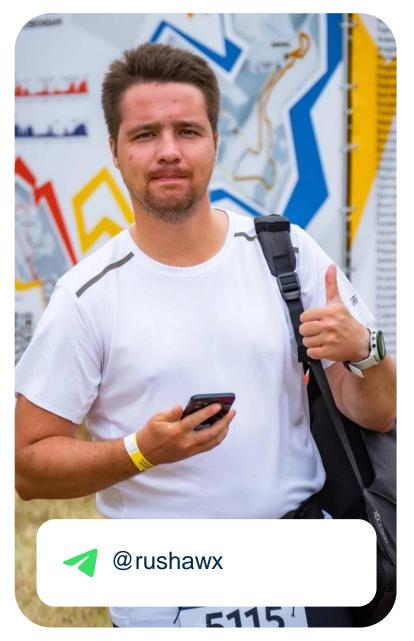
ozon{ech

Apache Airflow







Антон Шишков



Разработчик информационных систем

- Пишу сервисы на Python, Golang, SQL
- ✓ Из инструментов также использую Docker, Airflow, Postgres, Clickhouse, Redis, Elasticsearch, Superset, Grafana, Kafka, Spark
- В свободное время занимаюсь бегом и образованием

REPOSITORY



AGENDA

Что такое Airflow Что такое DAG Архитектура Airflow Подготовка окружения Написание приложения на FastAPI Создание инстансов Postgres, Clickhouse и Minio

Рассмотрение базовых концепций Airflow

Написание базового DAG

Написание продвинутого DAG



ozon{ech

Что такое Airflow



Airflow –

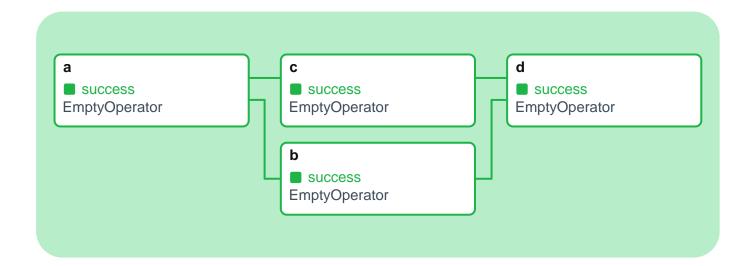
фреймворк для построения конвейеров обработки данных

- Ә Ключевая особенность Airflow заключается в том, что он позволяет легко создавать конвейеры обработки данных, запускаемых по расписанию, с использованием гибкой платформы Python
- → А также предоставляет **множество строительных блоков**, которые позволяют объединить огромное количество различных технологий, встречающихся в современных технологических ландшафтах

Что такое DAG

Ориентированный ациклический граф (DAG) –

это основная концепция Airflow, объединяющая задачи вместе, организованные с зависимостями и связями, чтобы определить, как они должны выполняться



DAG сам по себе не заботится о том, что происходит внутри задач; он лишь занимается тем, как их выполнять — в каком порядке их запускать, сколько раз повторять попытки, есть ли у них ограничения по времени и так далее

https://airflow.apache.org/docs/apacheairflow/stable/core-concepts/dags.html

Архитектура Airflow

В минимальном варианте Airflow состоит из **трёх** компонентов:

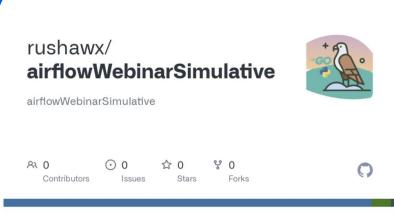
- **э** веб-сервер
- **э** планировщик
- база данных



- **Веб-сервер** и **планировщик** это процессы Airflow
- База данных отдельная служба, которую вы должны предоставить Airflow для хранения метаданных с веб-сервера и планировщика
- Папка с определениями ОАГ должна быть доступна планировщику

Подготовка окружения

- Инициализация Git-репозитория
- Ооздание виртуального Python-окружения
- Установка внешних зависимостей
- → Настройка pre-commit
- Oсздание Makefile
- Подготовка переменных окружения





GitHub rushawx/airflowWebin arSimulative: airflowWebinarSimulat ive

airflowWebinarSimulative. Contribute to rushawx/airflowWebinarSimulative development by creating an account on GitHub.



