

# 移动大数据管理平台实践



阎志涛

TalkingData研发副总裁

# 关于TalkingData

TalkingData创立2011年9月，是国内最大的数据管理、服务平台。TalkingData核心团队来自Oracle、IBM、HP等公司，长期从事分布式运算架构、海量数据处理、数据挖掘算法的研究工作。TalkingData深信数据本身蕴含巨大的价值，希望通过优秀的产品、完善的服务，将“大数据”落地，充分发挥数据的价值，用数据改变未来。

从基础数据分析、管理产品，到深度的数据咨询服务，TalkingData有着深厚的数据积累和应用经验。利用海量数据，不断实践科学计算领域内的各种算法、概念，不断尝试将数据与场景相结合，找到适合不同领域的数据模型，实现数据价值的最大化。

目前TalkingData产品及服务涵盖移动应用数据统计、移动广告监测、移动游戏运营、公共数据查询、综合数据管理、数据咨询服务等多款极具针对性的产品及服务。为超过80,000款应用、游戏提供数据统计、分析服务，覆盖超过13亿独立移动设备；为招商银行、中信银行、平安保险等大型企业提供全方位数据服务。

—— “用数据改变未来”



# TalkingData移动大数据平台



TalkingData移动大数据平台能够为客户提供基于移动互联网数据的全方位服务。

无论您是开发者，还是广告主，或是大型企业，都能找到适合自己需求的产品或服务。

TalkingData移动大数据平台产品及服务涵盖基础统计分析、游戏运营支持、移动广告监测、第三方数据管理平台、数据咨询服务，以及面向大型企业的综合数据解决方案。



# 移动互联网大数据特点

- 移动互联网大数据的4V

- Volume

随时随地都在产生数据，数据量更大

- Variety

随时随地联网的特性，使得移动互联网的数据更具有多样性。在移动侧可以有更为精准的位置数据，各种传感器数据。

- Velocity

对速度处理的要求性更高，很多的业务场景需要更实时的数据处理才能使得数据产生价值。

- Value

更多高价值的数据产生

- 万物皆可联网，数据方便人的生活

- IOT逐渐成为现实，万物都在贡献数据

- 各种智能硬件逐渐普及



# 移动互联网大数据处理流程



# 数据获取

- 获取哪些数据？

## 设备信息

- 设备ID
- 设备软硬件信息

## 数据业务信息

- 业务事件
- 会话信息

## 上下文信息

- 网络
- 位置
- 传感器



# 数据获取

- 如何获取数据？

## 存储转发

- 移动网络不稳定
- 移动应用不稳定

## 数据压缩

- 网络流量消耗
- 电池消耗

## 传输协议

- 数据安全



# 数据收集

- 数据收集器
  - 数据格式校验
  - 轻量级
  - 高并发处理
  - 无状态
  - 存储转发





# 数据存储

## 分布式文件系统



- ☐ 数据长久保存
- ☐ 数据冗余
- ☐ 离线计算服务

## NoSQL数据库



- ☐ 数据有时效性
- ☐ 为实时计算服务
- ☐ 缓存

## 关系型数据库



- ☐ 结果型数据
- ☐ 事务一致性保证
- ☐ 多表关联



# 数据计算

## 流式计算

- 实时指标
- 基于规则的标签

## 离线计算

- 批量统计
- 大时间尺度数据计算

## 数据挖掘

- 机器学习
- 迭代算法



# 数据服务

多维报表

数据可视化

数据服务接口



# 我们面临的挑战

- 业务发展的驱动，多个竖井

| App Analytics  | Game Analytics   | AdTracking   | DMP  | Insight  |
|--|--|--|--|--|
| <ul style="list-style-type: none"><li>• SDK</li><li>• Collector</li><li>• Data Store</li><li>• Compute</li><li>• Service</li></ul> | <ul style="list-style-type: none"><li>• SDK</li><li>• Collector</li><li>• Data Store</li><li>• Compute</li><li>• Service</li></ul> | <ul style="list-style-type: none"><li>• SDK</li><li>• Collector</li><li>• Data Store</li><li>• Compute</li><li>• Service</li></ul> | <ul style="list-style-type: none"><li>• Data Store</li><li>• Compute</li><li>• Service</li></ul> | <ul style="list-style-type: none"><li>• Data Store</li><li>• Compute</li><li>• Service</li></ul> |



