


Elasticsearch在爱奇艺用户画像系统的应用

杜益凡
爱奇艺



战略级赞助商  HUAWEI

钻石级赞助商  普翔

白金级赞助商  华夏博格

 神州数码
Digital China

金牌级赞助商  iDataAPI

合作伙伴  开源中国
oschina.net

 掘金

 知乎

 IT大咖说

 otpub

 Broadview
www.broadview.com.cn

 百格活动
bagevent.com

 MAXHUB
高效会议平台

Elasticsearch 在爱奇艺用户画像系统的应用

爱奇艺技术产品中心

杜益凡

2018.11



Elasticsearch在爱奇艺



集群数量：350+

物理机：500+

虚拟机：2500+

节点数：4700+

主要版本：2.3.2&6.0.0



Elasticsearch在爱奇艺的应用



日志分析

Venus日志收集分析系统

作为数据库

视频、音频检索，AI业务等



搜索引擎

MySQLIO等

数据分析

脸谱、画像系统等



用户画像

用户行为数据抓取、规则统计、分析预测获得的**用户特征集合**。
通常是将用户的社会属性、行为特征、潜在需求进行抽象提炼，
形成的动态**用户标签体系**。



汽车行业 用户画像					
用户基本属性	用户关联关系	用户兴趣偏好	用户价值信息	用户风险信息	用户营销信息
人口统计学 ▶ 姓名 ▶ 身份证号 ▶ 手机号 ...	生活关联关系 ▶ 家庭关系 ▶ 是否有子女 ▶ 同事关系 ▶ 朋友关系 ▶ 社区生活圈子 ...	车辆驾驶偏好 用户在本地的驾驶偏好 ▶ 目的地偏好 ▶ 驾驶习惯 ...	用户自身价值 用户自身的价值 ▶ 是否有房、位置 ▶ 房的大小、位置 ▶ 年收入区间 ▶ 是否企业高管 ▶ 是否有其他车辆 ▶ 是否有店外维修 ...	用户风险评价 从汽车金融维度对用户的风险进行评价 ▶ 综合征信评分 ▶ 信用风险等级 ▶ 洗钱风险等级 ▶ 综合授信额度 ▶ 信贷违约记录 ▶ 拖欠缴费记录 ▶ 还款能力 ▶ 违约概率 ...	近期需求信息 客户近期的需求 (包含汽车+非汽车) ▶ 近期是否准备结婚 ▶ 近期是否生小孩 ▶ 近期是否换工作 ▶ 近期是否出行 ▶ 近期是否想换车 ...
生活信息 用户基本生活类标签 ▶ 用水、用电 ▶ 天然气使用信息 ...	用车关联关系 用户在本地上车的关联 ▶ 驾驶关联 ▶ 维修关联 ▶ 车友关联 ▶ 乘客关系 ...	非车辆驾驶偏好 用户的兴趣爱好 ▶ 喜欢高尔夫 ▶ 经常看评测类文章 ...	用户对企业的贡献 用户在购买、使用和维修我品牌车辆时带来的贡献 ▶ 贡献排名 ▶ 客户综合价值 ▶ 钱包份额 ▶ 综合成本 ▶ 业务复杂度 ▶ 业务支持度 ...	黑名单信息 ▶ 信用卡逾期黑名单 ▶ 小贷逾期黑名单 ▶ 欠费用户名单 ▶ 车险退保用户名单 ▶ 最高失信人名单 ...	营销活动信息 用户对营销活动、以及企业各类产品服务的关系 ▶ 忠诚度 ▶ 用户满意度 ▶ 用户流失概率 ▶ 营销活动接受程度 ▶ 营销活动活跃度 ...
位置信息 ▶ 家庭、单位地址 ▶ 一般生活半径 ▶ 日常用车路径 ▶ 航空航班记录 ...	社交网络关联关系 用户社交网络图谱 ▶ 粉丝数量 ▶ 是否加V ▶ 微信朋友圈 ▶ 社交网络影响力 ...	经销商接触偏好 用户接触经销商的行为 ▶ 手机使用频率 ▶ 到店访问习惯 ...			
自定义信息 不同属性的自定义标签 ▶ 白领 ▶ 高收入人群 ...		非经销商接触偏好 用户全网渠道偏好 ▶ 上网习惯 ▶ 上网时长 ...			

用户画像的用途



爱奇艺的用户画像产品矩阵

画像数据层	脸谱 – 用户标签体系库	用户标签体系、标签质量管理、用户个人画像
画像服务层	画像分析接口服务	PC WEB 播放器 – 专辑流量分析 爱奇艺号 – 数据增值服务：自媒体受众画像 爱奇艺年度账单、偶像练习生粉丝画像、会员月刊 精准广告产品、广告DMP
	脸谱服务 – 精准广告人群	
	脸谱服务 – 个性化推荐	奇异果 长视频推荐
	脸谱服务 – RPC服务	商城商品推荐、奇秀主播推荐、爱奇艺文学...
画像应用层	达芬奇-爱奇艺画像系统	定制化人群、定制化画像分析
	后羿 – 精准营销系统	定制化人群、快速画像预览服务
	运营分析系统	爱奇艺指数、运营系统、内容消费指数

