



如何打造高性价比百亿级实时数据分析平台

——以精益创业理念修炼大数据分析平台

本产品保密并受到版权法保护

Confidential and Protected by Copyright Laws



郭炜 易观 CTO

郭炜先生2015年加入易观，担任易观CTO，构建易观技术团队完成易观大数据采集、平台、数据挖掘等技术架构与体系，从无到有完成易观混合云搭建、易观SDK升级并发布易观秒算实时计算平台，目前易观大数据平台日处理数据量30T，200亿条，月活用户3.58亿。

郭炜先生毕业于北京大学，加入易观之前，曾任联想研究院大数据总监，万达电商数据部总经理，并曾在中金、IBM、Teradata公司担任大数据方向重要岗位，对大数据前沿领域研究，包括视频、智能WIFI等大数据软硬数据一体技术有独特的见解。



- 精益化的修炼大数据
- 最小化的实时大数据分析闭环
- 新的技术框架迭代与扩展
- 高性价比的实时大数据框架

现在的易观的实时大数据平台

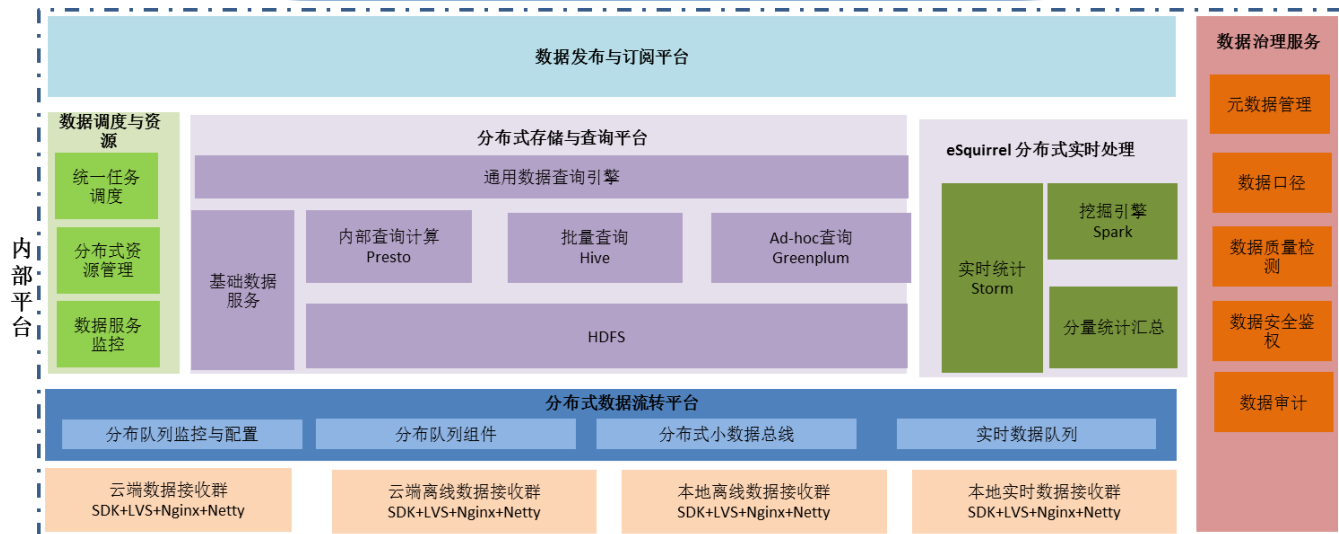


截止到2016年第4季度

产品展现与服务集群



大数据处理集群



数据采集与预处理

Android SDK



IOS SDK



价值

H5 SDK



微信小程序SDK



回顾 “精益创业的理念” 几个要点

- 最小化可行产品进行优化，而不是对其设定硬指标 v.s. 决策层说 “我们要建设大数据项目”
- 与最终客户与业务保持同步 v.s. “先有平台再加业务”
- 业务闭环，并形成针对大数据的数据分析 v.s. “管理层看到了Dashboard”
- 增速/转型/创新 ——最大的挑战，在于企业文化的改变

- 精益化的修炼大数据
- 最小化的实时大数据分析闭环
- 新的技术框架迭代与扩展
- 高性价比的实时大数据框架

最小可闭环的产品——“为什么需要实时大数据”



中生代技术
FRESHMAN TECHNOLOGY

大数据
大数据观

实时分析驱动用户资产成长

数据永远是“临时的”，分析永远是“有时效性的”

不适时的产品推荐 失效的活动分析与策略 过时的产品质量反馈



Dashboard
OLAP
Report
...

