

神盾开放通用推荐系统

雷小平

qq:106093647

email:leixp3636@qq.com

腾讯 社交网络事业群 数据中心

QCon

2016.10.20~22

上海·宝华万豪酒店

全球软件开发大会 2016

[上海站]



购票热线: 010-64738142

会务咨询: qcon@cn.infoq.com

赞助咨询: sponsor@cn.infoq.com

议题提交: speakers@cn.infoq.com

在线咨询(QQ): 1173834688

团 · 购 · 享 · 受 · 更 · 多 · 优 · 惠

7折

优惠(截至06月21日)
现在报名, 立省2040元/张

个人团队介绍

个人简介

雷小平 QQ大数据团队平台组组长

推荐系统/分布式计算

团队简介

QQ的基础数据挖掘系统

产品应用系统的研发和运营

<http://www.csdn.net/article/2014-07-03/2820520>

大数据解决方案

神盾

– 开放通用推荐系统

ADS

– 数据集市解决方案(第29届中国数据库学术会议)

COW

– 分布式流数据存储系统

R²

– 分布式实时计算系统(2014年全球互联网大会)

LBS云

– 提供统一的LBS云服务

其他

– hadoop、storm、spark



目录

1. 背景介绍
2. 架构介绍
 - 2.1 架构总览
 - 2.2 分布式计算
 - 2.3 数据引擎
 - 2.4 海量画像
 - 2.5 实时ABTEST
 - 2.6 立体监控
 - 2.7 通用开放
3. 运营情况



目录

1. 背景介绍
2. 架构介绍
 - 2.1 架构总览
 - 2.2 分布式计算
 - 2.3 数据引擎
 - 2.4 海量画像
 - 2.5 实时ABTEST
 - 2.6 立体监控
 - 2.7 通用开放
3. 运营情况

1 背景介绍(1) - 有哪些产品?



