



十年架构 成长之路

# SACC 第十届中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2018

2018年10月17-10月21日 北京海淀永泰福朋喜来登酒店



# 品友大数据分析平台的架构和演化

王晓鹏



第十届中国系统架构师大会  
SYSTEM ARCHITECT CONFERENCE CHINA 2018



# 议程

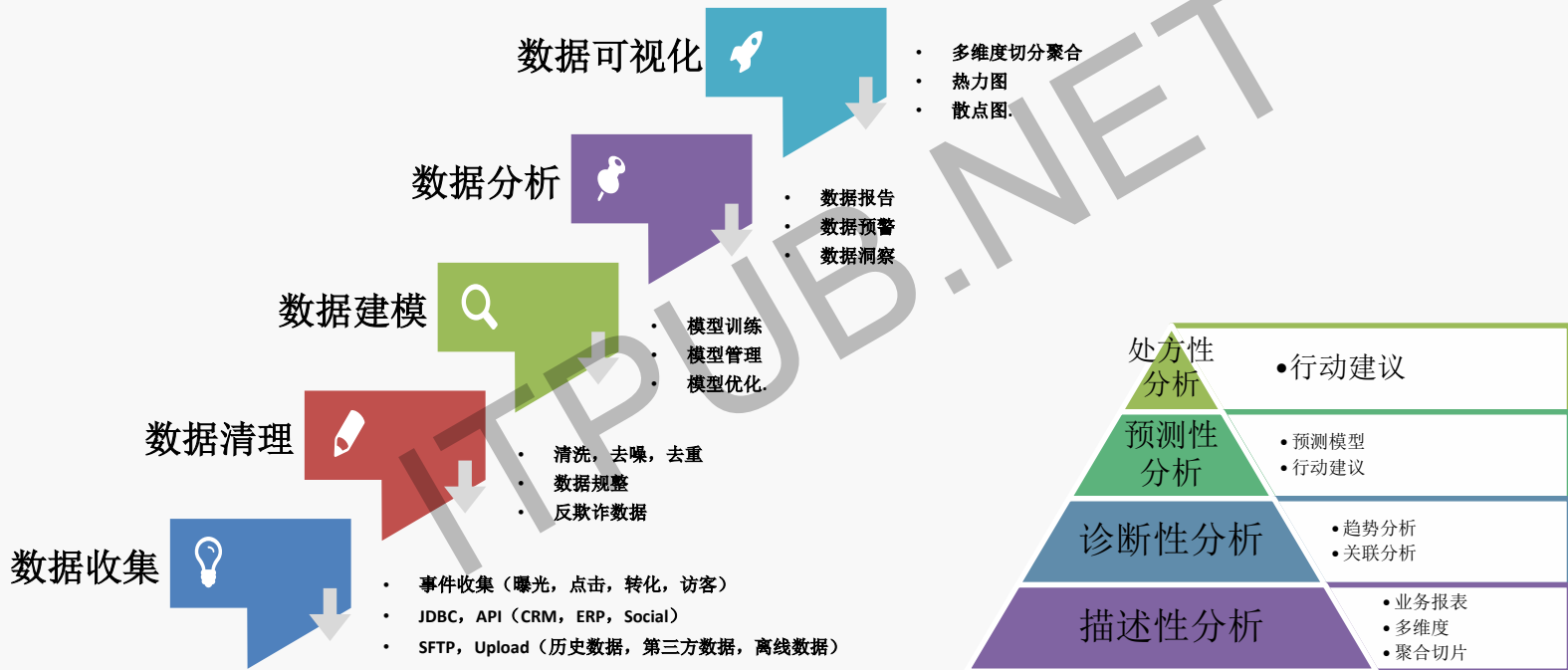
- 大数据分析的步骤
- 品友大数据分析平台的架构
- 投放分析平台的演进
- 数据管理分析平台的演进
- 工具的探索与实践
- 总结



