by the given device.
- Monitor (Constructing)**: A scalable distributed monitoring system for high-performance computing systems and common services.
- Data Sync (Constructing)**: Enable data sync across multiple IDC.
- Cloudera Manager**: The industry's trusted tool for managing Hadoop in production.



LEAP- Process



数据字典处理



LEAP- SQL

# + 示例：高效的图形化数据处理 workflow 配置

- ◆ 傻瓜式分布式计算开发；
- ◆ 可视化流程定制；
- ◆ 抽象多种大数据计算组件；
- ◆ 计算任务监控、预警；

计算任务管理

**组件集**

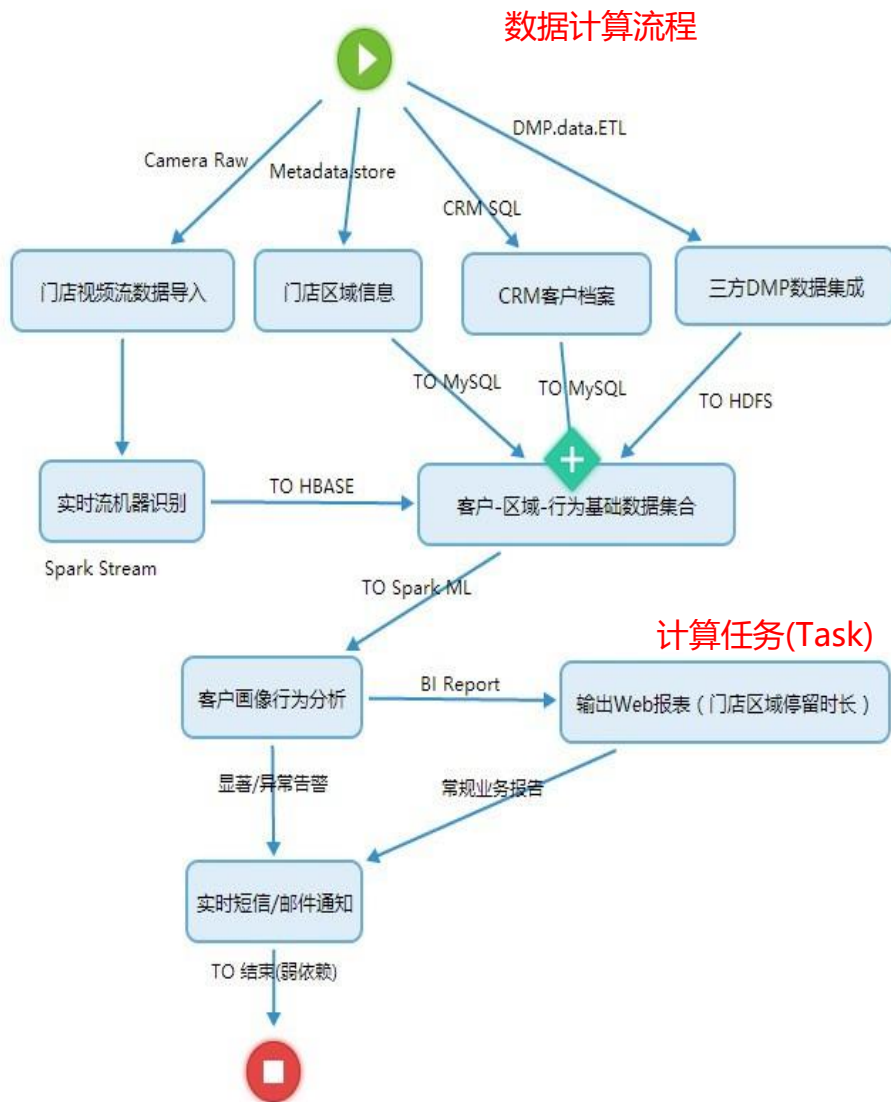
- 保存
- 选择
- 线条

**支持4类流程节点**

- 开始
- 结束
- 分支
- 合并

**支持11种计算任务**

- 导数
- Hdfs
- M/R
- Hive
- SparkJar
- SparkSql
- Mysql
- Oracle
- Shell
- Java
- Cache
- 依赖
- 通知



**属性**

流程ID:

流程名称: 门店区域停留时长

描述: 挖掘客户在门店关键区域的

有效期: 2018-05-31

优先级: 高优先级

应用: 请选择

调度类型: 天

调度日期:

调度时间:

最晚时间:

失败告警: 短信告警

告警对象:

流程变量:

应用变量:

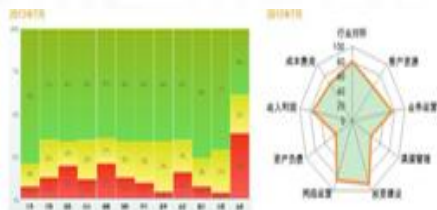
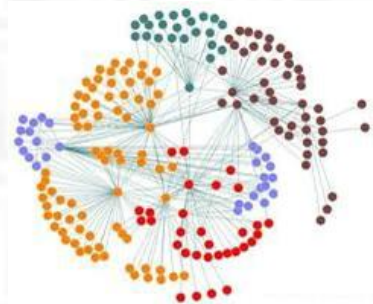
定义计算流程属性

流程调度参数设置

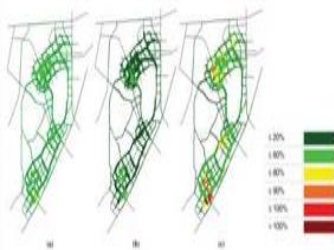
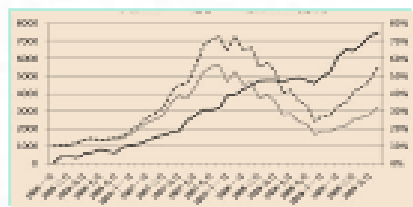
流程故障报警

runtime环境变量

# + 特性4：深度数据分析引擎，挖掘数据价值



ARIMA模型



## 模型分析

### 回归

- 简单线性回归 ( simple linear regression )
- 广义线性回归 ( generalized linear regression )
- 回归树 ( regression tree )
- 神经网络 ( neural network )
- 支持向量机 ( SVM )
- 随机森林 ( random forest )

### 分类

- 二进制响应的逻辑回归 ( logistic regression for binary response )
- 线性判别分析 ( linear discriminant analysis )
- KNN ( k nearest neighbors )
- 分类树 ( classification tree )
- 神经网络 ( neural network )
- 支持向量机 ( SVM )
- 随机森林 ( random forest )
- 朴素贝叶斯分类器 ( Naive Bayes classifier )

### 聚类

- k-means
- 层次聚类 ( hierarchical clustering )
- 基于密度的聚类 ( density-based clustering )

