



# DAMS

## 中国数据资产管理峰会

CHINA DATA ASSET MANAGEMENT SUMMIT

### 唯品会大数据实践



## CONTENT 目录

▶▶ 01 关于唯品会

▶▶ 02 数据平台建设

▶▶ 03 大数据应用建设

▶▶ 04 一些想法

## 数据平台实践

- 离线计算分析平台演化
- 实时计算平台演化
- 一些技术选型和经验

## 数据应用实践

- 系统开发和运营
- 业务和产品运营
- 恶意用户识别/风控系统
- 商品品牌推荐
- 个性化排序

## | 产品

数据仪表盘、数据魔方、比价系统、地图服务等

## | 算法

选品、分仓与预调拨

精准推荐

基础算法库

## | 数据

唯品会用户画像

细分人群

用户Lookalike

## | 系统

数据实时接入

离线计算平台

实时计算平台VRC

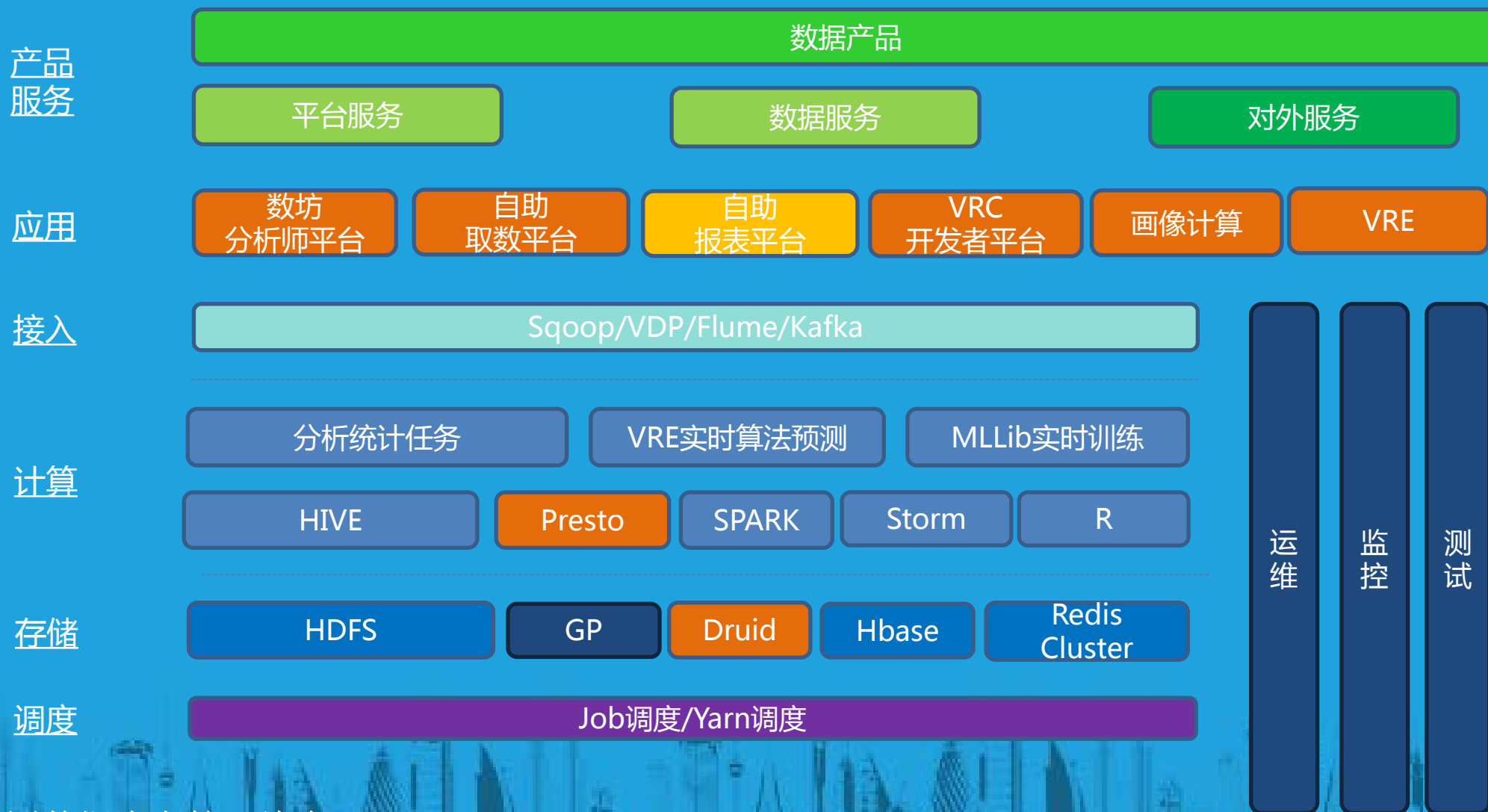
资源管理平台

运维

监控

测试

# 系统-大数据基础平台规划



# 数据平台的建设

- 离线计算分析平台选建设
  - 混合平台：Hadoop+Greenplum
  - 迁移策略和计划
  - daily job, hourly job, min job
  - 扩容，扩容，扩容
  - 离线和实时的混合
  - 开放平台
- 实时计算平台的建设
  - Binlog2Kafka →VDP
  - MySQL2Kafka
  - Spark vs Storm
  - Redis Challenge
  - 稳定性挑战
  - 开放平台
- 碰到的问题