日志分析
用户画像
推荐引擎
...

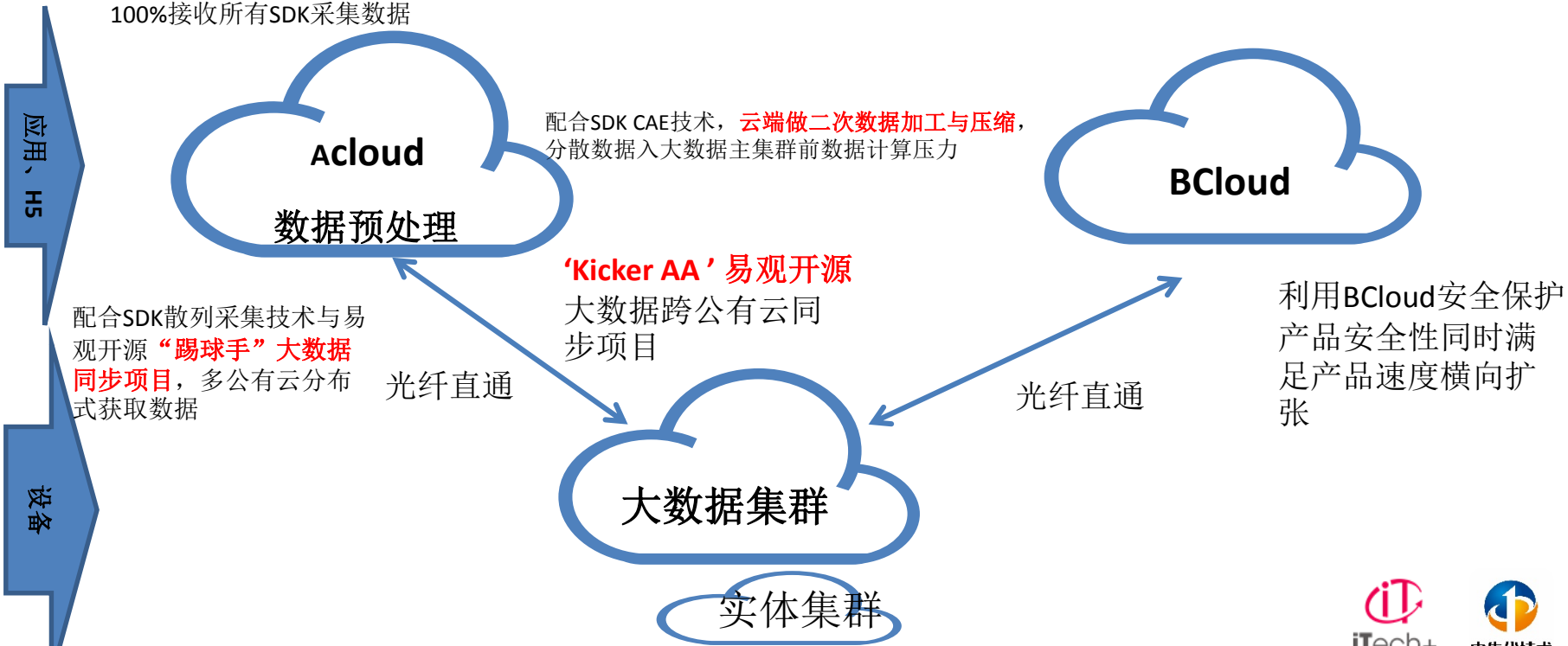
实时渠道分析
实时场景推荐
实时风险评估
...

智能机械
智能交互
智能策略
...