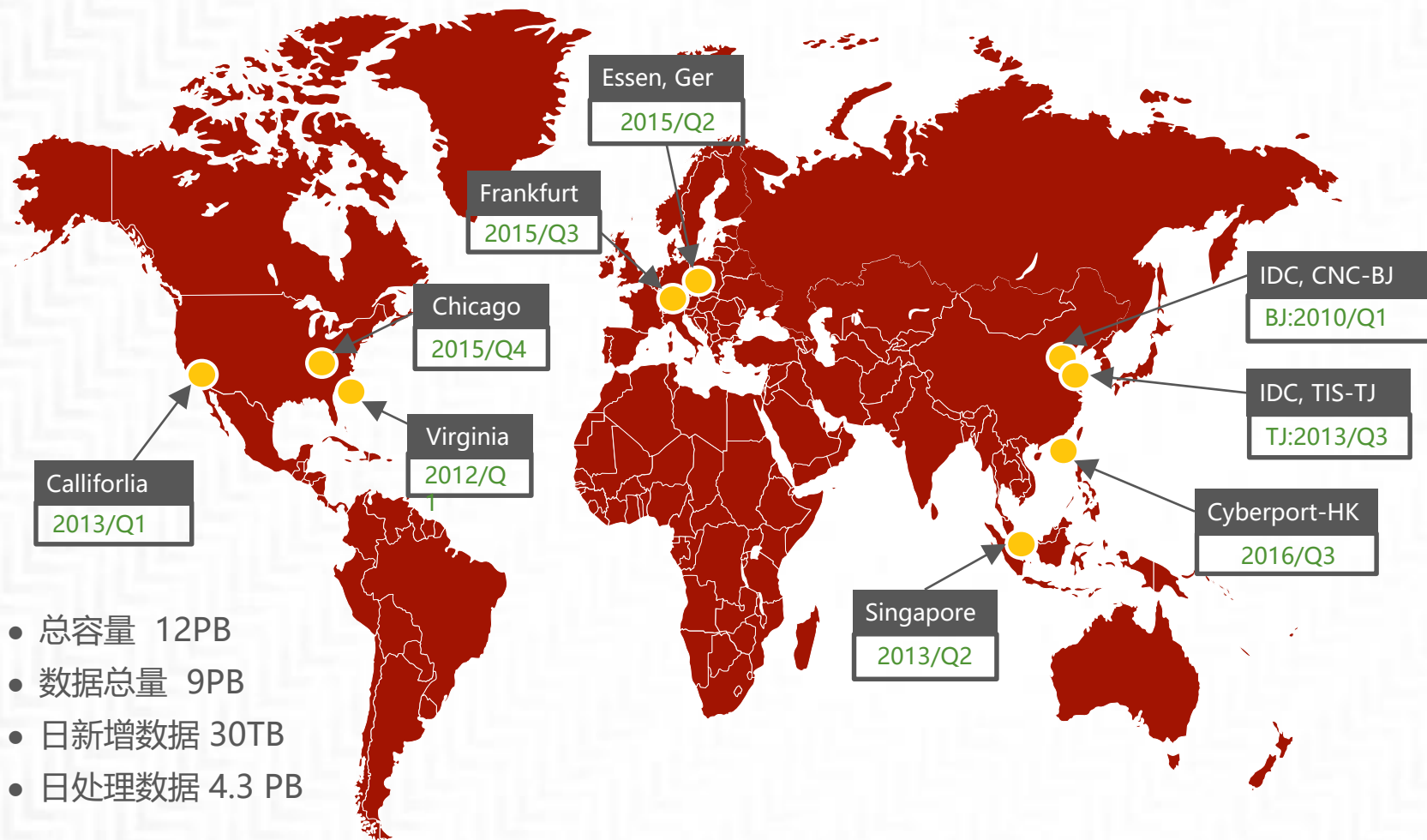
### 关联

- 关联规则 ( apriori algorithm )

### 时间序列

- 自回归积分滑动平均模型 ( ARIMA )
  - 移动平均模型 ( MA model )
  - 自回归模型 ( AR model )
  - 自回归移动平均模型 ( ARMA model )
  - 自回归积分滑动平均模型 ( ARIMA model )
- 广义自回归条件异方差模型 ( GARCH )

# + 特性5：全球超大规模实践，验证系统高可靠性



- 总容量 12PB
- 数据总量 9PB
- 日新增数据 30TB
- 日处理数据 4.3 PB

- 全球化多中心部署，2000台服务器，3000名操作用户
- 在实践中充分验证系统的高可靠性
- 企业数据本地化收集和存储
- 完全合规各国数据保护和隐私保护法律



