

The background of the slide is a vibrant, stylized illustration of an underwater environment. It features various types of coral in shades of orange, red, and purple, along with grey and blue rocky formations. Sunlight rays filter down from the top, creating a bright blue area at the top of the frame. The overall aesthetic is that of a modern, colorful digital artwork.

DataFunTalk

# 58同城大数据应用实践

2020.08.18-19



# 58房产数据服务之路

孙志文-数据架构师

1. 58房产数据服务的场景
2. 数据业务化面临的挑战和解决方案
3. 核心流程监控利器 – 智能机器人

# 58 58房产数据服务的场景

数据化运营

数据植入业务

数据业务化

# 58 数据植入业务



## 服务范围

服务公司: 上海中原物业顾问有限公司

主营商圈: 浦东-金桥、浦东-高行、浦东-航头

主营小区: 东陆新村六街坊、幸福小镇明丰绿都(公...、仁恒森兰雅苑(二期)

佣金范围: ≤0.5%

增值服务: 专车接送 新房代购 置换服务

## 用户评价

好评率50%, 2人评价

关注

微信

电话



二手房

新房

## 价格走势

近3月 近1年 近3年



上海二手房

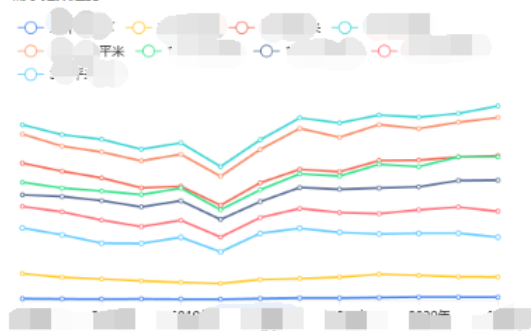
51513元/m²

关注

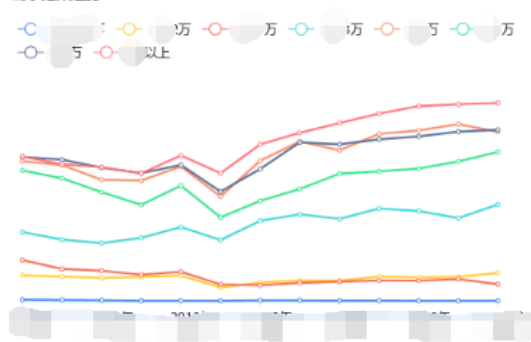
# 58 数据业务化



需求指数走势



需求指数走势





# 58房产数据架构

## 数据服务

数据化运营

运营报表

实时大屏

明细下载

Deep dive

任我查

销售决策

数据植入业务

房源Rank

房源推荐

楼盘榜单

用户画像

经纪人评分

房易通

数据业务化

楼盘洞察

经纪人参谋

门店参谋

公司参谋

锁客宝

房价地图

## 数据仓库

DWS

用户流量

收入

开发商

经纪人

市场渠道

销售行为

DWD

曝光

点击、连接

房源、楼盘

推广计划

会员

套餐、花费

ODS

浏览日志

微聊、电话

房源信息

经纪人信息

消耗明细

广告数据

## 数据平台

数据开发

元数据管理

数据监控

云窗

算法平台

实时计算平台

KV数据服务

Spark

Hive

Flink

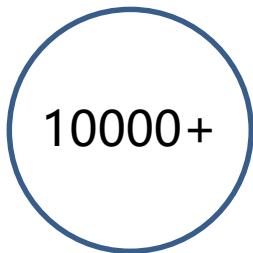
Tensorflow

DB 数据交换

KV 数据交换



# 58 数据服务相关统计



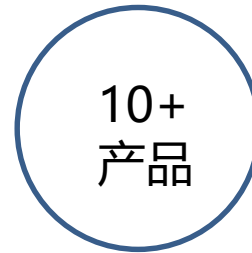
任务数



数据化运营



数据植入业务



数据业务化





# 数据业务化带来的挑战

数据量增加

- DB 无法满足



通用KV数据服务

数据质量要求更高

- 质量问题导致客诉



设计、流程改进

时效要求更高

- 8:00 点完成



智能机器人



# 通用KV数据服务

查询

负载均衡

查询条件转换

权限检查

存储

HBase

