



2016 杭州·云栖大会
THE COMPUTING CONFERENCE

云栖社区
yq.aliyun.com

基于数加的大数据仓库解决方案



The
Computing
Conference

主办单位:



战略合作伙伴:



署名：宁海元
职称：袋鼠云 CTO



扫码观看大会视频

公司简介

袋鼠云由多名前阿里云资深技术专家创立，核心员工来自**阿里巴巴**、**神州数码**等
做为阿里云战略级合作伙伴，专注于为企业客户提供**云计算和大数据技术服务及产品**

袋鼠云是数加平台**首个金牌合作伙伴**，也是阿里云认证的**区域服务商**和**云市场供应商**

袋鼠云总部位于**杭州**，在**北京**、**贵阳**、**苏州**等地设有分公司（办事处）

¥1000 万人民币

2015年11月20日注册成立，获得天使投资

\$450 万美元

2016年6月获得元璟资本领投盈动资本跟投的Pre-A轮融资



使命

让企业数据产生商业价值



愿景

做DT时代最好的云计算和大数据
服务提供商



目标

为10万+企业客户提供技术服务



目录

content

企业数据仓库现状

数加平台的优势

袋鼠云解决方案

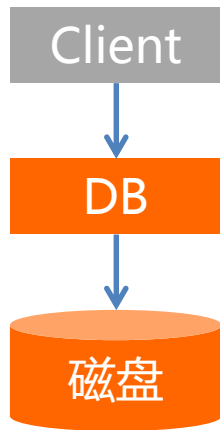
客户案例



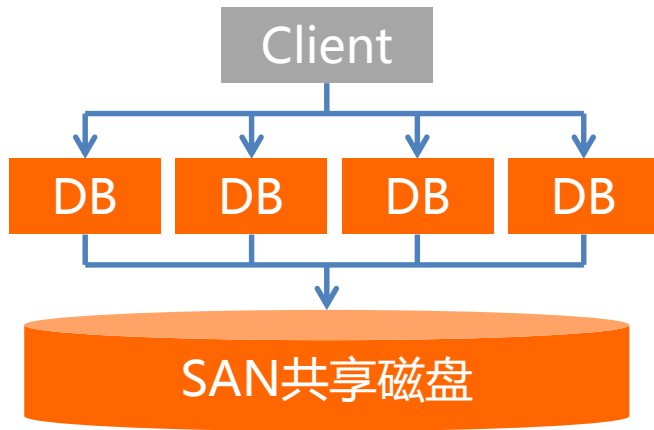
扫码观看大会视频

企业数据仓库现状

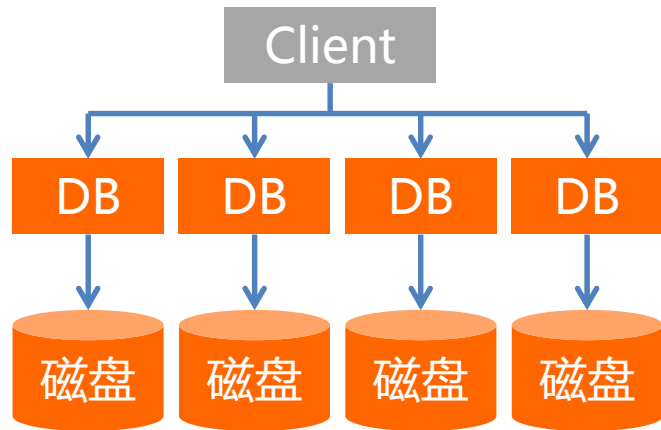




SMP
代表：小型机



共享磁盘
代表：Oracle RAC



MPP
代表：Greenplum, Teradata

- 存储能力受到制约，在数据日益增长的今天，集中式架构已不满足大数据的存储需求。
- OLTP与OLAP混合，导致平台管理已无法兼顾大数据的特殊要求(任务管理，元数据管理等)。
- 针对日益增长的大数据业务需求，需要提供大数据算法、组件支持。



存储

- 单机/共享存储扩展存在明显的天花板，大数据时代容易撞墙
- 价格高昂，无法保存全量数据

计算

- 单机、RAC、MPP的计算能力存在天花板，ETL时间越来越长，早上九点前无法完成全部计算任务
- 任务调度方案复杂，需要很高的学习成本

安全

- 数据仓库容灾成本高，很多企业的数据库没有做容灾，数据存在丢失的风险
- Hadoop的安全粒度较粗，无法满足多租户的隔离



数加平台的优势





拖拽式的操作界面

工作流调度

元数据管理

数据质量管理

一键式任务发布

可视化任务监控

完善的权限管控

弹性扩容

应用场景



数据仓库
海量数据加工



数据挖掘
海量数据挖掘



数据分析
Web excel报表



数据业务应用
智能推荐 | ...