十年架构 成长之路



# 大数据分析步骤



十年架构 成长之路



# 品友应用大数据分析的产品



## 投放分析平台

品友的广告投放平台  
用户为广告投放运营



## 企业数据管理平台

DMP, SaaS, In-House, Hybrid  
用户为企业市场, IT, 数据分析



目的不同



使用场景不同



数据不同



面向用户不同



十年架构 成长之路



# 品友投放数据及其分析需求

基础一方数据  
(用户资料, 用户标签):

$15\text{G}/\text{日} \times 365 + 40\text{G} \times 12\text{月} = 5\text{T}/\text{年}$   
考虑20%的业务增长率后为: **6T/年**

广告行为数据: 250T/年

考虑20%的业务增长率后为:  
**300T/年**

分析数据: 600G / 日

考虑20%的业务增长率后为:  
**272T/年**

## 数据存储的建议

- ❑ 鉴于用户换机周期为1.5年, 我们建议广告行为数据存储1.5年;
- ❑ 一方标签数据、分析数据、报表数据永久存储
- ❑ 建议分配存储: XX (与研发确认)

## 关键性指标

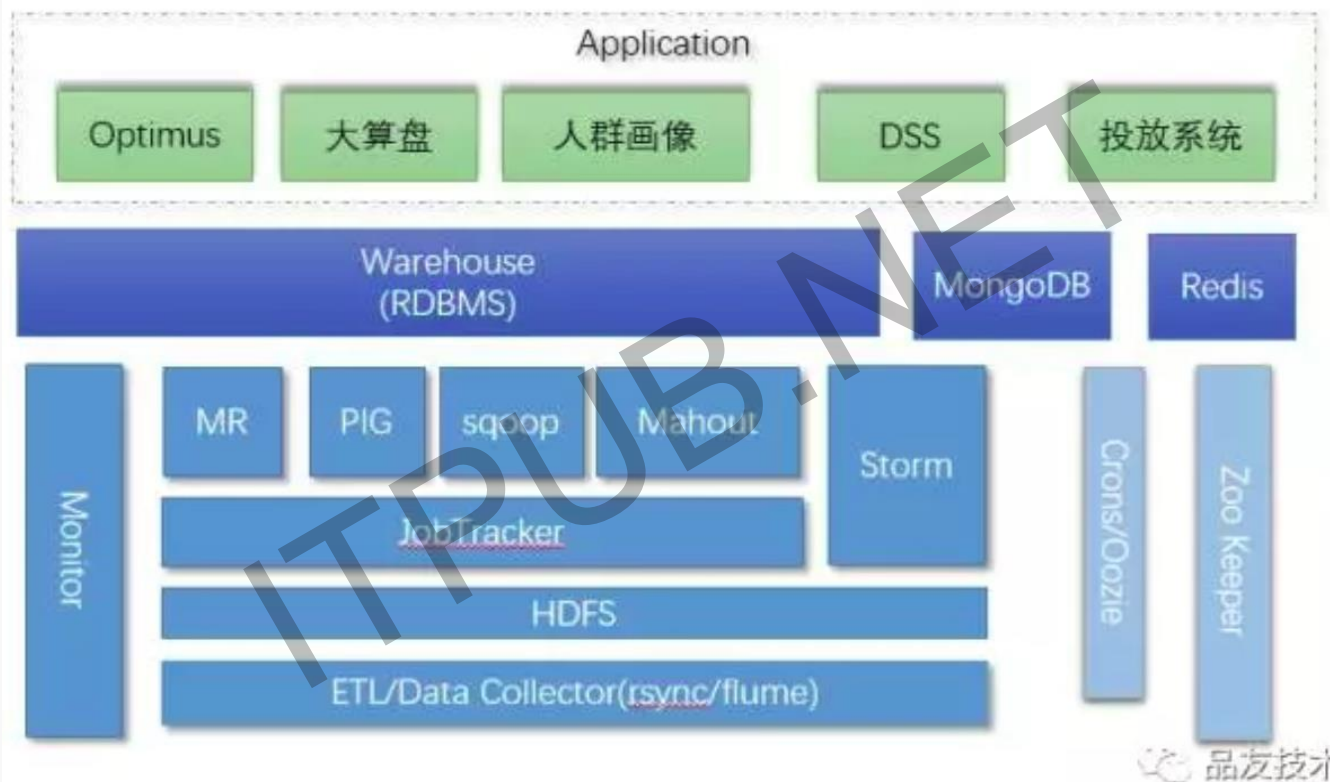
- ❑ 多维分析查询速度: 不高于10秒 (不含用户访问查询页面的时间);
- ❑ 人群预估响应时间: 秒级;
- ❑ 人群生成响应时间: 与人群规模有关;
- ❑ 并发查询数目: 1000请求/秒;