WTable

导入

HFile 生成

Protobuf 生成

数据过滤

Key 转换

元数据

运维

# 58 通用KV数据服务

Hive hdp\_anjuke\_dw\_db\_da\_ajk\_avg\_price\_zf\_room\_monthly  

集群  
58 Hive

数据库  
hdp\_anjuke\_dw\_c

表  
da\_ajk\_avg\_price

#	字段	注释	上传	RowKey
1	city_id	城市id	<input checked="" type="checkbox"/>	<input type="checkbox"/>
2	comm_id	小区id	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
3	areacode	区域代码	<input checked="" type="checkbox"/>	<input type="checkbox"/>
4	room_id	几室	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
5	room_desc	几室描述	<input checked="" type="checkbox"/>	<input type="checkbox"/>

Hive Where:

month\_id='\$cal\_dt'

例如: cal\_dt='\$cal\_dt', 表示取当天数据

调度 ID:

DWMS/58DP 中调度任务 ID

Rowkey Rules

服务名:  
hdp\_anjuke\_dw\_db\_da\_ajk\_avg\_price\_

# KEY

1 comm\_id ☒ 取模 ☐ MD5 ☐ 无

2 room\_id ☐ 范围

# 58 通用KV数据服务

## 基本信息

表名	hdp_anjuke_ana_db_sr	状态	ONLINE
Hive #1	hdp_anjuke_ana_db.sai	集群	58 Hive
用途	经纪人后台-三网推荐标签统一	负责人	

## 分客户端调用量

### 今天调用量

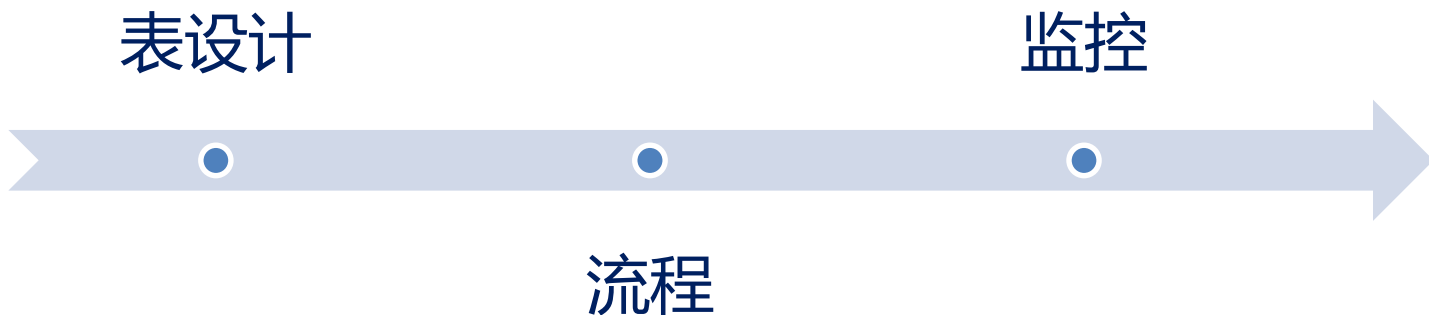


## 历史调用量

### 表结构

```
{
  "tableName": "hdp_anjuke_ana_db_sr",
  "columns": {
    "id": "bigint",
    "name": "string",
    "age": "int",
    "gender": "string",
    "city": "string",
    "phone": "string",
    "email": "string",
    "password": "string",
    "status": "int",
    "create_time": "timestamp",
    "update_time": "timestamp"
  }
}
```

# 58 设计、流程改进



1. 中间表设计
2. 应用层设计
3. 维表设计

1. 导出前质量检查
2. 设置数据信号

更全面的监控,  
提前发现问题

## 表设计

1. 中间层：针对内部的“大宽表”和对外服务的“小表”同时存在，中间层的解耦，确保对外数据流程的时效和稳定性
2. 服务层：根据业务特性（最新OR全部）决定是否加入分区键  
在满足业务需求的场景下，尽量减少存储压力
3. 维表：是否直接使用线上业务的维表？谨慎处理维度退化