阿里云大数据平台

开发层

数据
开发

数据
管理

数据
分析

数据
挖掘

管理
控制台

数据服务化接口：实时|离线

数据同步：实时|离线

计算层



ODPS离线计算引擎

统一
账号体系

统一
元数据服务



阿里云基础设施



扫码观看大会视频

实时流计算



功能强大

强大流计算引擎；丰富数据采集工具；深度整合各类云存储



性能优越

关键指标超越Storm的性能6到8倍，秒级乃至毫秒级延迟，单个作业吞吐量可做到百万级别



简单易用

支持流式SQL处理 (StreamSQL)；提供全流程流计算开发套件



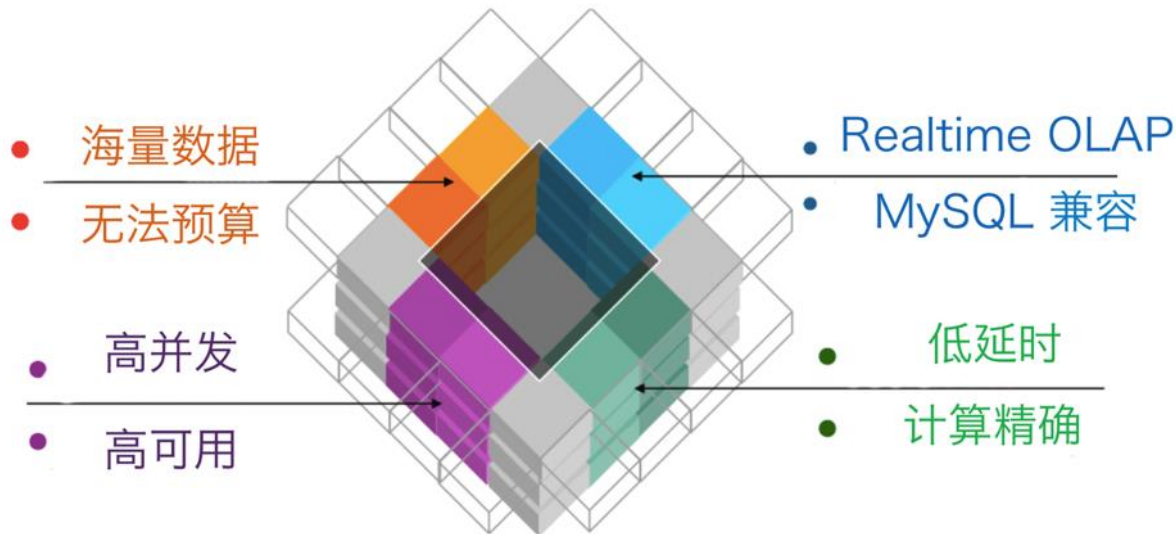
成本低价

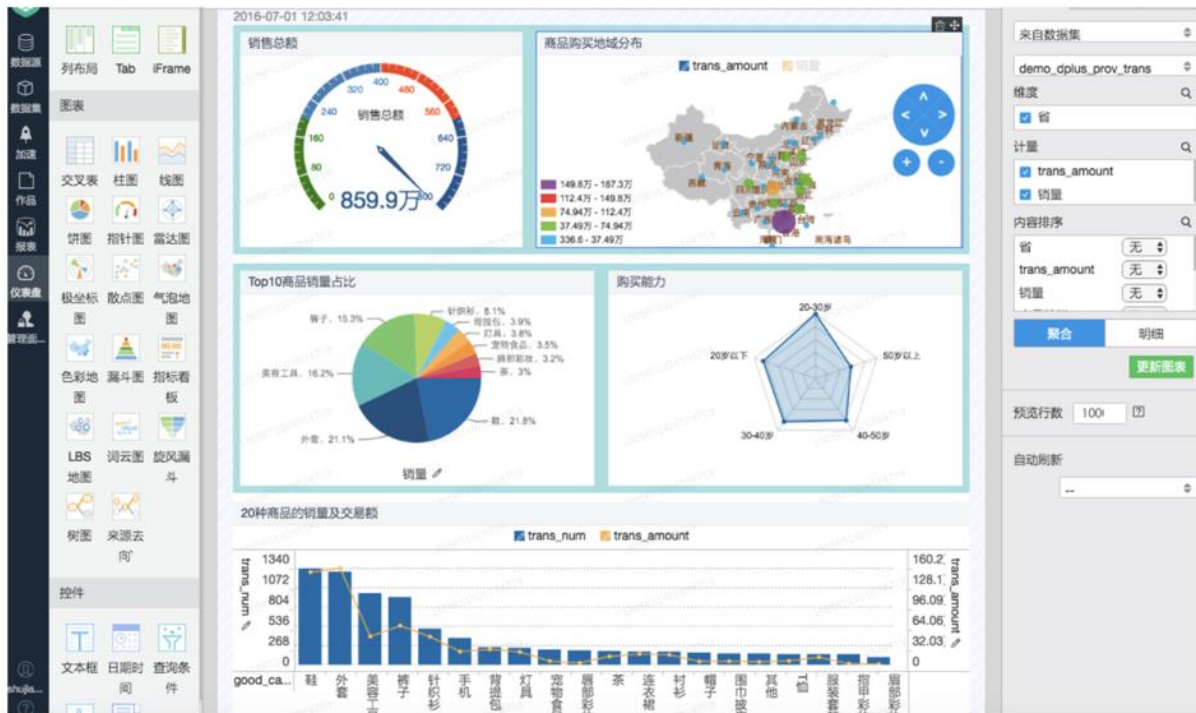
开发流程化、机器规模化；比Storm更省人力和机器成本



分析型数据库

海量数据实时高并发在线分析 (OLAP) 云计算服务, 百亿数据, 毫秒级计算。





Data IDE

阿里云大数据开发平台 dtstack_dev

数据开发 数据管理 运维中心 组织管理 项目管理 其他

输入关键词搜索

任务开发

- chengbw
- cxr_demo01
- junde
- qifeng
- shenhang
- 建表脚本
- cxr_test (汪禹 提交 2016-08-18)
- data_syn (汪禹 提交 2016-08-18)
- select_syn (汪禹 提交 2016-08-18)

脚本开发

节点组件

- 数据加工
 - OPEN_MR
 - OOPS_SQL
 - 机器学习
- 脚本
 - SHELL
- 控制节点
 - 虚节点

test_jiangfeng

```

graph TD
    start([start]) --> sync([sync])
    sync --> testsql[*testsql]
    sync --> testmr[*testmr]
    testsql --> sync_dest[*sync_dest]
    testmr --> sync_dest
    
```

流程图

阿里云大数据开发平台 数据开发 数据管理 运维中心 组织管理 项目管理 其他

root 中文

数据管理

- 全局概览
- 查找数据
- 数据表管理
- 权限管理
- 管理配置

总项目数: 2

总表数: 199

占用存储量: 287.67MB

消耗计算量: 0CU

项目血缘分布图

项目血缘概述

输入表所属项目: odps.dtstack_dev

输出表所属项目: odps.dtstack_dev

阿里云大数据开发平台 dtstack_dev

数据开发 数据管理 运维中心 组织管理 项目管理 其他

root 中文

概览

任务完成情况

已完成 8 失败 3 运行中 0 等待时间 0 等待资源 0 未运行 0

任务执行情况

任务执行时长排行 2016-08-18

| 任务名称 | 责任人 | 时长 |
|----------------|------|-------|
| test_jiangfeng | root | 12分7秒 |
| sync | root | 9分24秒 |