关系链

文本/音频/视频

APP

特权

1 背景介绍(2) - 有哪些场景？



QQ



企鹅FM



全民K歌



QQ音乐

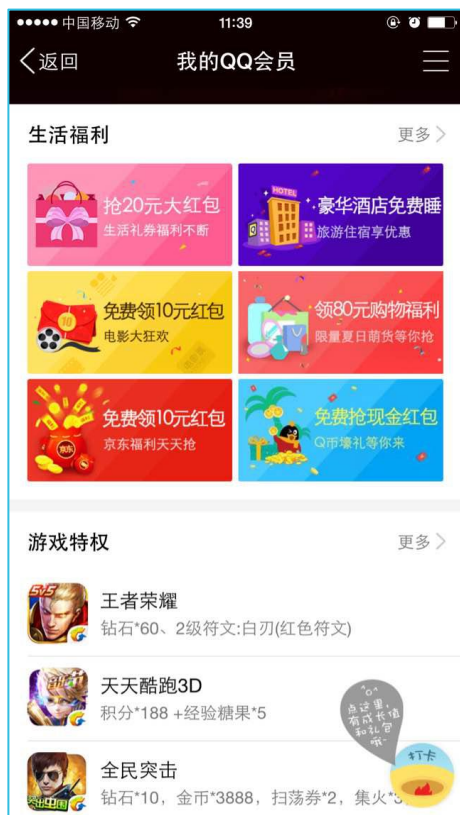
LBS推荐

相关推荐

社会化推荐

兴趣推荐

1 背景介绍(3) - 有哪些交互？



QQ会员



企鹅FM



QQ空间



腾讯课堂

用户浏览

产品push

用户反馈

交互式推荐



目录

1. 背景介绍

2. 架构介绍

2.1 架构总览

2.2 分布式计算

2.3 数据引擎

2.4 海量画像