# 离线平台的演化-1

- 2012 年底：CDC调度+GP10节点 系统稳定
- 2013 Q1：CDC调度+ETL Gp + Query Gp, Tuning
- 2013 Q2：
  - 自有调度平台开发 + 自有抽取系统+
  - Hadoop 流量开始迁移 +
  - GP交易数据 + Query GP
- 2013 Q3：
  - 自有调度平台+抽取迁移
  - Hadoop流量迁移结束 ( 70 ), 交易数据迁移开始
  - GP交易数据+Query GP
  - 核心数据小时级ETL
- 2013 Q4
  - 元数据管理系统，数据质量工具
  - ETL Gp完整迁移开始
  - Query GP扩容40节点
- 2014 Q1
  - 全部ETL@Hadoop
  - ~200 nodes cluster + 40 Ad-Hoc EDW
  - Hybrid node configuration

# 离线混合平台-2

- Referene:
  - Netflix, LinkedIn, eBay
- GreenPlum + Hadoop
  - 保护现有投资
  - Hadoop 海量数据分析
  - ETL复杂计算
  - 权限打通
- Greenplum :
  - GP擅长adhoc query速度快，
  - 分析师适应
  - 不足够scalable
  - 长期成本
- Hadoop
  - Massive scalable，但是单个查询慢
  - 海量ETL计算
  - Web查询

# 离线开放平台-3

- 开放平台
  - 自助ETL开发
  - 自助报表开发和展现
  - 自助取数分析
  - 成本breakdown, charge back
- 性能，实时，扩展性，成本
  - Presto
  - Druid

# 实时计算系统架构



# Hbase vs Redis

- 背景：
  - 个性化user profile, high QPS, very time sensitive
  - 用户信用体系user profile ,low QPS, non-critical
  - 用户实时浏览，订单历史，high tps, high qps
  - 都是海量数据
  - 看上去Hbase更加合适，但是不放心
- 选择：
  - Critical 的Redis
  - Non-critical 的Hbase
  - 积累经验，逐渐往Hbase dual write
  - 其实Hbase也不便宜，就是scale不动系统
  - Redis某种程度上也可以实现

# Redis

- Storm计算用redis保存中间和结果数据

- 流量一直增加
- 大促流量狂涨
- 计算复杂度一直增加
- 不停拆分。。。
- 每次改代码

- 怎么办？

- 逐个模块拆分
- 一开始就按模块写不同instance
- 一开始就Shard
- Twemproxy
- 优化数据结构
- Pipeline/Batch
- 不求100%准确hll log
- Redis Cluster

# Challenge

- 实时计算作为平台
- 离线和实时的融合
- 离线向实时的迁移成本

# 应用实践

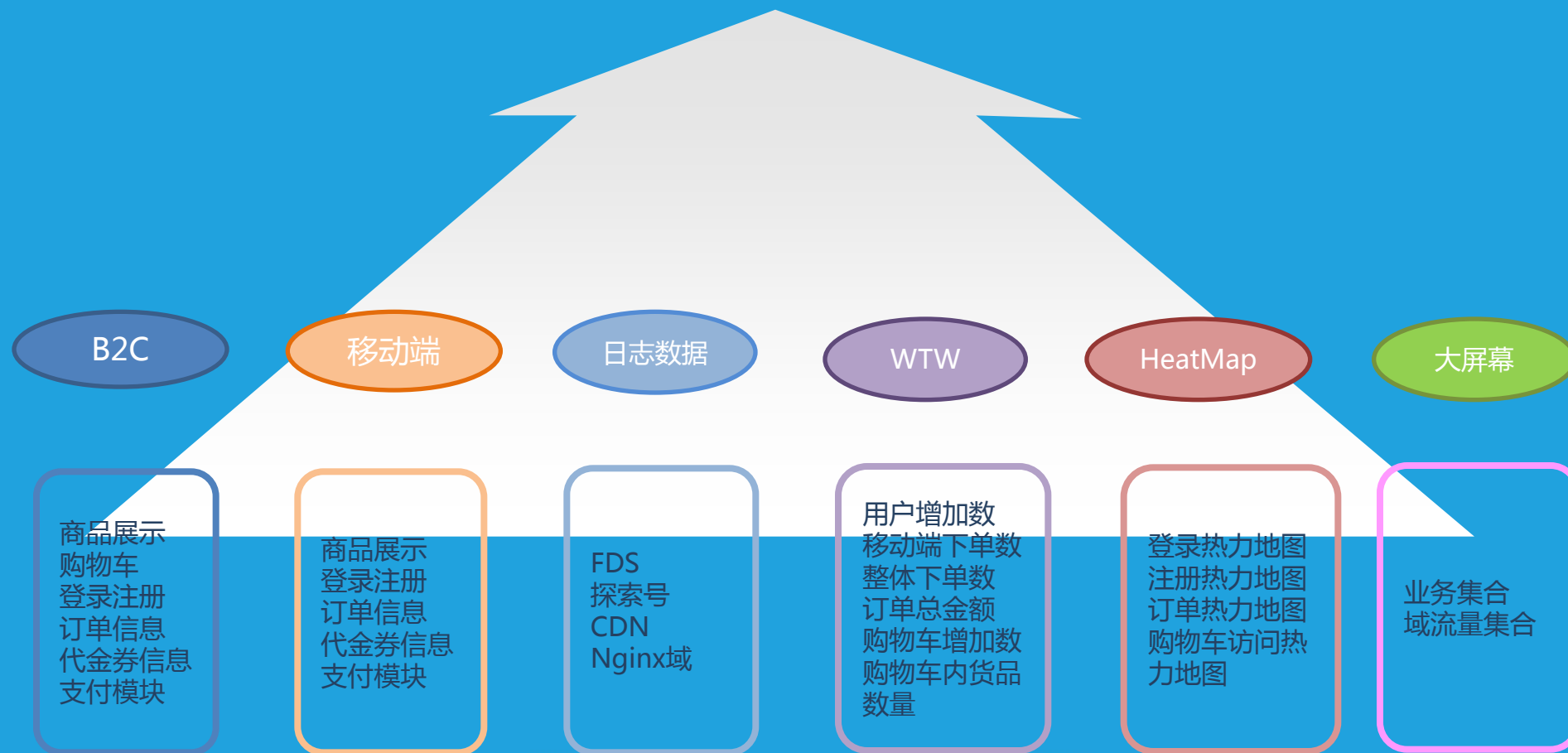
- 业务应用
  - 运营分析
  - 帮助公司买
  - 帮助公司卖
- 技术开发和运营
  - Telescope 业务监控(storm)
  - Logview/Titan 服务监控(spark)
  - Application logging(Spark)
  - CDN日志分析 (Hive)
  - Site speed分析(storm)
  - 安全审计分析(impala/storm)

