



Apache  
**Airflow**

# Назначение Apache Airflow

Apache Airflow – платформа для **создания, оркестрации, управления** расписанием и **мониторингом** Workflow-процессов загрузок данных.

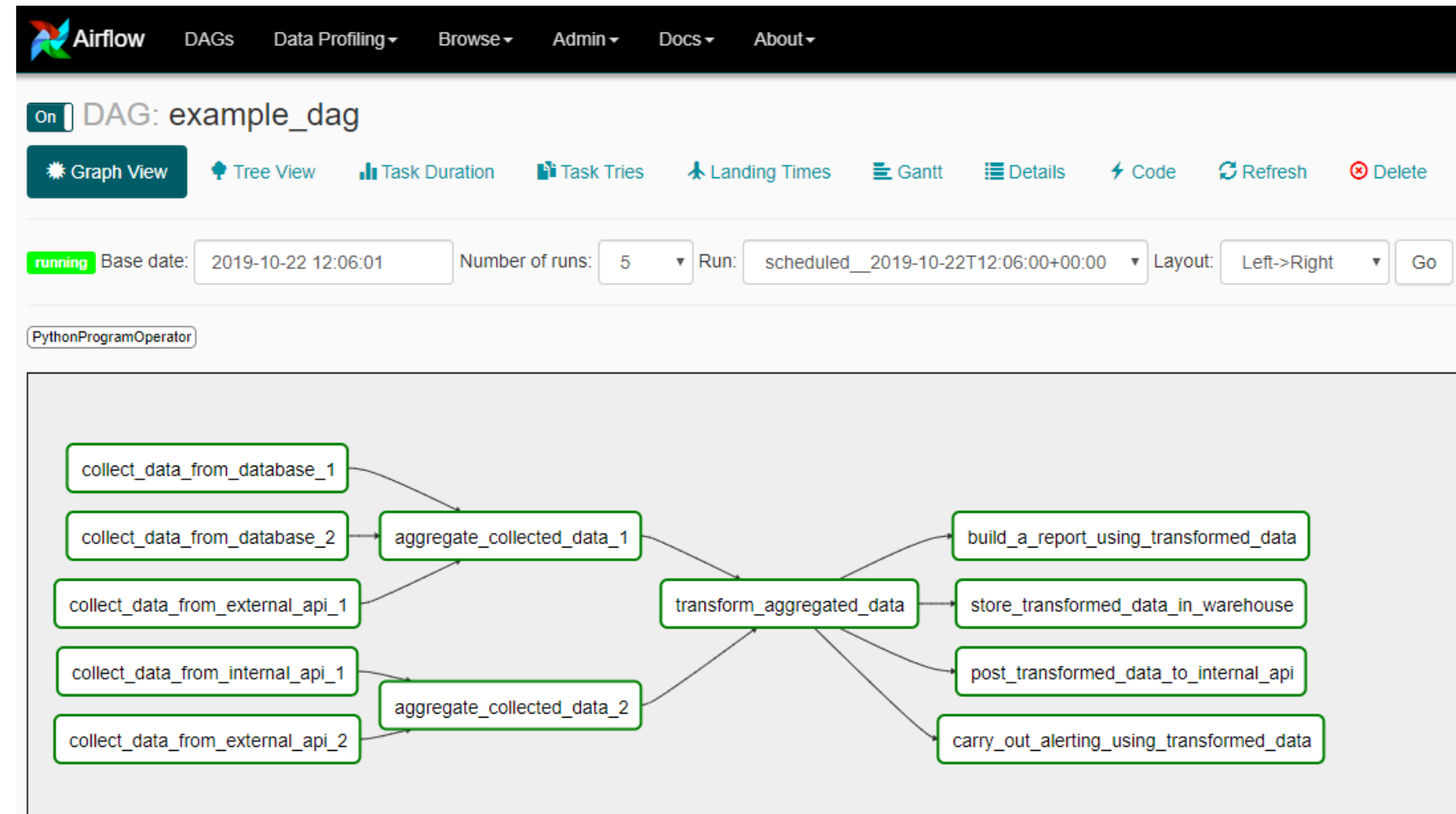
Основные сущности рабочего процесса на Apache Airflow:

- Направленные ациклические графы (DAG)
- Планировщик (Scheduler)
- Операторы (Operators)
- Задачи (Tasks)



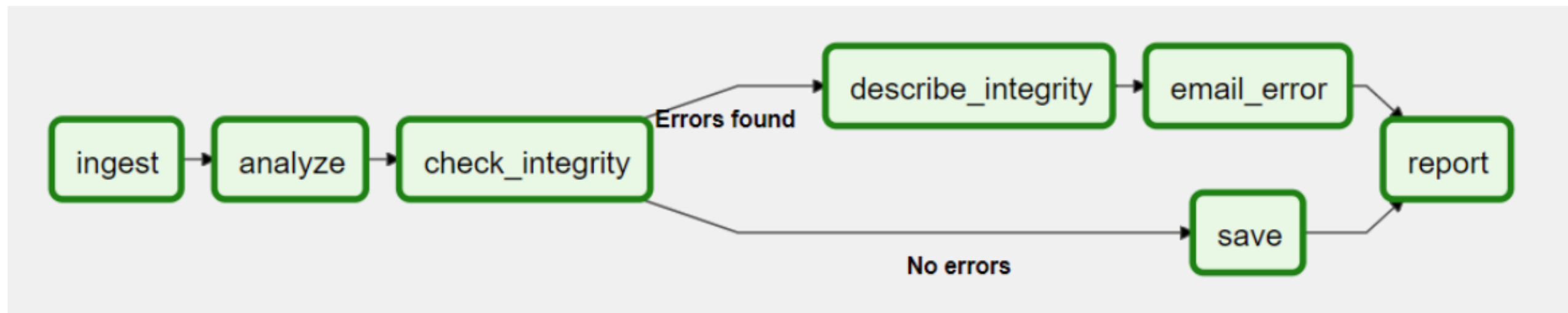
# Основные особенности

- Airflow используется для создания рабочих процессов в виде направленных ациклических графов (DAG) задач.
- Планировщик Airflow выполняет задачи в виде множества рабочих процессов, следуя указанным зависимостям.
- Пользовательский интерфейс позволяет легко визуализировать конвейеры, отслеживать ход выполнения и при необходимости устранять неполадки.
- Когда рабочие процессы определяются как код, они становятся более удобными в сопровождении, версиях, тестировании и совместной работе.



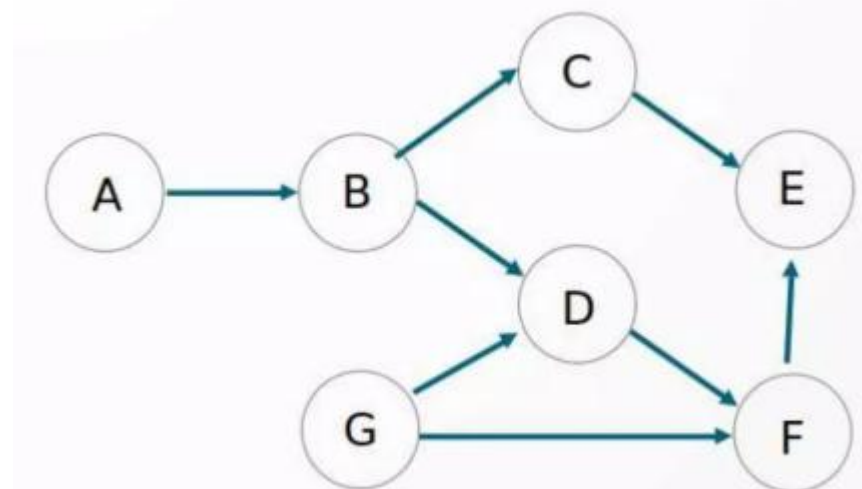
# Основная идея рабочих процессов в управлении данными

Airflow — это платформа, которая позволяет создавать и запускать рабочие процессы. Рабочий процесс представлен как **DAG** (ориентированный ациклический граф) и содержит отдельные части работы, называемые задачами, организованные с учетом зависимостей и потоков данных.



