

# ETL

**KARPOV.COURSES**

# БЛОК ETL

1. Введение в ETL
2. Знакомство с Airflow
3. Сложные пайплайны, часть 1
4. Сложные пайплайны, часть 2
5. Разработка своих плагинов
6. Установка и настройка Airflow

# ВВЕДЕНИЕ В ETL

**KARPOV.COURSES**

# ЛЕКЦИЯ «ВВЕДЕНИЕ В ETL»

1. Что такое ETL
2. Как правильно готовить ETL
3. Обзор планировщиков
4. Почему Airflow

# ЧТО ТАКОЕ ETL

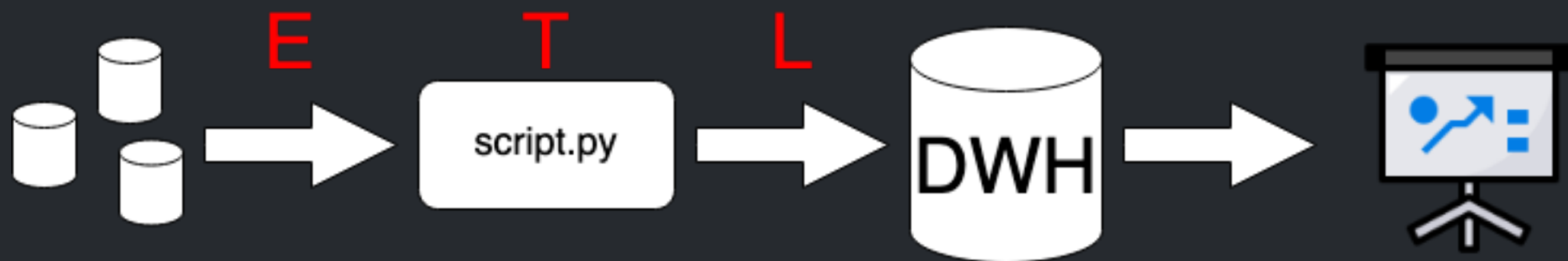
# ETL – ЭТО

процесс миграции данных из одного хранилища в другое

# ETL

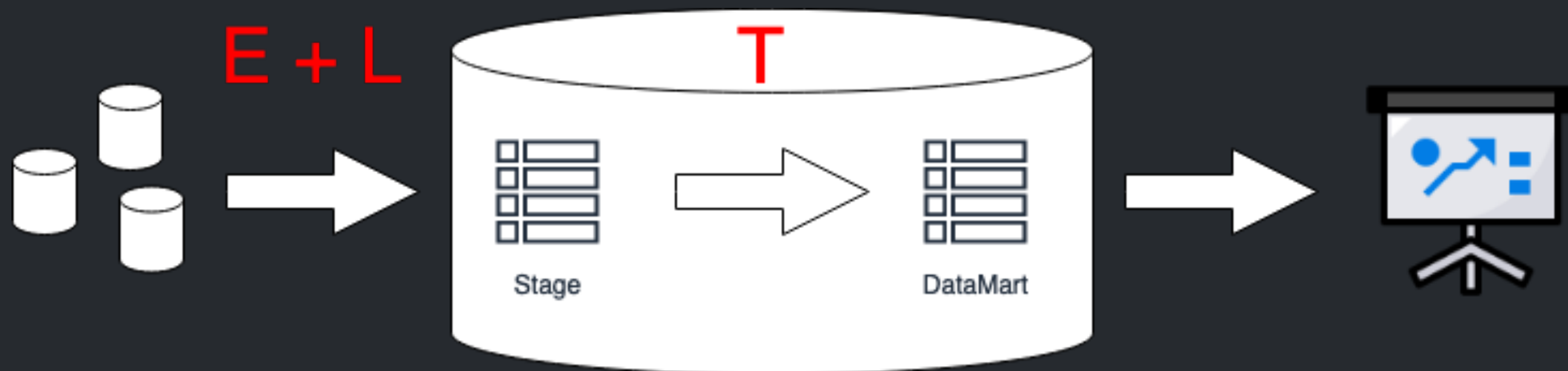
- Extract
- Transform
- Load

# ПРИМЕР ETL





# ПРИМЕР ELT



# КАК ПРАВИЛЬНО ГОТОВИТЬ ETL

# Принципы построения ETL

- Чистый код

# Принципы построения ETL

- Чистый код
- Простота

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой



# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)
- Партиционирование

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)
- Партиционирование
- Инкрементальный пересчёт витрин