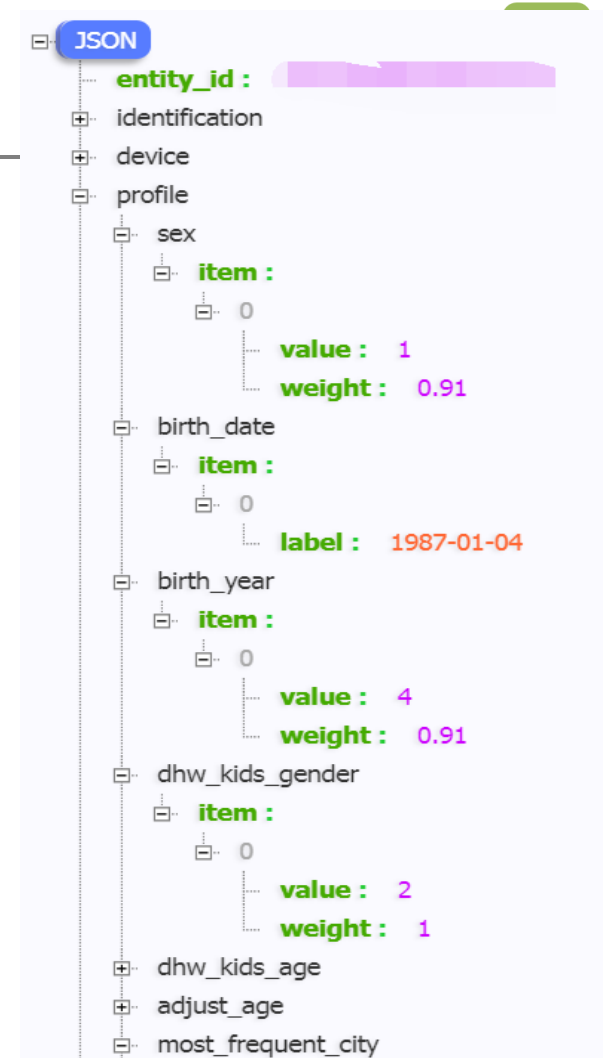
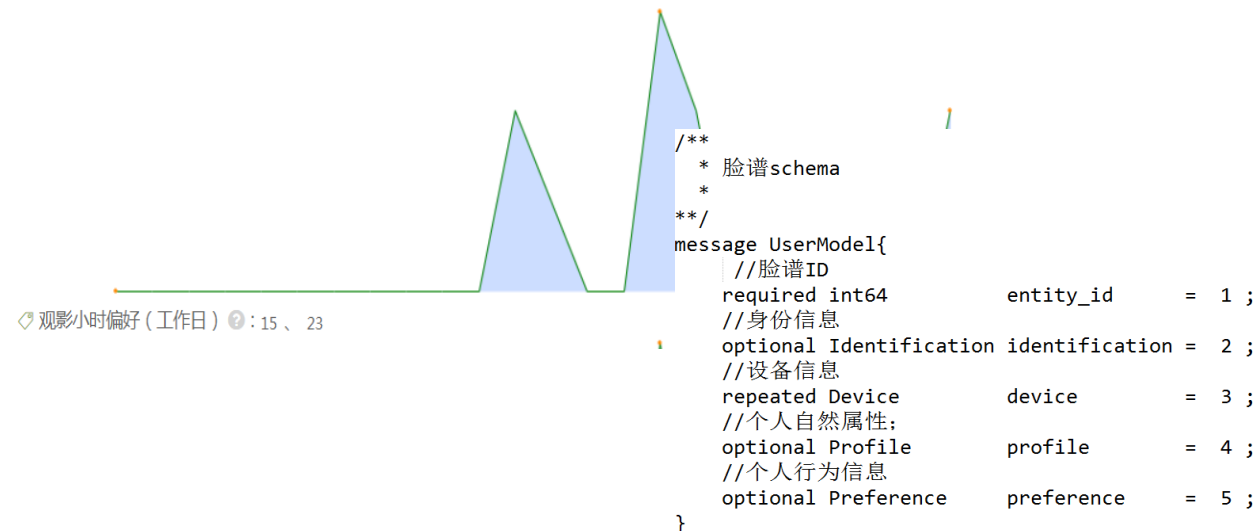
爱奇艺用户标签体系结构



爱奇艺用户标签数据结构

行为偏好

- 用户忠实度：7.50%
- 剧集偏好：画皮、奇葩说第5季、汪汪队立大功 第四季、爱情公寓、好戏一出、梦想改造家第5季、凉生我们可不可以不忧伤
- 明星偏好：马东、高晓松、蔡康永、李诞、薛兆丰、papi酱、施琰、骆新
- 视频tag偏好：爱奇艺出品(综艺)、奇葩说(综艺)、脱口秀(综艺)、马东(综艺)、高晓松(综艺)、蔡康永(综艺)、爱情(电影)、华语(电影)、7-10岁(少儿)、内地(电影)、...
- 频道偏好：综艺、电影、少儿、片花、娱乐
- 观影小时偏好 (General)：15、16、23、11、12



爱奇艺用户标签的生产



各产品的用户行为



用户标签



爱奇艺的用户画像数据流



爱奇艺及全网海量用户行为数据

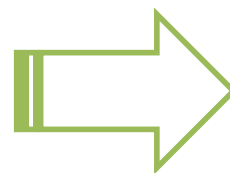
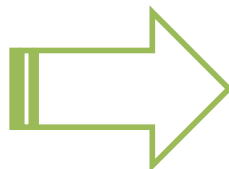
数据挖掘和建模，提炼用户标签

基于用户标签，输出定制化画像

6亿+ 月活用户行为数据
观影、直播、奇搜、会员
文学、ACG、商城 ...