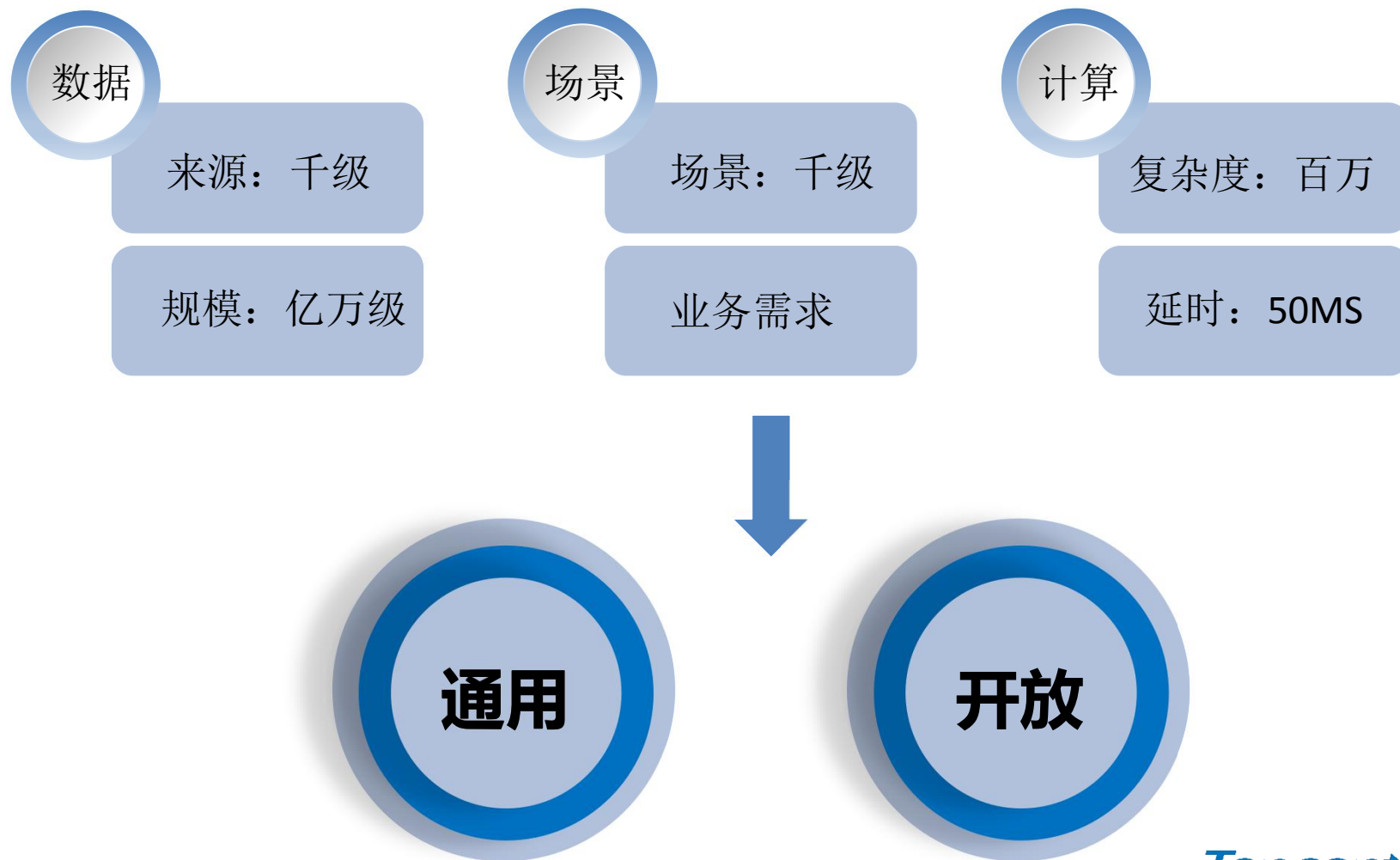
2.5 实时ABTEST

2.6 立体监控

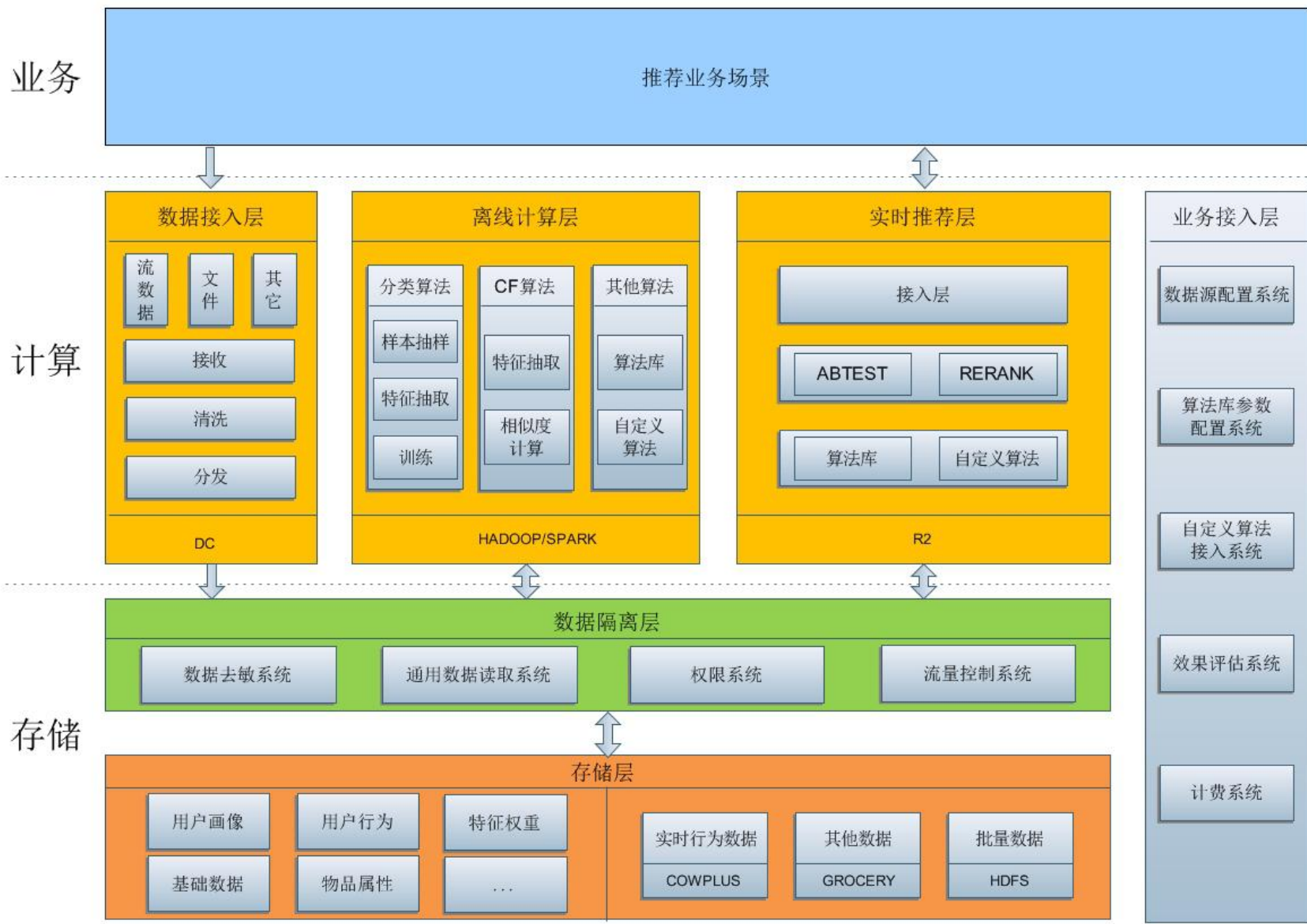
2.7 通用开放

3. 运营情况

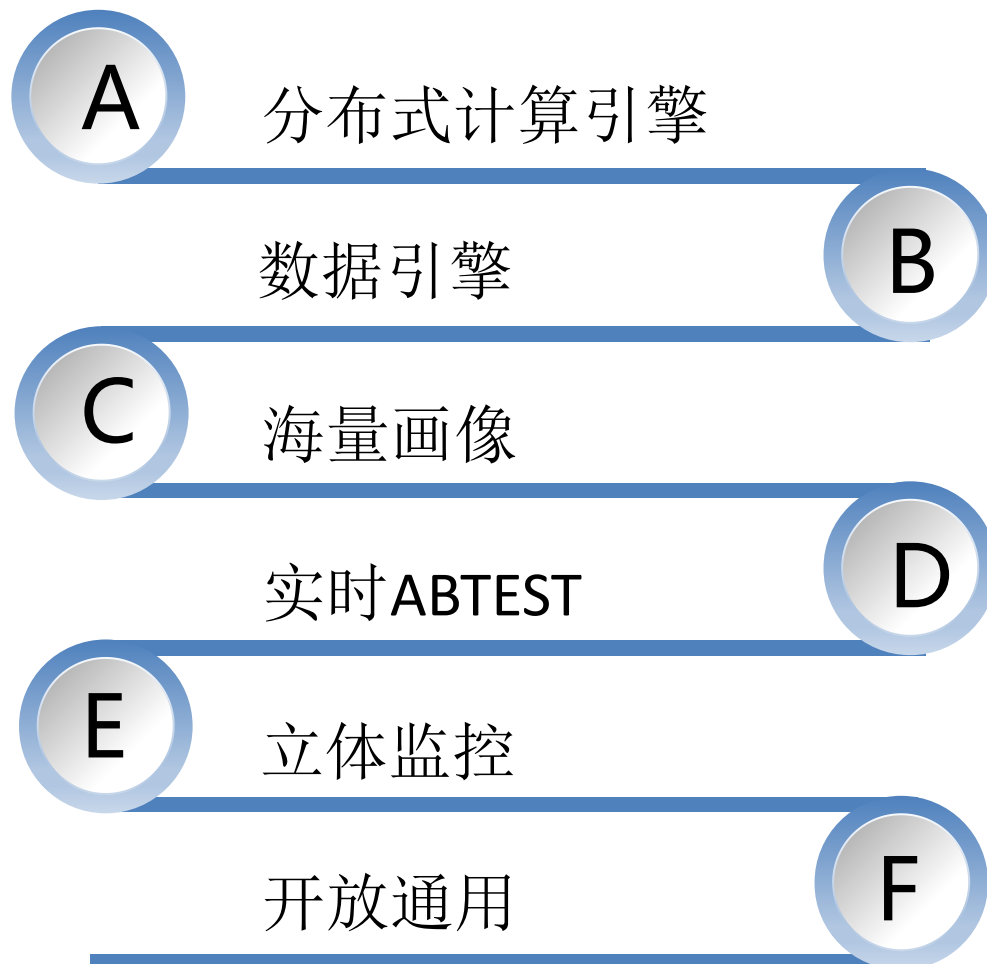
2.1 架构总览 - 系统要解决的问题



2.1 架构总览 – 架构图



2.1 架构总览 - 架构特点



2.2 分布式计算引擎(1) – 计算流

海选

- 获取有效池子
- 分人群等

初选

- 性别，年龄，行为等用户属性
- CF，热传导等
- 识别某些CVR高，但是用户不喜欢的badcase
