



混合异构数据的清洗、存储、 挖掘架构选型和设计策略

@卢亿雷 From AdMaster
johnlya@163.com



提纲

- 混合异构数据特点
- 混合异构数据分类
- 混合异构处理流程
- AdMaster混合异构数据平台架构
- AdMaster数据处理流程
- AdMaster混合异构数据分析
- Q/A

混合异构数据特点

- 不同的数据类型
- 不同的数据量级
- 不同的访问速度
- 不同的用户类型
- 不同的访问平台
- 不同的存储设备
- . . .

混合异构数据分类

在线数据		离线数据
数据内容	短周期数据	长周期（存档、归纳、 计算结果）
数据特性	字段固定	字段不固定
数据结构	高度结构化、复杂、适合操作计算	结构简单
使用频率	非常高（热数据）	一般（冷数据）
数据访问量	K B、MB级	GB、TB、PB级
响应时间	纳秒、微秒、毫秒级	秒、分钟、小时、天级

数据采集

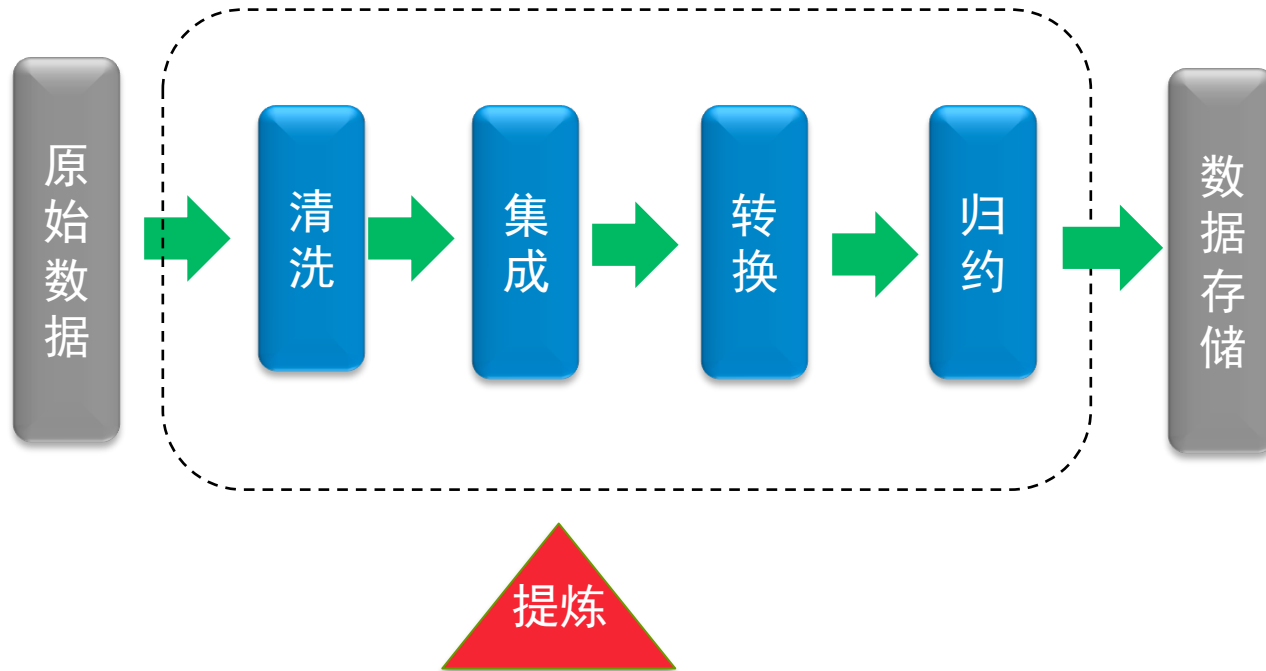
结构化数据

非结构化数据

Internet



数据预处理



数据分析

应用服务

Pig

Hive

Mahout

Flume

Sqoop

Oozie

Online
(HBase)

Batch
(MapReduce)

Streaming
(Storm, S4)

In-Memory
(Spark)

Interactive
(Tez)

HPC MPI
(OpenMPI)