十年架构 成长之路

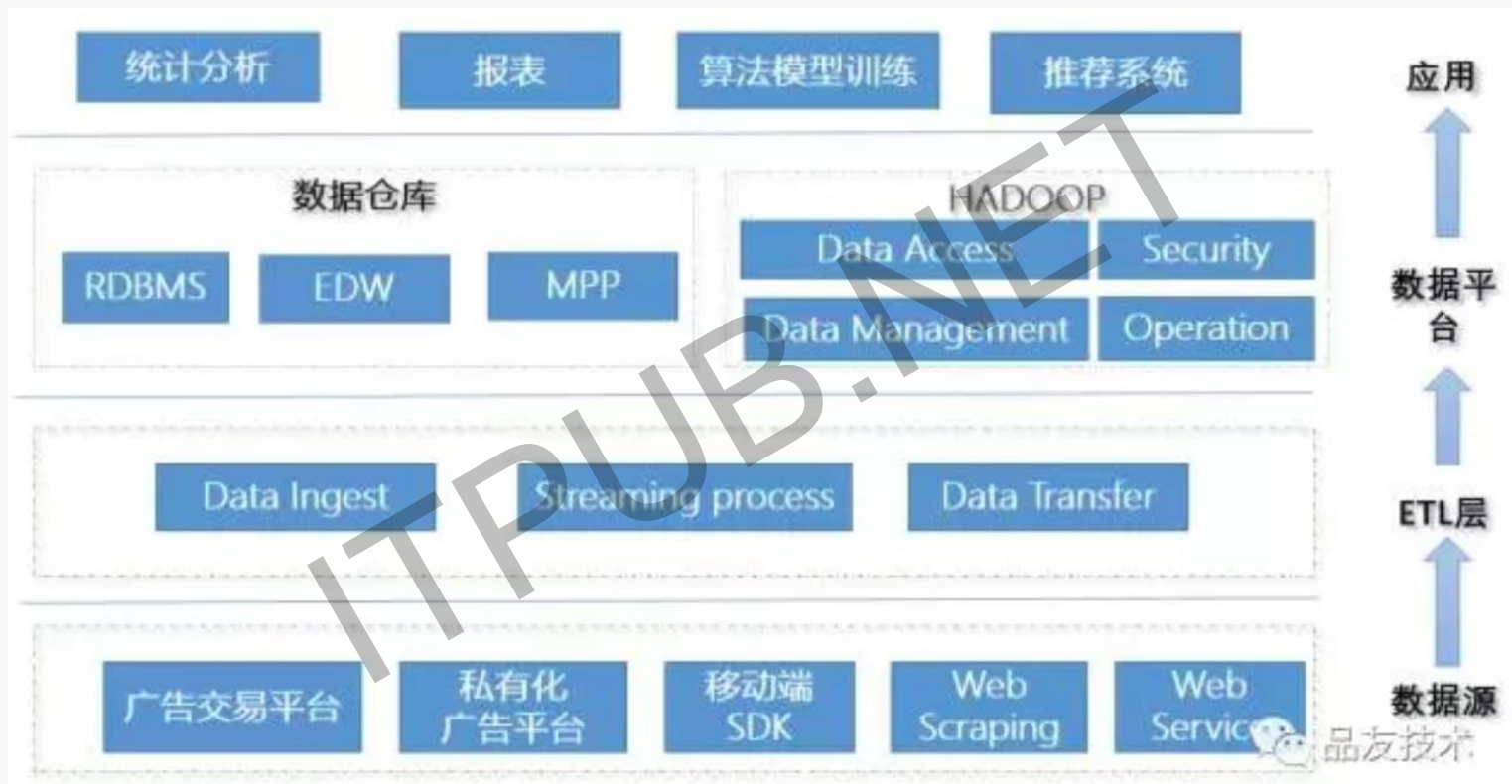


# 第一代大数据平台（2013-2014）



十年架构 成长之路

# 第二代大数据平台（2015-2016）

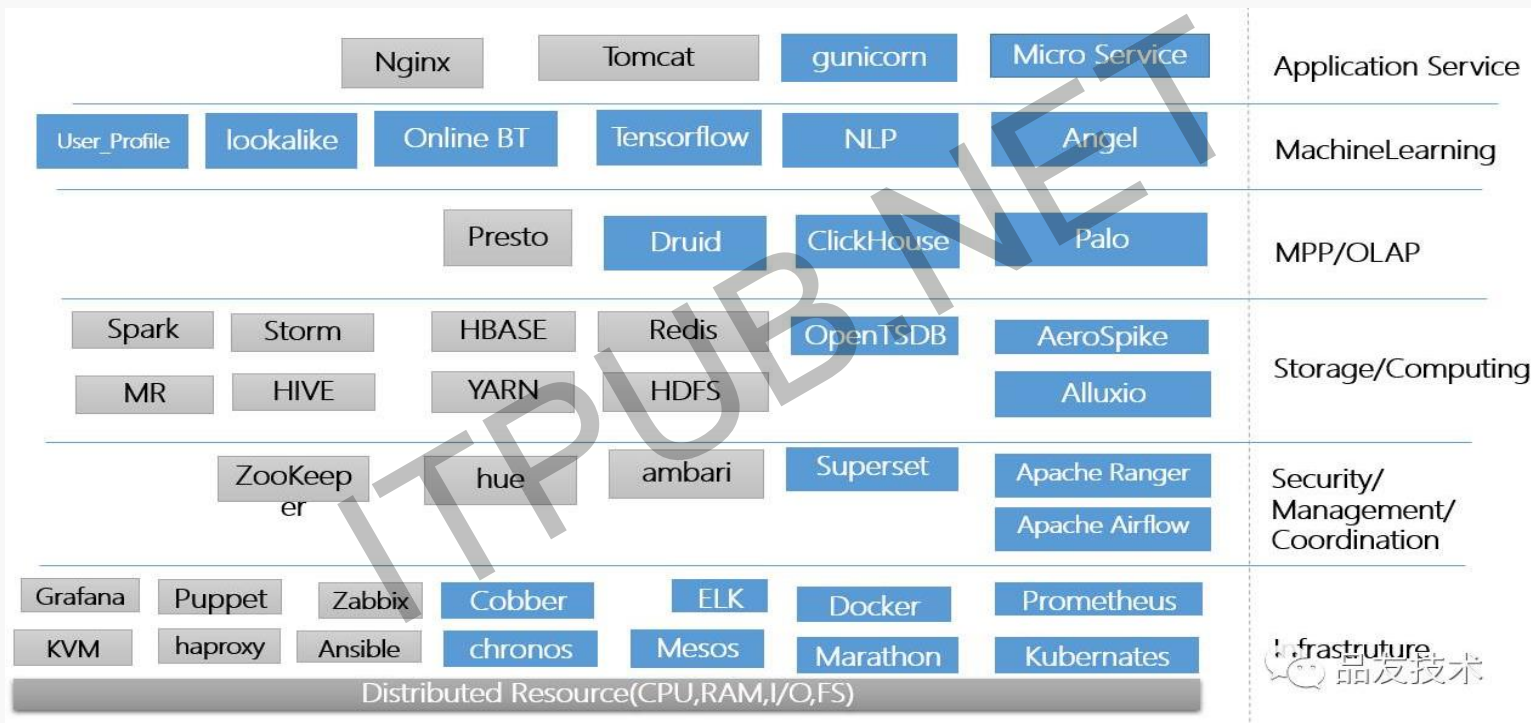


十年架构 成长之路