近30亿用户设备，近3亿 用户账号
5TB+特征数据，250维+ 用户标签
涵盖自然属性、设备特征、兴趣偏好及
21个业务线行为标签

支持50+维特征的定制化分析
支持10+垂线业务的交叉画像

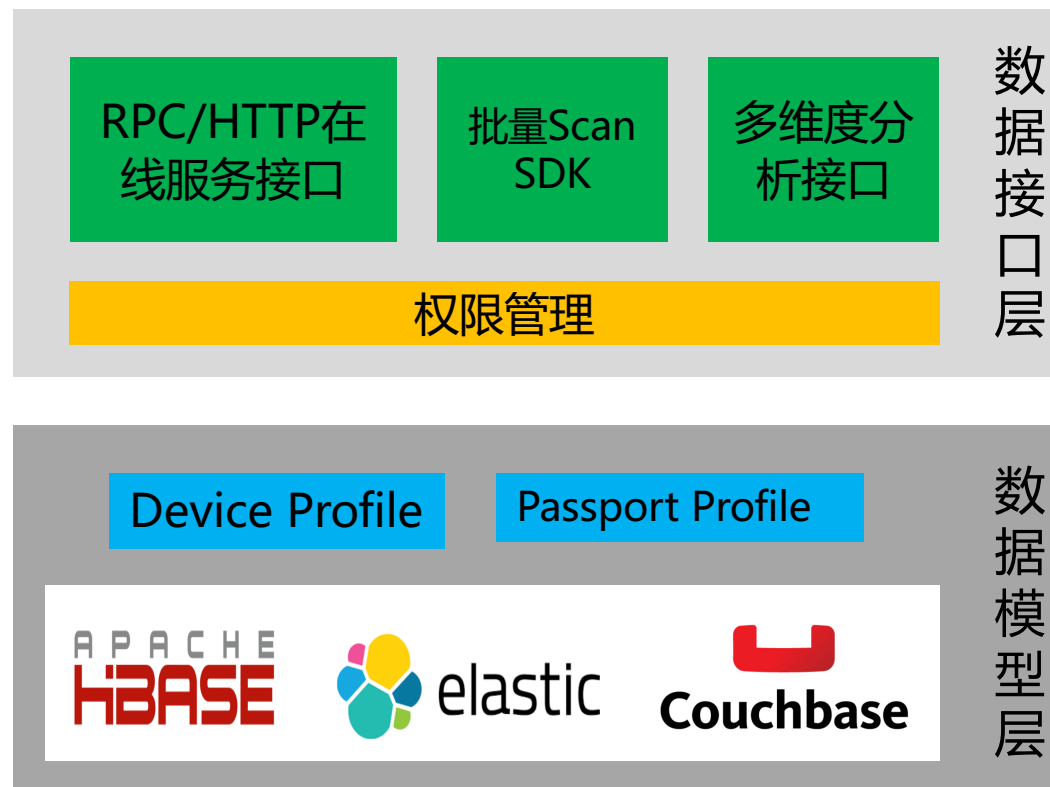


爱奇艺的用户标签体系库——脸谱

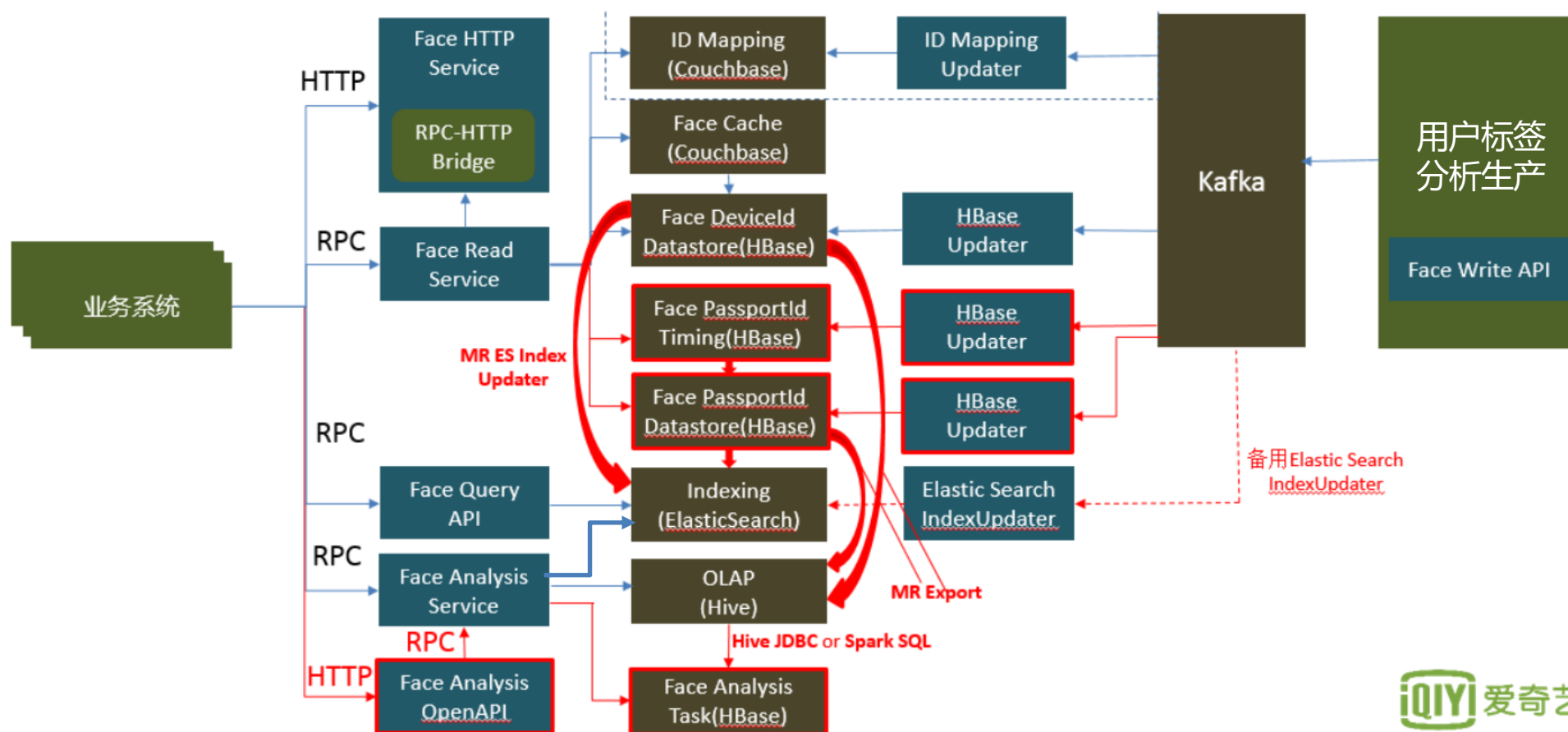


脸谱系统的架构

- 数据接口层
 - RPC/HTTP在线服务接口提供毫秒级接口服务用于支持个性化等在线服务
 - 批量Scan接口用于批量获取用户特征数据来进行用户画像
 - 多维度分析接口主要用于人群估量和简单的分布统计
- 数据模型层
 - 设备维度&Passport维度
 - 存储基于Hbase和Elasticsearch
 - ES用于服务常规业务
 - Hbase服务底层大数据量分析
 - 使用Couchbase作为实时数据缓存



用户标签数据写入脸谱



爱奇艺用户画像系统——达芬奇

- 基于爱奇艺全站用户行为数据、脸谱-用户标签体系的**通用人群画像分析平台**，为全公司各业务线提供**极速、精准、定制化**的画像分析服务

画像系统V1.9新功能介绍

查询画像数据 ...