# 我们面临的挑战

- 未来更多的数据业务
  - 竖井模式很难支持新业务的开展
- 更多的数据价值探索的需求
  - 竖井模式很难深入了解技术
- 更多的数据 ( Bigger than Bigger)
  - 竖井模式不利于资源的合理利用
- 没有统一的数据视图

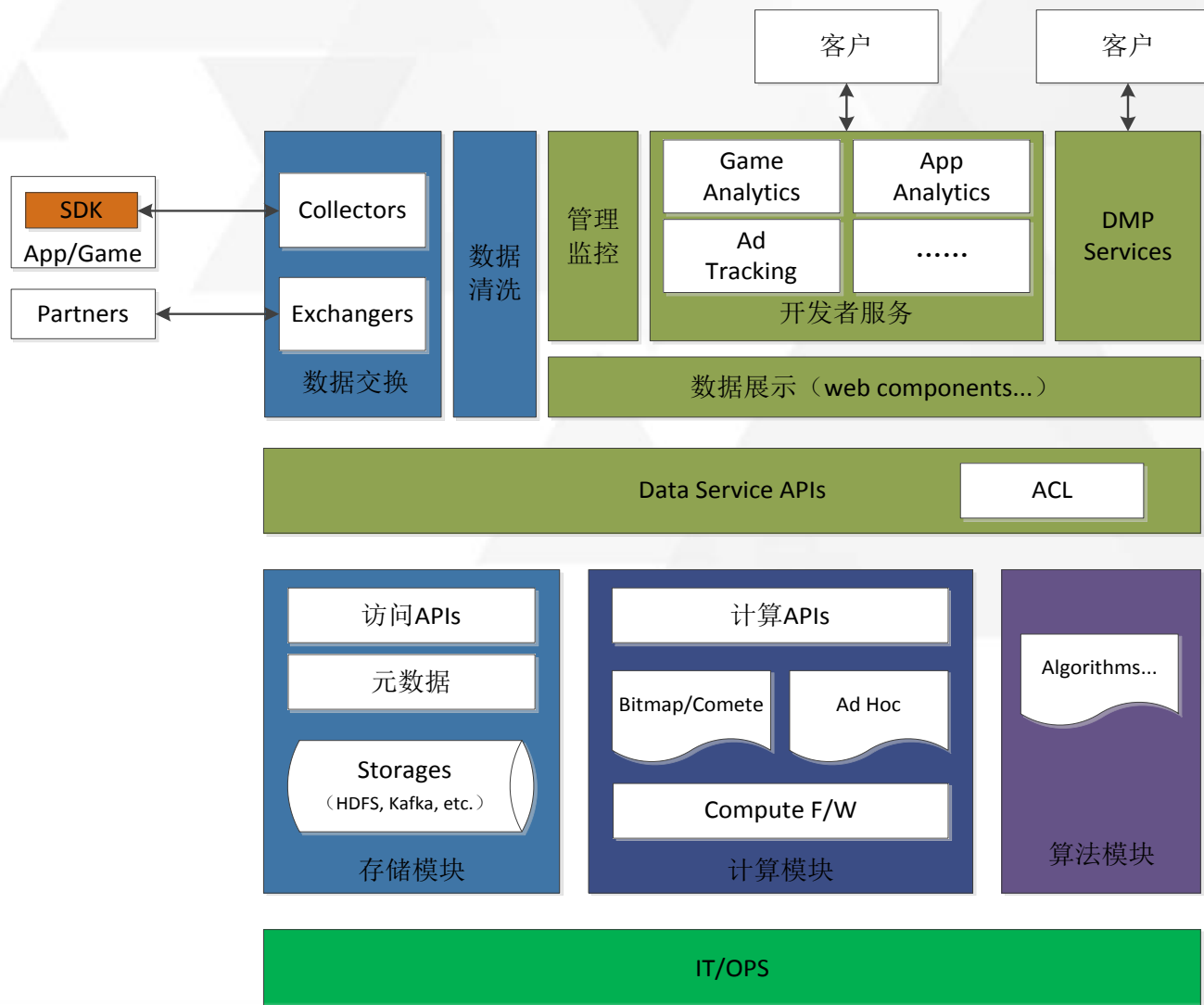


# TD移动大数据管理平台（ $\pi$ 系统）

- 整合多产品线的基础服务
  - 统一存储
  - 统一计算
  - 统一数据总线
  - 统一数据挖掘
  - 统一视觉呈现
  - 统一数据收集
  - 统一SDK
  - 统一监控和管理
- 提供更灵活高效的技术支撑
  - 产品能迭代速度更快
  - 研究成果加速流动



# $\pi$ 系统架构



# 统一SDK

- 统一SDK
  - 新的统一的数据收取框架
  - 业务层和基础层分离
  - 非阻塞模式
  - 处理各种异常
  - 高效存储格式





# 统一的数据收集

- 统一数据收集系统
  - 基于Node.js和C++开发
  - 支持分布式部署
  - 数据收集系统支持存储转发
  - 分布式收集节点和中心节点数据传输高压缩比



# 统一的数据总线

- 统一数据总线
  - 基于Kafka的数据总线
  - 规范不同业务线的topic命名规则
  - 统一的管理



# 统一存储

- 统一的离线存储(HDFS)
  - 数据域管理，多业务系统可以共享存储资源
  - 数据文件按照时间进行切片
  - 数据文件时效管理，中间数据可以自动删除
  - 数据自动归档
  - Parquet列式存储格式，方便数据计算
  - 计划支持数据EC(Eraser Coding)
  - 分布式缓存Tychyon



# 统一存储

- NoSQL数据库
  - 开发Bitmap存储，bitmap基本运算下沉到存储层，底层基于RocksDB
  - MongoDB 3.0(WireTiger引擎)，基于SSD
  - Redis



# 统一存储

- 关系型存储
  - MySQL Cluster