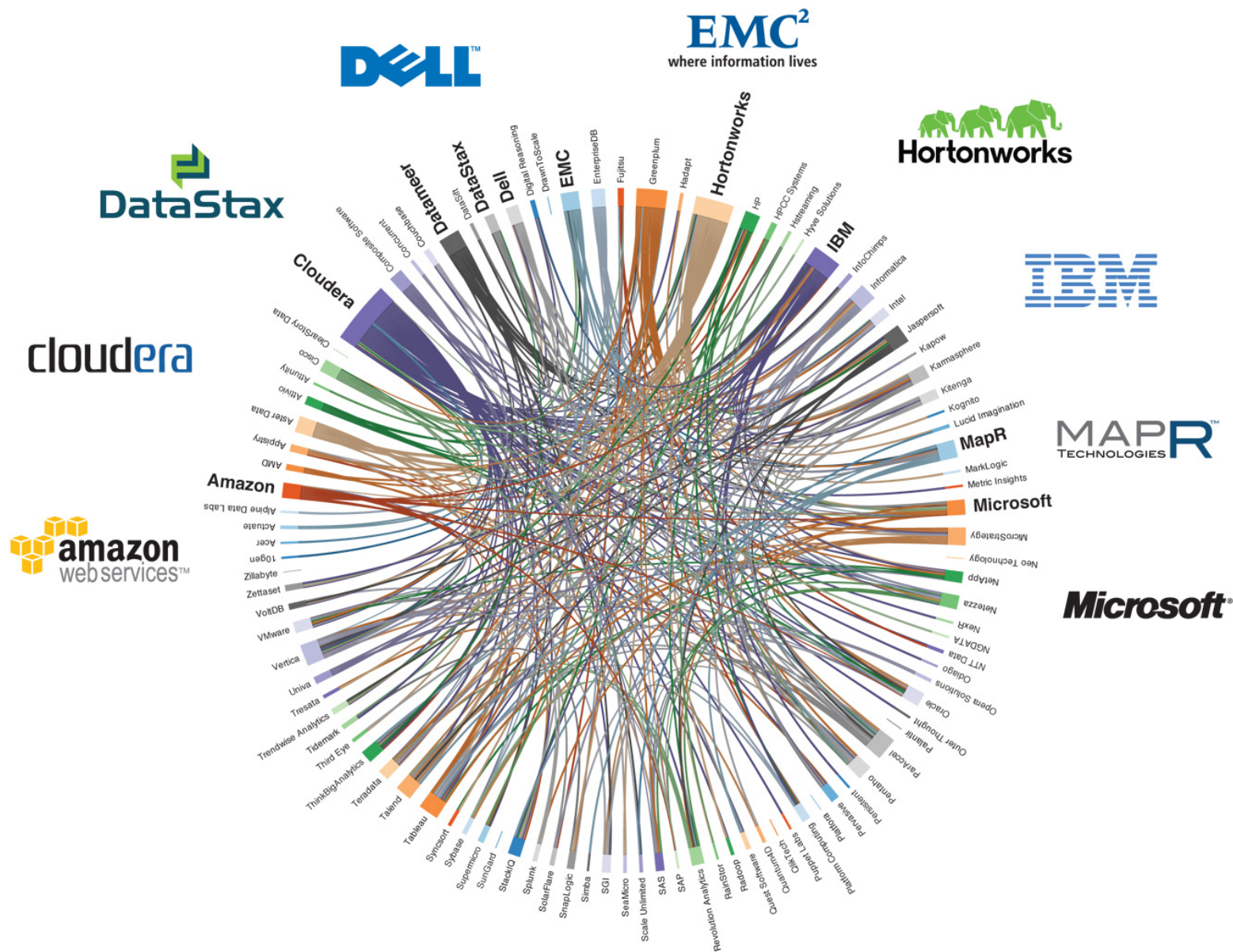
Zookeeper

YARN Cluster Resource