Написание приложения на FastAPI

fakerApi. Dockerfile

```
# Dockerfile
      FROM python:3.12-slim
 3
 4
 5
      WORKDIR /app
 6
      COPY requirements.txt .
      RUN pip install --no-cache-dir -r requirements.txt
 8
 9
     COPY . .
10
11
      CMD ["uvicorn", "app.app:app", "--host", "0.0.0.0", "
12
          --port", "8000"]
```

Dockerfile

fakerApi. Models

```
# fakerApi/app/models/person.py
     import datetime
     import uuid
     from pydantic import BaseModel
     from typing import Optional
10
     class PersonResponse(BaseModel):
          id: uuid.UUID
12
          name: str
13
          age: int
          address: str
          email: str
16
          phone_number: str
          registration_date: datetime.datetime
18
          created_at: datetime.datetime
19
          updated_at: datetime.datetime
20
          deleted_at: Optional[datetime.datetime] = None
21
```

Models

fakerApi. Handlers

```
# fakerApi/app/handlers/person.py
import datetime
from fastapi import APIRouter
from faker import Faker
from app.models.person import PersonResponse
router = APIRouter(prefix="/person", tags=["person"])
faker = Faker(locale="ru_RU")
@router.get("/", response_model=PersonResponse)
async def get_person():
    person = PersonResponse(
        id=faker.uuid4(),
       name=faker.name(),
        age=faker.random_int(min=18, max=99),
        address=faker.address(),
        email=faker.email(),
        phone_number=faker.phone_number(),
        registration_date=faker.date_time_between(start_date="-1y", end_date="now"),
        created_at=faker.date_time_between(start_date="-1y", end_date="now"),
        updated_at=datetime.datetime.now(),
        deleted_at=None,
    return person
```

Handlers

fakerApi. Main

```
# fakerApi/app/app.py
     import uvicorn
     from fastapi import FastAPI
      from app.handlers.person import router as
          person_router
     app = FastAPI()
      app.include_router(person_router)
10
12
     @app.get("/")
13
      async def root():
          return {"message": "Hello, Simulative!"}
15
16
     if __name__ == "__main__":
17
          uvicorn.run(app, host="0.0.0.0", port=8000)
18
19
```

Main



Поднимем Postgres, Clickhouse, Minio и fakerApi

Postgres



```
# docker-compose-services.yaml
     services:
       pg:
         image: postgres:latest
         container_name: pg
         env_file:
         - ./.env
         environment:
          - POSTGRES_PASSWORD=${PG_PASSWORD}
           - POSTGRES_USER=${PG_USER}
           - POSTGRES_DB=${PG_DATABASE}
         ports:
          - "5432:5432"
         healthcheck:
           test: /usr/bin/pg_isready
           interval: 10s
           timeout: 10s
           retries: 5
         restart: unless-stopped
21
         volumes:
22
           - ./init/pg:/docker-entrypoint-initdb.d
23
```

Postgres

Clickhouse



```
# docker-compose-services.yaml
services:
  zookeeper:
   image: zookeeper:latest
   container_name: zookeeper
   hostname: zookeeper
  ch:
    image: clickhouse/clickhouse-server:latest
   container_name: ch
   hostname: ch
   ports:
     - "8123:8123"
     - "9000:9000"
   volumes:
     - ./data/clickhouse/node1:/etc/clickhouse-server
     - ./data/clickhouse:/docker-entrypoint-initdb.d
   depends_on:
      zookeeper
   healthcheck:
      test: wget --no-verbose --tries=1 http://127.0.0.1:8123/ping ||
          exit 1
      interval: 10s
     timeout: 10s
      retries: 5
```

Clickhouse

Minio



```
# docker-compose-services.yaml
      services:
        minio:
          image: quay.io/minio/minio
         command: server /minio_data --console-address ":9001"
         ports:
           - "9001:9000" # Remap MinIO API port
           - "9002:9001" # Remap MinIO Console port
         environment:
           - MINIO_ROOT_USER=minioadmin
           - MINIO_ROOT_PASSWORD=minioadmin
         volumes:
           - minio_data:/minio_data
         healthcheck:
            test: ["CMD", "curl", "-f", "http://localhost:9002/minio/health/live"]
           interval: 30s
21
           timeout: 20s
           retries: 3
24
```

Minio

fakerApi



```
# docker-compose-services.yaml
      services:
        #
        faker-api:
          build: fakerApi
          ports:
10
            - "8000:8000"
11
          depends_on:
12
            pg:
              condition: service_healthy
13
14
            ch:
              condition: service_healthy
15
16
```

fakerApi

Базовые концепции

- DAG
- DummyOperator
- BashOperator
- PythonOperator
- SQLSensor
- ExternalTaskSensor