最小化的可执行的闭环架构——基础架构

数据的实时采集 → 数据的实时接收 → 数据的实时计算 → 数据的实时查询 → 数据的实时挖掘 → 数据的实时服务

100%接收所有SDK采集数据



最小化的可执行的闭环架构——为什么用混合云

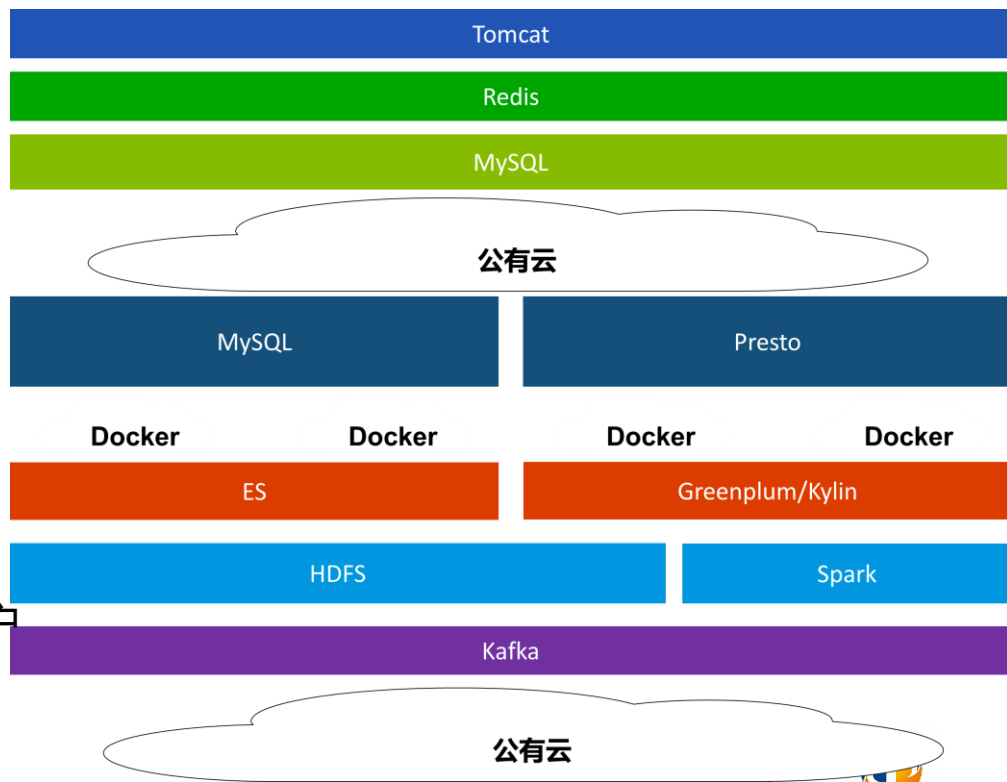
接收端公用云弹性扩展

- 网络带宽
- 接收性能
- 安全防控

处理端平台

- 独占性能
- 技术迭代迅速
- 投入TCO可控

经过1年多的检验，每日去重活跃3000万用户，600多个合作伙伴，每日200亿条数据



最小化的可执行的闭环架构——可执行的接收端



中生代技术
FRESHMAN TECHNOLOGY

IT 大数据观

实时分析驱动用户资产成长

当大数据有批量传输变为实时传输时发生了质的变化：

- 更类似于并发的交易系统
- 关注于数据流的疏导大于处理

技术关键点：

- 良好的扩展网络架构
- 云+端的控制策略