新建画像 任务管理 人群管理

热门画像 业务画像 频道画像 公共画像 更多>

社交 (泡泡) 画像

文学 画像

全站观影 画像

奇秀基线 画像

游戏 画像

电商 画像

头条 画像

奇搜 画像

爱奇艺号 画像

奇秀APP 画像

传片助手 画像

电视果 画像

电影票 画像

奇巴布 画像

我的任务 常看任务 更多>

预约节目 测试 0919 09-19

其他业务线画像查看 09-19

人民的名义iPhone 09-19

人民的名义安卓 09-19

测试 - TGI 修复 0919 09-19

延禧攻略iPhone画像 09-19

延禧攻略安卓画像 09-19

达芬奇画像解决的问题



《奇葩说》第四季的用户是什么样的？

2线城市女白领们在关心什么？

喜欢看《延禧攻略》的用户特征是什么

喜欢看《中国新说唱》的男性还看什么综艺

奇秀与漫画的用户重合情况

搜吴亦凡的更容易买VIP么？

《奇葩说》第四季的用户喜欢什么商品？

年卡半价吸引的这波会员都有什么特点？

华为和小米用户有什么不同？

纳逗APP的竞品都有谁？



达芬奇系统的核心功能



达芬奇系统的核心功能



登录账号人群

设备ID人群

通用

大陆会员

台湾会员

文学

漫画

头条

奇巴布

电影票

爱奇艺号

电视果

基本属性

行为偏好

消费特征

业务线行为

筛选方式

非重合用户

重合用户

交叉人群示意

业务线受众

业务线专属人群圈定

最终选定人群

访问行为

观影

奇搜

泡泡

观看弹幕

发送弹幕

文学

游戏

游戏SDK

游戏中心

漫画

头条

电商

奇秀基线

奇秀APP

游戏直播

奇巴布

电影票

电视果

VR

热点

收起>>

最后访问时间-奇秀基线

距今

小于等于