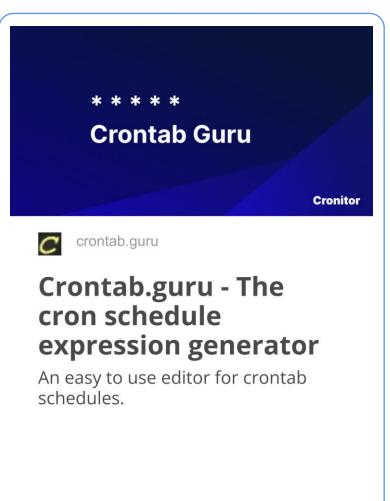
DAG

```
import datetime
      from airflow.models.dag import DAG
      from airflow.operators.dummy import DummyOperator
      from airflow.decorators import task
      from airflow.sensors.sql import SqlSensor
      from airflow.sensors.external_task_sensor import ExternalTaskSensor
10
      with DAG(
11
          dag_id="basic_dag",
          schedule="0 18 * * *",
13
          start_date=datetime.datetime(2025, 1, 1),
14
          catchup=False,
          tags=["simulative"],
15
16
      ) as dag:
17
18
          #
```

DAG

Cron





DummyOperator

```
from airflow.operators.dummy import DummyOperator
dummy_task = DummyOperator(task_id="dummy_task")
```

DummyOperator

BashOperator

```
from airflow.decorators import task
@task.bash
def bash_task():
    return "echo https://airflow.apache.org/"
bash_task = bash_task()
```

BashOperator

PythonOperator

```
from airflow.decorators import task
3
     @task
    def python_task():
5
         print("Hello Python!")
6
8
     python_task = python_task()
```

PythonOperator

SqlOperator

```
from airflow.sensors.sql import SqlSensortask
3
    sql_sensor = SqlSensor(
         task_id='sql_check',
6
         conn_id='postgres',
         sql="""SELECT 1""",
8
```

SqlOperator

ExternalTaskOperator

```
from airflow.sensors.external_task_sensor import ExternalTaskSensor
     external_sensor = ExternalTaskSensor(
          task_id='external_check',
          external_dag_id='simulative_example_basic_dag',
          external_task_id='print_hello',
          execution_date_fn=lambda dt: dt,
         dag=dag
10
```

ExternalTaskOperator

executionOrder

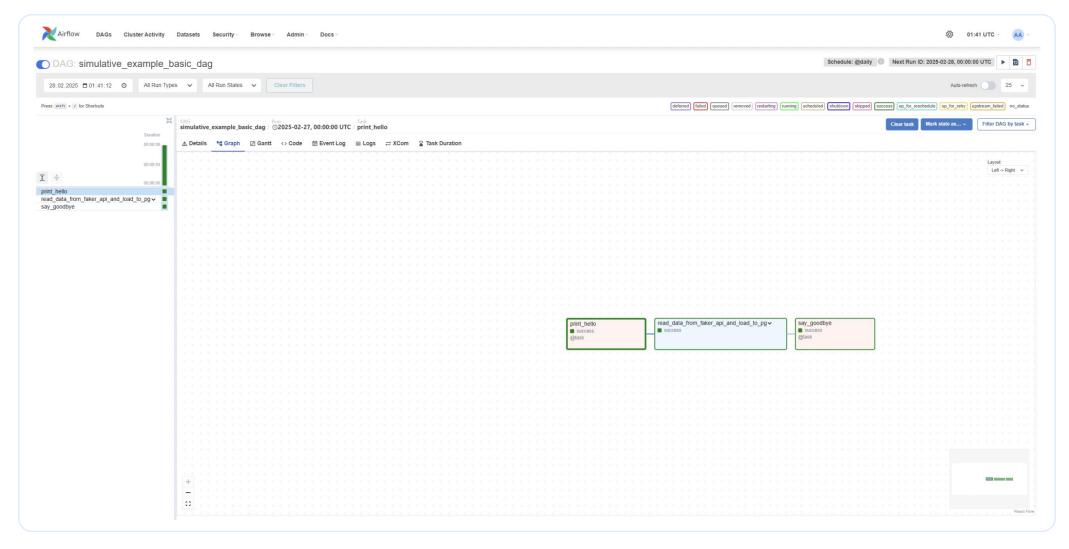
```
bash_task >> python_task >> [sql_sensor, external_sensor] >> dummy_task
```