  - WebScaleSQL?



# 统一存储

- 统一接口封装
  - 存储层对计算层通过接口提供数据
  - 存储对于计算完全透明



# 统一存储

- 元数据管理
  - 基于Hcatalog进行二次开发
  - 支持不同数据源
  - 支持json,protobuf等数据格式
  - 支持版本



# 统一计算

- 统一的计算框架和接口
  - 基于Yarn进行计算资源调度
  - 基于Spark的并行计算框架
  - 基于预先生成Bitmap的OLAP解决方案
  - 利用Spark Streaming进行流式计算
  - 自行开发的任务调度系统
  - 统一的计算查询接口





# 统一的数据挖掘

- 数据挖掘服务化
  - 基于统一计算框架
  - 针对Spark,自行实现了LR,DT等数据挖掘算法库
  - 将数据挖掘服务化，变成统一计算的一种能力



# 统一的视觉呈现

- 统一的视觉呈现
  - 视觉呈现组件化
  - 支持各种自定义报表
  - 支持各种数据可视化效果



# 统一监控和管理

- 统一监控
  - 基于Zabbix开发
  - 支持CPU、内存、硬盘、网络以及进程运行状态等等的监控
  - 支持短信、邮件、微信报警



# 带来的好处

- 更方便的增加新的数据业务
- 工程师可以更深入的了解技术
- 资源可以更合理的进行配备



# 未来

- 进一步优化存储能力
  - 热数据、冷数据、归档数据的合理分层
  - 硬盘、SSD、内存的合理使用
  - 基于latency的存储提供
- 进一步优化计算能力
  - 更好的支持即时分析
  - 更细粒度的资源调度能力(Myriad?)
- .....





Contact:

阎志涛 Tony Yan

Email: [tony.yan@tendcloud.com](mailto:tony.yan@tendcloud.com)

Wechat: [zhitao\\_yan](#)

Weibo: <http://weibo.com/ztyan>

<https://www.talkingdata.com>