# 第三代大数据平台（2016- 现在）

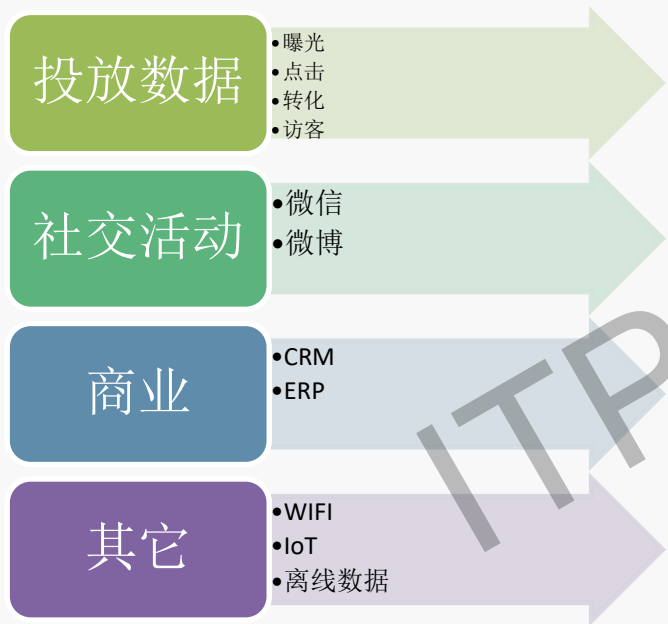


十年架构 成长之路



# DMP的数据及分析需求

## 数据来源



## 用户角色及需求



十年架构 成长之路

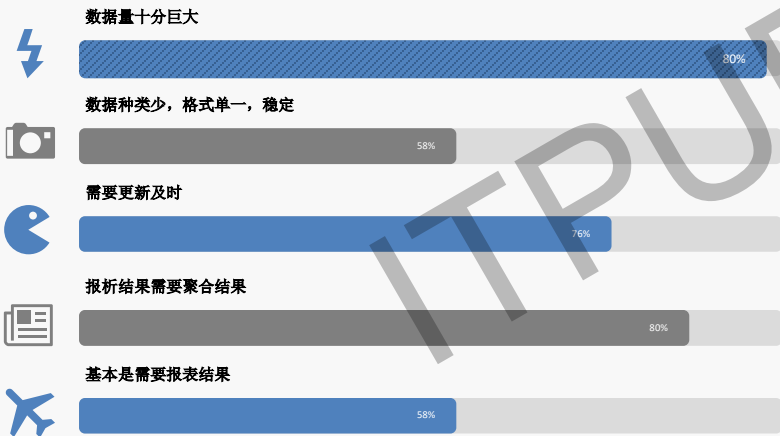