30

天

自定义日期

开始日期

~

结束日期

最后访问时间-游戏直播

距今

小于等于

30

天

自定义日期

开始日期

~

结束日期

所选条件

通用

业务线行为：

X 圈选方式

重合用户

X 访问行为

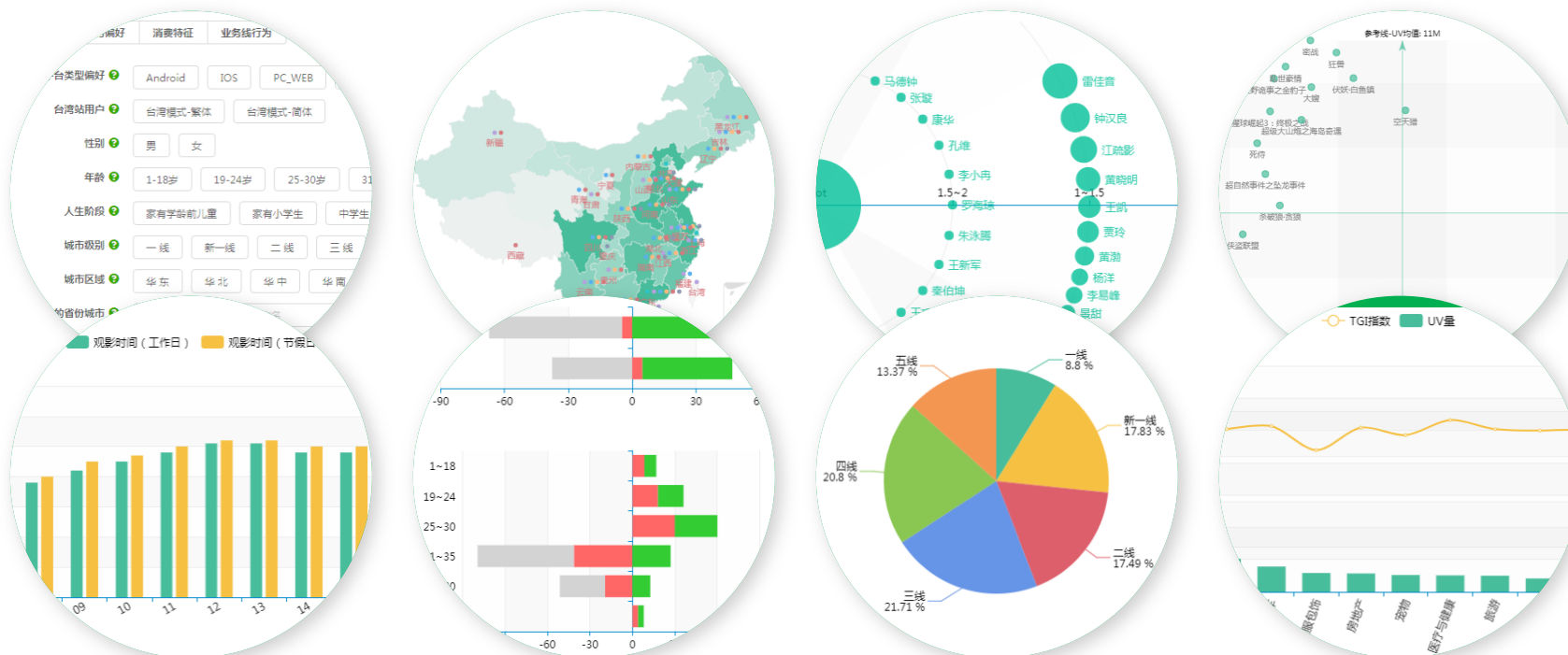
小于等于30天访问奇秀基线，小于等于30天访问游戏直播

人群实时预估

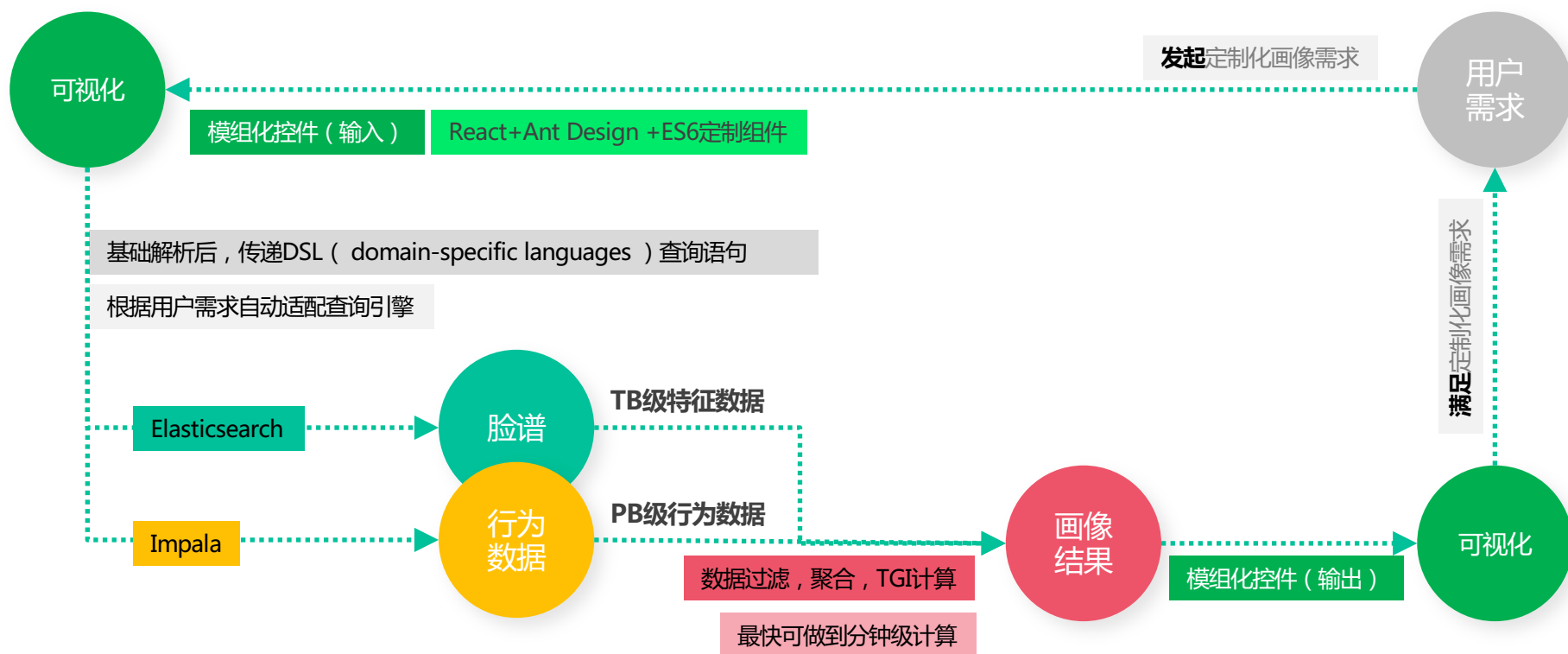
预估UV：



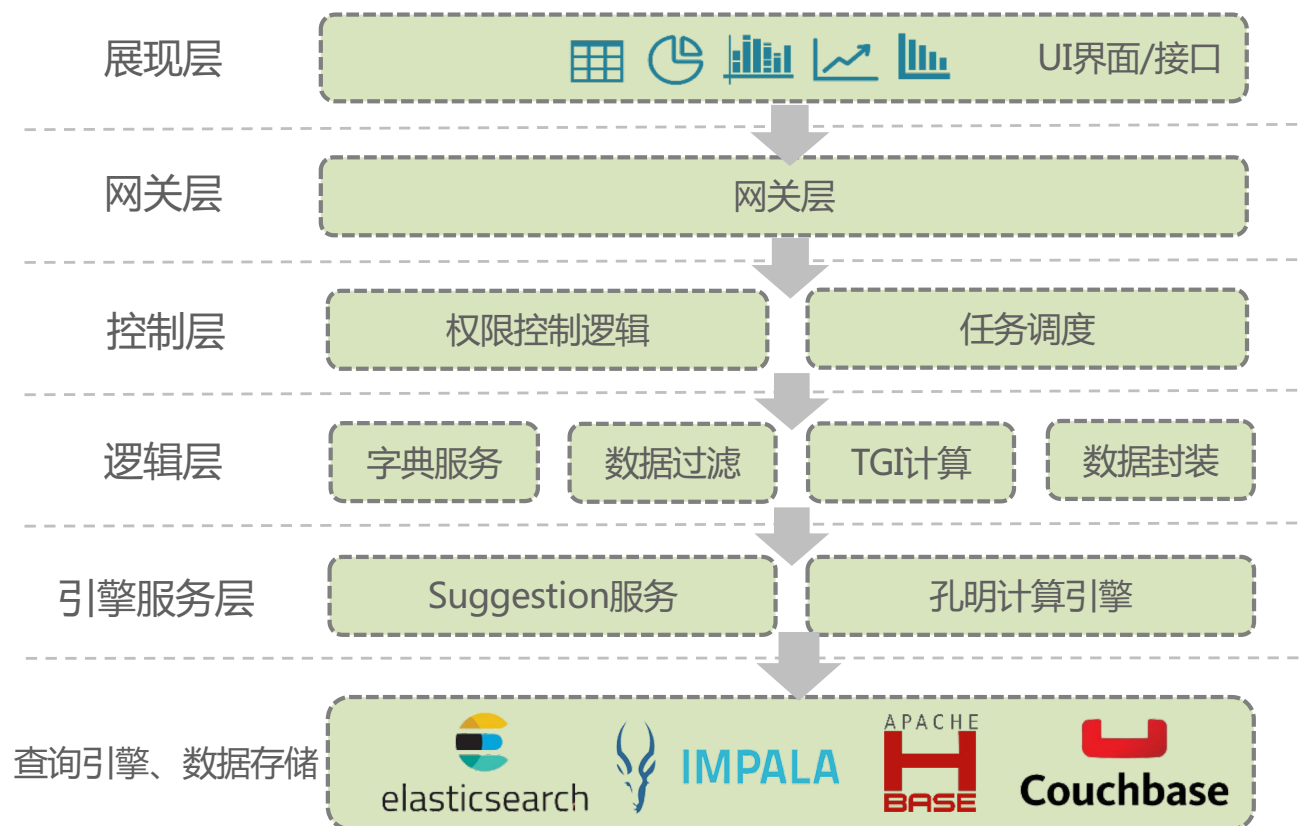
达芬奇系统的画像数据结果



达芬奇系统的技术逻辑

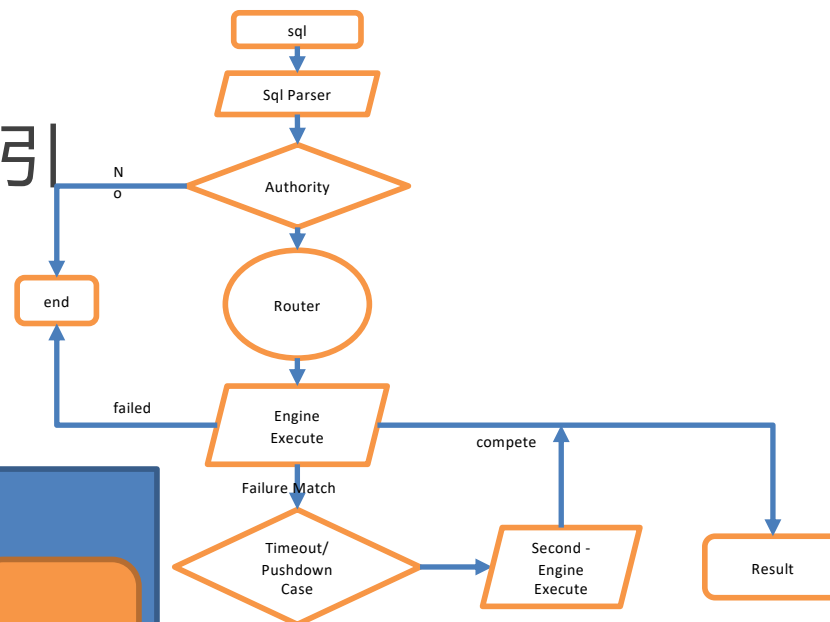
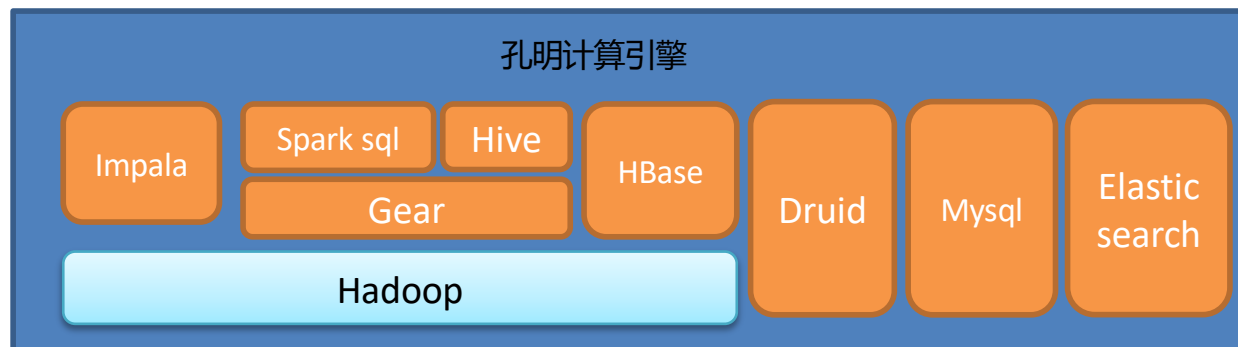