DAG определяет зависимости между задачами и порядок их выполнения и выполнения повторных попыток

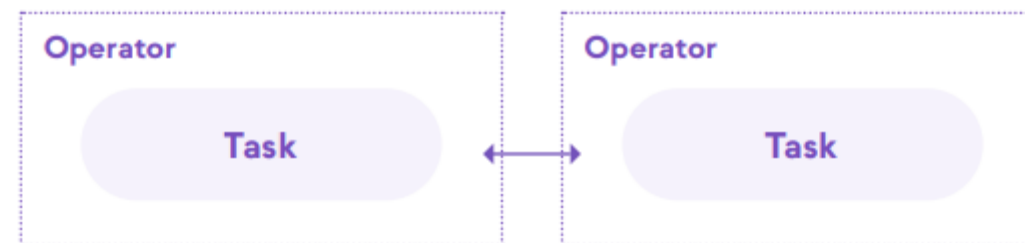
Сами задачи описывают, что нужно делать, будь то получение данных, запуск анализа, запуск других систем, скриптов ML, проверки качества данных, e-mail рассылка и т. д.



# DAG

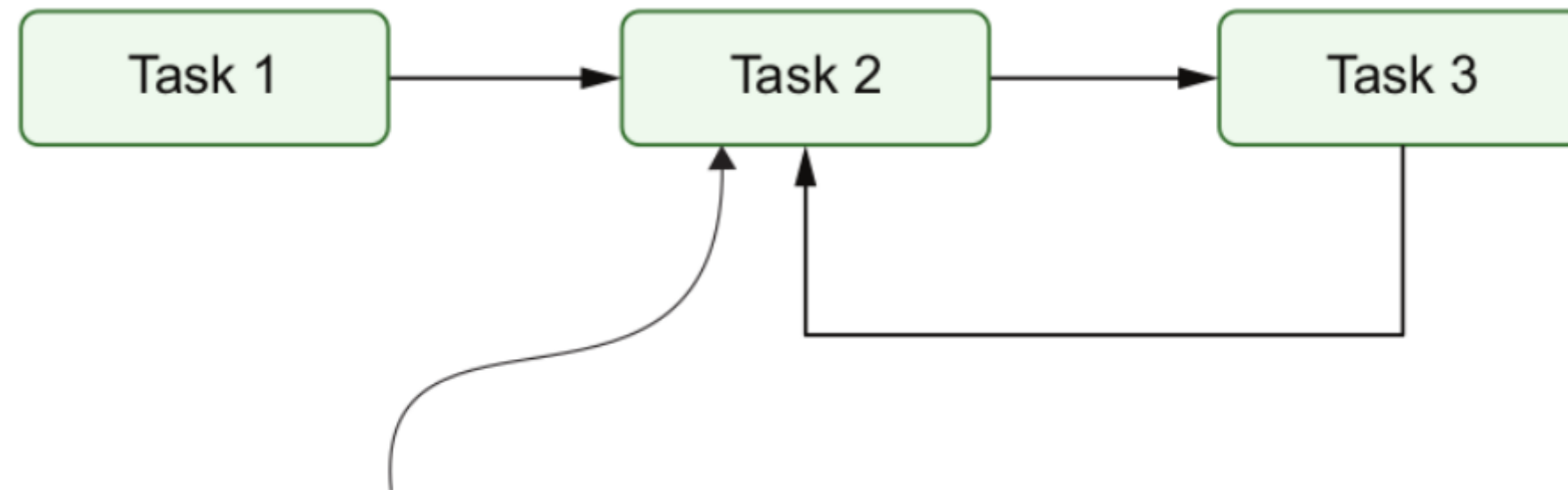
- **Зависимости в DAG** гарантируют, что ваши задачи обработки данных каждый раз выполняются в одном и том же порядке, это создает надежный процесс для вашей повседневной инфраструктуры обработки данных.
- **Графический компонент DAG** позволяет визуализировать зависимости в пользовательском интерфейсе Airflow.
- **Поскольку каждый путь в DAG является линейным**, легко разрабатывать и тестировать конвейеры данных на соответствие ожидаемым результатами.