# 没有银弹 (No Silver Bullet)



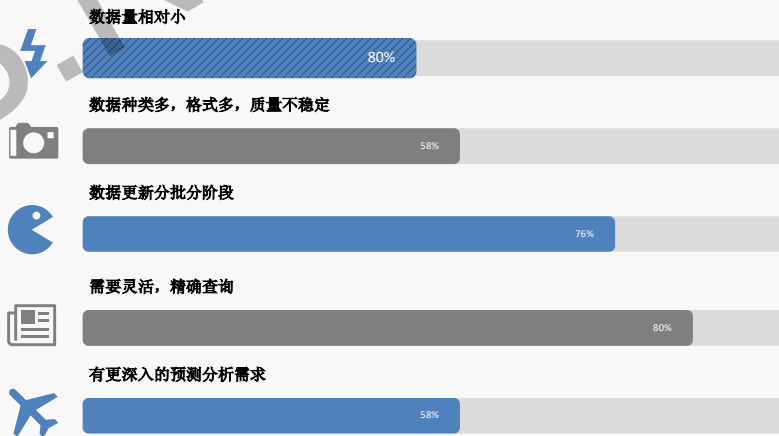
## 投放数据

来自广告投放及网站访问事件



## 第三方数据

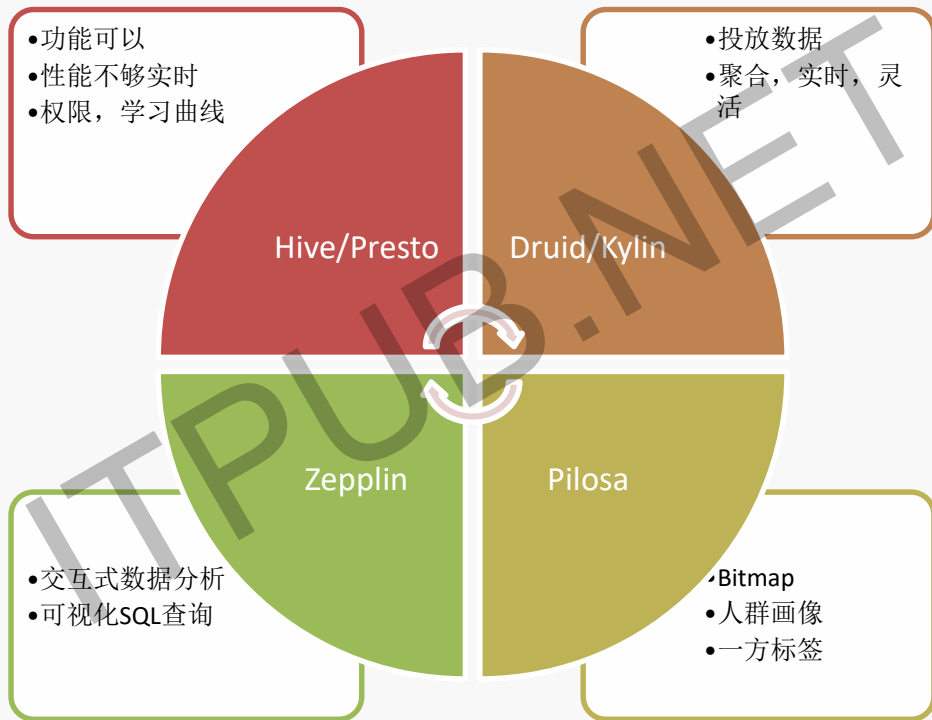
属于客户第一方的多种数据



十年架构 成长之路



# DMP数据分析的全面开花



十年架构 成长之路