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)
- Партиционирование
- Инкрементальный пересчёт витрин
- Загрузка всего без ограничений

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)
- Партиционирование
- Инкрементальный пересчёт витрин
- Загрузка всего без ограничений
- Избавляться от неактуального

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)
- Партиционирование
- Инкрементальный пересчёт витрин
- Загрузка всего без ограничений
- Избавляться от неактуального
- Идемподентность

# Принципы построения ETL

- Чистый код
- Простота
- Единообразие
- Время выполнения пайплайна
- Меньше сетевого трафика
- Работа с репликой
- Оптимизация забора данных
- Способы загрузки данных (SCD)
- Партиционирование
- Инкрементальный пересчёт витрин
- Загрузка всего без ограничений
- Избавляться от неактуального
- Идемподентность
- Аудиторский след



Будьте готовы

# Будьте готовы

— Отсутствие целостности

# Будьте готовы

- Отсутствие целостности
- Сетевые проблемы

# Будьте готовы

- Отсутствие целостности
- Сетевые проблемы
- Незапланированные изменения

# Будьте готовы

- Отсутствие целостности
- Сетевые проблемы
- Незапланированные изменения
- Пайплайны задерживаются

# Будьте готовы

- Отсутствие целостности
- Сетевые проблемы
- Незапланированные изменения
- Пайплайны задерживаются
- Данные в разных системах противоречивы

# ОБЗОР ПЛАНИРОВЩИКОВ

# CRON

## Плюсы

— Максимально простой

## Минусы

— Максимально простой



# JENKINS / GITLAB CI

— Инструмент для CI/CD



GitLab

НАПИСАТЬ СВОЙ

# HH.RU (2021-10-03)

— Airflow	685
— SSIS	296
— Oracle Data Integrator	252
— NiFi	138
— Oozie	67
— Talend	60
— Informatica PowerCenter	55
— Luigi	50
— SAS DIS	44

# Платные шедулеры

- Дорогие
- Нет доступа к коду
- Но есть поддержка
- Визуальный редактор

# MS SQL SERVER INTEGRATION SYSTEM

- Интегрирован в MS SQL Server
- Транзакционность из коробки
- Визуальный редактор
- <https://www.microsoft.com/en-us/sql-server/sql-server-2019>

# ORACLE DATA INTEGRATOR

- Входит в экосистему Oracle
- Визуальный редактор
- <https://www.oracle.com/middleware/technologies/data-integrator.html>

# INFORMATICA POWER CENTER

- Визуальный редактор
- Встроенные утилиты контроля качества данных
- <https://www.informatica.com/products/data-integration/powercenter.html>

# SAS DATA INTEGRATION STUDIO

- Бесплатное обучение
- Визуальное построение пайплайнов
- Язык SAS Base
- SAS Access — плагины для источников



