



# Day11 让编排变的简单-数据湖工厂

HUAWEI TECHNOLOGIES CO., LTD.

[www.huawei.com](http://www.huawei.com)



# 大数据处理

大数据技术，是指从各种各样类型的数据中，快速获得有价值信息的能力。

大数据处理的过程就像做饭一样。

买菜



洗菜



配菜



炒菜



数据集成

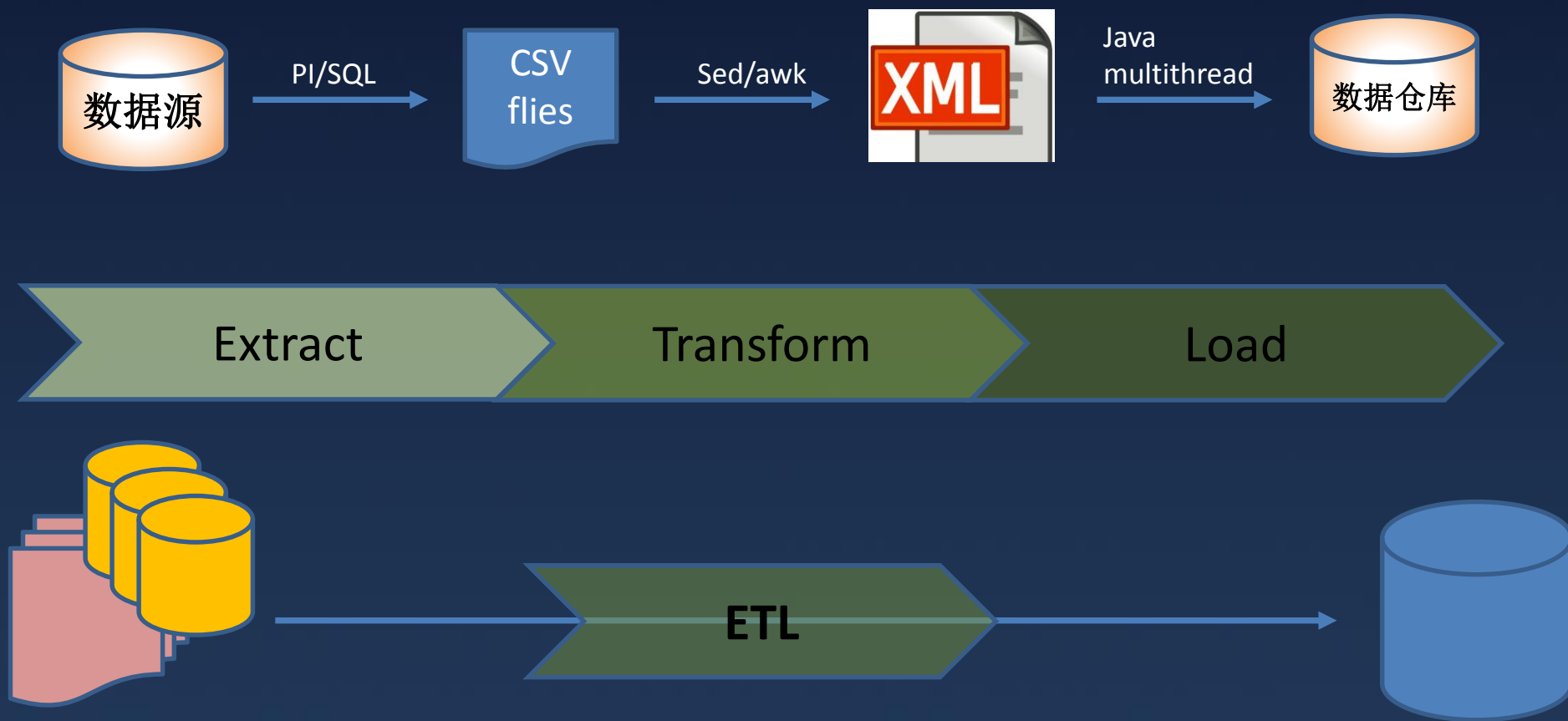
数据清洗

数据建模

数据加工

# 大数据处理工具——ETL

**ETL**，是英文 Extract-Transform-Load 的缩写，用来描述将数据从来源端经过萃取（extract）、转置（transform）、加载（load）至目的端的过程。



## 使用ETL工具会遇到的一些问题.....

- 数据没有按时处理完成怎么办？
- 数据处理失败了怎么办？
- 数据处理过程中其中一个步骤失败怎么办？
- 如何对数据处理的过程进行运维和监控？
- 如何控制数据处理工作流的调度？