- 应用宝广告推荐提升超过20%

精排

- LR/GBDT/CF等
- 多维度特征

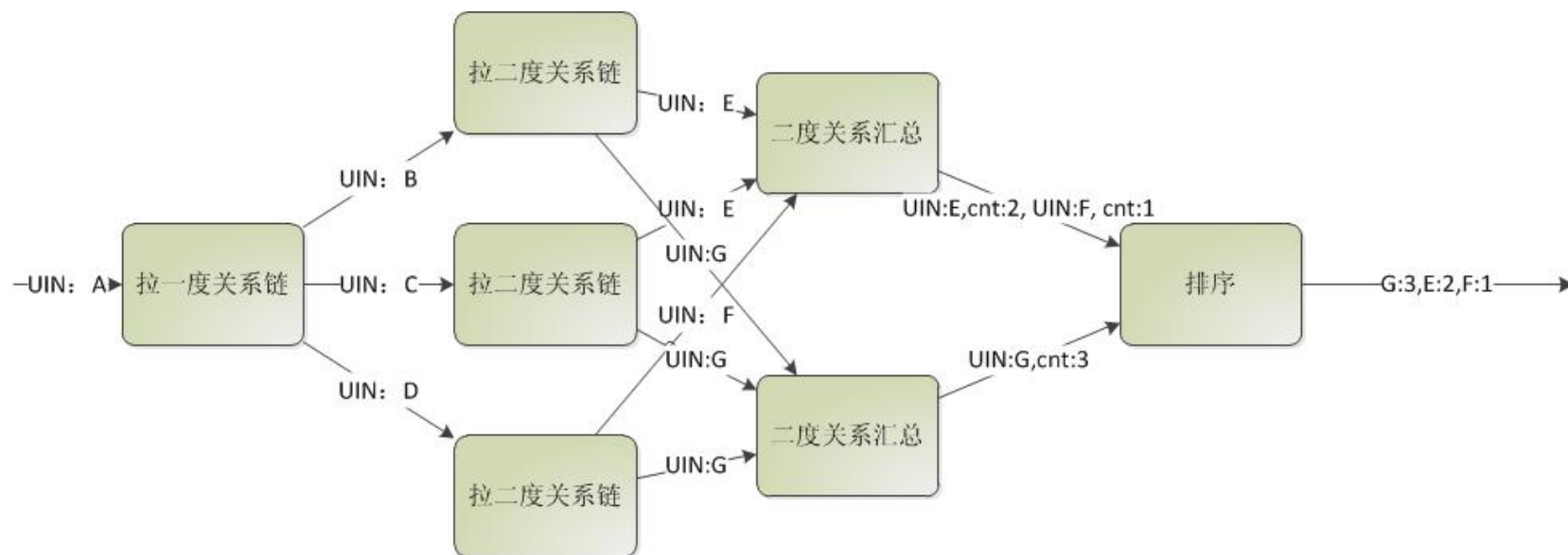
重排

- 曝光衰减
- 产品规则
- 其他badcase

2.2 分布式计算引擎(2) – 复杂计算

要解决的问题：复杂计算VS低延时要求

以QQ好友推荐的抽象逻辑为例：



2.2 分布式计算引擎(3) – why not storm

- 面向服务
 - 稳定性
 - 流量控制
- 性能
 - 数据拷贝
 - 资源
- 运维
 - 动态负载均衡
 - 自动扩缩容
- 开发
 - 多线程调试