达芬奇系统的技术架构

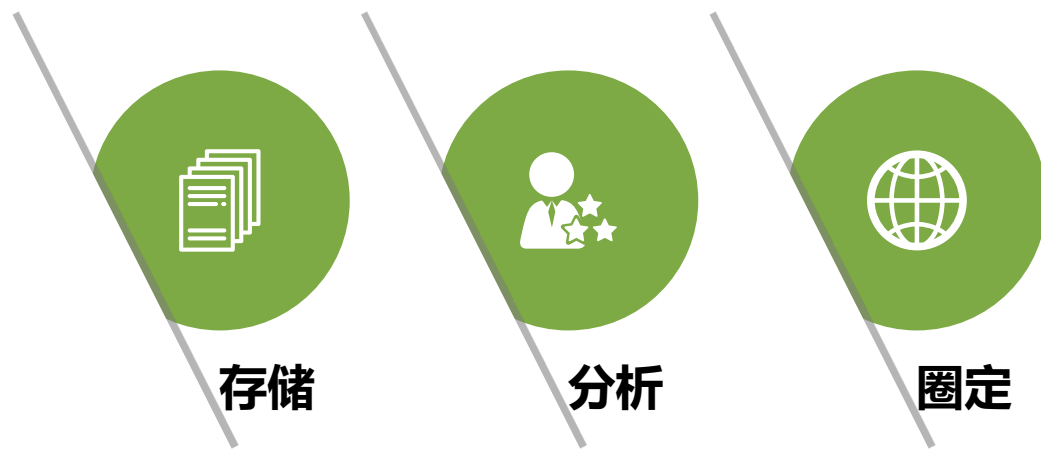


孔明计算引擎

- 统一查询接口的计算引擎
- 智能路由，选择最合适的实际执行引擎
- 智能下沉，确保查询可靠性



Elasticsearch在达芬奇系统中的作用



存储

存储特征标签并建立索引,来进行快速检索和 Suggestion

分析

使用检索统计功能对人群进行快速分析

圈定

使用聚合功能根据条件圈定人群



Suggestion服务



频道偏好 偏好程度: 普通 电影 电视剧 纪录片 动漫

视频偏好 偏好程度: 普通 xhs

明星偏好 偏好程度: 普通

标签偏好 请搜索具体视频标签, 例如: 爱情

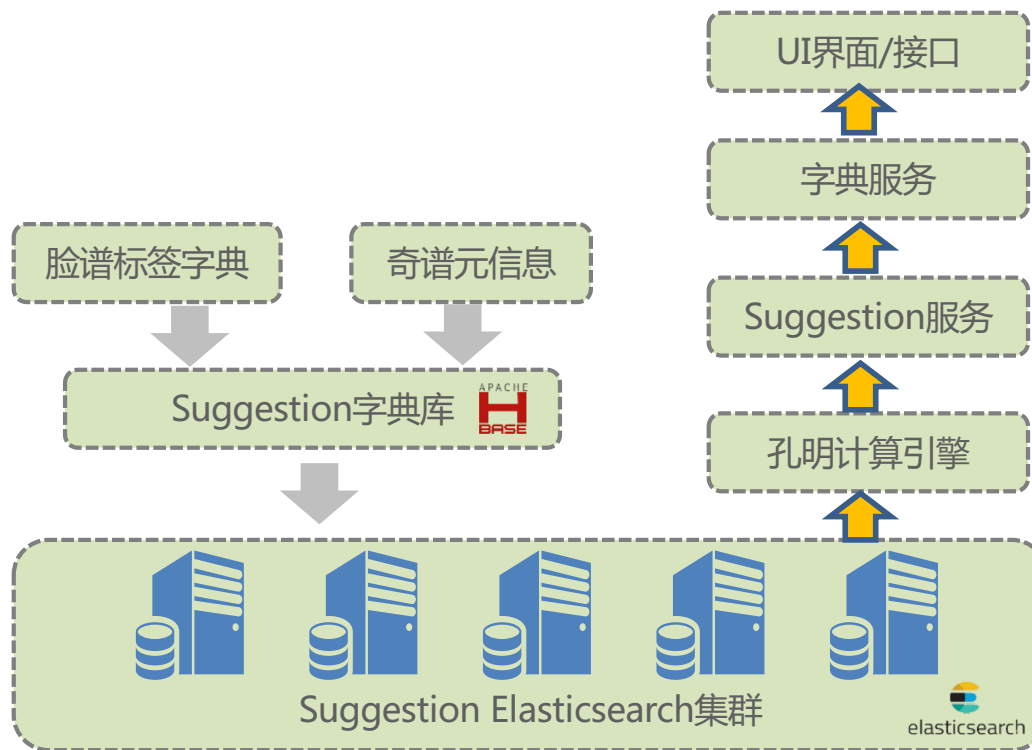
视频偏好 请输入视频名称或专辑ID, 支持最

西虹市首富(电影)

一禅小和尚(动漫)

生死狙击少云解说新号上传老粉丝请订阅(游戏)

在地下城寻求邂逅是否搞错了什么(动漫)



人群圈定与分析

基本属性：

X 性别 男

行为偏好：

X 视频偏好 偏好程度：较高，西虹市首富(电影)

X 明星偏好 偏好程度：极高，沈腾

人群实时预估

预估UV：



```
// build statement
FaceTermBuilder faceTermBuilder = FaceTermBuilder.newBuilder()
    .feature(UserModel.PROFILE_FIELD_NUMBER)
    .field(Profile.BIRTH_YEAR_FIELD_NUMBER).value(1).build(); // 年龄在0-18岁
FaceTermBuilder faceTermBuilder2 = FaceTermBuilder.newBuilder()
    .feature(UserModel.PROFILE_FIELD_NUMBER)
    .field(Profile.INCOME_FIELD_NUMBER).value(6).build(); // 收入在5001-8000
FaceBoolBuilder faceBoolBuilder = FaceBoolBuilder.newBuilder()
    .must(faceTermBuilder).must(faceTermBuilder2).build(); // and
// 关系，还支持or(should)以及not(mustNot)