# 大数据对于技术运营

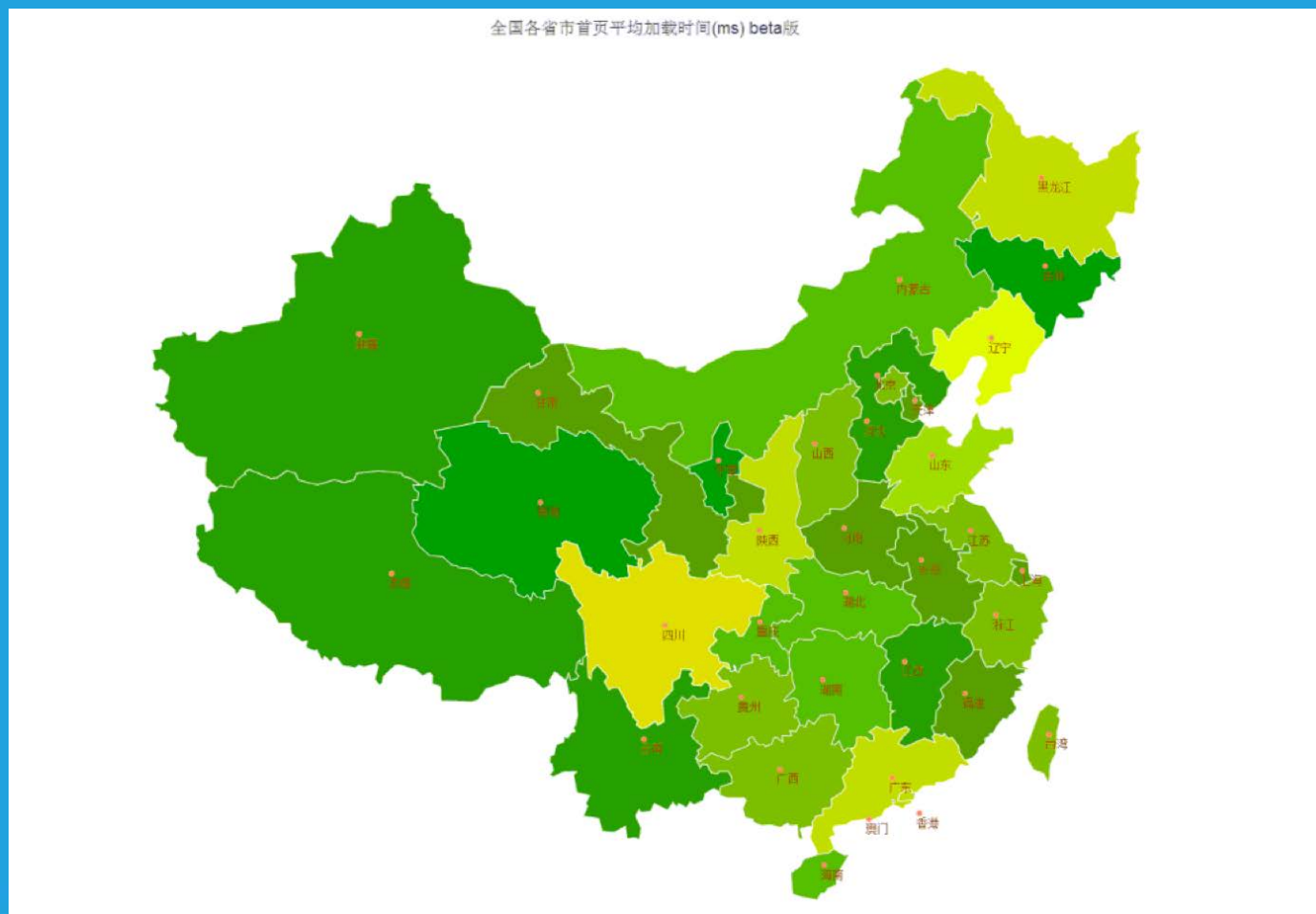
# 实时业务监控

## ➤ 现有平台

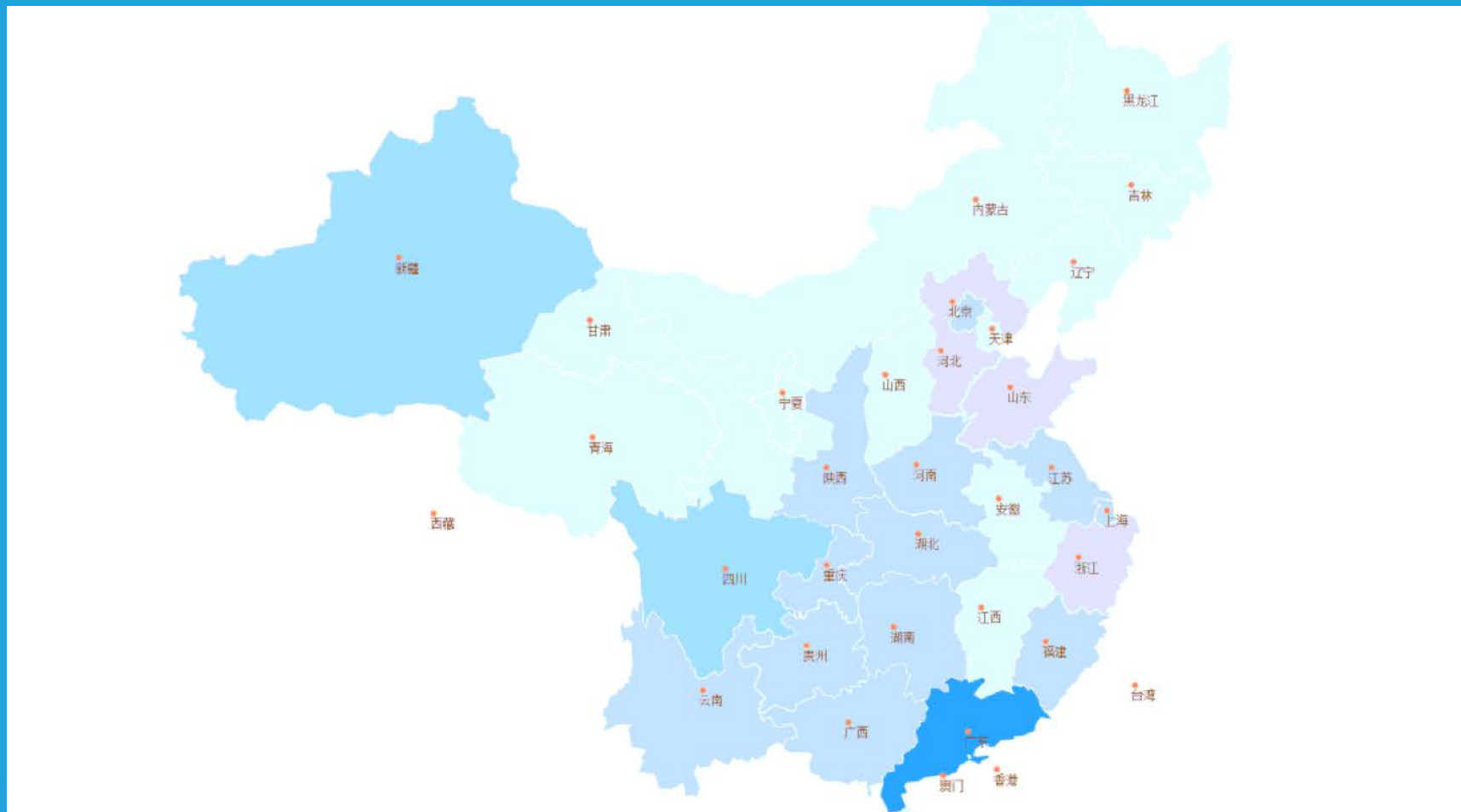
访问地址：xxxx.vipshop.com



# 实时页面加载时间监控



# 实时PV分布监控



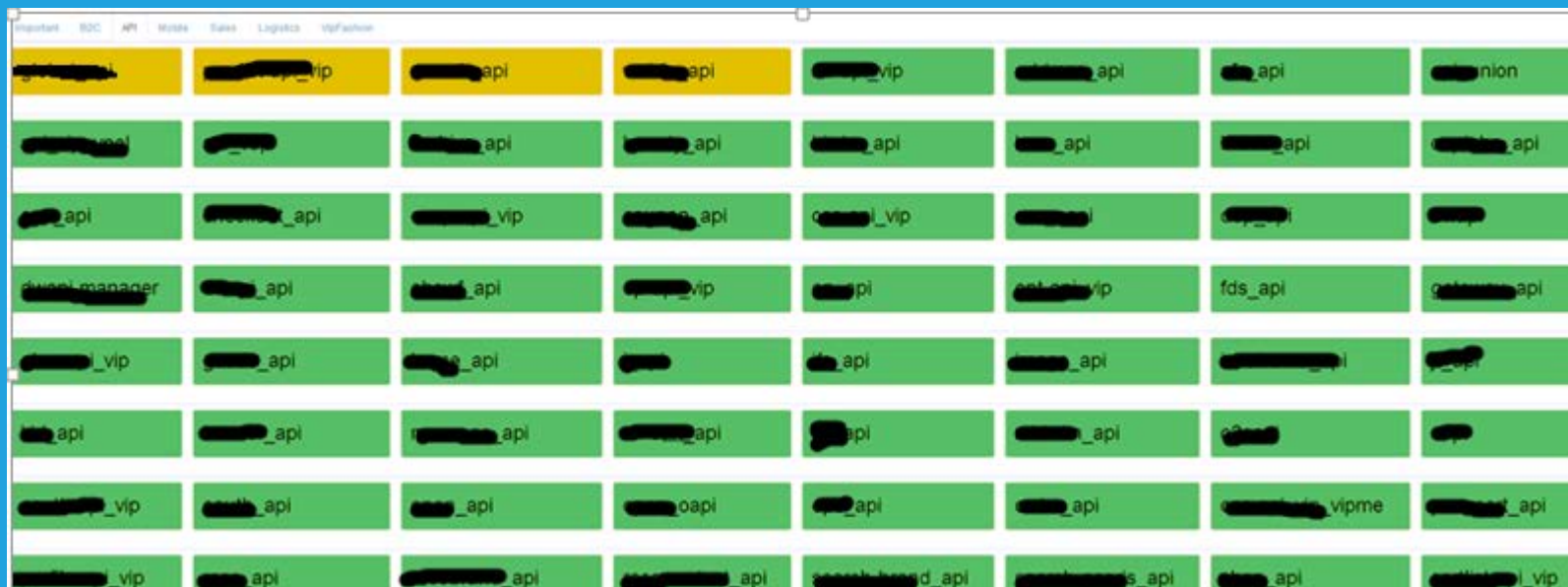
# 商业CDN质量分析

厂商	域名	总访问量	可用性(%)	平均响应时间	大于1s比例	404比例	历史数据查看
chinanetcenter	ipshop.com	289399266	99.9447	0.0695	1.0708	0.0376	>>
chinanetcenter	ipshop.com	134545159	99.7986	0.2984	4.0261	0.1781	>>
chinanetcenter	ipshop.com	70041140	99.9858	0.0858	0.6454	0.0052	>>
chinanetcenter	ipshop.com	49061096	99.9687	0.1596	2.1460	0.0004	>>
chinanetcenter	ipshop.com	21211867	99.9840	0.1576	1.1194	0.0046	>>
chinanetcenter	ipshop.com	16802384	99.9824	0.1718	1.2501	0.0020	>>
