

优酷土豆大数据架构和数据平台发展

--合一数据，唾手可得

杨大海
2016年1月

合一数据平台现状

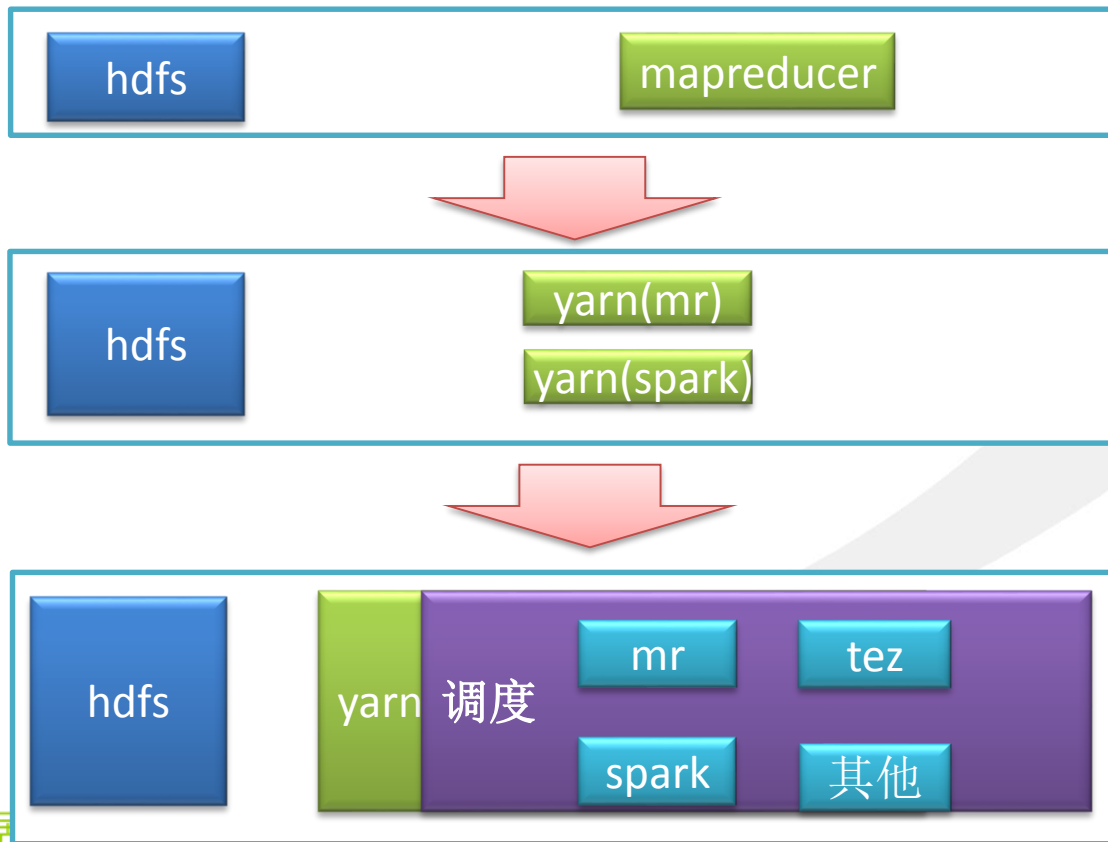
Hadoop体系发展介绍

数据平台整体架构

实时计算流程,产品

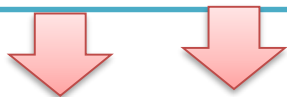
数据产品架构介绍