FieldAnalysisBuilder fieldAnalysisBuilder = FieldAnalysisBuilder
    .newBuilder().feature(UserModel.PROFILE_FIELD_NUMBER)
    .field(Profile.SEX_FIELD_NUMBER).build(); // 相当于group by 性别
FaceAnalysisBuilder faceTermAnalysisBuilder = FaceTermAnalysisBuilder
    .newBuilder().faceBoolBuilder(faceBoolBuilder)
    .addField(fieldAnalysisBuilder).build(); // 相当于 select sex,
// count(entity_id) as
// count_id from
// gipw_face_table
// where birth_year=1
// and income=6 group
// by sex order by
// count_id;
// RPC request and response
FaceESAnalysisService faceAnalysisService = FaceESAnalysisService
    .newBuilder().dataCenter(DataCenter.JYLT).username("abc")
    .password("abcde").timeoutMs(10000) // 根据复杂程度设置超时时间
    .build();
```


画像系统所使用的Elasticsearch集群

- Suggestion服务集群+画像服务集群
- 5台+8台物理机
- 系统：CentOS6.4
- 内存：192G
- CPU：E5-2650 V3
- 硬盘：900G SSD
- JAVA版本：1.8
- ES版本：2.3
- 每日数据量：50T

一些应用实践

优先使用Java8，客户端和服务端使用相同版本JVM

需要创建大量连接就使用传输客户端

只是少数持久的对象连接到集群，客户端节点可以更高效，但是会造成耦合

线程池的线程个数根据CPU核数设置，最多两倍，再多也浪费

主分片数设置为data节点数的1倍

能用short就不用long

读写分离场景关闭sniffer





elastic 中文社区

专业、垂直、纯粹的 Elastic 开源技术交流社区

<https://elasticsearch.cn/>