2.2 分布式计算引擎(4) – 解决方案R2

● 解决的问题

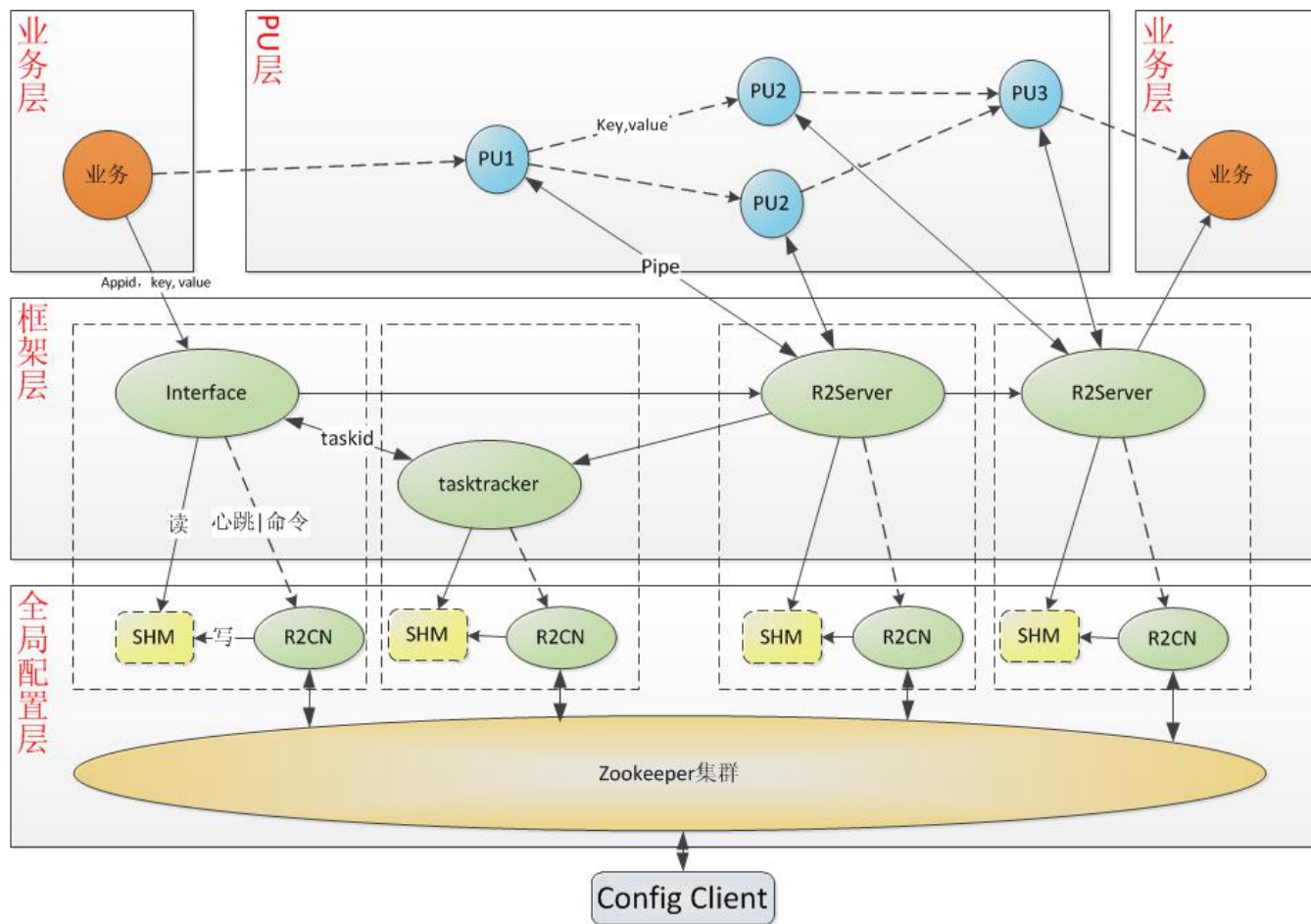
- 大计算量与低延迟
- 即时数据/处理/应用

● 适用业务

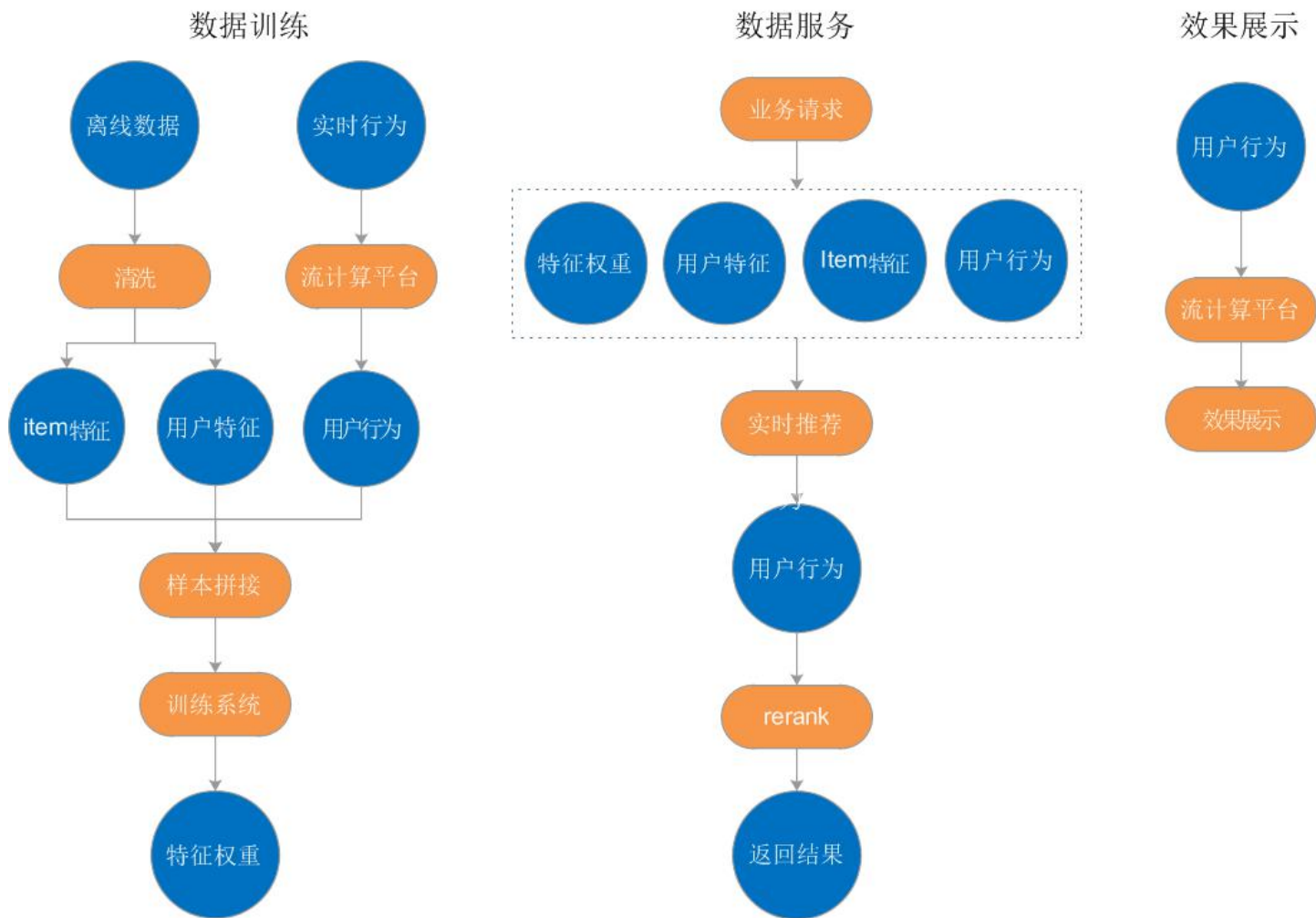
- 计算量比较大
- 对低延迟要求较高
- 实时数据计算&处理

● 计算模式

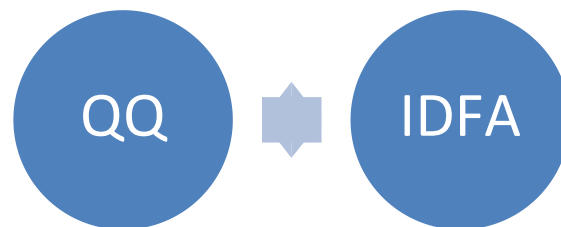
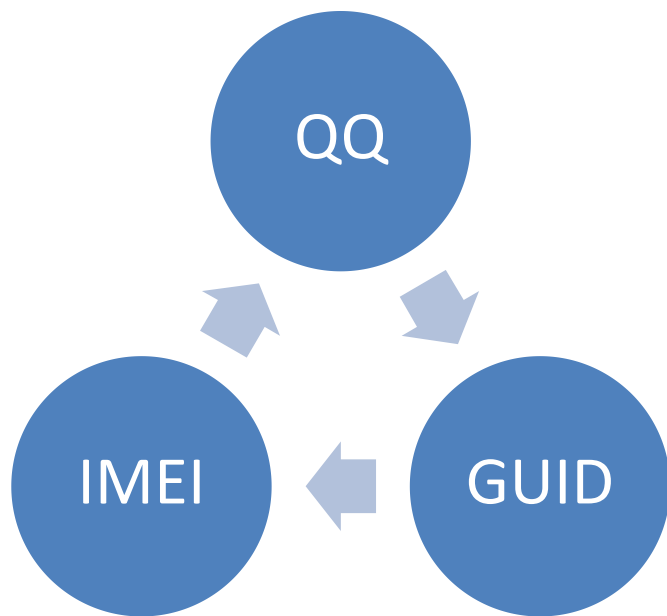
- 类map-reduce
- 实时化，多层化