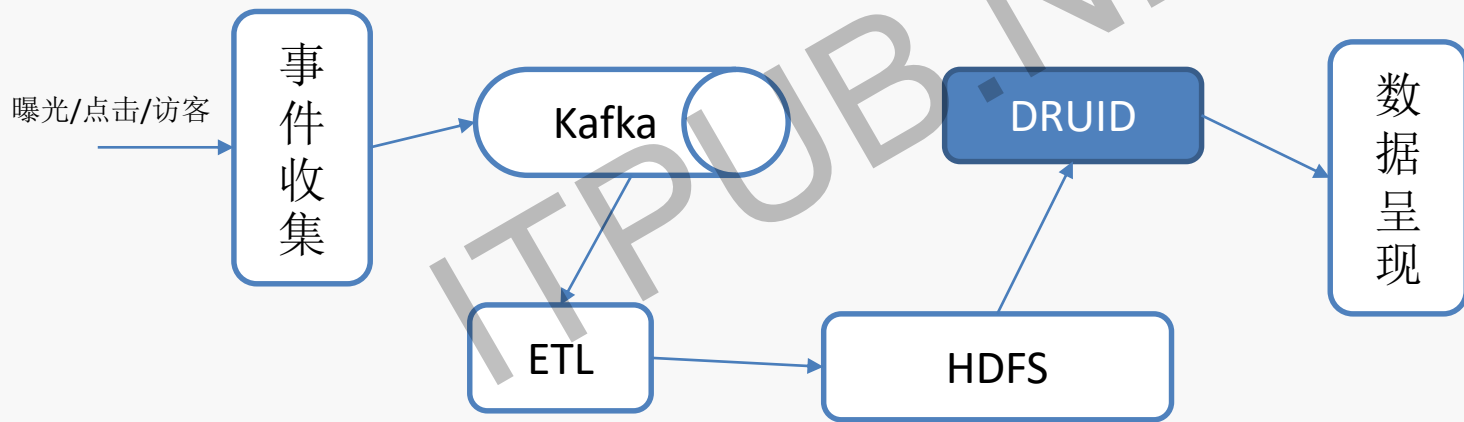
# DRUID

- 高性能的，分布式列存储的MOLAP框架
- 特点
  - 亚秒级查询
  - 实时数据注入
  - 可扩展的PB级存储
  - 支持多种数据源：hadoop，spark，kafka，storm和samza等
- 缺点
  - 只有聚合结果，没有明细