hadoop架构演变



hadoop规模发展

2012.一个团队在用（50台左右规模），裸奔无任何控制

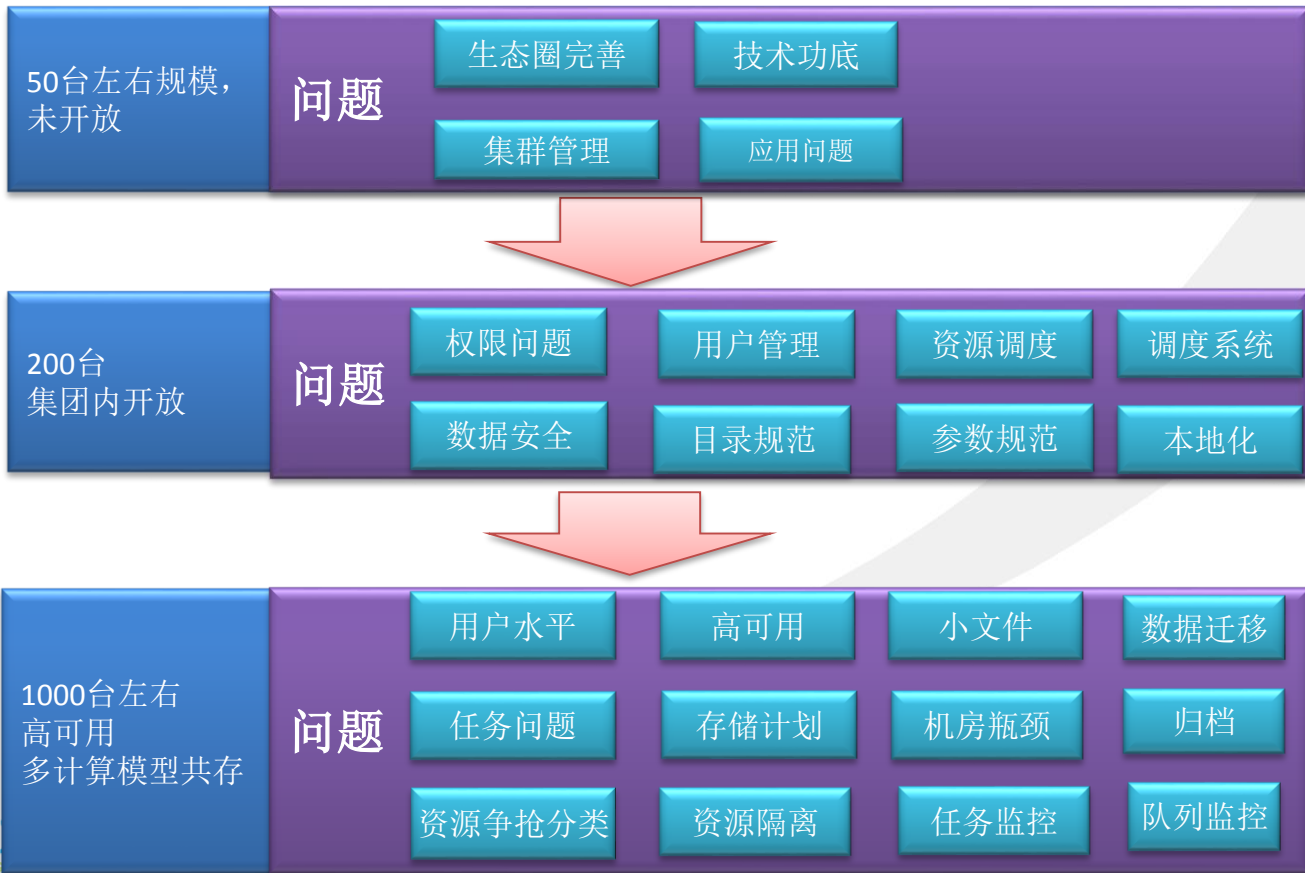


2013.平台集团内开放（200台左右），加入规范和权限控制

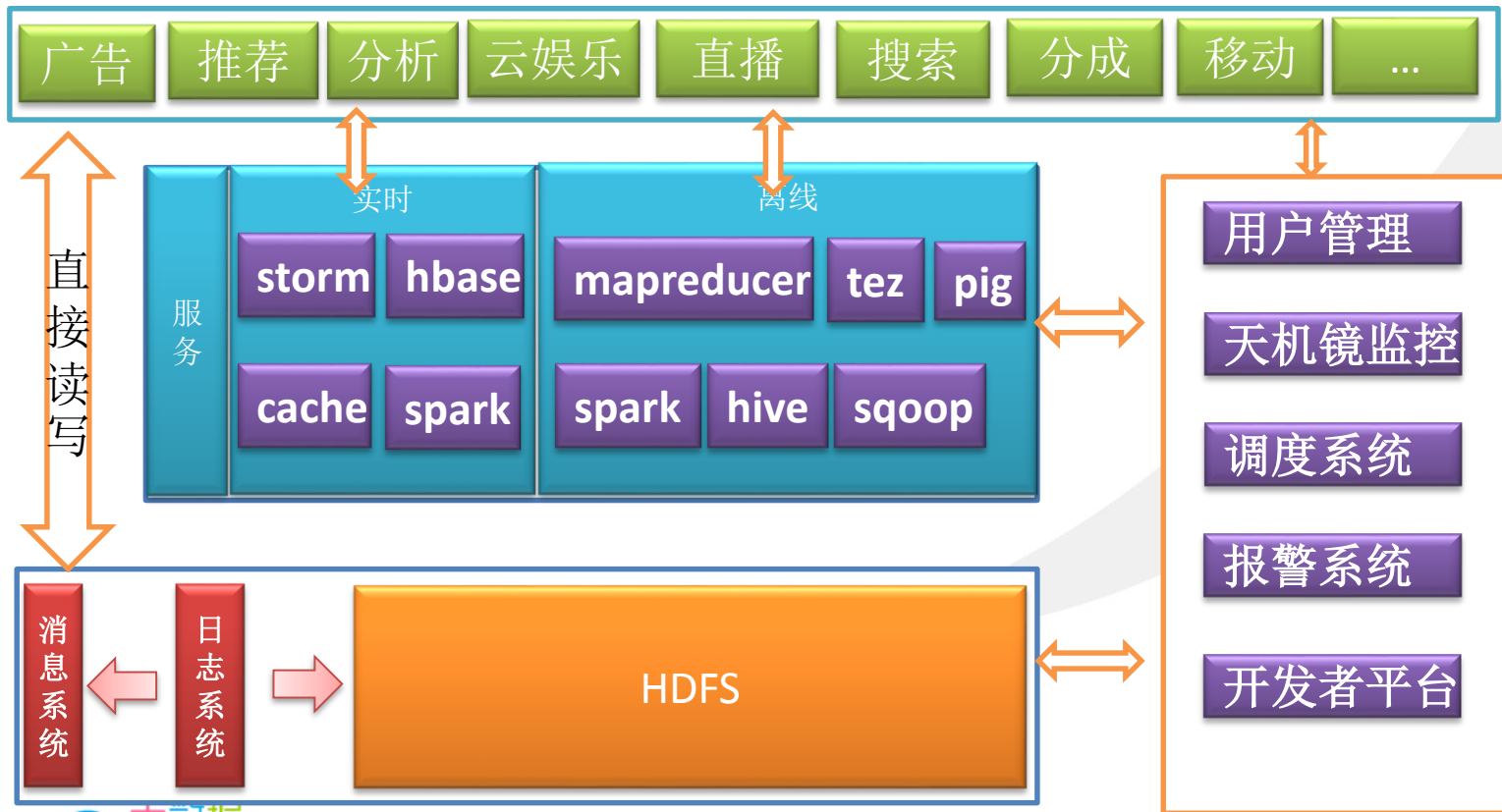


2016.多元化计算模型，几百个用户和1000台左右节点