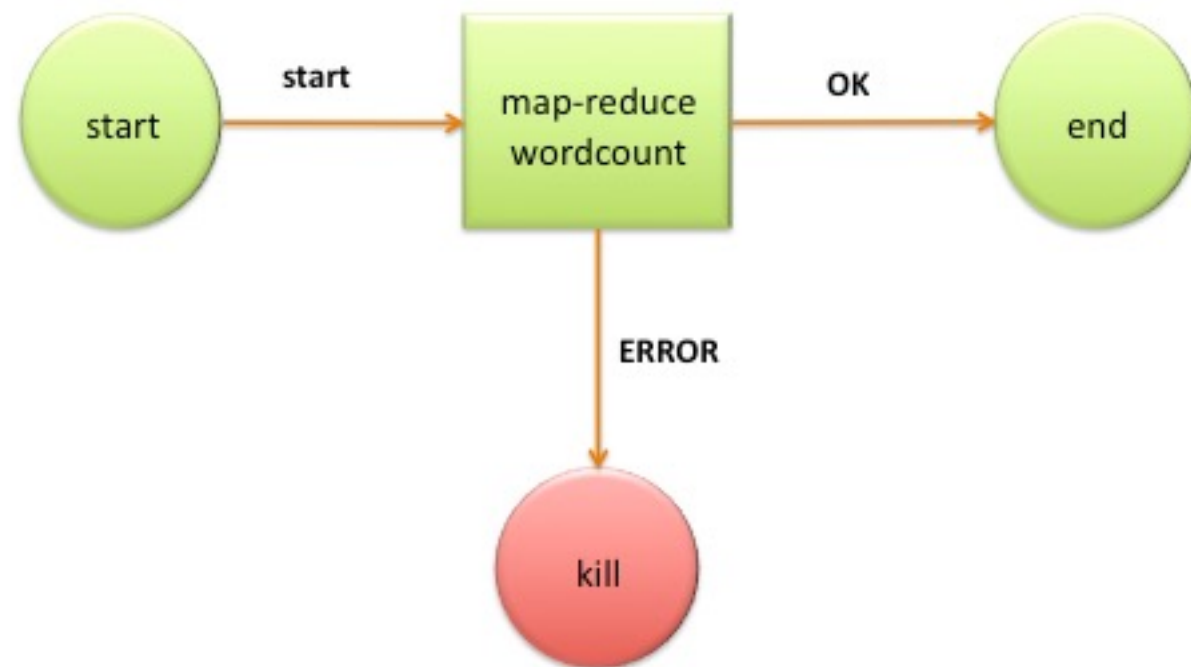
# Open Source

- Бесплатные
- Можно посмотреть в код
- Можно контрибьютить

# APACHE OOZIE



- Планировщик для Hadoop
- Пайплайны на xml
- <https://oozie.apache.org/>

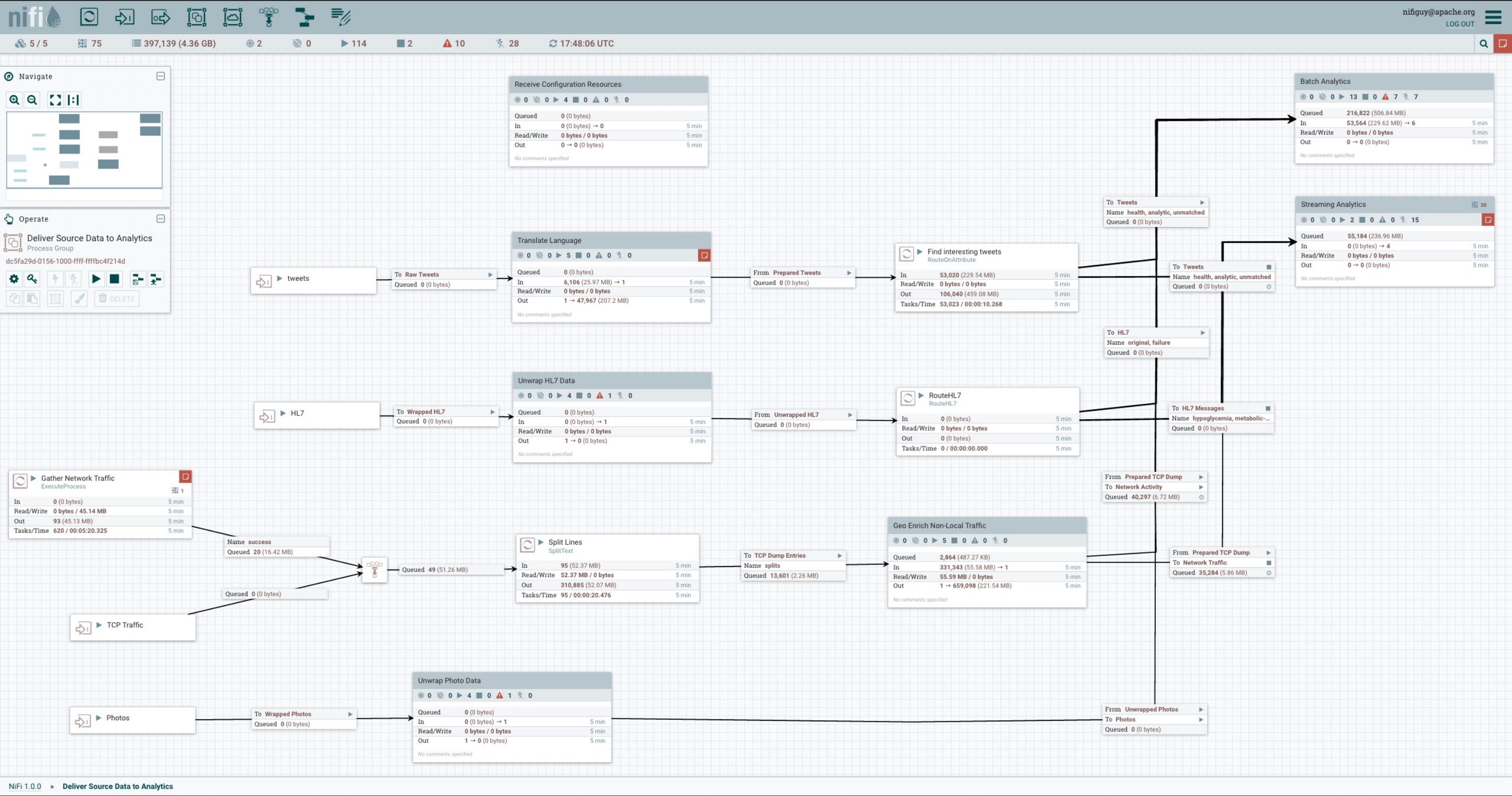


# APACHE NIFI

- Open Source
- Работает внутри Hadoop
- В том числе с внешними источниками
- Поточковая и батчевая обработка
- <https://nifi.apache.org/>



# APACHE NIFI



# TALEND

- Open Source
- Визуальный редактор
- Потoki между задачами
- Требуется знания Java
- <https://www.talend.com/>



# TALEND

**Pipeline Designer**

**Load Salesforce Accounts**

SELECT A RUN PROFILE

**Field Selector 1**  
Field Selector

**Configuration** Info

**SELECTORS** **NEW ELEMENT**

- .Id ID
- Name Account
- .BillingStreet Address
- .BillingCity City
- .BillingState State
- .BillingPostalCode PostalCode
- .Industry Industry

**Data preview - Field Selector 1**

Input Both Output

**Input** 50 records

(record) 198

- Id: 0010H00002Qt57yQAF (string)
- IsDeleted: false (boolean)
- MasterRecordId:
- Name: Macys (string)
- Type: Customer - Direct (string)
- ParentId:
- BillingStreet: 800 EAST CARPENTER (string)
- BillingCity: SPRINGFIELD (string)
- BillingState: IL (string)

**Output** 50 records

(record) 179

- ID: 0010H00002Qt57yQAF (string)
- Account: Macys (string)
- Address: 800 EAST CARPENTER (string)
- City: SPRINGFIELD (string)
- State: IL (string)
- PostalCode: 627690001 (string)
- Industry: Entertainment (string)