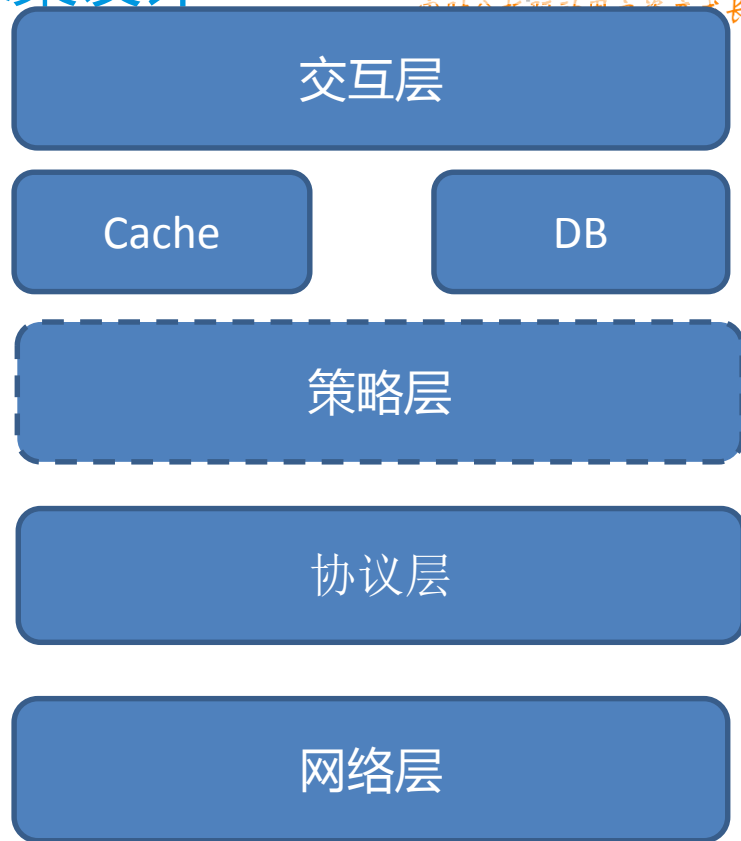


麻雀虽小五脏俱全——最小化的实时采集设计

无论是IoT还是设备，一般都需要五层实时获取方式设计：

- 交互层
- 存储层
- 策略层
- 协议层
- 网络层

建议：行为数据传输采用短连接,http协议，长连接用于心跳保持



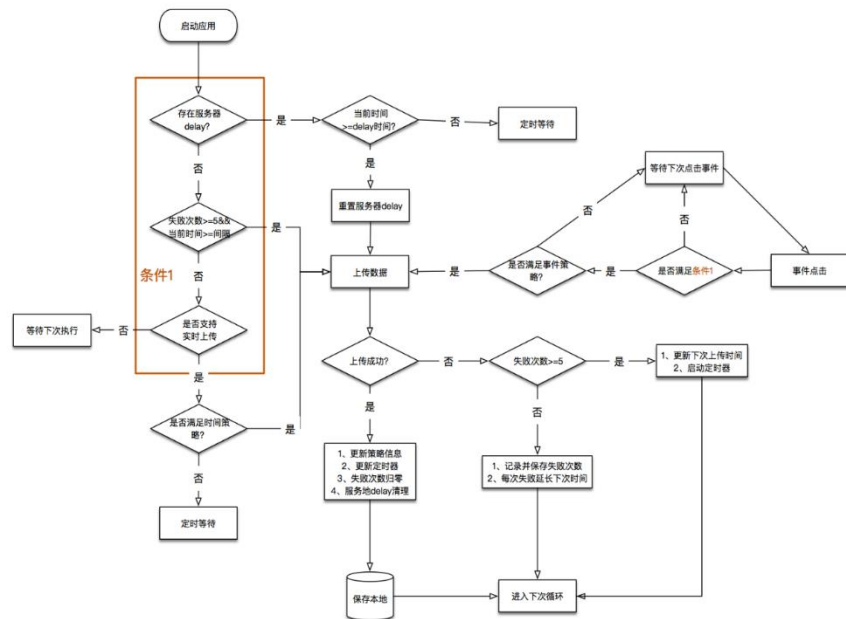
麻雀虽小五脏俱全——最小化的策略设计

云端的主要策略：

- 时间间隔
- 失败策略
- 清洗策略
- 分流策略

设备端主要策略：

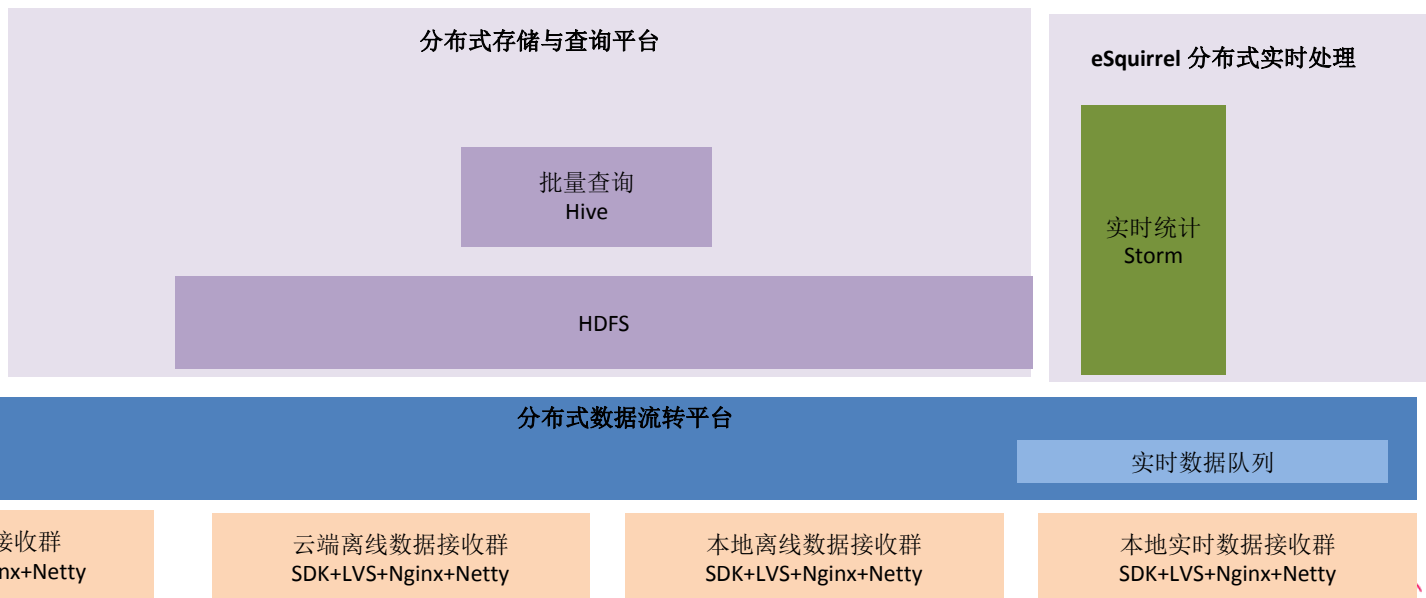
- 失败策略
- 更新策略
- 保活策略



目前验证可以分不同合作伙伴，即时调整策略从5秒到6小时，也可以屏蔽或分流问题设备

最小化的实时计算框架

内部平台

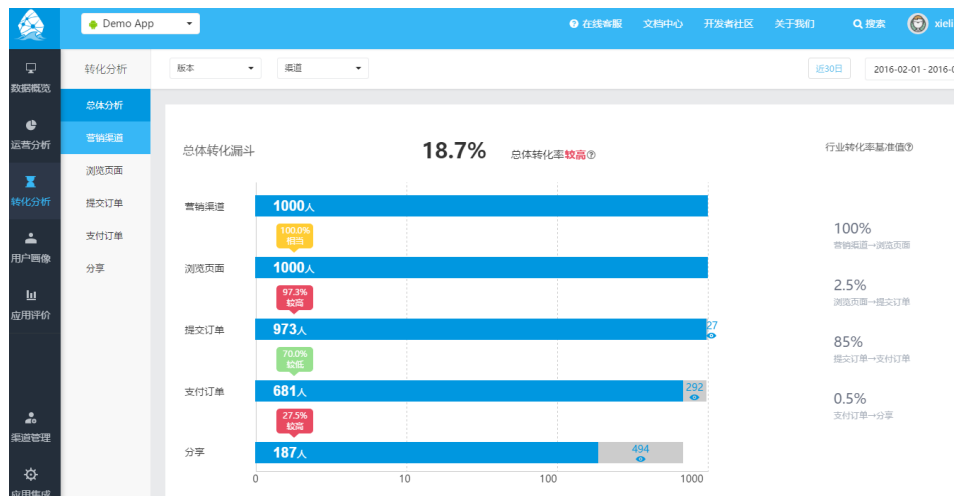


最小化的实时计算框架

	Spark Streaming	Storm
处理模型及延时	微批处理,秒级时延	一次一条数据,亚秒级时延
准确性保证	Exactly once	At least once(也可以实现exact once(使用storm trident的事务状态更新机制))
容错性	checkpoint , 血统进行恢复	ack/fail机制标记跟踪每条数据
易用性	很多内置函数, 可以使用sql/dataframe/datasets等	一些功能需额外编码
编程语言api	scala、java、python	Java、Python等
开发语言	Scala	Clojure
生产支持公司	Cloudera、MapR、Databricks	Hortonworks
资源管理	基于Spark , Mesos和Yarn	Mesos和Yarn
成熟度	中上, 正在积极发展中	相对成熟稳定
社区	社区非常活跃, bug fix及时	社区相对不太活跃