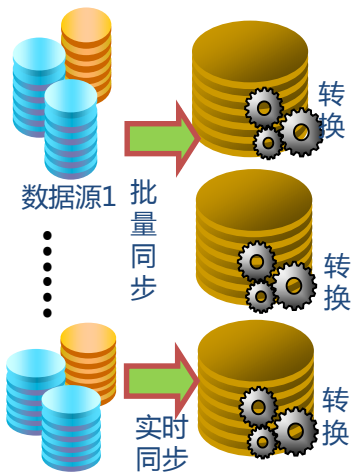
数加平台：完整的数据仓库平台

• 数据如何进入数据仓库

实时/批量数据同步+数据清洗转换

数据同步

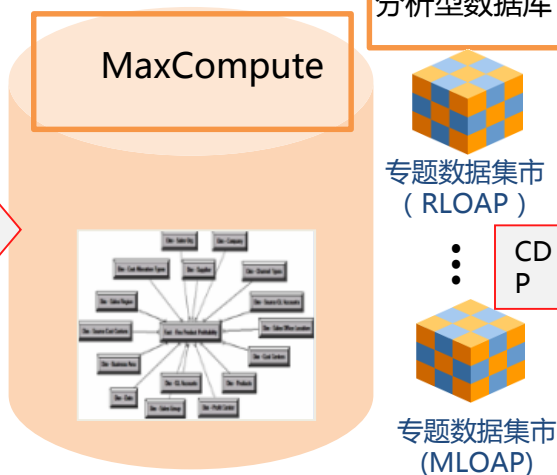
数据转换



• 数据如何管理,性能和空间

数据仓库+主题数据集市

中央数据仓库



• 用户如何使用数据提升管理思想

商务智能分析平台+分析主题

信息展现平台



数据仓库技术平台+企业级数据仓库模型



扫码观看大会视频

存储费用 1TB < 4000/年

| 基础价格 | 大于100GB部分 | 大于1TB部分 | 大于10TB部分 | 大于100TB部分 | 1PB以上部分 |
|--------------|--------------|--------------|--------------|-------------|-----------|
| 0.0192元/GB/天 | 0.0096元/GB/天 | 0.0084元/GB/天 | 0.0072元/GB/天 | 0.006元/GB/天 | 请通过工单联系我们 |

计算费用 10CU 15000/年起卖

| 资源定义 | 内存 | CPU | 售价 |
|------|-----|-------|--------|
| 1 CU | 4GB | 1 CPU | 150元/月 |



