

ETL - метод обработки

Extract, Transform, Load

Витрузка

Преобразование

Загрузка в БД

DWH-БД запраив быстро раб.
всеструктурированно
Data Lake - сырым виде,
лучшего копнуть можно

отправка данных на потту (отправка результата)

DAG - контейнер для задач.

task - это экземпляры задач.

- это набор задач которые выполняются в какойто последовательности
от начала и до конца.

является Class.

которые создаются с помощью операторов.

Рекомендации писать DAG

1. Весь код декларативный

императивный код

= 2 + 2 - знаем результат

x + y = x3

2. Идентичность

x + 0 = 0 0 + 1 = 1

Это при повторном процессе операции будет 1 табл, а не 2

+ 3 есть обязательные аргументы DAG

Оператор - нужны для выполнения действий
основ. компоненты, из которых состоит DAG
Python Operator
Email Operator

виды Трансформеры - нужны передать данные из одного источника в другой

• sensor - проверит наличие чего-либо, где-либо
x есть ли такой файл.

• Action - выполнять конкр. задачу (запуск скрипта, отправка на почту)

Зависимости

1 Параллельные t1 → [t2, t3] ① → ② ③

2. Последовательные t1 → t2 ① → ②

stepic of:
3.592 с windows
01.02.24

vmorkin@dnail.com
981017006

task - r

так называемый оператор, который позволяет выполнять ту или иную операцию

Hooks - позволяют взаимодей с внешними приложениями, например БД postgres

- 1 docker-compose down -v
- 2 examples меняем на false
- 3 docker-compose up airflow-init
4. docker-compose down -d ир

Из того состоит Airflow

Сканирует

Scheduler - рекурсивно сканирует папку с DAGами, когда нужно запустить

↓
после того как решил, что пора запускать задачу входит

Запускает задачу

Executor - отвечает за запуск задач, и они отправляются в IPSC сервер

Витаскивает задачи из очереди

↓, появляется

Worker - он витаскивает задачу из очереди и начинает исполнение

↓
для операторов Sensor возможно ситуация повторения задачи, которая именуется как Rescheduler, это позволяет обратно в очередь пока не достигли нужного события.

Статус задачи Success - выполнен

Failed - задача падает или повторяется.

2 типа

так локал

2 + ПК

sequential Executor - использует sqlite бд. не умеет запускать параллельно, только

Local Executor - требует локального инстанса Postgres

Celery Executor - реализован в кластере

Kubernetes Executor - позволяет запускать Airflow в кластере - задача задачи будет каждая в своем собственном контейнере.

Pools - это механизм ограничения параллелизма задач. Если очень много будет задач, сервер улетит, по этому нужно ограничение Pools - 10 и budget max 100 задач

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

Контекст

это информация доступная во время выполнения Оператор
она помогает работать ^{задачами} вместе друг с другом

Контекст содержит след. инфу.

1 Execution Date - дата выполнения (не дата запуска)

2 Task Instance - экземпляр задачи

3 Parameters - параметр - год параметри

4 Маблон - Templates - это перемен или знат, которые можно шлоуз. в коде задачи
с помощью выражение jinja

Jinja - это шаблонизатор, это код напрямую втект.

Макроси - спец. команда, вида $\{\{ds\}\}$ - для упрощ. выполн. задач.

Поддержка - использ. для отправки калоту остатке выполняе задачи

Платини - позволяет расширить возможности AirFlow + оператори, сенсори и тд.

Sensors Оператор - они проверяют условие в течение определенного времени

Если условие выполнено, то DAG будет идти дальше, если нет то sensor

ждет след. проверки и так в течение определенного кол-ва времени, или бесконечно

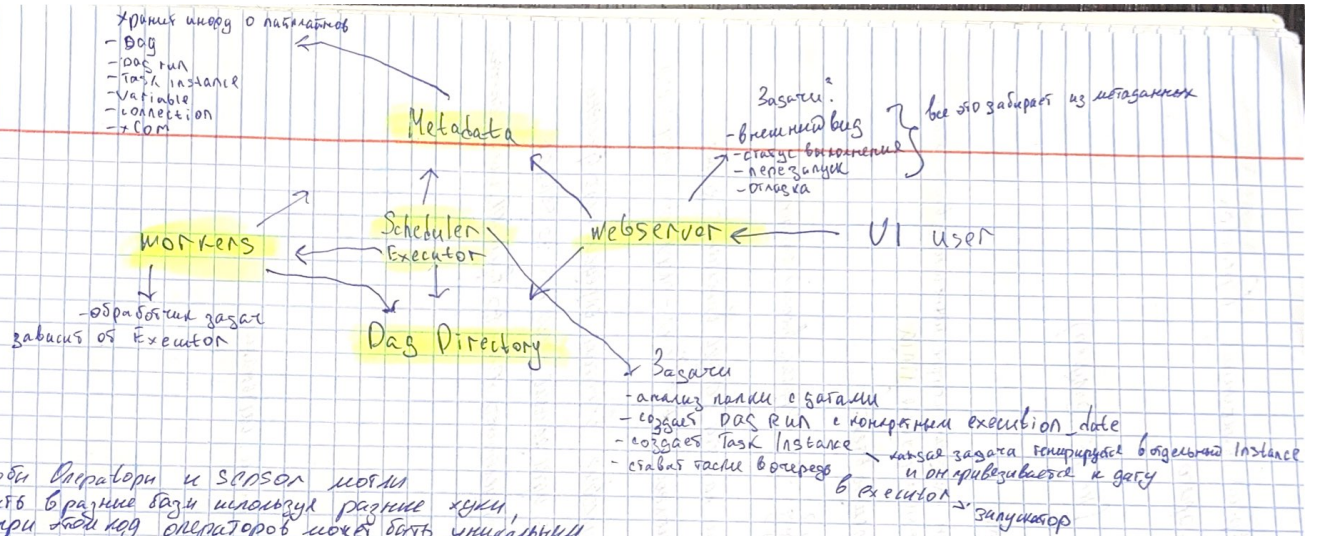
→ ветвления: if else

Branch Python Operator

ShortCircuit Operator

Hook - это интерфейс для взаимоз. с внешн. системами

→ уменьшает дубликаты кода



Хуки - чтобы операторы и sensor могли
хотать разные базисы и использовать разные хуки,
но при этом код операторов может быть уникальным
один код, который берет разные хуки

pool - ограничивает потребление ресурсов.

Хуки, операторы, сенсоры

Connections - для подключения к различным базам

хук - это интерфейс для соединений - PostgreSQL Hook

Операторы - шаблоны для задач - Postgres Operator

Сенсоры - охватывают моменты наступления какого-либо события

параметры → ExternalTaskSensor

timeout

soft-fail

poke-interval

mode - poke, reschedule

Ветвление

Branch Python Operator

Шаблоны Jinja

`{{ ds }}`

...

Макросы - для упрощения задач.

xsom - механизм позволяющий задачам взаимодействовать друг с другом

↓
методы

xsom_push - передает аргумент, который мы хотим получить в другой задаче

xsom_pull - забирает эти параметры

Taskflow API

↓
каждая задача взаимодействует друг с другом напрямую, т.е. результаты работы одного task являются входными параметрами для другого task.

Sub Dags - дочерние даги

↓
периспользование кода

визуальная группировка

Task Group - цели одинаковые

- группировка задач

- периспользование кода

Динамическое создание дагов