最小的业务闭环

漏斗的转化 V.S. 实时Dashboard



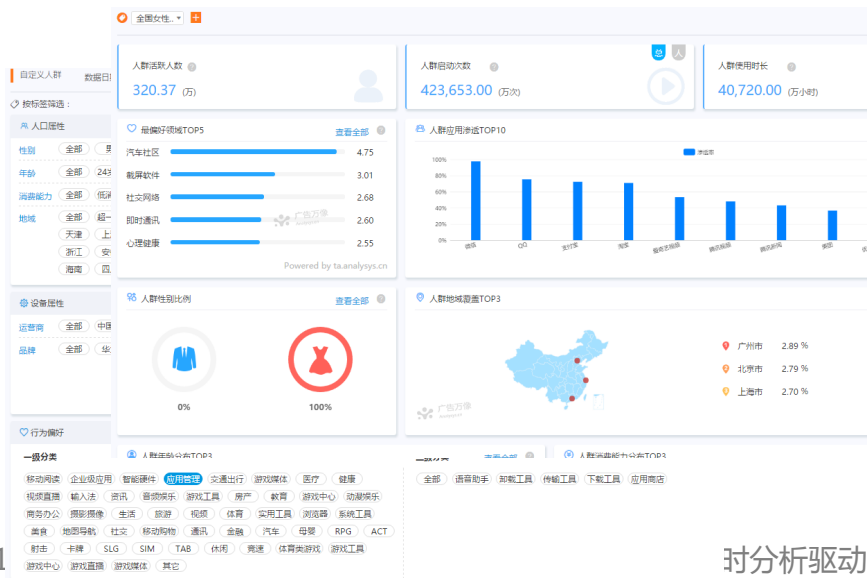
- 精益化的修炼大数据
- 最小化的实时大数据分析闭环
- 新的技术框架迭代与扩展
- 高性价比的实时大数据框架

开始新一轮扩展

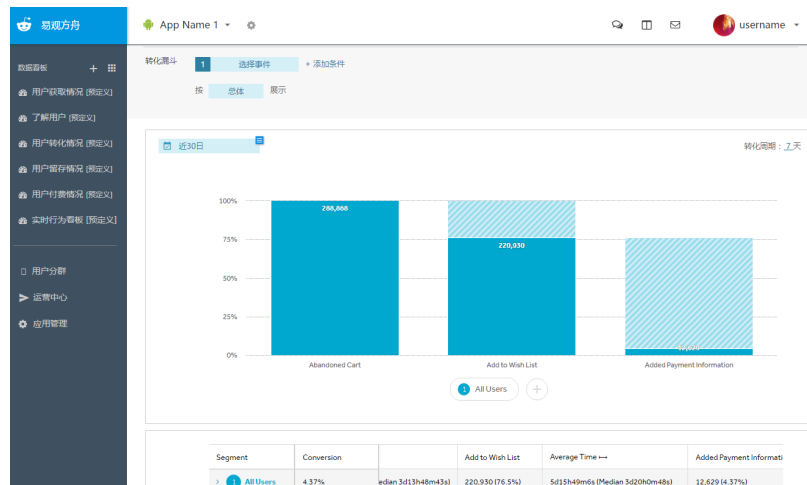
先有鸡？ V.S. 先有蛋？

——先有“鸡窝”