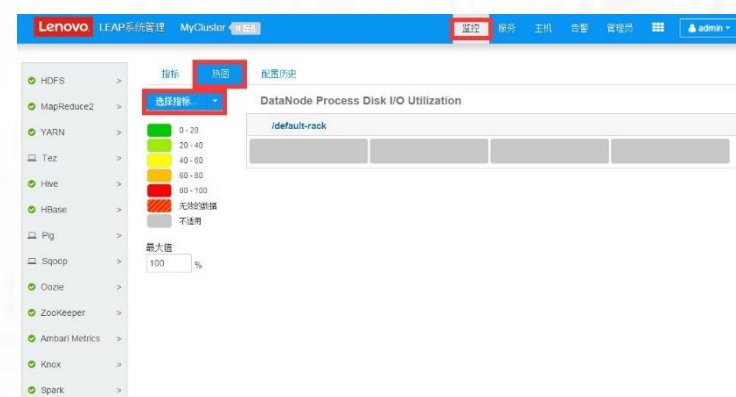
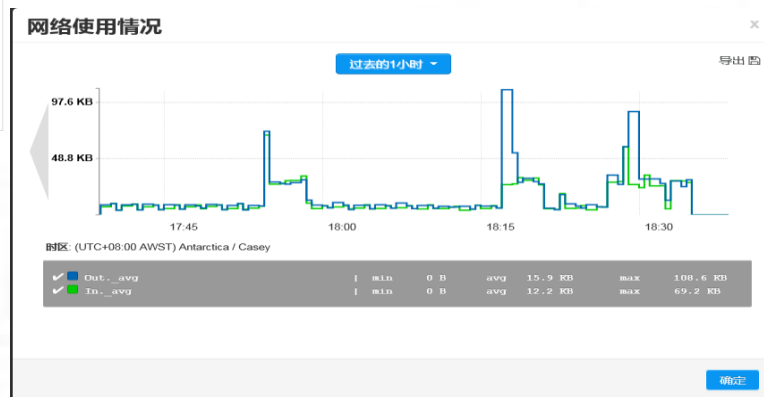
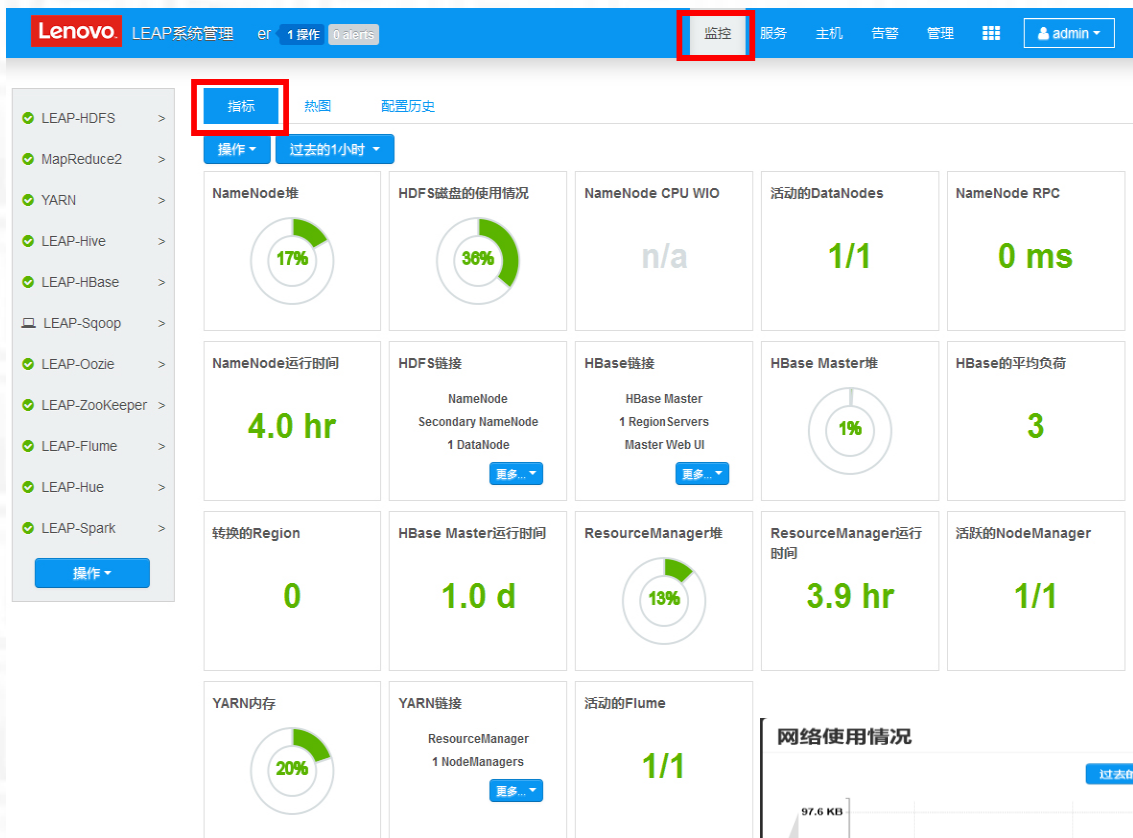
## + 特性6：自动化智能运维，易用易管理

提供全面丰富的平台运维功能，各类性能指标和功能故障监控功能，结合运维知识库，利用自动化策略，实现平台的智能运维。



# + 各类监控管理---指标监控

平台提供集中监控管理平台，丰富的监控管理指标体系使使用者可以直观便捷的掌握平台信息。



# + 自动化部署

支持Web图像化界面和快速向导，帮助用户短时间内部署一个或者多个集群，大大减轻工作量。

集群安装向导

- 开始安装
- 选择版本
- 安装选项**
- 确认主机
- 选择服务
- 指派Masters
- 指派Slaves and Clients
- 定制服务
- 检查
- 安装，启动和测试
- 概要

集群安装向导

- 开始安装
- 选择版本
- 安装选项
- 确认主机
- 选择服务**
- 指派Masters
- 指派Slaves and Clients
- 定制服务
- 检查
- 安装，启动和测试
- 概要

## 安装选项

输入要包含在集群中的主机列表，并提供你的SSH密钥。

### 目标主机

输入使用完全限定域名（FQDN），每行一个主机列表。或者使用模式表达式

host names

### 主机注册信息

提供您的 SSH私钥 自动注册主机

未选择任何文件

ssh private key

SSH用户帐户

root

执行 手动注册 在主机上 并且不使用SSH

[← 返回](#)

[注册和确认 →](#)

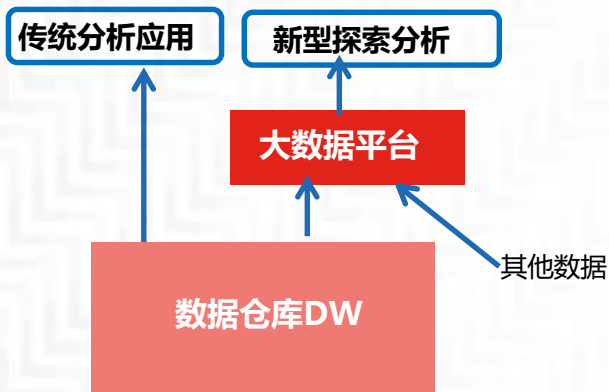
## 自动部署服务

选择您想要安装集群上哪些服务。

<input checked="" type="checkbox"/> 服务	版本	描述
<input checked="" type="checkbox"/> HDFS	2.7.1.2.3	Apache Hadoop 分布式文件系统
<input checked="" type="checkbox"/> YARN + MapReduce2	2.7.1.2.3	Apache Hadoop NextGen MapReduce (YARN)
<input checked="" type="checkbox"/> Tez	0.7.0.2.3	Tez是一个在YARN顶级上的下一代Hadoop查询处理框架
<input checked="" type="checkbox"/> Hive	1.2.1.2.3	数据仓库系统即席查询、大型数据集和表分析及存储管理服务
<input checked="" type="checkbox"/> HBase	1.1.1.2.3	非关系型分布式数据库
<input checked="" type="checkbox"/> Pig	0.15.0.2.3	分析大型数据集的脚本平台
<input checked="" type="checkbox"/> Sqoop	1.4.6.2.3	一个用于在Apache Hadoop和结构化数据存储批量转换数据的工具，比如相关的数据库之间

# + 大数据平台的部署形态

## 专业探索分析集市



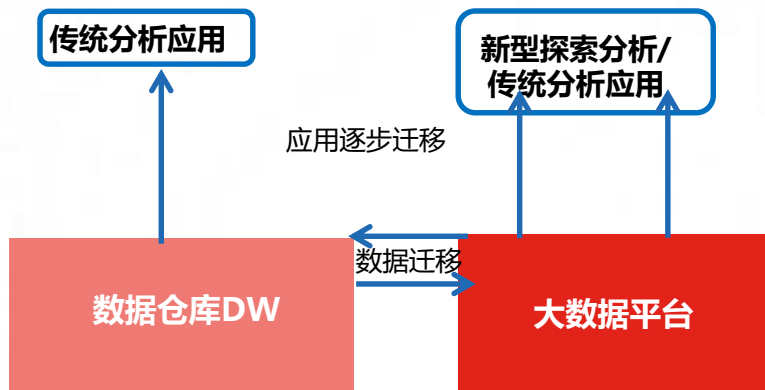
### 大数据平台定位：

- 是数据仓库平台的一个补充系统，主要面向新型数据和部分仓库数据的存储和处理，通过数据挖掘算法等，发现隐性的数据规律和价值。

### 特点：

- 小：系统规模小，使用人员少（以专业研究分析人员为主）
- 快：针对特定专题快速分析；支持实时处理和分析
- 灵：专用平台，灵活响应和尝试
- 深：专业深挖，挖掘算法、模式分析、图分析、文本分析等