chinanetcenter	ipshop.com	10873474	99.9646	0.1880	1.3357	0.0080	>>
dnion	ipshop.com	9930950	99.7703	0.5040	4.5622	0.0208	>>
chinanetcenter	ipshop.com	9044900	99.9840	0.0928	0.7620	0.0000	>>
chinanetcenter	ipshop.com	7135595	98.8645	0.2824	3.8467	1.0292	>>
chinacache	ipshop.com	5900711	99.7132	0.0495	0.2472	0.2640	>>
chinanetcenter	ipshop.com	4419063	99.9821	0.0695	0.5381	0.0000	>>
chinanetcenter	ipshop.com	4080247	99.5972	0.2397	2.1786	0.1266	>>
chinanetcenter	ipshop.com	3250084	99.9488	0.0687	0.8412	0.0121	>>
chinacache	ipshop.com	995471	99.9244	0.0074	0.0548	0.0309	>>
chinacache	ipshop.com	552782	98.5955	0.5161	6.7341	1.1858	>>
chinacache	ipshop.com	289916	99.6189	0.5581	9.0795	0.0000	>>
chinacache	ipshop.com	141017	99.7298	0.2177	2.7564	0.1163	>>
chinacache	ipshop.com	11023	97.7683	0.6400	8.1012	2.0321	>>

蓝汛 qos:

服务器IP	总访问量	错误访问量	错误比例	平均响应时间
219	731576	195	0.03	0.0871
	2198805	141	0.01	0.0192
61	560964	124	0.02	0.0969
	561022	122	0.02	0.084
18	386613	121	0.03	0.1635
11	1550742	117	0.01	0.0734
	775691	114	0.01	0.2394
11	1445115	109	0.01	0.0282
61	47359	108	0.23	0.2718
	657436	104	0.02	0.0326
	863049	102	0.01	0.047
1	338099	96	0.03	0.049
1	414780	90	0.02	0.0496
1	1074185	90	0.01	0.0297
	1317888	89	0.01	0.0252
	783104	86	0.01	0.0217

# App Service Quality



- 进去可以看到每个pool，每个服务器，每个url的请求次数，响应时间，错误率，在过去两周的各个维度的统计数据 and 曲线；
- 可以看到pool之间的互相调用关系；
- 全无入侵，应用上线即插即用；