例子1：易观万象



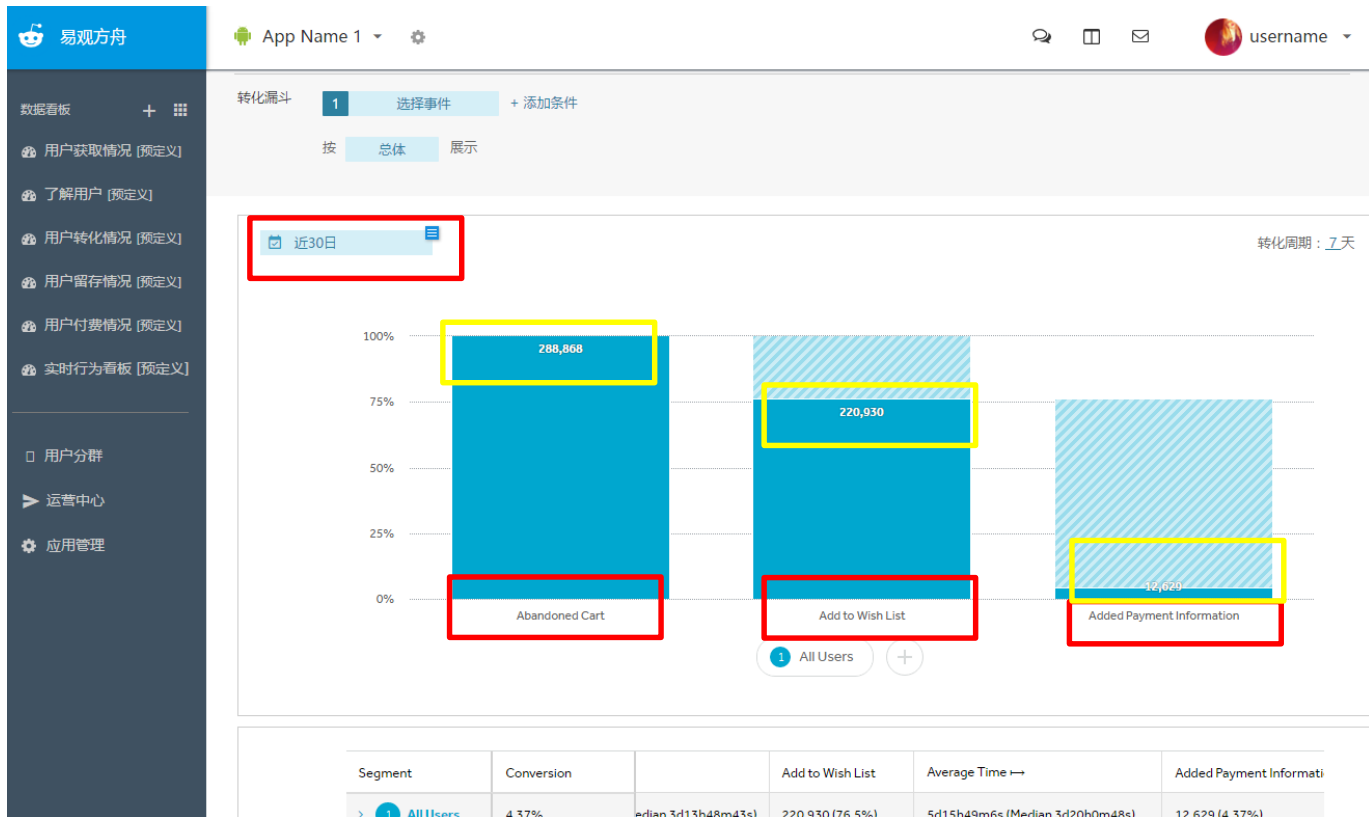
例子2：易观方舟转化



例子1——易观万象场景，基于用户画像+行为的明细查询



例子2——易观方舟场景，有序漏斗查询



有序漏斗查询——抛砖引玉

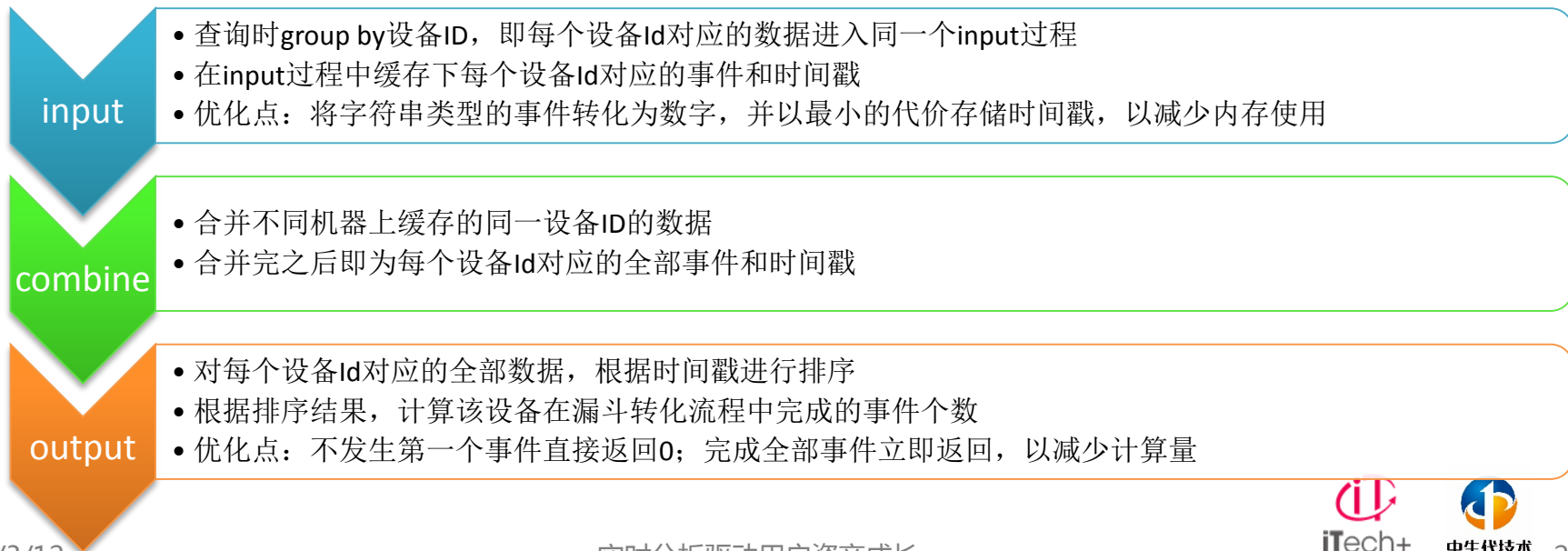
设备ID	时间戳	事件名称	事件属性	应用标识	日期
1111	1490970560539	搜索商品	{.....}	3zprjdg7yyp5	20170227
2222	1490965747548	浏览商品	{.....}	3zprjdg7yyp5	20170227
3333	1490972189107	提交订单	{.....}	3zprjdg7yyp5	20170227
.....