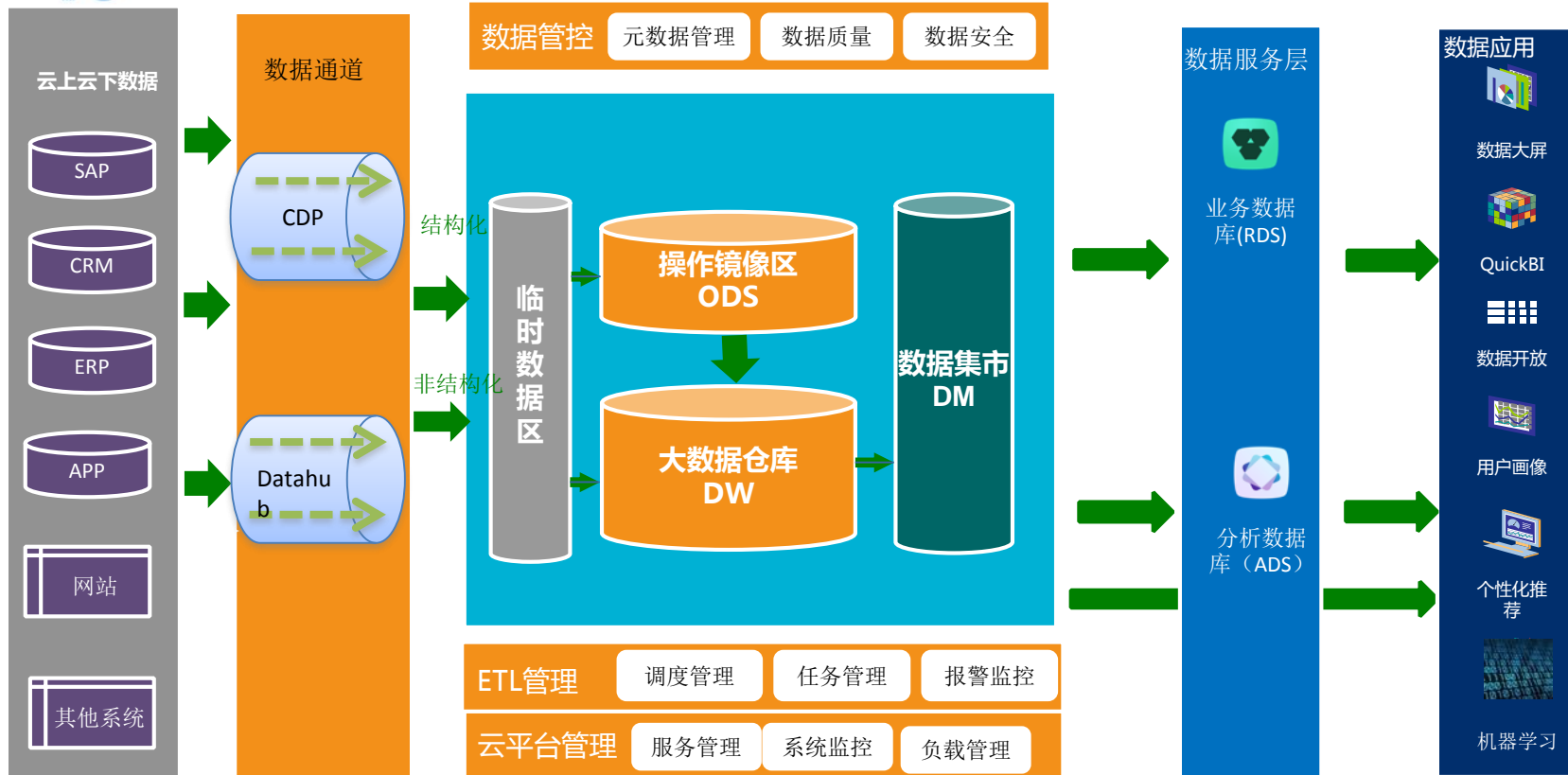
袋鼠云解决方案





2016.6.20 袋鼠云成为数加首个签约的金牌合作伙伴







用户画像



精准营销



推荐引擎



可视化大屏



智能语音



图像识别



公众趋势

基于数加的大数据仓库(MaxCompute/StreamCompute/Analytic DB)



扫码观看大会视频



袋鼠云解决方案



某知名服装企业

客户介绍

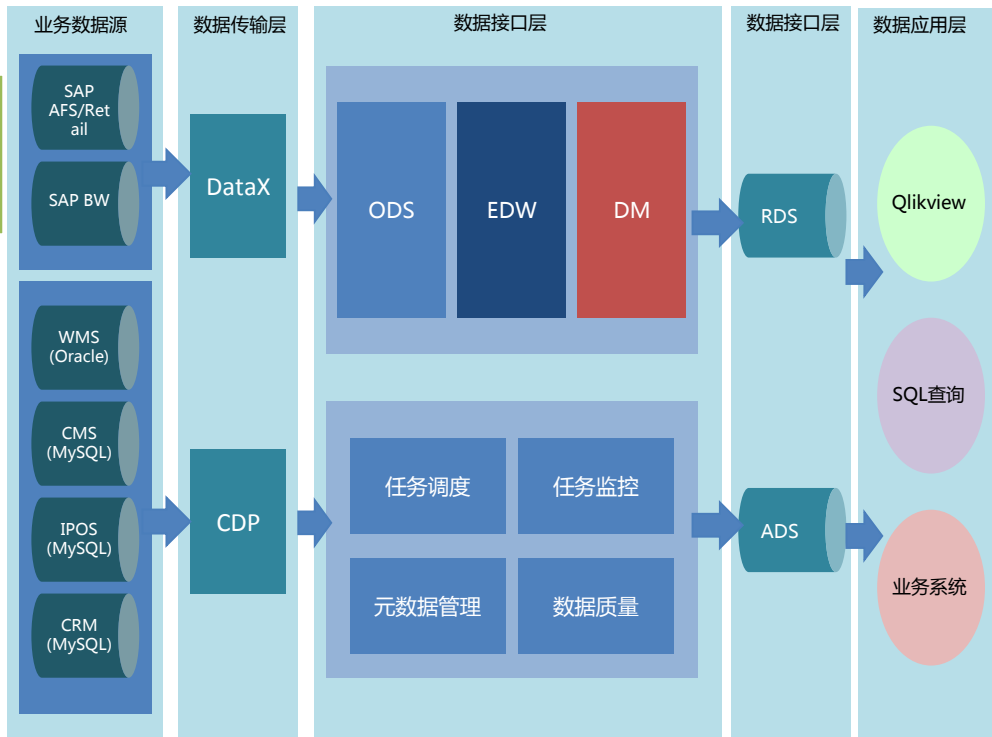
某知名品牌在亚太地区全权运营机构
负责在亚太地区的品牌运作业务,主要经营品牌女装

客户需求

1. 随着数据量增长,即将达到单机数仓的内存天花板
2. 不支持弹性扩容,ETL计算任务执行超过9小时
3. 系统、工具较多,运维管理成本高

袋鼠云服务

1. 基于DataX二次开发服务,支持SAP数据源
2. 基于阿里云大数据仓库架构,ETL时间优化到2小时内
3. 传统数仓无缝升级到云平台,开发运维可视化统一管理





浙报传媒
ZHEJIANG DAILY MEDIA



百联集团

GQY



天弘基金
TIANHONG
ASSET MANAGEMENT



亿享金服
ESHARE.CN

JIAQUWEN®
甲骨文超级码



轻松筹

Leq 乐具
ee

美购

全球精品超市



优速快递



扫码观看大会视频

2016 The
Computing
Conference
THANKS