## 混搭双中心



### 大数据平台定位：

- 是数据仓库平台的重要并列系统，
- 分担DW系统的存储和计算压力，提高处理效率、降低成本
- 传统应用逐渐迁移到大数据平台
- 通过数据挖掘算法等，发现隐性的数据规律和价值。

### 特点：

- 地位重要，承载的作用更大
- 支持新型分析方法和传统应用
- 系统可靠性、支撑能力要求更高
- 数据仓库DW的重要性下降

## 企业级大数据中心



### 大数据平台定位：

- 企业级大数据中心，采集全企业层面的各类内部数据及相关外部数据，并对这些结构化/非结构化海量数据进行整合、加工、处理，完成信息的深加工，逐步形成数据资产，为公司进行企业决策管理和生产一线的营销服务等工作提供完整、及时、准确、科学的信息支撑。

### 特点：

- 一个中心承载各类数据，进行各类分析应用，服务企业内外部各类用户
- 系统可靠性、系统稳定性、系统开放性、支撑能力要求很高

# + 联想与其它HDADOOOP厂商的测试数据对比



## 基础功能

Hadoop平台功能性最为完善，**超过**其他友商等



## 查询性能

HBase查询性能最好，在清单查询等典型场景中优势明显



## 数据处理性能

Hadoop的数据处理性能最优，与客户的典型应用场景最匹配



## 复杂计算

面向未来的复杂计算条件下



## 抗压能力

在计算资源紧张的大压力条件下

### 联想VS友商

附表.某制造行业客户组织Hadoop平台各厂商测试详细结果

厂商	功能得分	Hbase查询	两文件关联	五文件关联	汇总	压力
		主索引查询	26G gz 包含踢重、上传	26G gz 包含上传	47G gz 包含上传	大文件加载
联想	86	24.9"	28'11	15'04	50'37	58'56
XXX	81	39"	47'41	31'28	58'49	69'
XXX	85	29"	32'	17'12	61'04	61'06