Management

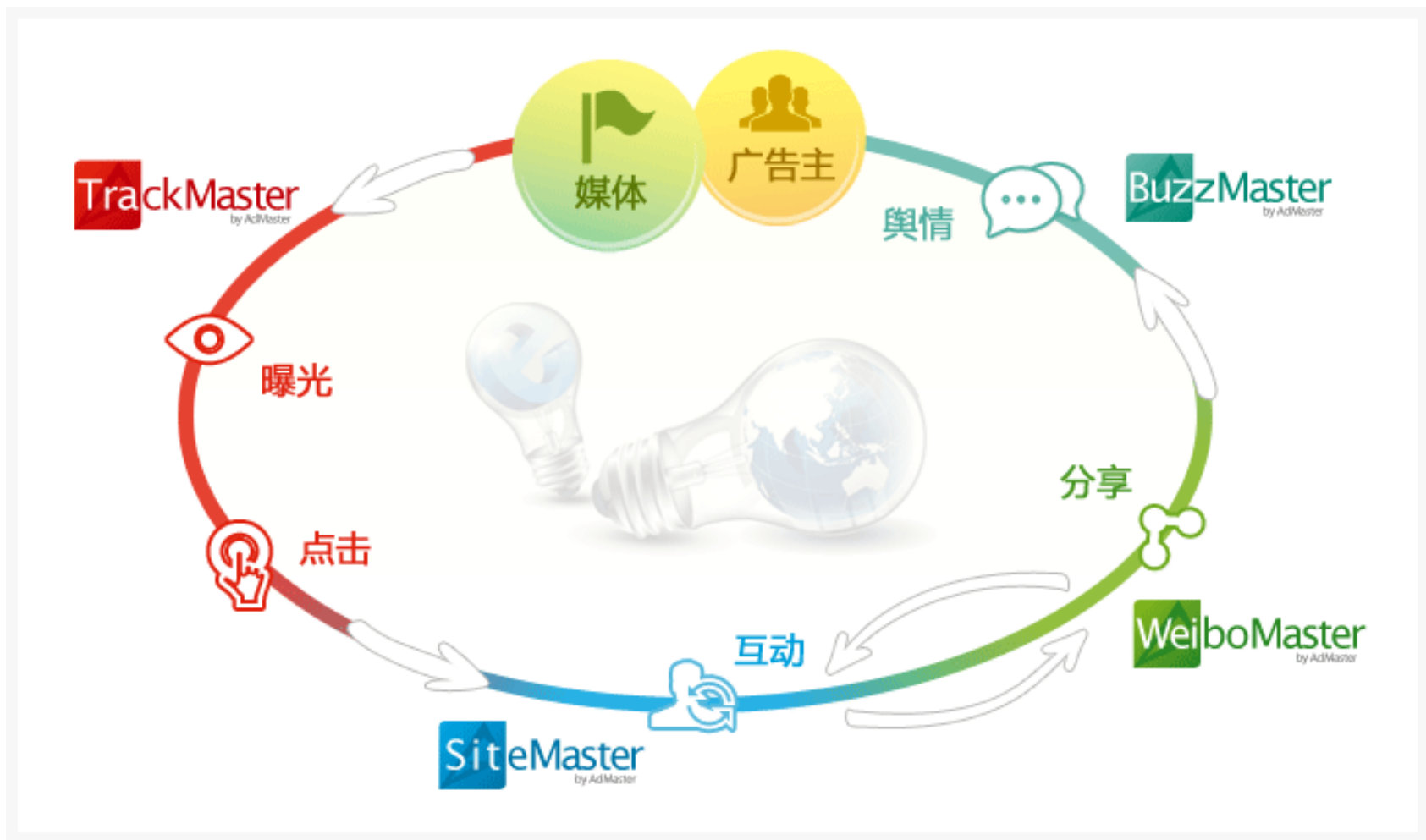
HDFS

OS (操作系统)