## 流程

1. 导出数据到线上服务时，先做质量检查；异常时，人工介入
2. 导出数据完成后，写数据信号，方便线上业务切换到新的数据

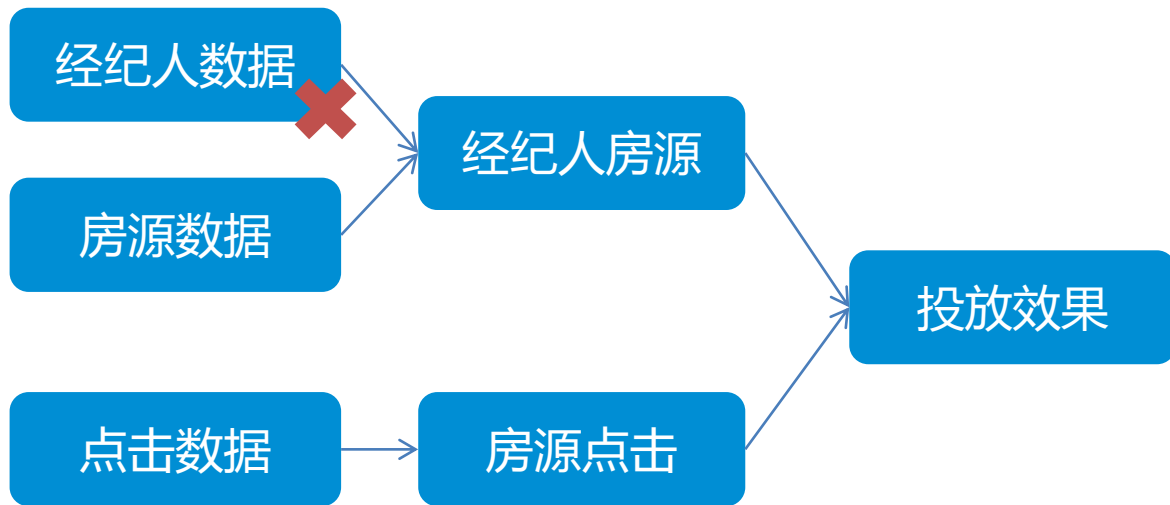


## 监控

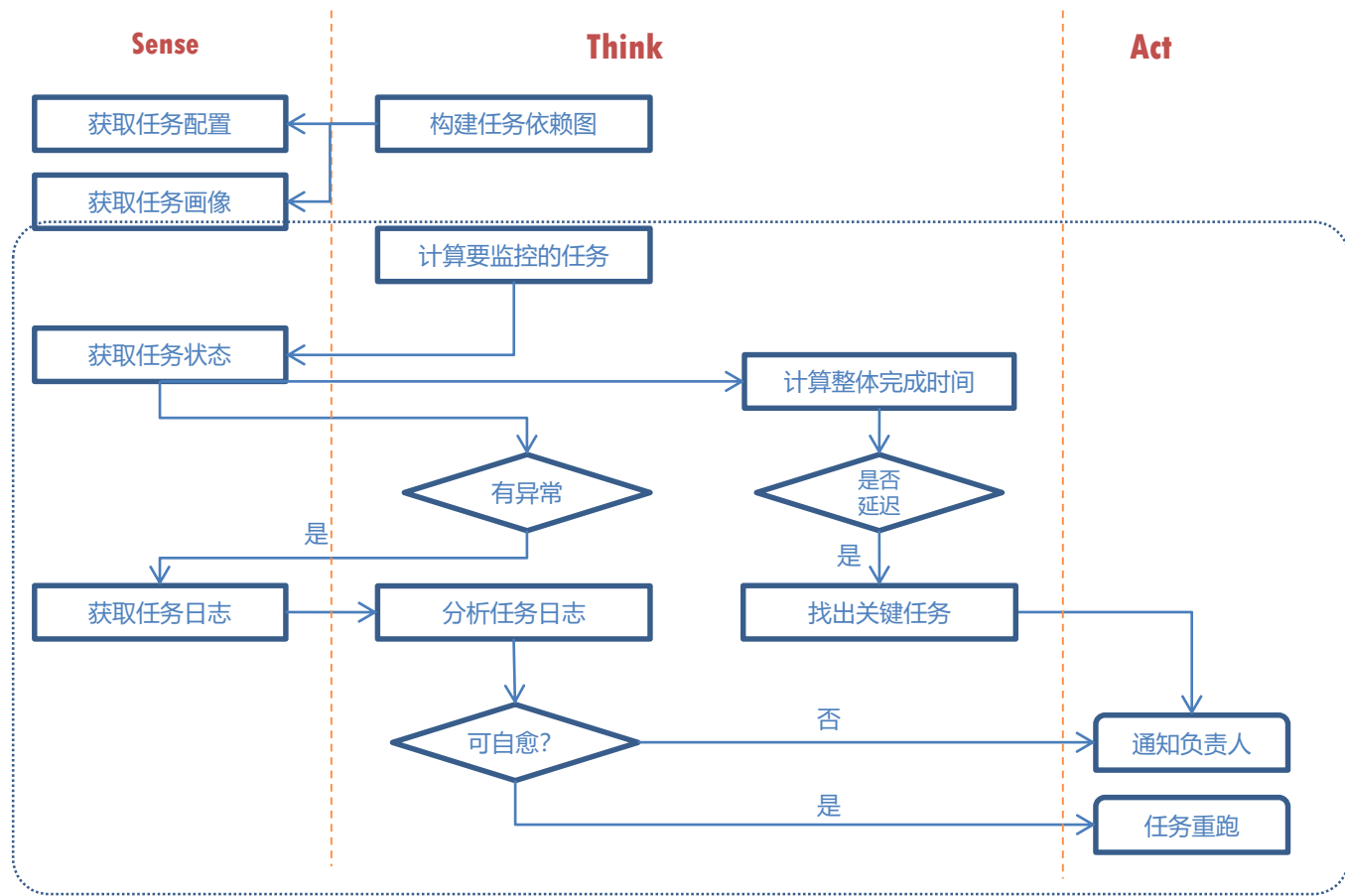
任务监控从点到面

# 58 单个任务监控方案的缺陷

1. 监控片面
2. 延迟发现晚
3. 任务诊断难



# 58 智能机器人架构



## See

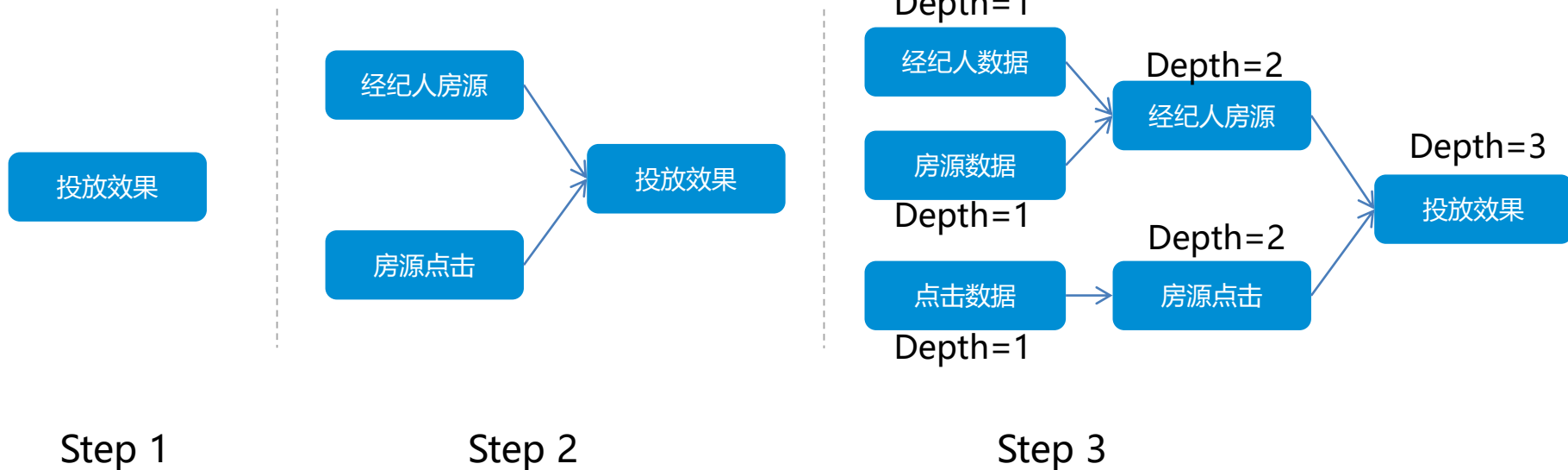
- 任务依赖关系
- 任务运行状态和日志
- 任务画像
- 集群状态

## Act

- 任务重跑
- 通知负责人

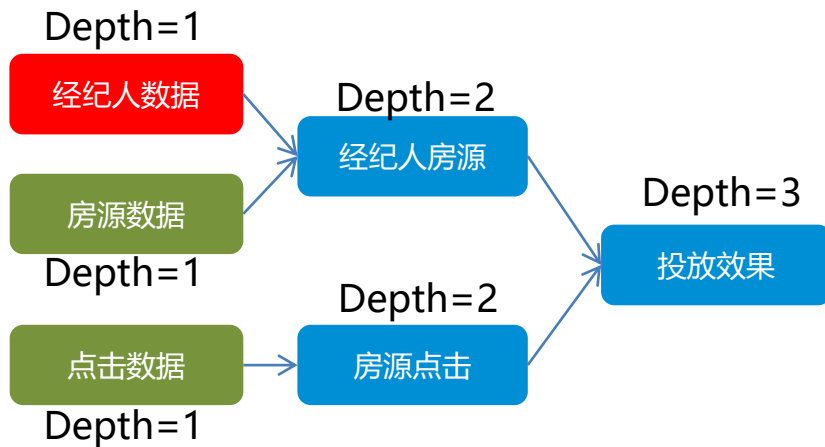
## Think

1. 构建 job 依赖图
2. 设置 job 深度



## Think

1. 监控当前深度下的未完成任务，若均完成，深度 + 1
2. 若发现异常任务，则采取对应动作



## Think

通过日志等信息，将异常分类

### 不可重试

1. 语法错误
2. 相关配置错误
3. 缺少权限
4. Mapjoin 异常
5. 表不存在