（1000台左右）总存储：28P使用21P，集群cpu使用80%以上

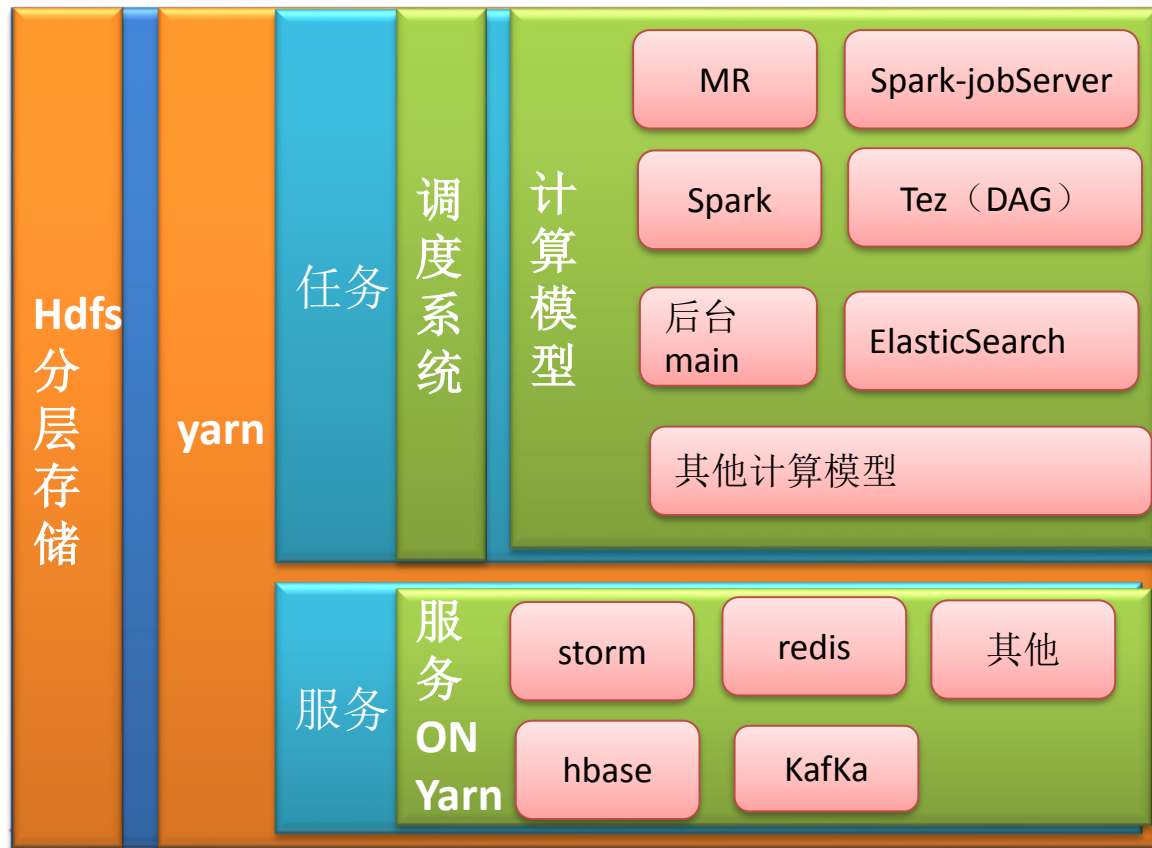
hadoop问题演变



数据平台现状



hadoop平台未来



hadoop平台挑战

Hdfs

Namenode内存瓶颈

日志系统对日志大小控制

节点块操作Api

多机房方案