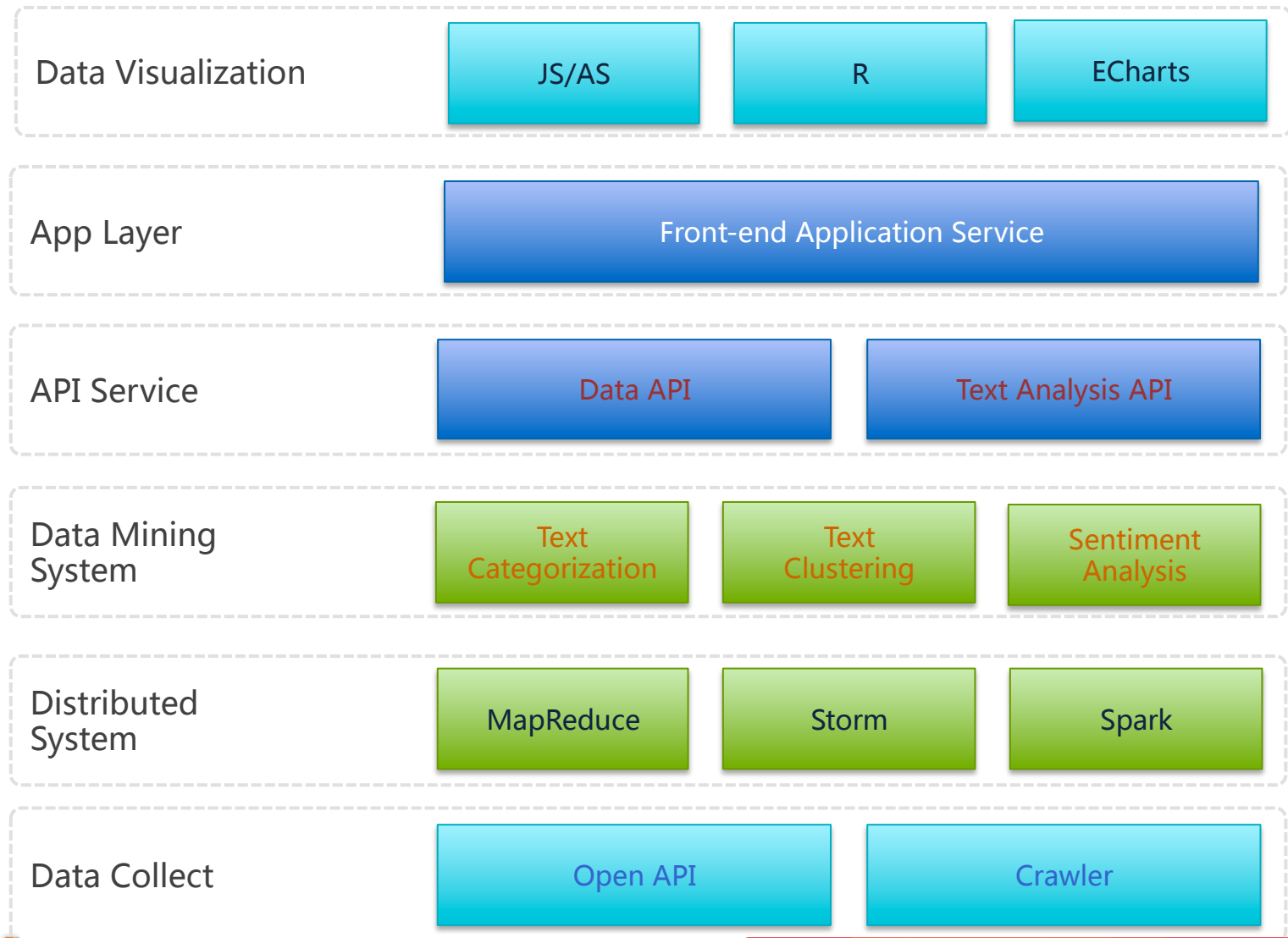
数据展示



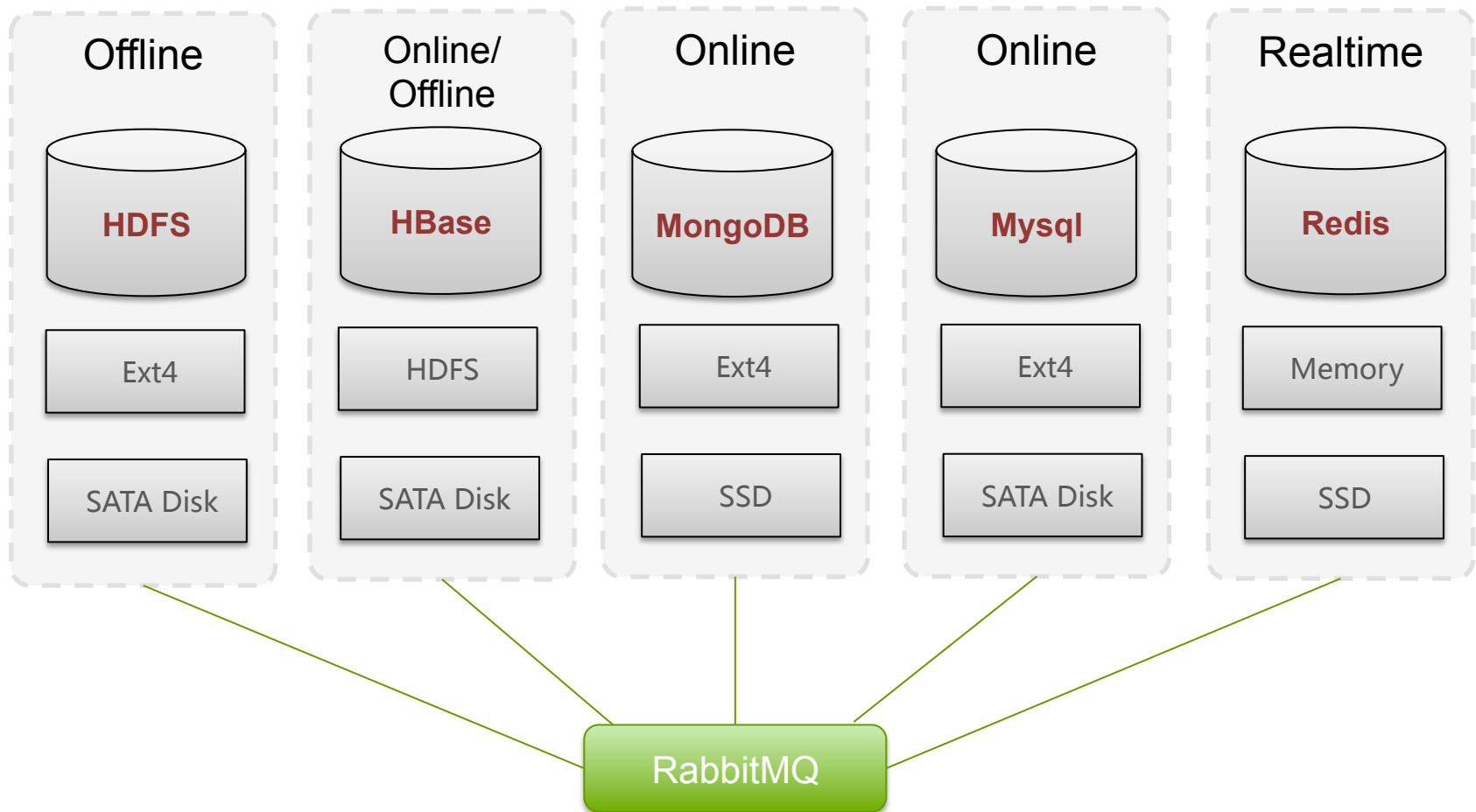
互联网广告监测全流程



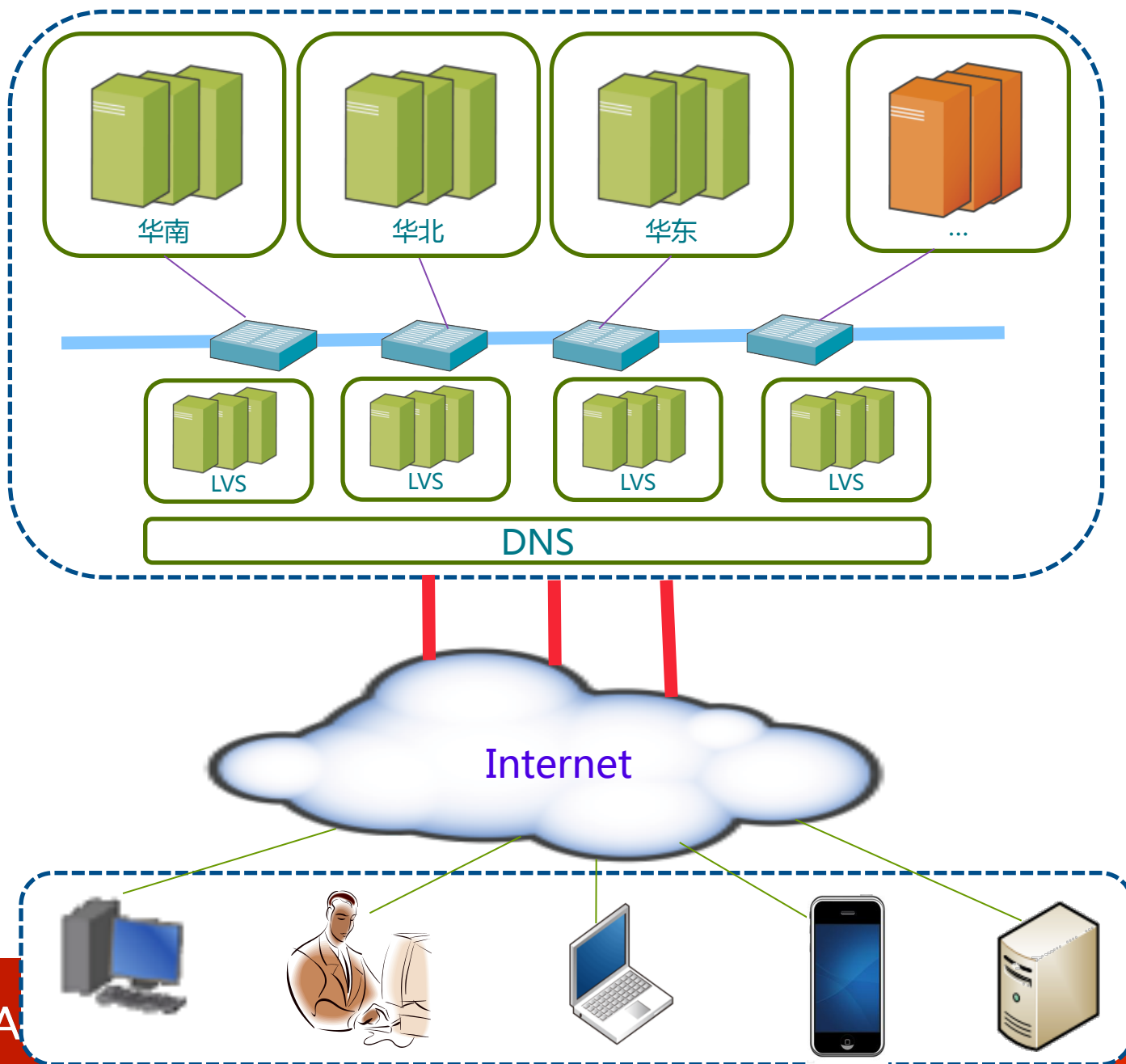