# Druid在品友的实践

- 使用场景：广告实时统计分析
- 数据：投放数据，20亿/天



十年架构 成长之路



# Zepplin, 数据分析师的心头好

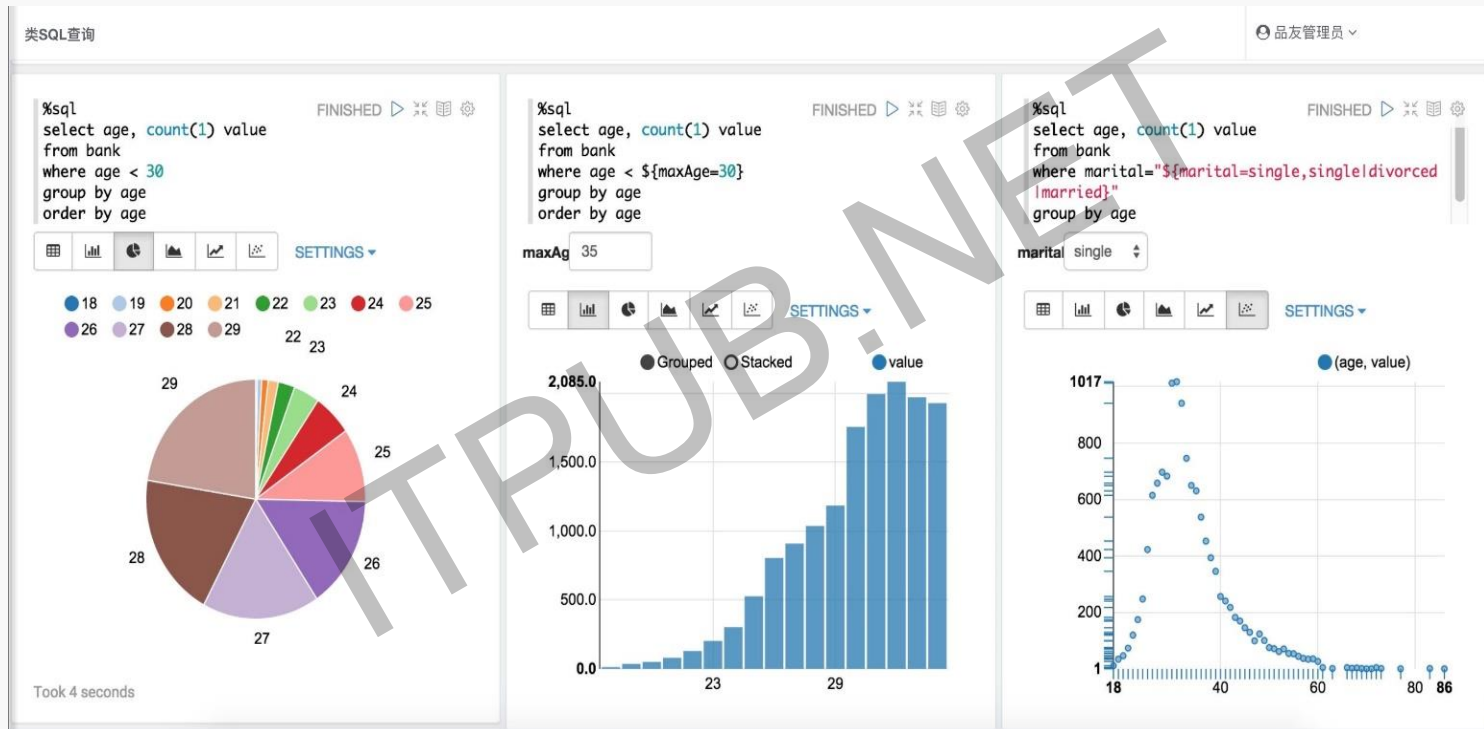
- Apache Zepplin是可视化框架
- 应用于交互式数据分析，七牛云，
  - 支持多种语言，默认是scala(背后是[Spark](#) shell),  
SparkSQL, Markdown 和 Shel
- 功能
  - 数据可视化
  - 用SQL来进行可视化查询