集群规模太大namenode性能瓶颈

Yarn

调度个性化分类

资源隔离

数据仓库的必要性

基于标签调度完善

更加强大的监控平台

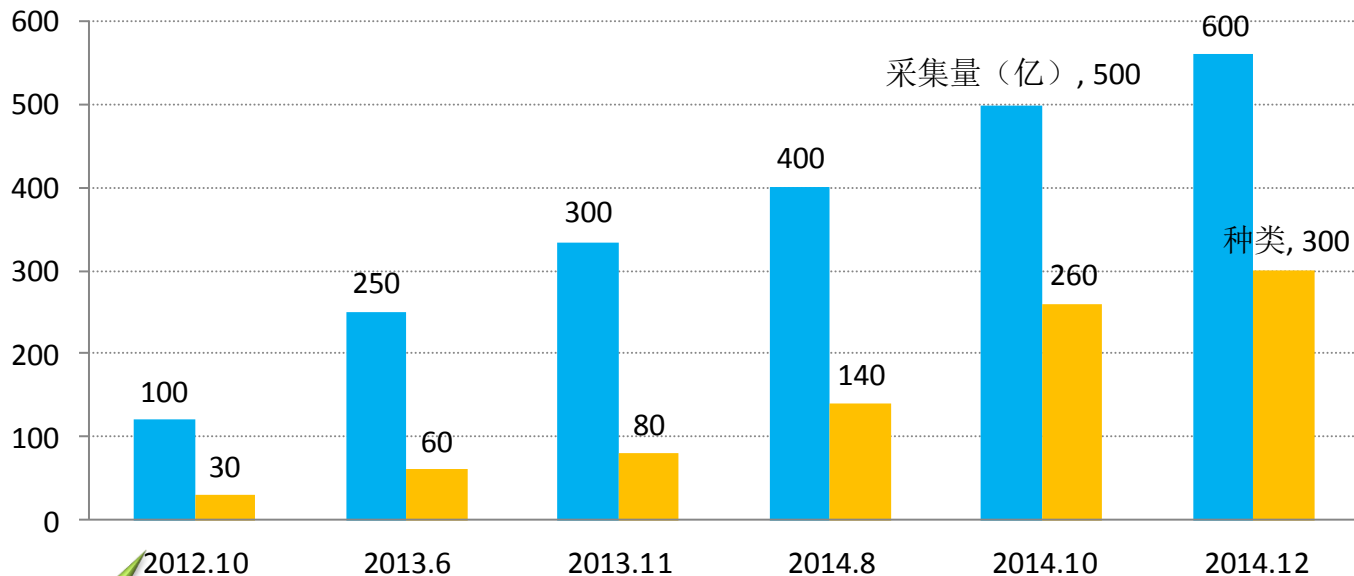
Client

Docker统一管理

配置问题

业务依赖升级问题

日志系统：采集量两年6番



300种600亿
日增20T;
周增3种;