(record) 175

RESET SAVE

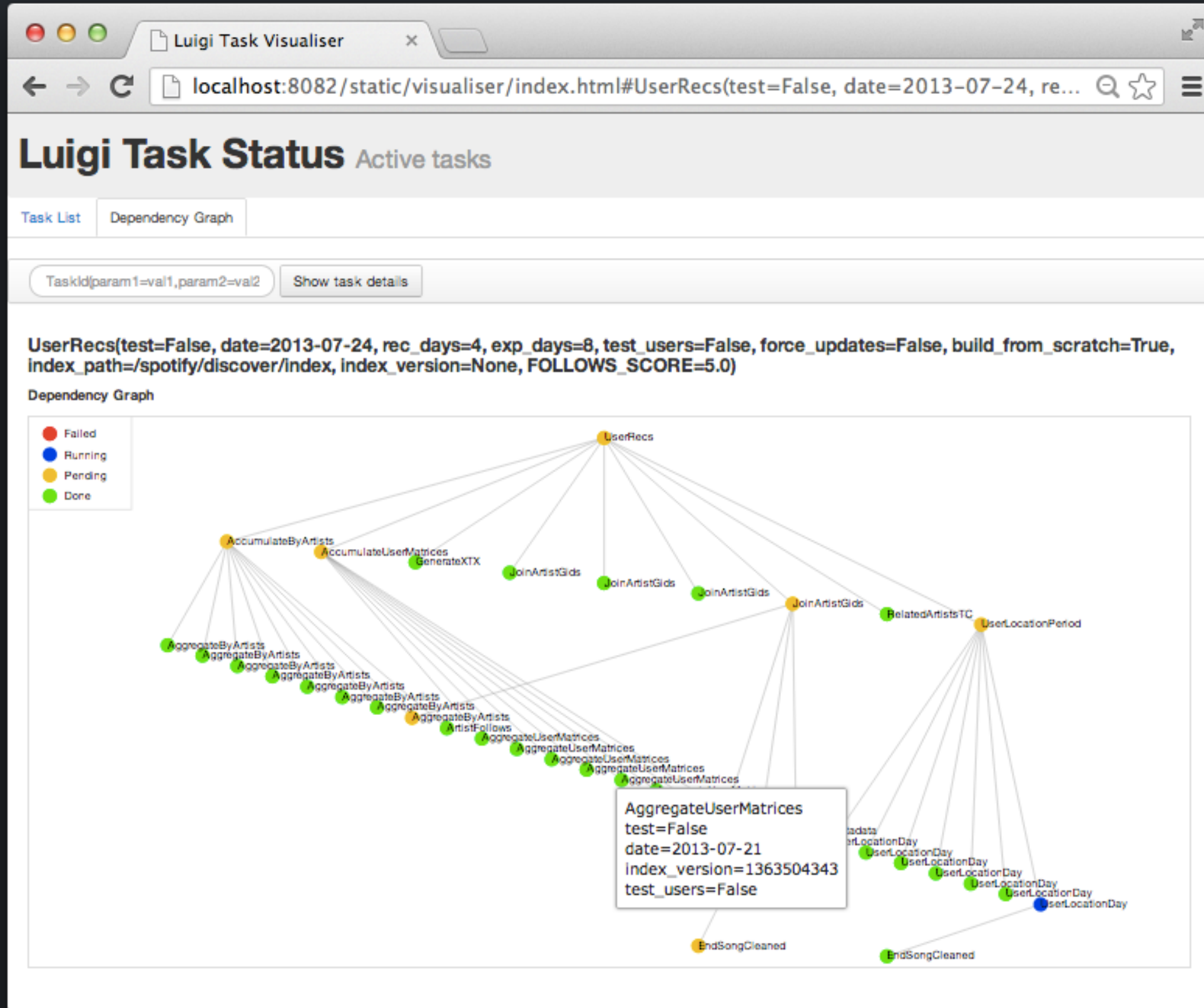
# LUIGI

- Open Source
- Код на Питоне
- Работает поверх артефактов
- Нет встроенного планировщика
- Почти не развивается
- <https://github.com/spotify/luigi>





# LUIGI





ПОЧЕМУ AIRFLOW

# Преимущества Airflow

- Open-source

# Преимущества Airflow

- Open-source
- Отличная документация

# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне

# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне
- Удобный веб-интерфейс

# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне
- Удобный веб-интерфейс
- Алертинг и мониторинг

# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне
- Удобный веб-интерфейс
- Алертинг и мониторинг
- Интеграция с основными источниками

# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне
- Удобный веб-интерфейс
- Алертинг и мониторинг
- Интеграция с основными источниками
- Кастомизация



# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне
- Удобный веб-интерфейс
- Алертинг и мониторинг
- Интеграция с основными источниками
- Кастомизация
- Масштабирование

# Преимущества Airflow

- Open-source
- Отличная документация
- Простой код на Питоне
- Удобный веб-интерфейс
- Алертинг и мониторинг
- Интеграция с основными источниками
- Кастомизация
- Масштабирование
- Большое комьюнити

# СПАСИБО



**ДИНА САФИНА**