十年架构 成长之路



# Zepplin在品友的实践



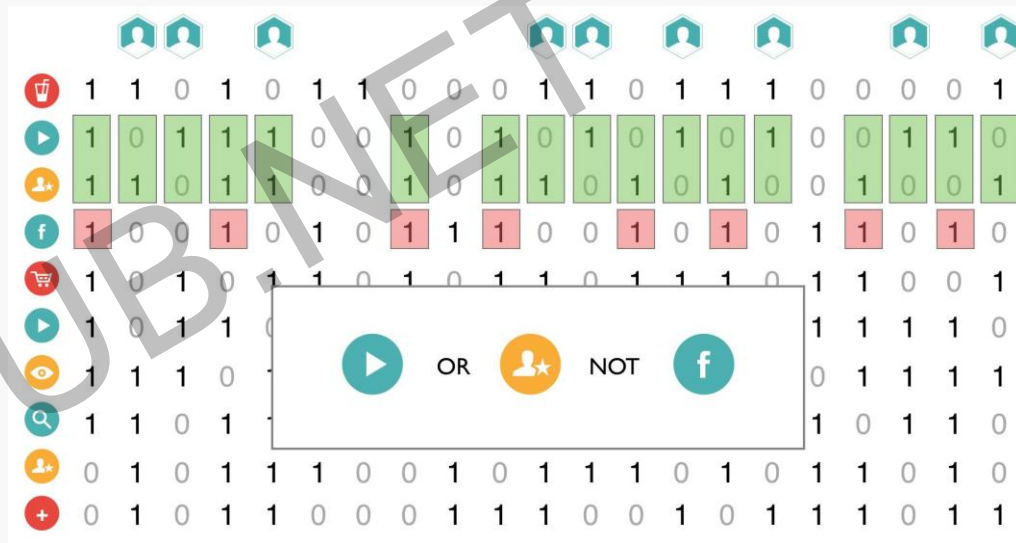
十年架构 成长之路





# Pilosa在品友的实践

- Bitmap对海量用户进行标签
  - 0/1来代表有某标签
- 人群画像速度提升
- 人群查询方便
- 易于扩展



十年架构 成长之路



# Palo vs ClickHouse

|               | Palo    | ClickHouse |
|---------------|---------|------------|
| Count         | 89.7秒   | 0.23秒      |
| Group By(单线程) | 0.14秒   | 0.47秒      |
| Group By(多线程) | 略优      | -          |
| 实时导入          | 不支持     | 支持         |
| 写入速度          | 2.5万行/秒 | 14.2万行/秒   |
| 存储空间          | 8.7G    | 9.5G       |

品友技术



十年架构 成长之路



# 选择。 。 。 选择。 。 。

- ClickHouse vs Palo
- Druid vs Kylin
- GreenPalm vs Elastic Search



十年架构 成长之路



# 自己动手，丰衣足食

- 数据Console
- 机器学习平台
- 加速数据分析 Alluxio
- Knime



十年架构 成长之路



# 品友数据分析平台总结

- 根据数据量，使用者角色设计设计分析平台很重要
- 对工具的选择来说，没有银弹
- 走工具+自我开发的道路



十年架构 成长之路





A network diagram consisting of several blue circular nodes connected by thin blue lines, forming a web-like structure across the top half of the image.

# THANKS



A large, light gray watermark text "ITPUGS.NET" is oriented diagonally across the center of the image, partially overlapping the word "THANKS".



Abstract geometric shapes in the bottom right corner, including overlapping triangles and curved bands in shades of pink, orange, yellow, and light blue.