2.3 数据引擎(1) – 数据流



2.3 数据引擎(2) -账号体系打通

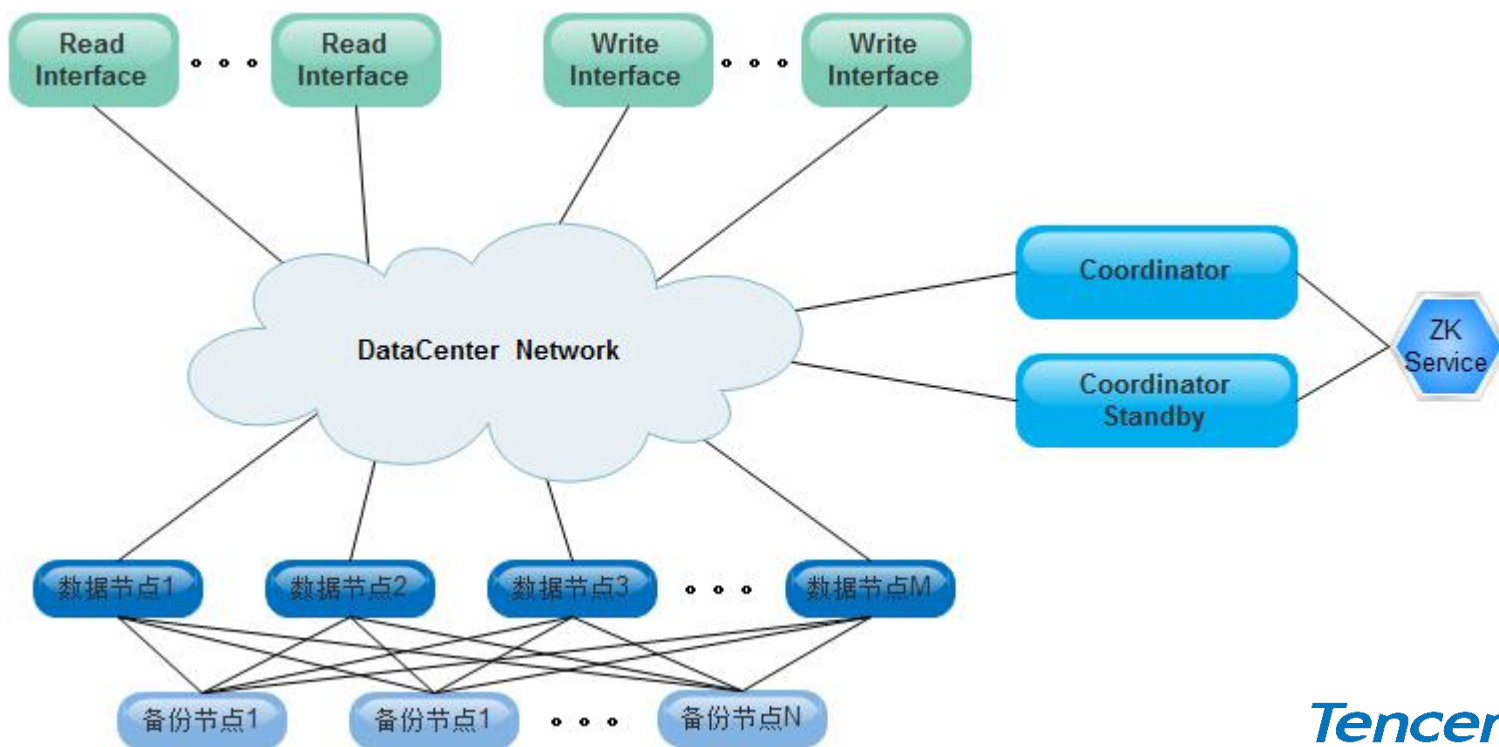


2.3 数据引擎(3) – 行为数据存储cowplus

- 数据量：300T+
- 访问量：百亿级访问
- 多场景，多维度读取



- 时间维度淘汰
- 多索引查询
- 简单计算

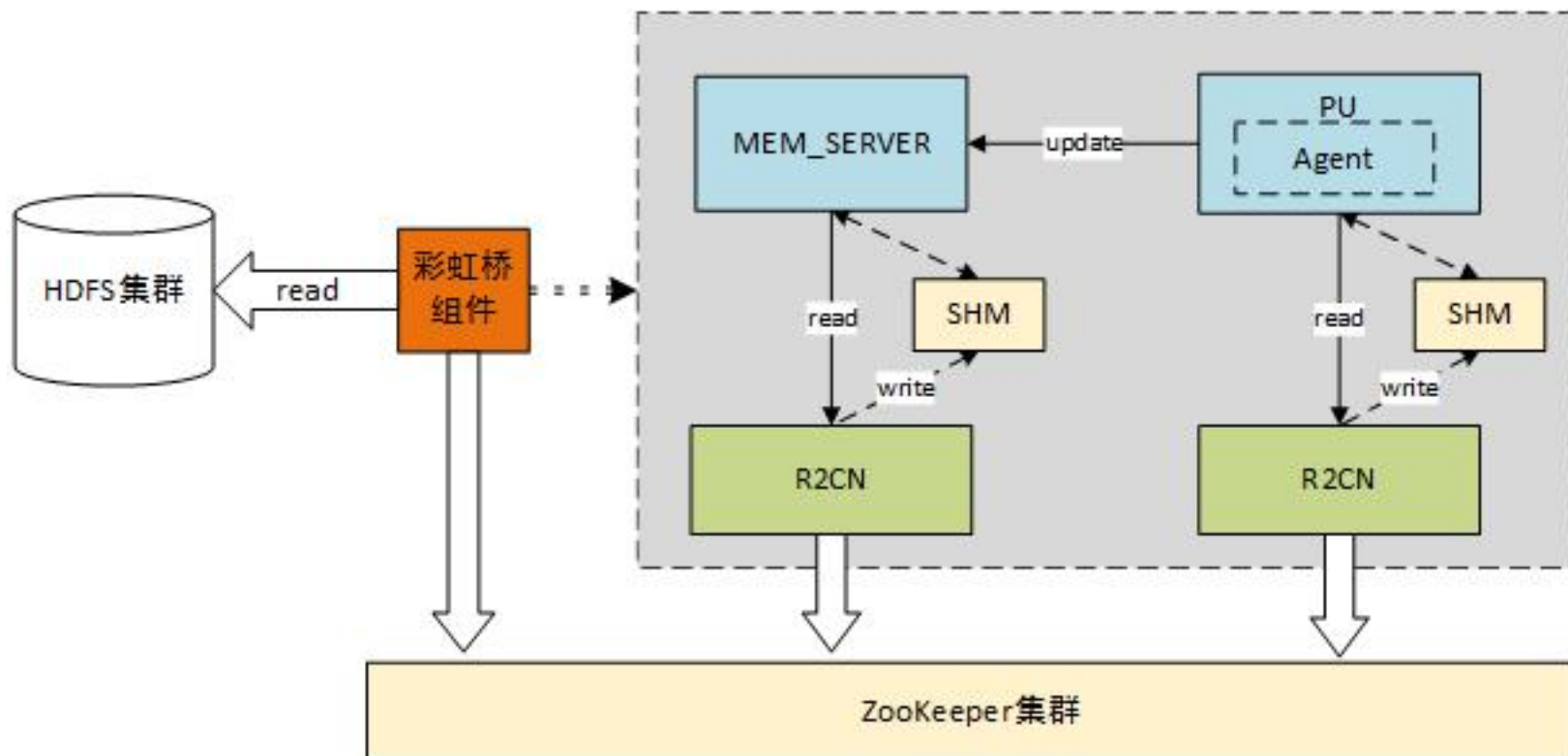


2.3 数据引擎(4) – 高速cache

- 每秒千万级读写
- 随时更新，秒级生效
- 数据一致性
- 通用数据格式



- Client按需拉取
- 数据一致性：MD5+ VERSION
- 数据更新不影响服务：主备切换



2.3 数据引擎(5) – 多数据源读取优化

- 读取多种存储系统
- 同步编码：开发快，性能差
- 异步编码：性能高，开发慢



- 协程
- 多数据源同步编码，异步执行
- 通用接口

- STEP1: 同步编码获取数据源

```
//同步编码 -- 获取数据源1
get_source1()
{
    //打包
    pack_1(...);
    //收发包API
    send_and_recv();
    //解包
    unpack_1(...);
}
//同步编码 -- 获取数据源2
get_source2()
{
    pack_2(...);
    send_and_recv();
    unpack_2(...);
}
```

- STEP2: 异步并行拉取数据源

```
//将获取数据源1的实现函数放到集合
closures.push(&get_source1);
//将获取数据源2的实现函数放到集合
closures.push(&get_source2);
//异步并行执行多数据源读取
closures.action();
```