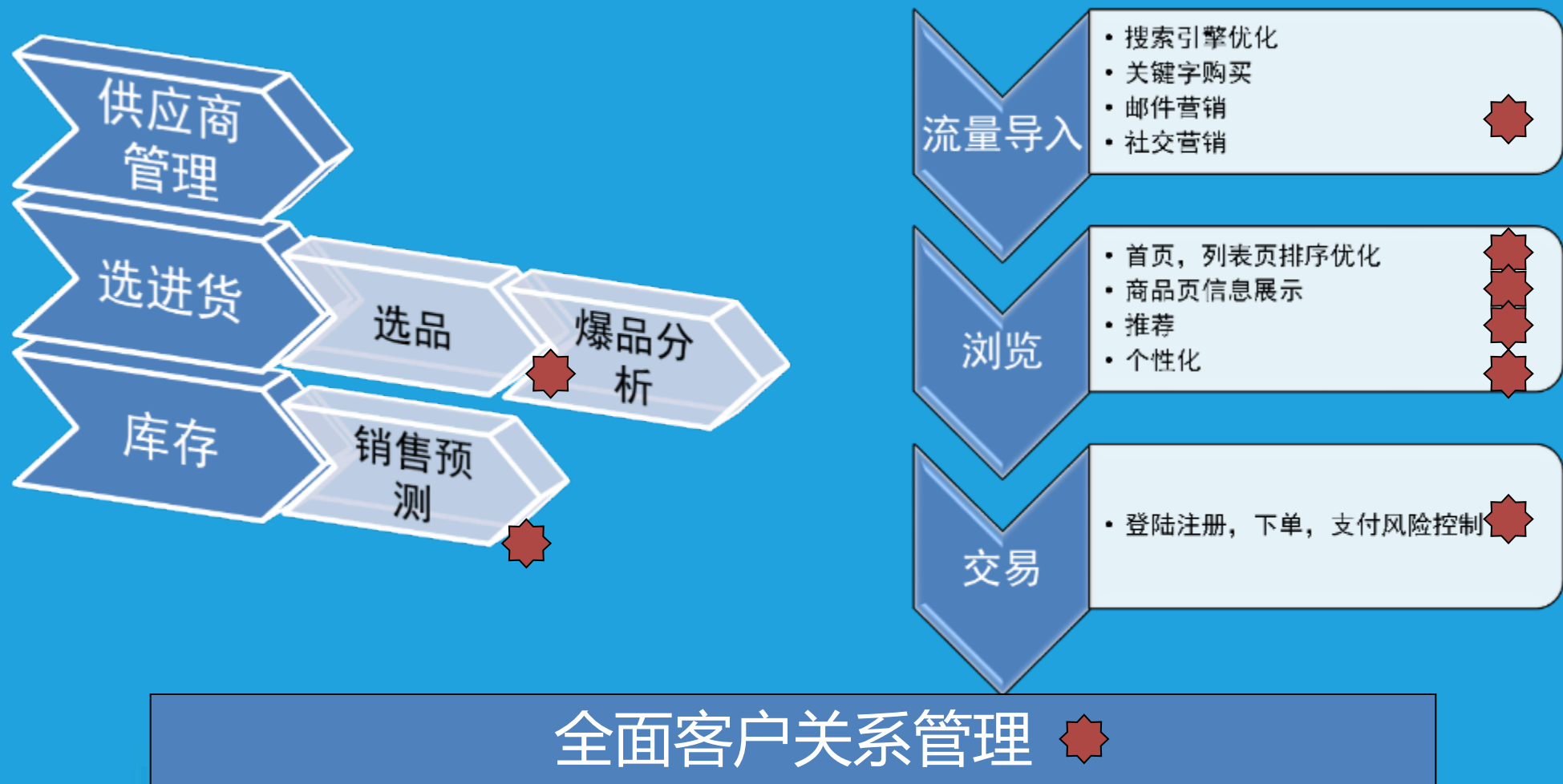
# Data Service Quality



# 大数据在唯品会特卖模式的业务价值

# 大数据对于数据化运营

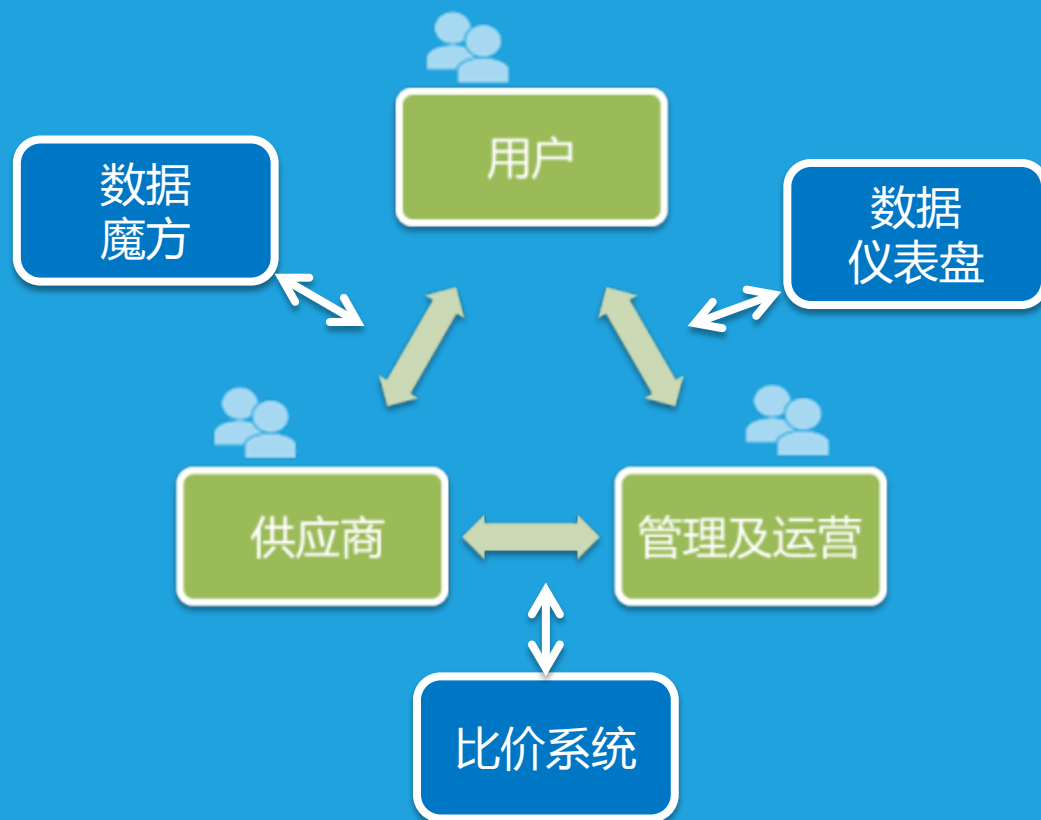
# 应用于唯品会



# 数据化运营-数据产品

- 对外：
  - 供应商：数据魔方
- 对内：
  - 高管：手机数据仪表盘，经营分析
  - 商务：选品，比价
  - 物流：分仓，预调拨
  - 产品/运营：指导产品分析和决策，经营分析，效果评估，产品优化
  - 金融：供应商贷款，
  - 消费者：个性化推荐，唯品白条
  - 营销：个性化EDM，个性化Push，CRM
  - 业务安全：风控

# 产品-数据产品及服务



## 数据仪表盘

打法一：  
数据从按天更新向实时化转变  
丰富数据可视化交互方式

## 数据魔方

打法二：  
合规前提下，开放更多数据给供应商  
丰富数据接口格式及实时性

## 比价系统

打法三：  
实时比价与价高告警  
比价数据与销售转化率数据关联分析

# 系统架构



# 挑战

- 用户
  - 数据稀疏，有效反馈少
  - 长尾严重
  - 用户体验，50ms返回
- ITEM
  - 冷启动
  - 特征难抽取，比如图片素材
- 场景
  - 缺少上下文
  - 没有明显意图，不同于“搜索”

# 底层数据

## 品牌

- 历史和实时销售数据
- 价格，品类，颜色尺码风格，季节
- 品牌相似性

## 商品

- 商品profile的长期开发
- 历史和实时商品信息（库存，销售，转化）

## 用户

- 用户点击浏览，购物车，购买，收藏行为
- 按品类，风格，价位，性别，尺码