executionOrder

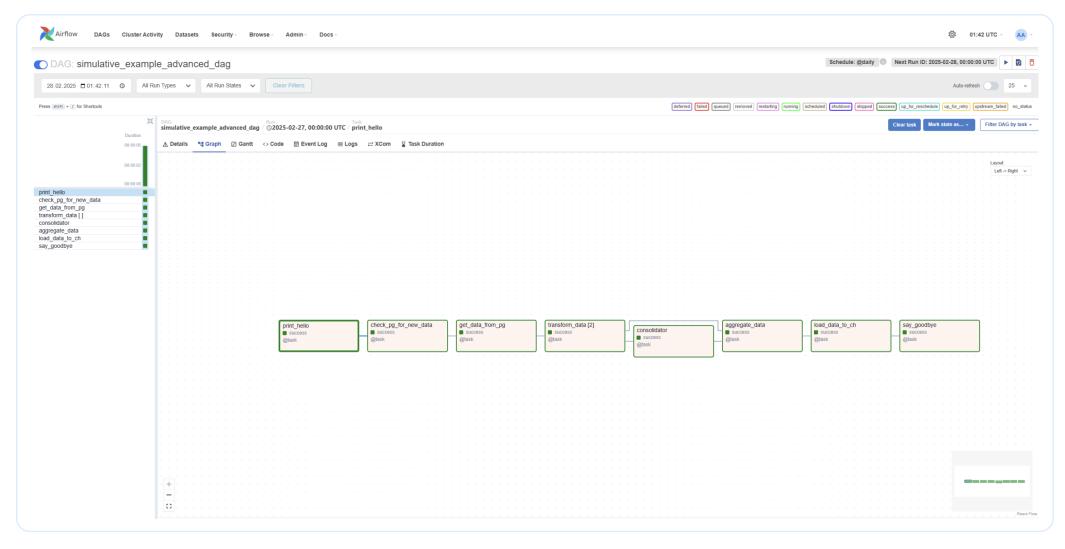


LiveCoding. Airflow

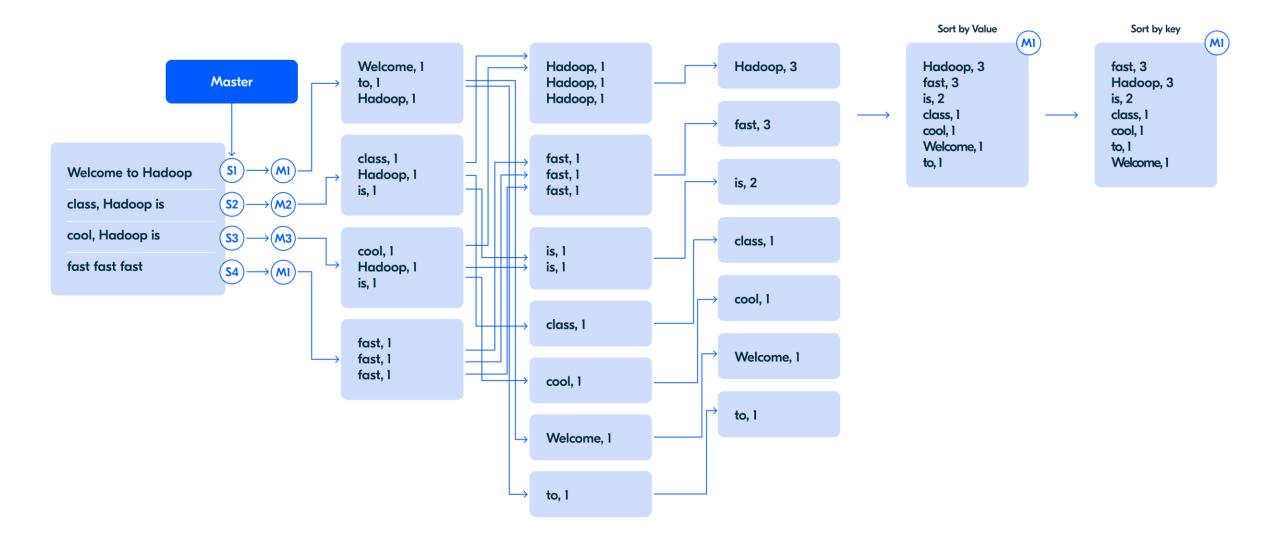
LiveCoding. Airflow. Basic DAG



LiveCoding. Airflow. Advanced DAG



MapReduce



Заключение

За сегодняшний вебинар мы с вами:

- ⇒ познакомились с Airflow
- ⇒ узнали что такое DAG
- ⇒ создали веб-сервис на FastAPI
- → подняли Airflow, Postgres, Clickhouse, Minio
- ⇒ настроили ETL процессы HTTP-SQL-NoSQL



Data Analyst Intern



ozon{ech

Спасибо за внимание!



Антон Шишков

Разработчик информационных систем

ashishkov@ozon.ru / @rushawx