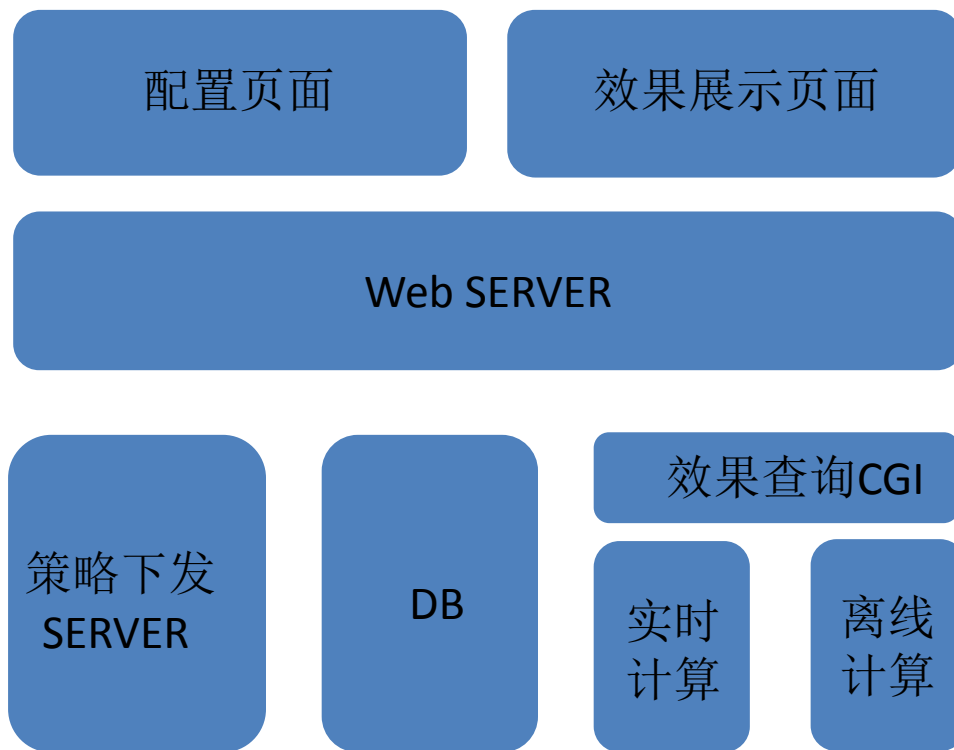
2.4 海量画像

- 10多亿用户, 几十款产品
- 上千种画像维度
- 多渠道准确度评估



2.5 实时ABTEST(1) – 整体架构

- 分钟级别更新
- 小时维度展示



2.5 实时ABTEST(2) – 用户配置

● ABTEST策略配置

输出策略Id	策略名称	策略条件	策略参数
1005	LR-社会化	{"tail":{"begin":"0","end":"39"}}	{}
1004	LR-GBDT特征	{"tail":{"begin":"40","end":"79"}}	{}
1001	首页参考系-热度	{"tail":{"begin":"80","end":"99"}}	{}

● 效果指标配置

指标配置			
操作	序号 ∨ ^	指标来源	
 	<input type="radio"/> 1	曝光用户数	usrs101
 	<input type="radio"/> 2	曝光次数	cnt101
 	<input type="radio"/> 3	点击用户数	usrs102
 	<input type="radio"/> 4	自定义	[usrs102]/[usrs101]
 	<input type="radio"/> 5	点击次数	cnt102
 	<input type="radio"/> 6	自定义	[cnt102]/[cnt101]
 	<input type="radio"/> 7	报名人数	usrs103
 	<input type="radio"/> 8	自定义	[usrs103]/[usrs102]
 	<input type="radio"/> 9	上课时长过去7天	times104

2.5 实时ABTEST(3) – 效果展示

- 恶意数据过滤
- 效果平滑处理
- 效果波动告警



2.6 立体监控(1) – 线上服务监控

● 业务总体监控

➤ 请求量 失败量 平均延时 ...

四级模块: [通用推荐接入层][接入] 主调模块: 神盾接入服务层测试 被调模块: REQUIRE_APPID_GUAIJIAN 接口: REQUIRE_APPID_GUAIJIAN	6715	400	7.35	99.94%	6.77ms
四级模块: [通用推荐接入层][接入] 主调模块: 神盾接入服务层测试 被调模块: FM 手Q资料卡推荐 - 默认图标 接口: RSQUIRE_APPID_DEMOGRAPHIC_HOT	1171	11	3.851	99.95%	7.3ms
四级模块: [通用推荐接入层][接入] 主调模块: 神盾接入服务层测试 被调模块: 企鹅fm 接口: RSQUIRE_APPID_DEMOGRAPHIC_HOT	142382	3869	16.37	99.95%	7.31ms

● 重点函数调用次数监控

➤ 历史对比

➤ 波动告警



2.6 立体监控(2) – 其他监控

● 计算过程重放

实时流水查询										
接口ID			环境设置	<input checked="" type="radio"/> 正式 <input type="radio"/> 测试		索引启动时间: 2016-04-05 15:03:53	索引状态: 运行中	关闭索引查询	查看接口配置	
查询日期	2016-04-14		开始时间	14	时 50	分	结束时间	15	时 50	分
uin							最大记录数	100	查询	导出excel
									切换表头	
#	reporttime	clientip	uin	imei	seq_no	rt	test_id	r2app_id	ruleid_out	output_value
1	2016-04-14 14:50:05	.62			1460616605		200106	60200	1009	rt: 2001%0Drulei...
2	2016-04-14 14:50:05	.00.19			1460616605		200106	60200	1009	rt: 2001%0Drulei...
3	2016-04-14 14:50:05	.62			1460616605		200106	60200	1009	rt: 2001%0Drulei...
4	2016-04-14 14:50:05	.233			1460616605		200106	60200	1009	rt: 6001%0Drulei...
5	2016-04-14 14:50:05	.100			1460616605		200106	60200	1009	rt: 2001%0Drulei...

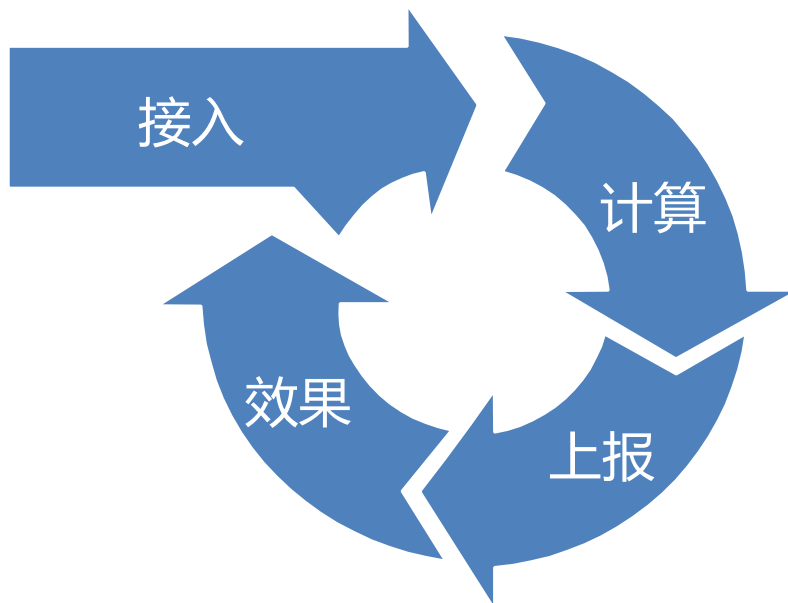
● 推荐效果监控

- 恶意数据过滤
- 效果平滑处理
- 效果波动告警

● 数据监控

- 数据依赖
- 失败告警/重传

2.7 通用开放(1) – 通用接入



- 接入
 - 系统注册
 - 接入API
- 计算
 - 算法选择
 - 参数配置
- 上报
 - 上报字段配置
 - 上报API
- 效果
 - 效果展示字段设置
 - 效果查看