# + LEAP同业对标

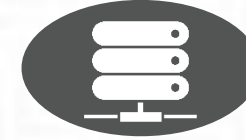
## LEAP



## 某产品化方案



## 开源方案



### 可靠性

统一监控，故障及时拉起  
全部组件支持HA  
支持热切换，对业务与数据影响为零

统一监控，故障及时拉起  
部分组件支持HA  
支持热切换，运行任务需中断重跑

原生态分布式，被动拉起  
部分组件支持HA  
不支持热切换

### 安全性

统一集中的用户账户管理  
限制匿名访问  
HDFS数据加密存储

集中用户账户管理  
限制匿名访问  
HDFS数据加密存储

分散用户账户管理  
原生无限制匿名访问  
无知识库

### 易用性

基于知识库的开发助手  
丰富的开发组件库  
基于机器学习的平台知识库

开发文档  
部分代码样例  
无知识库

开发文档  
无开发代码或样例  
无知识库

### 易管理性

导向式部署，一键式升级  
图形化管理与监控，提供优化建议  
丰富的调优参数配置

集成部署环境，手工升级  
图形化管理界面  
提供调优参数

手工部署与升级  
无图形化管理界面  
人工配置调优

### 可持续性

300+核心研发团队全球技术共享  
根植Apache社区，回馈社区  
组件迭代周期平均1个月

100+研发团队  
封闭内核  
组件迭代周期平均3个月

无固定人员投入  
基于Apache社区  
无固定组件迭代周期

# + 目录

**01 联想大数据概述**

**02 联想大数据平台-LEAP**

**03 行业大数据解决方案**

# + 联想大数据在行业中的应用



## 制造业

现代化的制造生产线安装有数以千计的小型传感器，来探测温度、压力、热能、振动和噪声。利用大数据将所有信息收集分析，对生产计划进行实时调整与预测。



## 政府及公共事业

政府通过运用大数据，将多渠道的数据采集和快速综合的处理分析，可提升治理社会的能力，实现政府公共服务的管理创新和服务模式创新。



## 交通运输业

交通运输每时每刻都产生大量的数据。利用大数据的海量处理和分析能力，为交通运输行业的决策和服务带来新的解决思路。



## 能源行业

通过运用大数据，及时呈现能源使用情况，发现低效与浪费，并提出优化建议，帮助合理规划能源的生产和使用，从而为全社会的绿色节能开拓新领域。



## 零售业

零售业通过整合各类消费数据，利用大数据分析各类顾客的消费行为，商品间的销售关联性，以快速精确支持营销。



## 医疗行业

通过大数据分析的公共卫生数据，提高疾病预报和预警能力。整合基础健康相关数据，提高危机探测能力。



## 金融行业

运用大数据技术对交易信息、调查报告、业绩报告及消费者研究信息等数据的分析，可实现精准的营销、风险的管控、精细化的管理。同时带来金融服务和产品的创新。

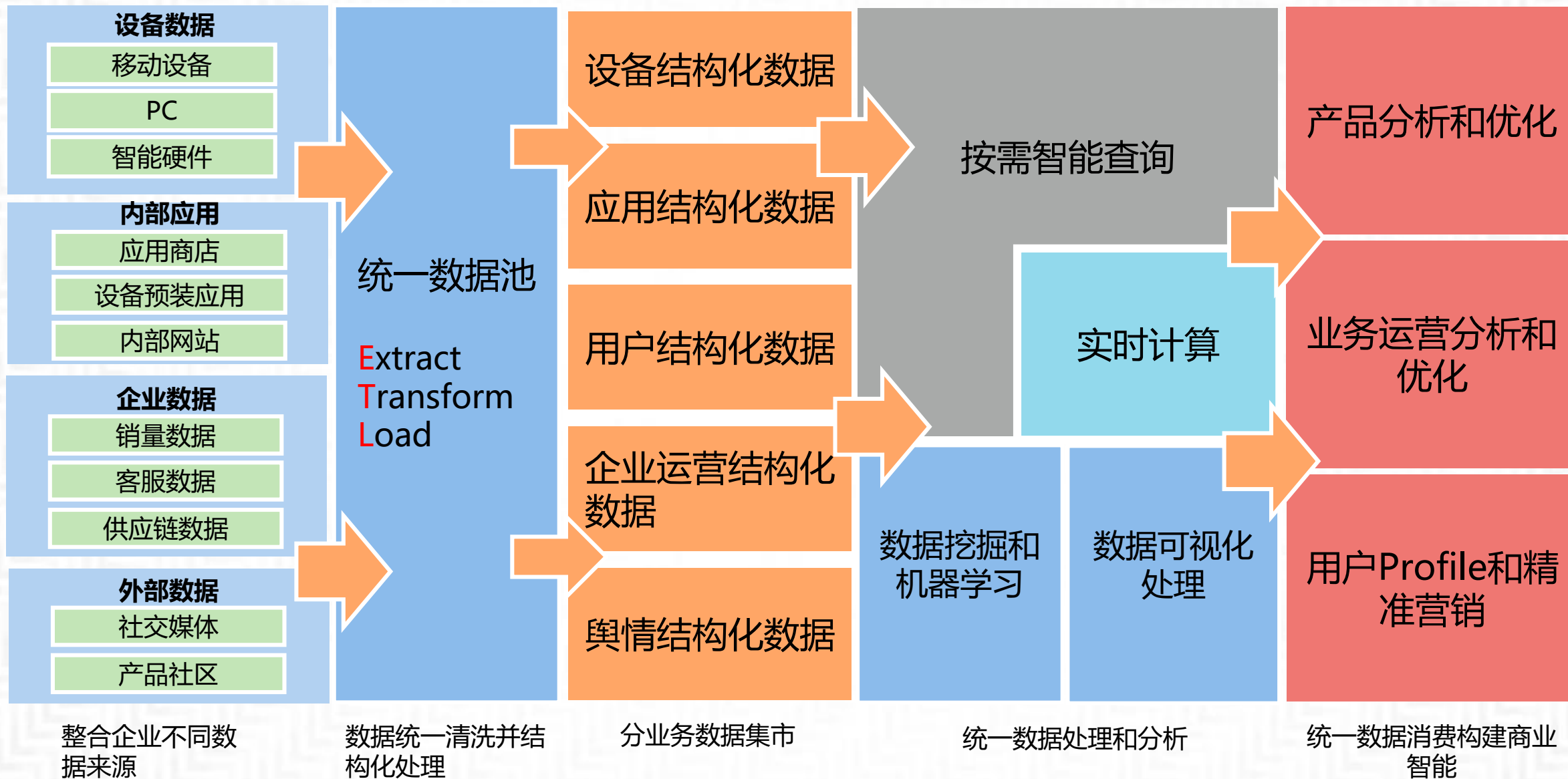


## 通用行业

通过大数据实现对各类网络海量信息的爬取、分析及管理。随着供应链越来越复杂，通过大数据分析可为供应链提供精准的需求预测、敏捷的资源获取、优化的库存等多方面能力。



# + 联想内部大数据分析平台



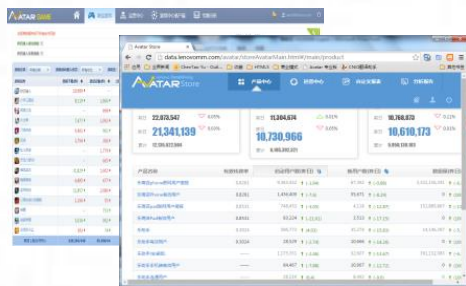
# + 联想大数据的业务分析应用示例

## 设备分析 ( Avatar Device )



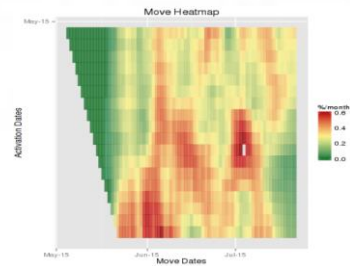
设备激活, 用户行为跟踪, 系统优化, 设备画像

## 应用分析 ( Avatar apps )



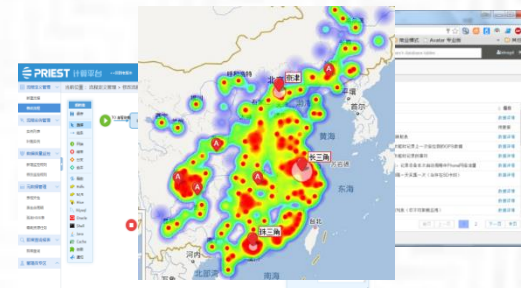
应用数据分析&服务经营

## 手机质量管理 ( MQM )



设备质量预警, 备件预测

## 店面热力图



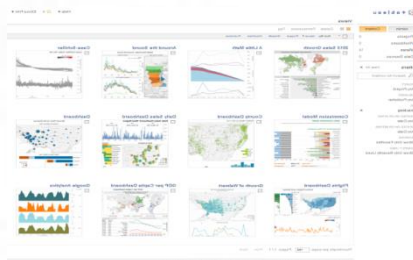
基于店面camera采集数据, 分析店面布局、客流量、客户对产品的关注度

## 舆情分析



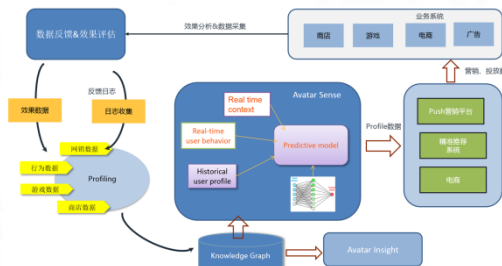
市场矩阵, 渠道优化, 公众情绪  
从各个评论网站获取手机的评价  
正向评价/反向评价