# DLF解决这些问题

大数据处理一般较为复杂，数据建模、数据集成、脚本开发、作业调度、运维监控需要使用多种工具和服务组合，使用门槛较高。

## 大数据门槛高

- 涉及多个大数据服务，要串起来一个业务技术要求高
- 每个服务存储格式，访问接口各不相同

## 业务开发难

- 打通本地开发工具与服务配置复杂，安全性差
- 每个服务开发工具各不相同

## 运维专业性强

- workflow 缺乏可视化编排与调度
- 业务执行过程各自监控，运维难度高

## 开发难



## 运维难



# DLF服务介绍

DLF是EI大数据核心服务MRS、DWS、DLI、CloudTable、MLS的配套服务，基本上各种EI大数据的业务场景都可以使用，它是简化用户使用华为云大数据服务的工具。



## 产品优势



一站式IDE平台：**一站式**建设云上数仓，无需切换多个工具



数据湖开发：管理**多种**大数据服务，可实现**跨服务**作业编排和调度



简单易用：在线SQL/Shell脚本编辑调试；预设**10多种**任务类型，拖拽式 workflow 编排



调度稳定高效：丰富的调度配置策略，**百万级别**的作业调度能力

# DLF服务介绍

DLF为serverless多租户架构，即开即用，而且目前是免费的。DLF使用的流程是按照用户使用大数据业务的流程来划分自身的功能。



## DLF目前支持6种数据连接：

- 数据湖探索DLI
- 数据仓库DWS
- MapReduce服务Hive
- MapReduce服务SparkSQL
- 表格存储服务CloudTable
- 云数据库RDS

## 数据集成：

- CDM
- DIS
- CloudStream

## 脚本开发功能：

- 支持立即执行
- 支持格式化SQL
- 支持导出执行结果
- 支持参数设置

## 工作流开发功能：

- 支持多种节点
- 支持拖拽式开发
- 支持立即执行做调试
- 支持配置运行参数
- 支持添加便签

## 作业调度与监控功能：

- 支持启用、暂停或停止调度
- 支持查看调度执行历史
- 支持查看执行日志
- 支持可视化调度依赖关系
- 支持对单个原子重试、暂停、跳过、强制成功等操作



# DLF workflow编排界面

数据湖工厂

总览 数据管理 数据集成 数据开发 运维调度 配置管理

节点库

数据集成

- CDM Job
- DIS Channel
- DIS Dump

计算&分析

- CloudSearch
- CloudStream
- MapReduce
- MachineLearning
- Shell Script
- Spark
- SparkSQL
- DLI Spark

数据库操作

- DWS SQL

保存 测试运行 停止调度 清空画布 全屏 导入 导出

输入节点名称查询

start\_calc

top\_rating\_mo...

top\_active\_mo...

调度方式 \*

☐ 单次调度 ☒ 周期调度 ☐ 事件驱动调度

调度属性 ▾

生效时间 \* 2018/10/31 至 2018/10/31

☐ 从不

调度周期 \* 分钟

开始时间 \* 00 时 00 分

间隔时间 \* 5 分钟

结束时间 \* 23 时 59 分

依赖属性 ▾

依赖作业 ? 请输入作业名称

作业参数配置

咨询·反馈

日志

[INFO] [2018/10/31 15:47:50 GMT +08:00] : 作业开始运行...

[INFO] [2018/10/31 15:48:01 GMT +08:00] : 节点start\_calc运行完成。(开始时间-结束时间:2018/10/31 15:48:42 ~ 2018/10/31 15:48:45)

[INFO] [2018/10/31 15:48:01 GMT +08:00] : 节点top\_rating\_movie开始运行...

- 节点库
- workflow操作按钮
- workflow画布
- workflow调度配置
- workflow运行日志





# 典型场景1：BI报表开发和自动化调度

## 典型的报表工作流：

- 1、使用CDM服务例行将数据集成到数仓
- 2、使用多个SQL脚本进行报表数据生成
- 3、数据生成后再用可视化BI软件展现报表结果

## DLF在BI报表生成场景下的价值：

- 1、支持报表脚本在线开发和调试
- 2、支持工作流编排，从数据集成到多个BI报表脚本
- 3、支持管理员方便进行运维，比如单个节点重试，暂停，跳过，强制成功等



# 典型场景2：机器学习流程自动化（如销售预测）

## 典型的机器学习流程：

- 1、将数据导入到大数据存储
- 2、使用数据处理服务进行数据清洗，整理
- 3、使用机器学习任务进行数据挖掘

## DLF在机器学习场景下的作用：

- 1、支持跨服务的混合流程编排
- 2、机器学习流程端到端可视化



# 典型场景3：电商商品推荐

## 典型的推荐系统流程：

1. 将电商商品推荐的信息统一导入OBS
2. DLI SQL对OBS上的数据进行预处理
3. 使用分析推荐系统对OBS上预处理过的数据进行分析，输出推荐列表

## DLF在推荐系统场景下的作用：

- 1、支持多数据源定时导入
- 2、调用推荐系统API定时完成推荐



数据湖工厂帮助文档详情查看：

<https://support.huaweicloud.com/dlf/index.html>



# Thank You.

**Copyright©2016 Huawei Technologies Co., Ltd. All Rights Reserved.**

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.

# 华为云数据湖工厂服务DLF