数据说明：

- (1) 时间戳精确到毫秒，并能保证同一个设备在同一时间戳内，不会同时发生两个事件。
- (2) 事件个数有限，不会超过50个，但事件属性无法固定。
- (3) 每天的数据量大概有： $100\text{万DAU} * 5\text{事件} * 2\text{次/事件} = 1000\text{万}$ 。

有序漏斗查询——抛砖引玉

- (1) 底层存储用HDFS
- (2) 建立Hive表，并以应用标识、日期、事件名称为分区
- (3) 查询自定义UDAF，或者利用Spark core自定义相同逻辑



针对每个点持续，不断的优化...

市面最小的SDK Android SDK=66K iOS SDK=927K，H5 SDK 4.57K，TV SDK

- 目前具有的独特特性：
 - **事件防火墙技术**：避免开发者埋点错误造成数据洪峰，在本地利用防火墙技术阻挡异常埋点
 - **云端互动旋钮技术**：本地采集策略，可以通过云端配置参数控制本地情况，甚至控制停止数据采集
 - **CAE (Computer at the edge) 预计算技术**：根据特殊模型在本地事先计算，减少云端计算载
 - **H5混合APP事件采集技术**：面对越来越H5混合APP出现，采用混合APP采集技术，支持H5与APP在同样APP中做事件采集与分析
 - **混合云散列分享采集技术**：易观面对大量头部应用，采用散列采集，与混合云结合允许跨多个公有云多地区，针对不同应用分布式采用接口
 - **代码融合压缩技术**：所有功能通过代码极致优化与压缩，Android 66K
 - ...
- 目前具有的通用功能：
 - **应用数据采集**：采集打开关闭信息，统计PV，UV等
 - **自定义事件采集**：根据开发者定义事件，用于页面事件分析
 - **应用安装列表采集**：采集安装列表，用于用户画像
 - **应用打开关闭采集**：采集应用打开关闭信息，用于用户画像
 - **地理位置与传感器信息**：SSID，GPS等传感器信息采集，感知用户所处环境

- 精益化的修炼大数据
- 最小化的实时大数据分析闭环
- 新的技术框架迭代与扩展
- 高性价比的实时大数据框架

高性价比的架构——多一分则肥，少一分则瘦

对外服务

易观千帆

易观方舟

易观博阅

易观万象

内部平台

数据发布与订阅平台

数据治理服务

数据调度与资源

分布式存储与查询平台

分布式实时处理

统一任务调度

通用数据查询引擎

元数据管理

分布式资源管理

基础数据服务

内部查询计算
Presto

批量查询
Hive

Ad-hoc查询
Greenplum

数据口径

数据服务监控

HDFS

实时统计
Storm

挖掘引擎
Spark

数据质量检测

分量统计汇总

数据安全鉴权

分布式数据流转平台

数据审计

分布队列监控与配置

分布队列组件

分布式小数据总线

实时数据队列

云端数据接收群
SDK+LVS+Nginx+Netty

云端离线数据接收群
SDK+LVS+Nginx+Netty

本地离线数据接收群
SDK+LVS+Nginx+Netty

本地实时数据接收群
SDK+LVS+Nginx+Netty

如何打造高性价比百亿级实时数据分析平台

——以精益创业理念修炼大数据分析平台

温室中再美丽的花，也不如野外的一根草

微博：William-郭炜



微信：guodaxia2013



实时分析驱动用户资产成长



- 易观千帆
- 易观万像
- 易观方舟
- 易观博阅