## 服务效率(Service)



客服中心优化, 部件预测优化, 新产品  
引进优化

## 手机市场分析(Accelerator)



从手机出货到渠道-》渠道到客户-》  
客户开机激活的监控

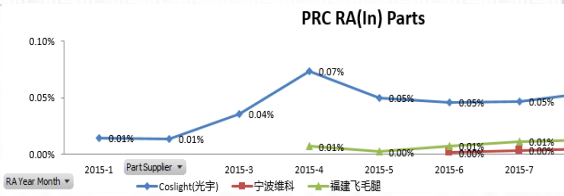
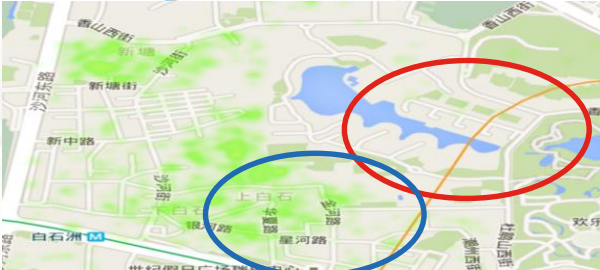
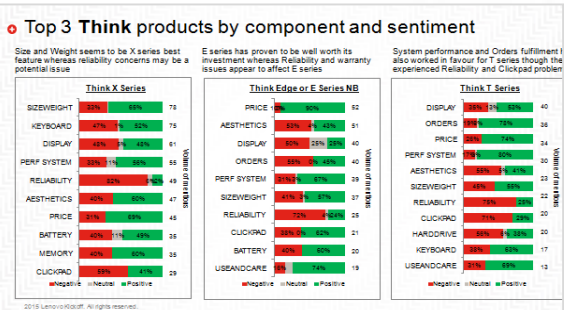
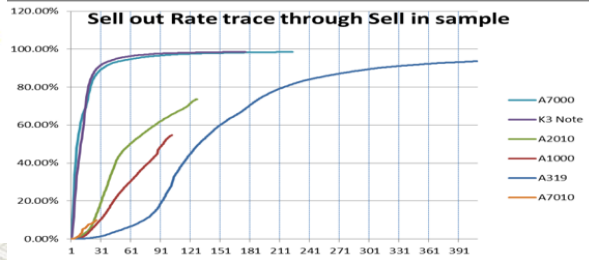
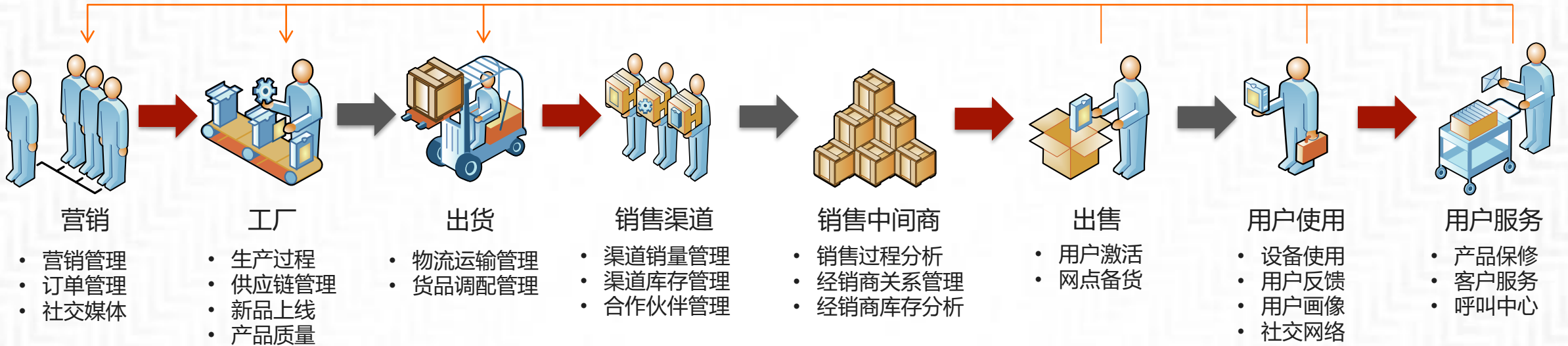
## 用户洞察