AdMaster混合异构数据平台架构



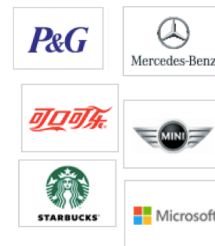
AdMaster混合异构数据平台架构



AdMaster数据采集



广告数据采集



Buzz Resource

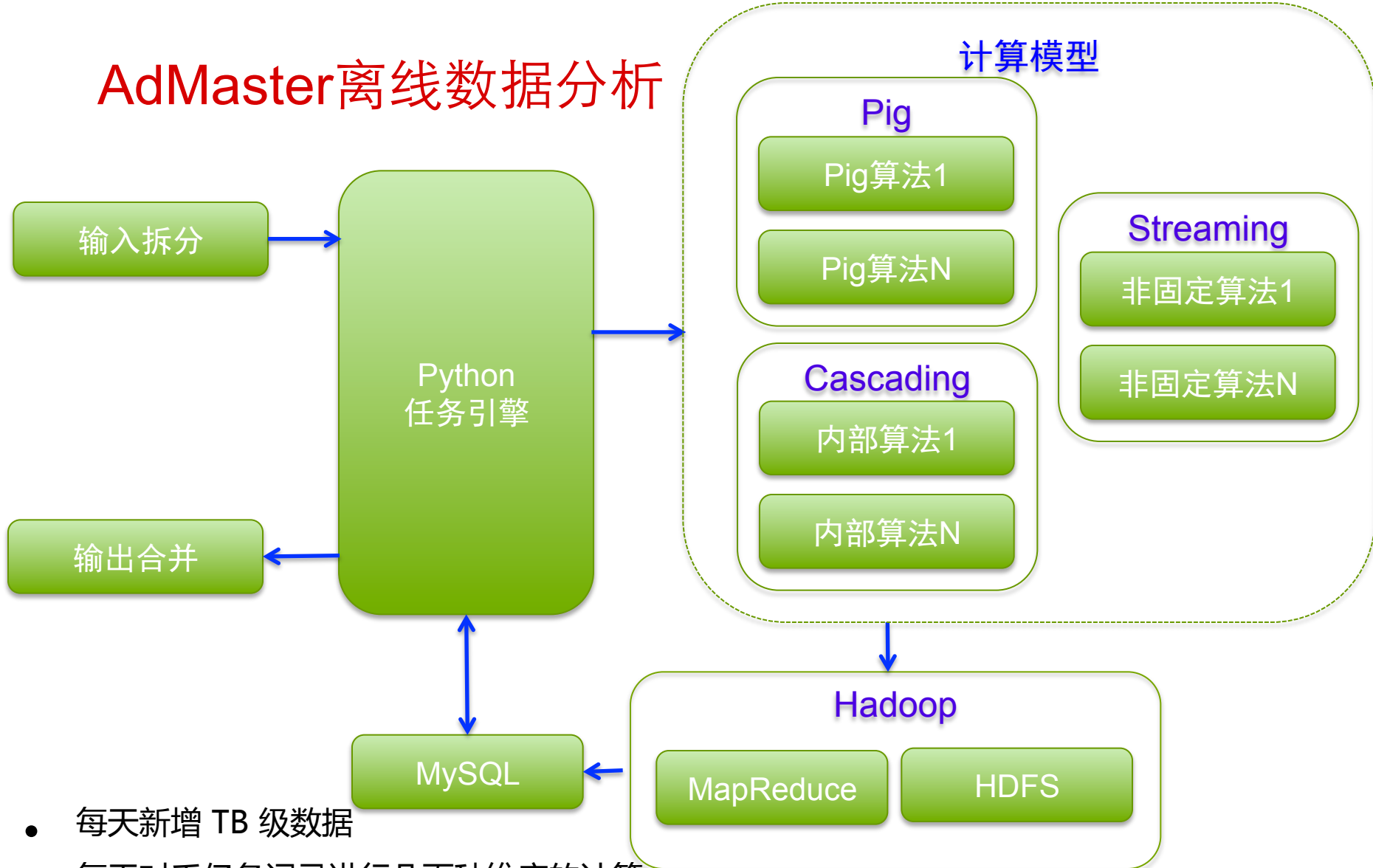


Others ...

AdMaster数据采集

- `cat /proc/sys/net/ipv4/tcp_mem`
- `cat /proc/net/sockstat`
- `cat /proc/sys/net/ipv4/tcp_max_orphans`
- `net.netfilter.nf_conntrack_max`
- `net.netfilter.nf_conntrack_tcp_timeout_established`

