

# SACC

卓越 5周年 变迁

SeoueMedia  
盛拓传媒

ITPUB  
www.itpub.com

ChinaUnix

ITPUB

## 2013中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2013

大数据下的IT架构变迁

# 美团数据仓库的演进

美团网数据组 刁士涵

# 美团是一家数据驱动公司

- 每天运行 job 20000+
- 每天新增数百GB原始数据
- 1000多个ETL流程
- 十几条业务线，2000多个业务指标
- 几十个专职数据分析与研究人员



# 演进过程

- **Pre**数据仓库
- 引入ETL
- 构建完整的数据仓库
- 开放和协作

# Pre数据仓库

- 工程师写一段PHP或者Shell统计脚本
- 自己连接业务DB，提取数据
- 在内存中完成统计计算
- 将结果写入报表DB
- 写一个PHP页面作为报表给需求方

# Pre数据仓库

- 很多重复劳动和代码
- 中间数据缺失，中间结果不能共享
- 程序语言五花八门，方法各异很难管理
- 清洗和转换没有统一方法，容易出错
- 不同数据源的数据很难综合使用

# 演进过程

- Pre数据仓库
- 引入**ETL**
- 构建完整的数据仓库
- 开放和协作

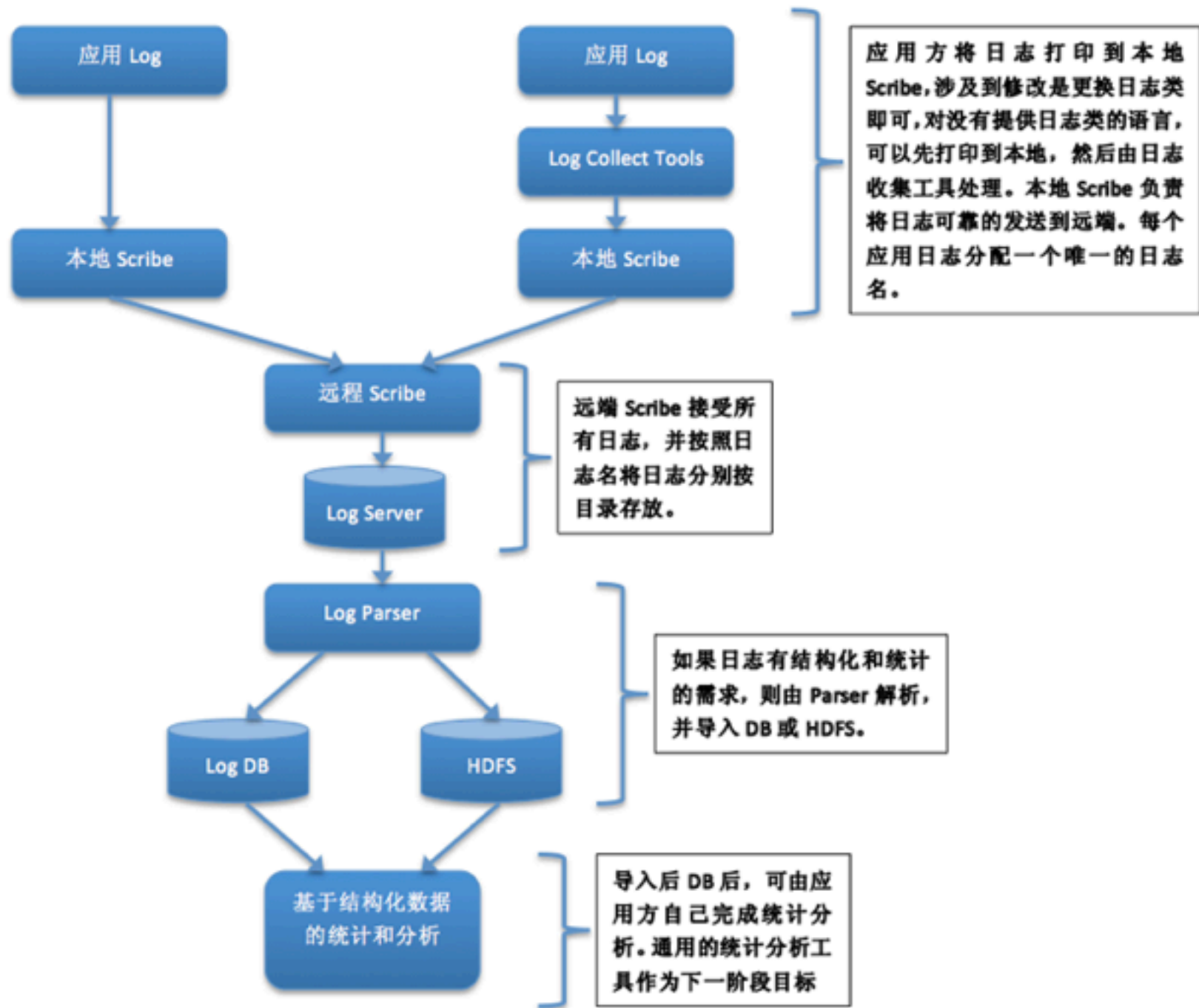
# 引入ETL





# 引入ETL

- 使用独立DB，集中数据，复用中间结果
- 以ETL作为数据处理的核心，简化操作
- 用数据表示逻辑
- 规范数据命名和组织方式
- 进入数据仓库时完成清洗
- 独立出日志收集系统