- 用户实时行为路径

# 我们走过的路

- 2013Q4-2014Q1:基于人群分组和人工排序的个性化运营尝试
  - 人群划分
  - 首页人工排序
  - 列表页人工规则自动排序
  - 无效果。。。
- 2014Q2:开始有机会在小流量新版首页尝试技术主导
  - 机器学习+业务规则
  - 首页动态生成个性化推荐模块
  - 首页动态生成个性化排序页面
  - 提高了首页到列表页转化率，降低了跳出率，提高了销售

# 我们走过的路

- 2014 Q3-Now: 首页和列表页的个性化排序
  - 机器学习train model
  - Hadoop 生成 user profile/brand profile
  - Storm 计算实时转化销售数据，用户实时行为和意图
  - 实时排序首页和列表页
- 下一步
  - 更多引入个性化因子(feature)
  - 细化user/brand profile ,更多数据
  - 引入更多其他算法，做到算法可以灵活替代
  - 不但个性化排序和推荐，还可以有更多

# 个性化推荐下一个阶段

- 实时，实时，再实时
  - 实时计算商品品牌信息，用户profile
  - 实时推荐
  - 实时算法迭代更新
  - 实时A/B test verify
- 个性化，个性化，个性化
  - 移动天然是个个性化的好场所
  - 更多的个性化因子
  - 更加全面的数据：用户画像建设，曝光数据的收集...

# 个性化阶段性成果

## PC端

- 推荐：
  - 10%~12% PC销售占比
- 首页个性化排序
  - ~4%销售金额提升

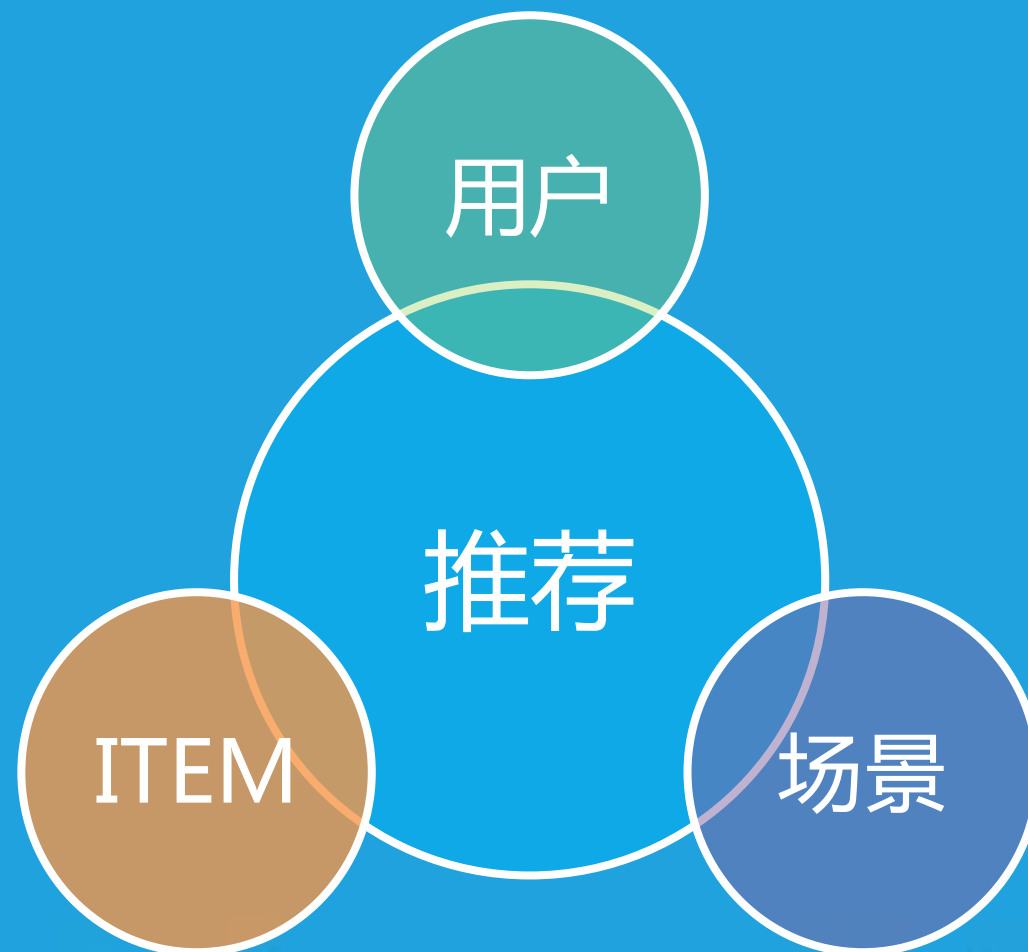
## 移动端(2014/12)