### 可重试

1. 连接数过多
2. 读取数据超时
3. 拿不到锁



## Think

3. 估算所有任务完成时间
4. 若最终完成时间超出设定阈值，找出慢的任务，并告警

$$\textit{End time} = \textit{Max}(\textit{start time} + \textit{execute time}, \textit{current time})$$

$$\textit{Start time} = \textit{Max}(\textit{schedule time}, \textit{parent job end time}, \dots)$$

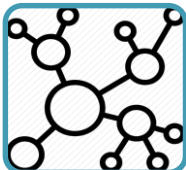
$$\textit{Execute time} = \textit{Max}(\textit{Median}(\textit{Execute time}), \textit{current execute time} + R)$$

# 58 智能机器人效果



## 常规巡检

- 减轻负责人凌晨压力



## 全局可控

- 任何影响到核心流程的异常都会被监控



## 快速定位

- 快速定位异常点，更快恢复流程



**70%**

减轻执勤压力



**98%**

准确率

1. 实时数仓
2. 数据治理

An underwater scene with a large black rectangular box in the center. The background is a vibrant blue and green underwater environment. In the foreground, there are various coral reefs, including a large orange branching coral on the left and a purple and orange coral on the right. The floor is covered in smaller corals and rocks. In the background, there are large blue and green rock formations. The lighting is bright, suggesting sunlight filtering through the water. The overall style is a colorful, stylized illustration.

THANKS !