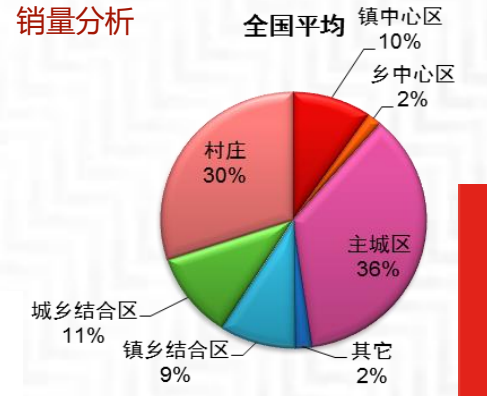
用户画像, 设备画像, 联想分析, 用  
户反馈

# 联想大数据支持全集团业务的生命周期管理优化

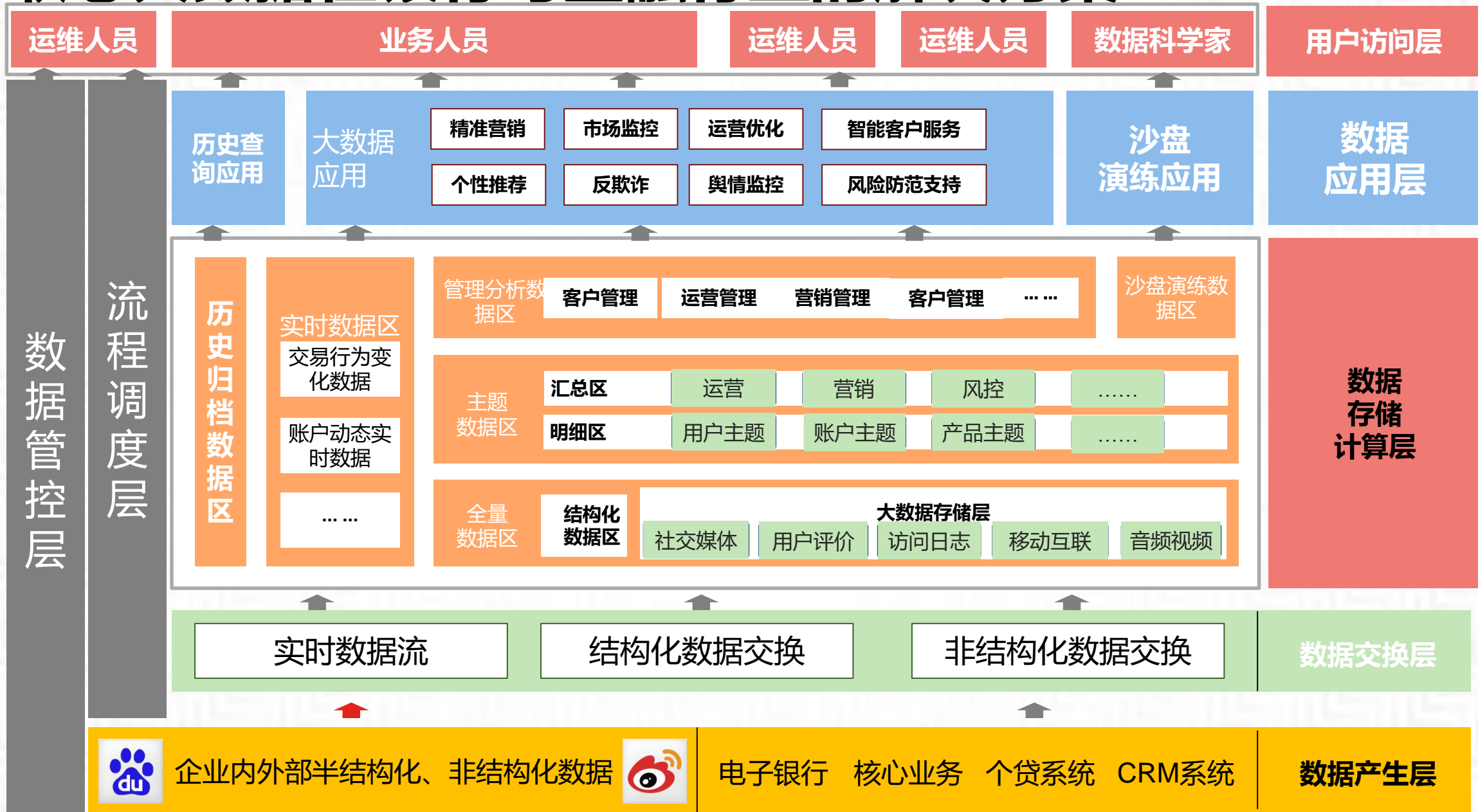
## 全生命周期的数据分析和产品管理



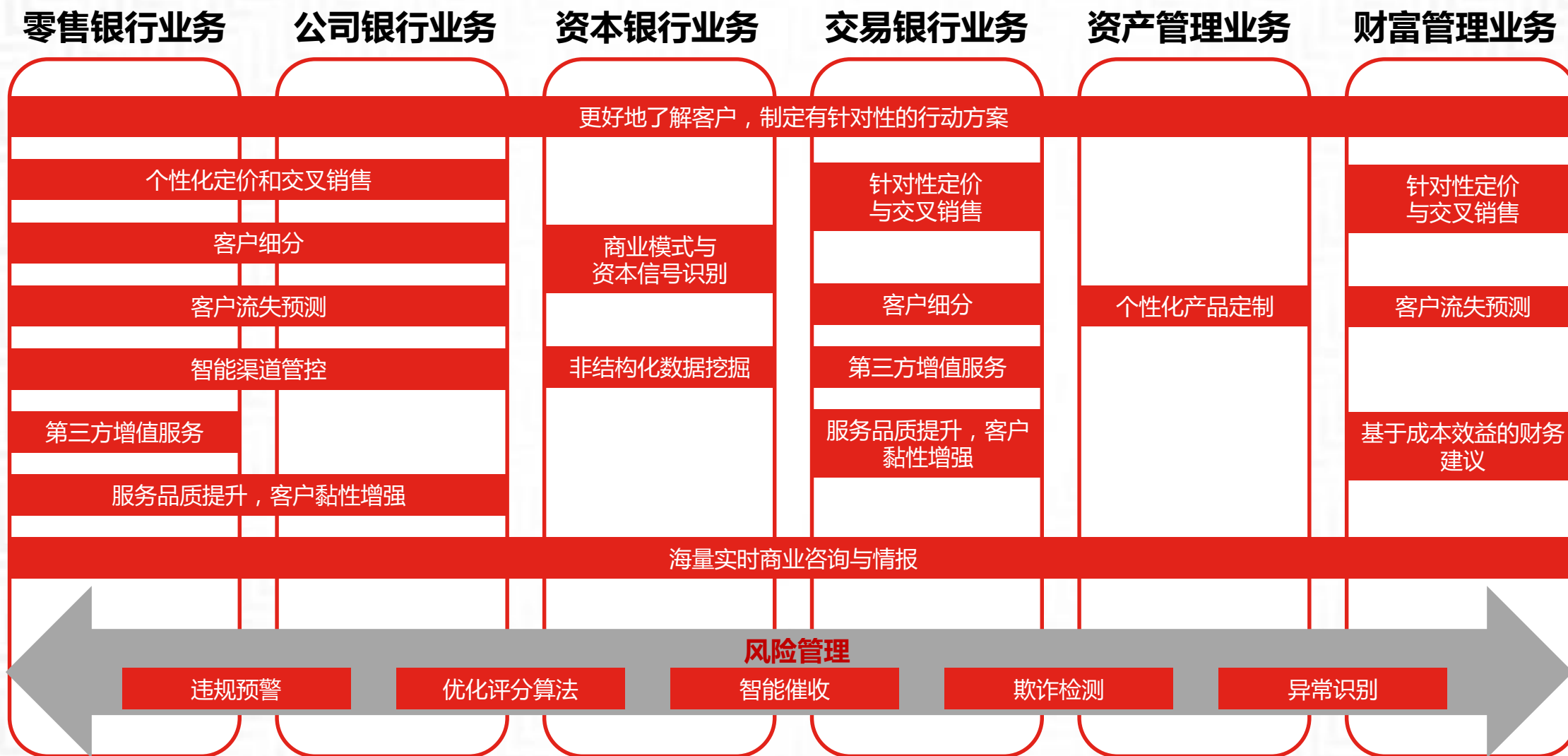
### 销量分析



# + 联想大数据在银行与金融行业的解决方案



# + 银行大数据应用实践

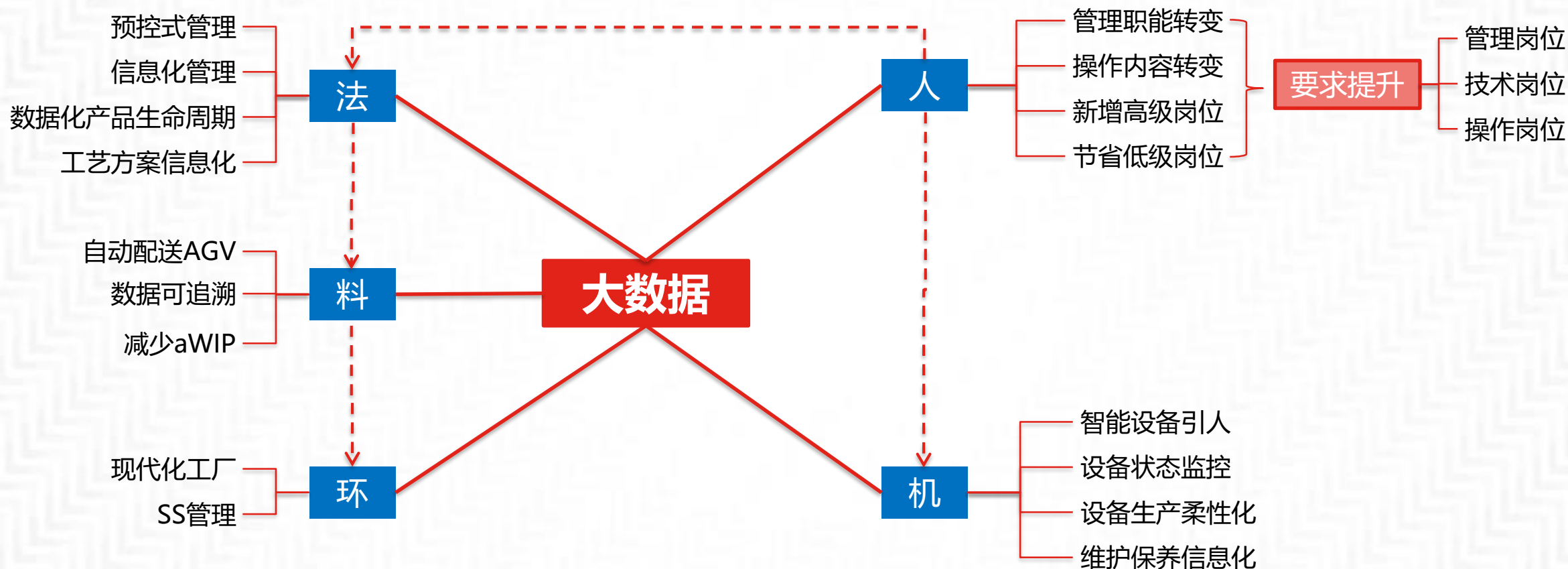


# + 联想大数据在制造业的解决方案



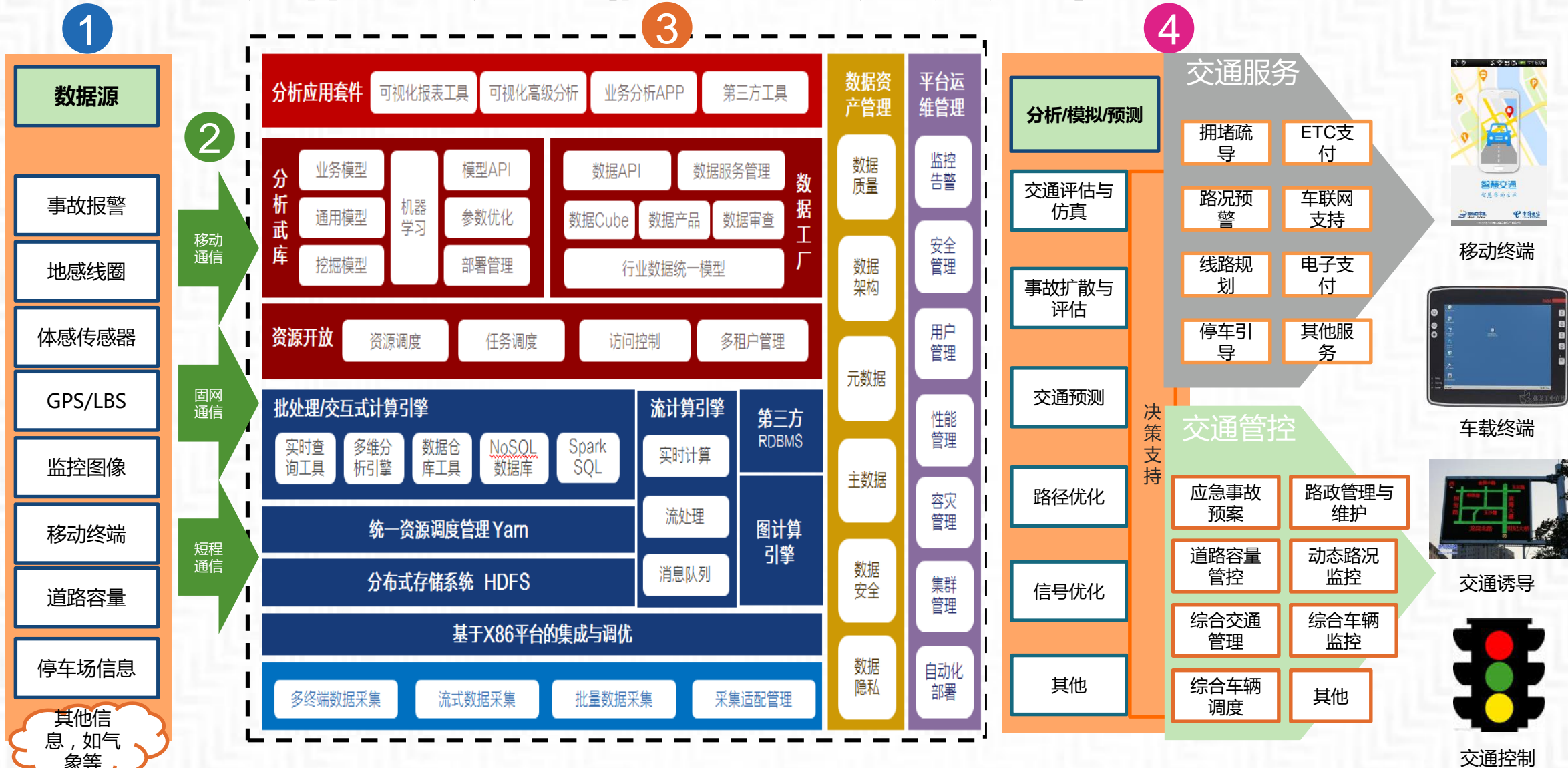
# + 制造行业大数据应用实践

未来的制造将围绕大数据平台构建智能化生产体系，将人，机，法，料，环链接起来，实现多维度数据融合，为企业的运营提供预见性的支撑与指导。



制造业大数据的侧重点在于将所有人，机，法，料，环等信息有效整合起来，加以分析并应用于整个工业生产过程，对整个生产链条进行监控、调整、管理。从而形成高度灵活、个性化、网路化的产业链。大数据是实现工业4.0的关键。

# + 联想大数据在交通运输行业的解决方案



1 感知层    2 网络层    3 大数据分析平台    4 业务应用层