新版上线,
统一优土
采集

开发完成连接
日志采集, 时长
采集更精确

发布v2版本,
性能提高,
多机房部署

来疯、移
动部分日
志上线

线上广告日
志全部切到
日志系统

移动API相关
日志正在切

实时计算架构流程



数据分析产品背景

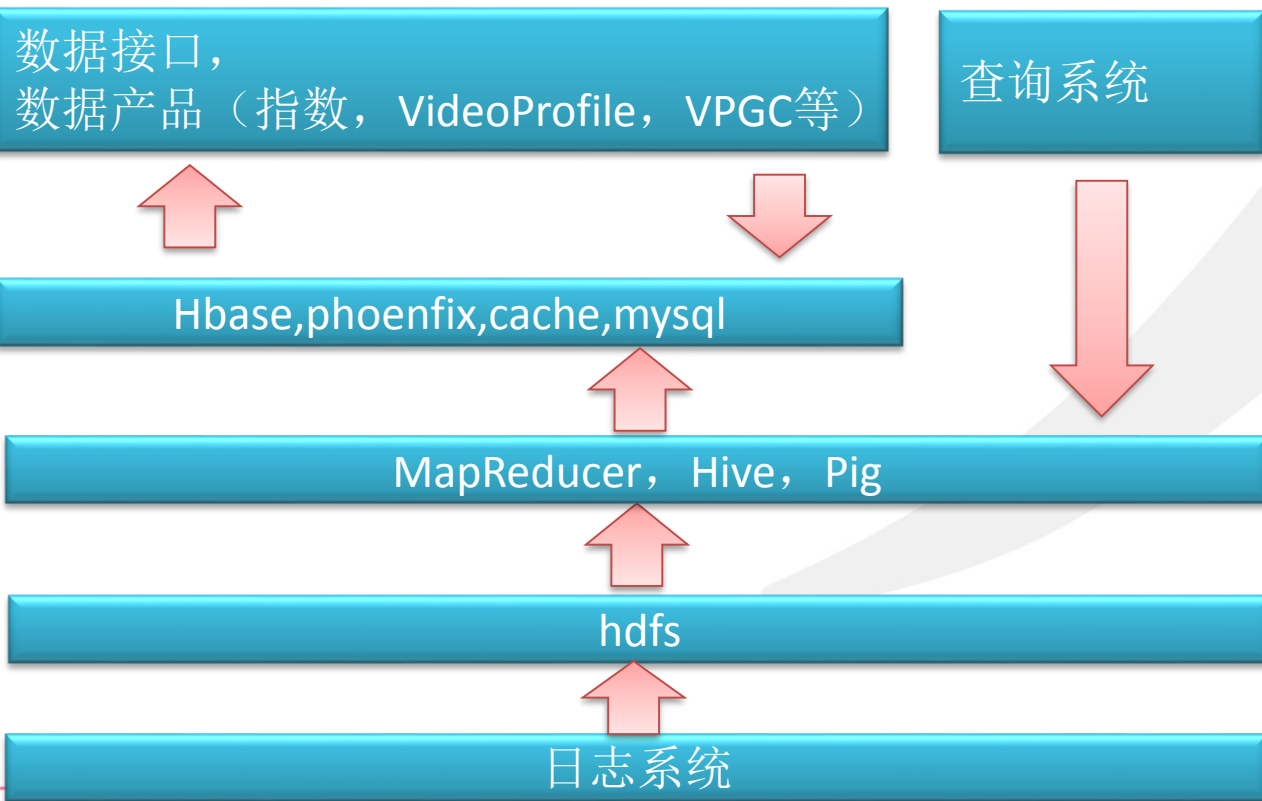
指数:index.youku.com全网视频榜单

VIdeoProfile 全站视频播放行为分析

Vpgc用户行为分析

查询系统，个性化用户需求处理

数据分析产品架构



数据分析架构-问题

1.缺少数据仓库：缺少数据仓库导致数据重复扫描集群占用资源后期无法控制的情况

2.数据接口 平台应统一包装：数据接口应有数据平台统一包装开放，业务团队只用申请使用，做到数据接口和数据产品分离（好处是数据接口共享，数据接口可以作统一的流量请求监控等）



2016

大型企业信息运维高峰会

Big Enterprise Information System Maintenance Summit

企业信息运维的互联网+

主办单位：黑龙江嘉和安泰商务服务有限公司

协办单位：国网黑龙江省电力有限公司信息通信公司 黑龙江省电力调度实业公司