# 管理工具

- Crontab
- Shell脚本

# 演进过程

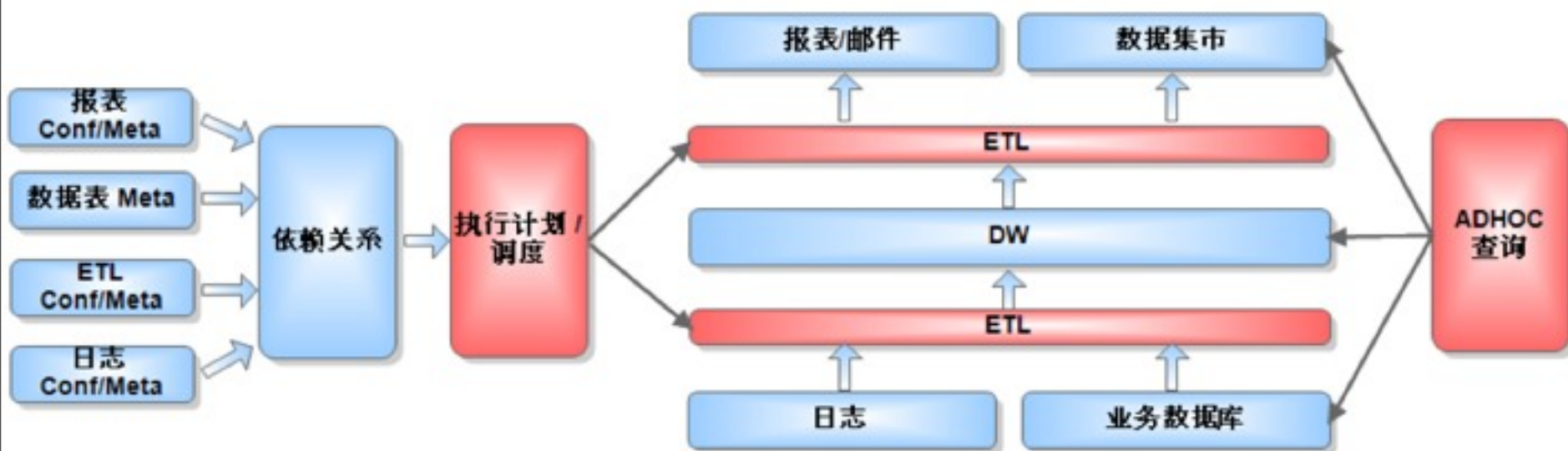
- Pre数据仓库
- 引入ETL
- 构建完整的数据仓库
- 开放和协作

# 遇到了新问题

- 有多少流程在线？执行了多久？
- 很多依赖关系，A必须在B后执行
- 并行开发，出现冲突
- 指标多，概念相近，解释麻烦
- 手写报表效率低

# 针对问题开发工具

- 流程注册、管理、查看工具
- 流程依赖关系解析，画出依赖关系图
- 开发调度系统，根据关系图调度ETL执行
- 抽象报表工具，屏蔽报表页面开发
  - 报表 = SQL + 配置
- 建立数据字典，解释概念和指标计算过程



- 数据仓库是一套完整的软件环境，包括数据抽取、存储、计算、查询、展示，以及管理这些过程的工具



# 演进过程

- Pre数据仓库
- 引入ETL
- 构建完整的数据仓库
- 开放和协作

# 新挑战

- 美团的高速发展对数据的需求也高速发展
  - 数据提取和分析需求增长
  - 数据分析人员的增加
  - 数据分析复杂度增加
- 数据团队疲于应付，大量重复性工作
  - 迫切需要需求方自助获取数据并分析

# 开放数据平台

- 建立主题表，方便分析人员
- 开发自助查询工具
- 开放报表工具
- ETL编辑环境Web化，开放给RD
- 建立数据集市

# 开放带来的问题

- 资源分配与管理
- 版本、审核、测试、部署
- 监控：数据质量、锁、进度、资源使用
- 调度策略改进，持久化、多队列、并发重导
- 权限与事后审计
- 日志系统升级

# 美团数据仓库后续的挑战

- 实时性
- 多机房
- 数据安全
- 一致性
- 集成平台

# Thanks!

SequeMedia  
盛拓传媒

IT168.com  
www.it168.com

ChinaUnix.net

ITPUB