# + 交通运输行业大数据应用实践

关键商业过程	公路交通管理的业务优化方向				
规建运维管理	公路规划	工程建设	质量和验收	公路养护	设备资产
路网效率管理	交通流量监控	拥堵预测	拥堵原因分析	分流管理	资费杠杆管理
安全与应急管理	道路灾害管理	天气灾害预测	路段行驶安全预警	公共安全	行政协作
公民服务管理	客流分析	货流分析	出行线路优化	配套设施改进	数据公开
财务与结算管理	公路投资管理	公路运营成本分析	跨省结算	收入滴漏管理	经营企业收益管理
行政监管与合规管理	路网运营监控	公路收费政策管理	部门职能管理	人员与流程管理	合规管理
数据资产管理	数据资产整合	数据逻辑模型	数据质量改进	数据标准管理	数据与分析成熟度评估

# + 联想大数据在能源行业的价值框架



# + 联想大数据与各行业客户携手共赢



## 技术能力

提供企业级大数据分析平台、各类数据工具、以及管理服务

## 业务价值

实现大数据分析 with 业务价值交付，释放大数据资产的生产力

## 商业生态

与各界共建跨业大数据技术与商业合作的全价值生态圈

“

# THANK YOU

DAKUJEM DANK BEDANKT MERCI TAKK 谢谢  
ありがとう СПАСИБО GRACIAS DZIĘKUJĘ DANKE  
OBRIGADO БЛАГОДАРЯ GRAZIE תודה GRACIAS