2.7 通用开放(2) – 训练开放

特征构造

算法选择

特征选择

训练

<input type="checkbox"/> 全选	特征库名	特征主键	特征副键	来源	来源（归一化）（可以为空）
<input type="checkbox"/>	pendant	cross_try_gender_dis	gender_dis	$\${base:gender_dis_cnt_4w_gender_dis:.*}$	$\${base:gender_dis_cnt_4w_gender_dis:.*}$
<input type="checkbox"/>	pendant	cross_try_vip	vip	$\${base:vip_cnt_4w_vip:.*}$	$\${base:vip_cnt_4w_vip:.*}$
<input type="checkbox"/>	pendant	cross_try_vip_dis	vip_dis	$\${base:vip_dis_cnt_4w_vip_dis:.*}$	$\${base:vip_dis_cnt_4w_vip_dis:.*}$
<input type="checkbox"/>	pendant	set	.*	f_rcmd_fea_norm_score_d	f_rcmd_fea_norm_score_d
<input type="checkbox"/>	pendant	set_cnt_4w_age	.*	f_rcmd_fea_norm_score_d	f_rcmd_fea_norm_score_d

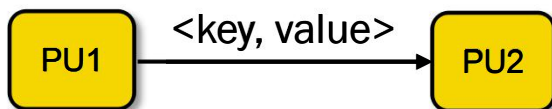
<input type="checkbox"/> 全选	业务ID	业务名称	算法ID	算法名称	算法类型	样本表(可以为空)	抽样规则(可以为空)	更新周期(可以为空)	建模方式(1-按场景)
<input type="checkbox"/>	4	挂件推荐	40004	LR-新用户	LogisticRegression	f_rcmd_pendant_sample_new_d	随机	w	1
<input type="checkbox"/>	4	挂件推荐	40005	LR-活跃用户	LogisticRegression	f_rcmd_pendant_sample_active_d	随机	w	1
<input type="checkbox"/>	4	挂件推荐	40006	LR-流失用户	LogisticRegression	f_rcmd_pendant_sample_losing_d	随机	w	1

<input type="checkbox"/> 全选	算法ID	特征模块	特征类型	特征ID	取数模式	更新时间
<input type="checkbox"/>	50200	yyb	dis_app_copair	.*	2	2016-04-16 13:50:00
<input type="checkbox"/>	50034	yyb	app_filter_540_adpos1_cvr_day	.*	2	2016-04-16 12:56:07
<input type="checkbox"/>	50034	yyb	cross_career_app_filter_540_adpos1_cvr_day_cross_discretization_dl_count	.*	2	2016-04-16 12:56:07



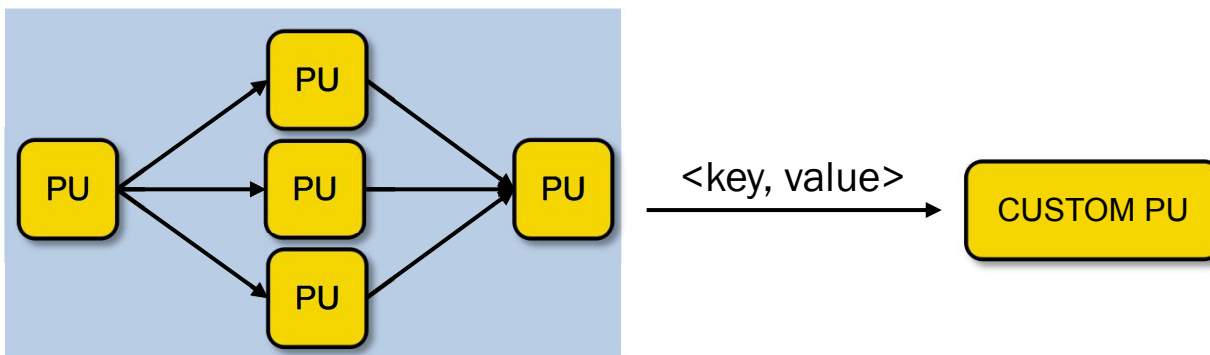
2.7 通用开放(3) – 计算开放

- 数据传递方式



- 计算逻辑

➢ 用户自定义custom pu逻辑



2.7 通用开放(4) -自定义计算

●说明

- R2提供库文件，由用户引用
- 用户只用实现代码逻辑函数“proc”即可。（其余功能由系统完成）

●步骤

- step1 继承基类PuPB，实现proc函数

```
class PbApp:public r2::pu::PuPB
{
public:
    PbApp(){}
    ~PbApp(){}

    virtual int proc(const std::string & key, const std::string & value)
    {
        std::string key_out;
        std::string value_out;

        key_out=key+"_hello";
        value_out=value+"_world";
        send_pack("CLIENT", key_out, value_out);

        return 0;
    }

private:
};
```

- step2 make

```
[tashanji@PLT_Kdc_DE ~/R2_proj/trunk/app/per_test/L1PU]$ make
g++ -c -I ../../example/include/ -o l1_1.o l1_1.cpp
g++ -o l1_1 l1_1.o -L ../../example/lib/ -lappframe -lpthread
```



目录

1. 背景介绍
2. 架构介绍
 - 2.1 架构总览
 - 2.2 分布式计算
 - 2.3 数据引擎
 - 2.4 海量画像
 - 2.5 实时ABTEST
 - 2.6 立体监控
 - 2.7 通用开放
3. 运营情况

3 运营情况(1) – 系统

- 海量

- 日调用：数十亿
- 扩散量：数百亿

- 实时

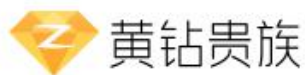
- 90%以上请求在20MS之内
- 所有请求处理延时50ms之内

- 稳定性

- 系统稳定性99.95%
- 确保每个消息的可靠传递

3 运营情况(2) – 运营

- 典型业务



- 运营

- 场景：近200个
- 业务接入：2人

- 应用宝推荐效果

- 推荐分发占可控分发超过80%
- 为应用宝带来超过60%的收入提升

- 手游推荐

- Banner场景为业务带来530%的效果提升
- 每天带来的下载超过百万

欢迎加入腾讯QQ大数据团队