- 首页个性化排序
  - ~4%销售金额提升
- 列表页排序优化
  - ~15%销售金额提升
- Overall: ~17%

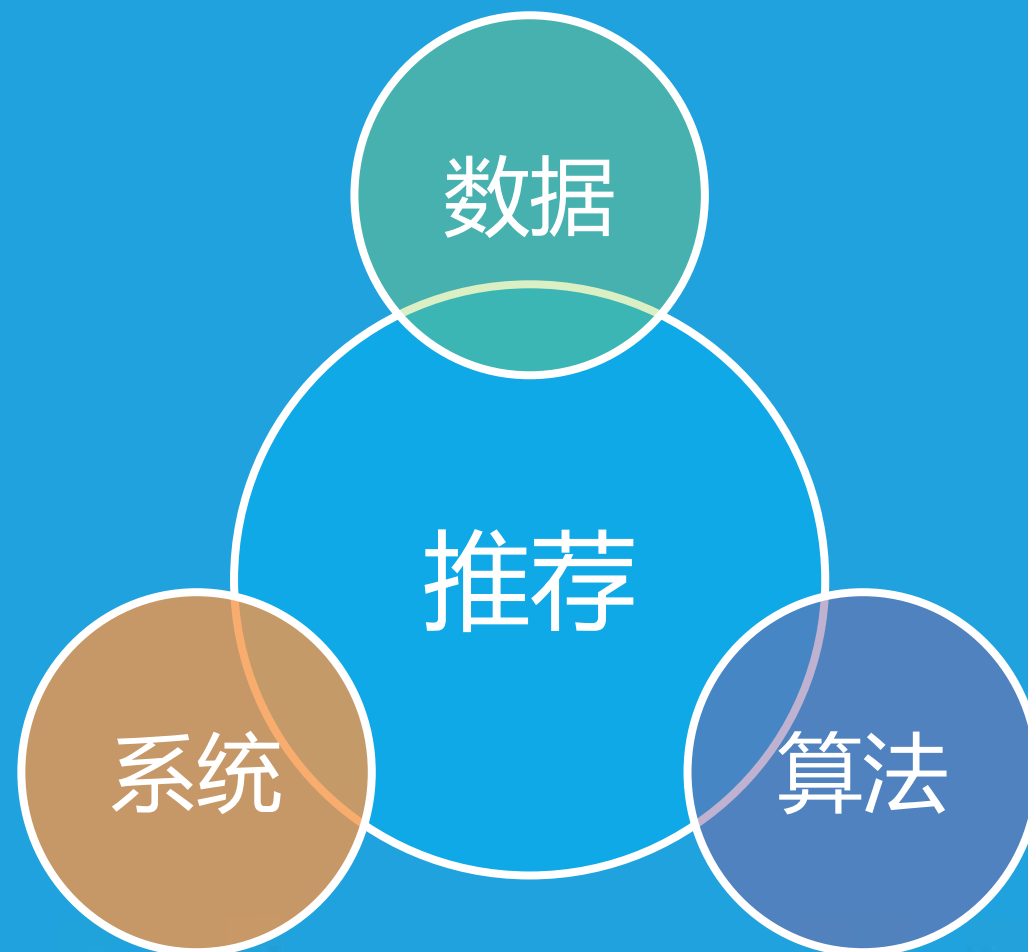
# 一些小结

- 技术选型：
  - 业界标准best practice
  - 成熟技术: 技术本身的成熟度，和我们队这个技术的把控力
  - reference customer/implementation
  - 用最合适的技术，而不是最先进的技术
- Don' t reinvent the wheel
  - 框架
  - 算法
- 基础架构/数据很重要
  - 模块化
  - 通用化
- Things Change  
半年前不用，可能现在用；  
Spark , Hbase

# 推荐关键点



# 解决之道





# DAMS

## 中国数据资产管理峰会

CHINA DATA ASSET MANAGEMENT SUMMIT

# THANK YOU