AdMaster离线数据分析

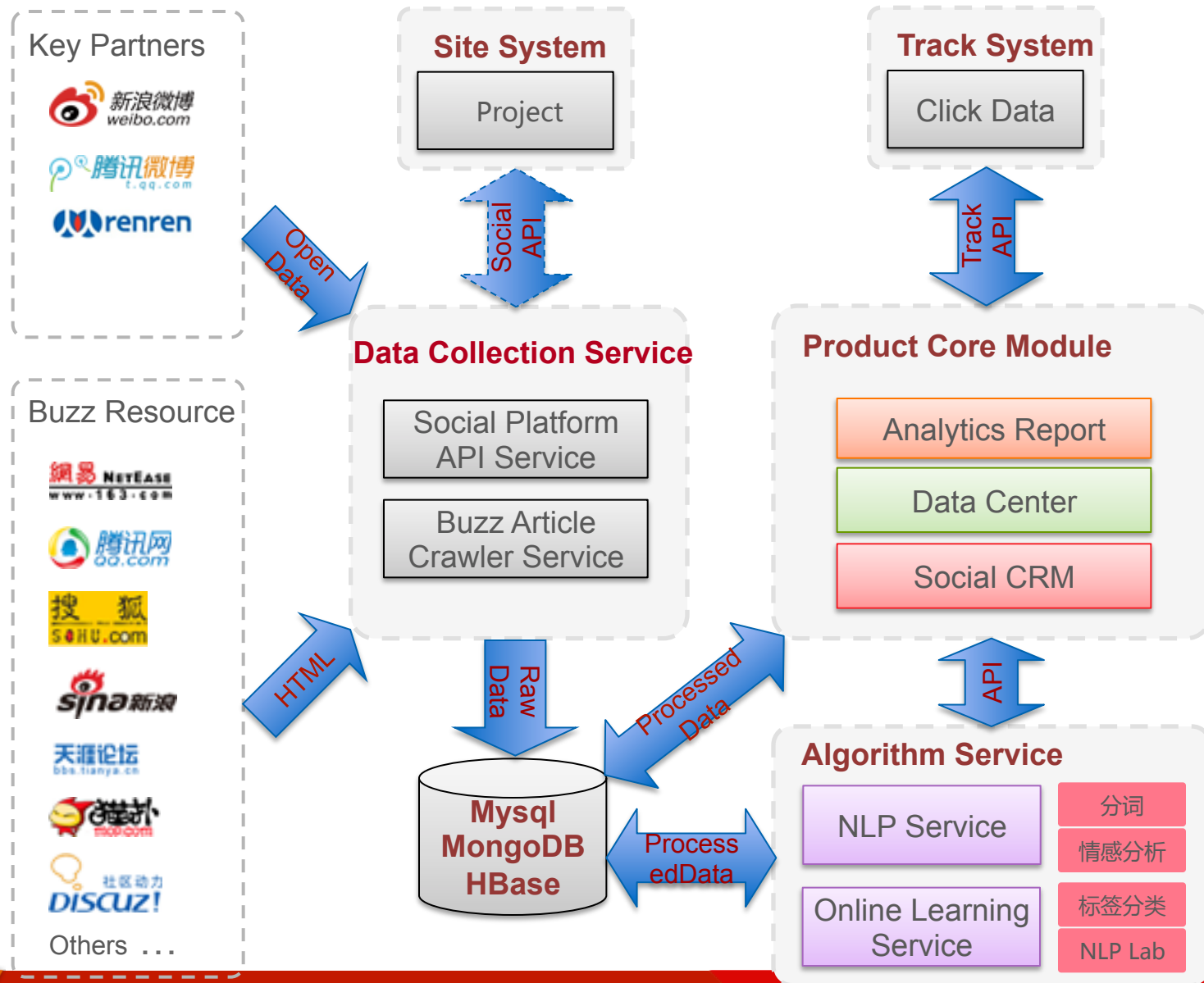


- 每天新增 TB 级数据
- 每天对千亿条记录进行几百种维度的计算

AdMaster离线数据分析

- /sys/kernel/mm/
redhat_transparent_hugepage/enabled
- /sys/kernel/mm/
redhat_transparent_hugepage/defrag
- dfs.socket.timeout
- dfs.datanode.max.xcievers
- dfs.datanode.socket.write.timeout
- dfs.namenode.handler.count

AdMaster 在线数据分析



AdMaster 在线数据分析

- Kafka & Tail
- HBase & MongoDB
- Storm & Rsync
- Spark & Hana

AdMaster 数据可视化



- 主题 (配色、品牌名、品牌logo)
- 轮播信息 (Screens、Slides、标题)
- 权限 (用户、用户组)

配置信息



- Social数据源
- Site数据源
- Track数据源

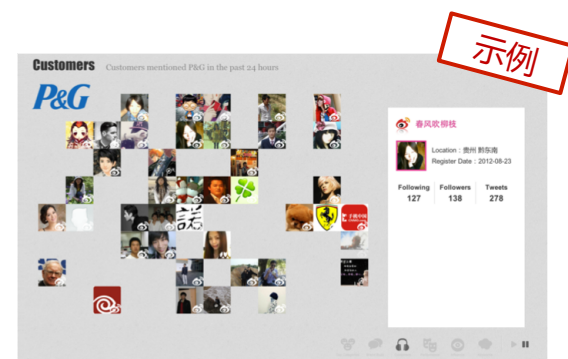
数据源



AdMaster 数据可视化



- 数据呈现方式
- 展示层与数据层松耦合，多种数据源接入
- 极高的可靠性和容错机制



Q & A