## DAG



# А почему не циклические графы?

Ориентированный циклический граф привел бы к дэдлоку



Видим, что задача 2 элементарно не может выполняться, т.к. зависит от результата задачи 3, а та, в свою очередь, зависит от задачи 2

# Поток управления

- DAGs предназначены для многократного запуска, и несколько их запусков могут выполняться параллельно.
- DAGs параметризуются, всегда включая интервал, для которого они «выполняются» (интервал данных), но также и с другими необязательными параметрами.
- Задачи имеют зависимости друг от друга. Вы увидите это в DAG либо с помощью операторов >> и <<:

```
first_task >> [second_task, third_task]  
fourth_task << third_task
```

- Или с помощью методов `set_upstream` и `set_downstream`:

```
first_task.set_downstream([second_task, third_task])  
fourth_task.set_upstream(third_task)
```

- Эти зависимости составляют «ребра» графа и то, как Airflow определяет, в каком порядке выполнять задачи. По умолчанию задача будет ждать, пока все ее вышестоящие задачи не будут выполнены успешно, прежде чем она запустится, но это может быть настроено также с помощью функций

# Передача данных между задачами

## XComs («кросс-коммуникации»)

система, в которой вы можете задавать задачи для передачи и извлечения небольших данных

Загрузка и скачивание больших файлов из общей системы хранения

(либо запущенной вами, либо части общедоступного облака)



Airflow отправляет задачи для выполнения на рабочих процессах по мере освобождения места (пула ресурсов), поэтому нет гарантии, что все задачи в вашем DAGe будут выполняться на одном Worker или на одной и той же машине.

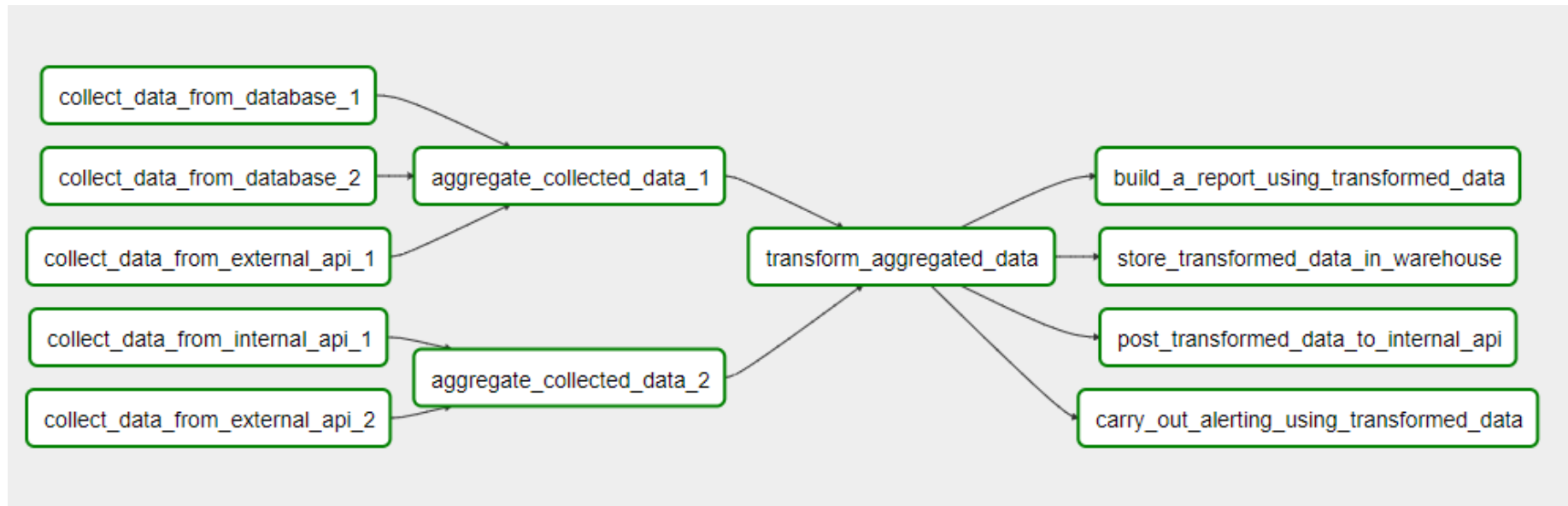


# Передача данных между задачами

По мере того, как вы создаете свои DAG, они становятся очень сложными и большими, поэтому Airflow предоставляет несколько механизмов для того, чтобы сделать это более устойчивым и удобным для восприятия:

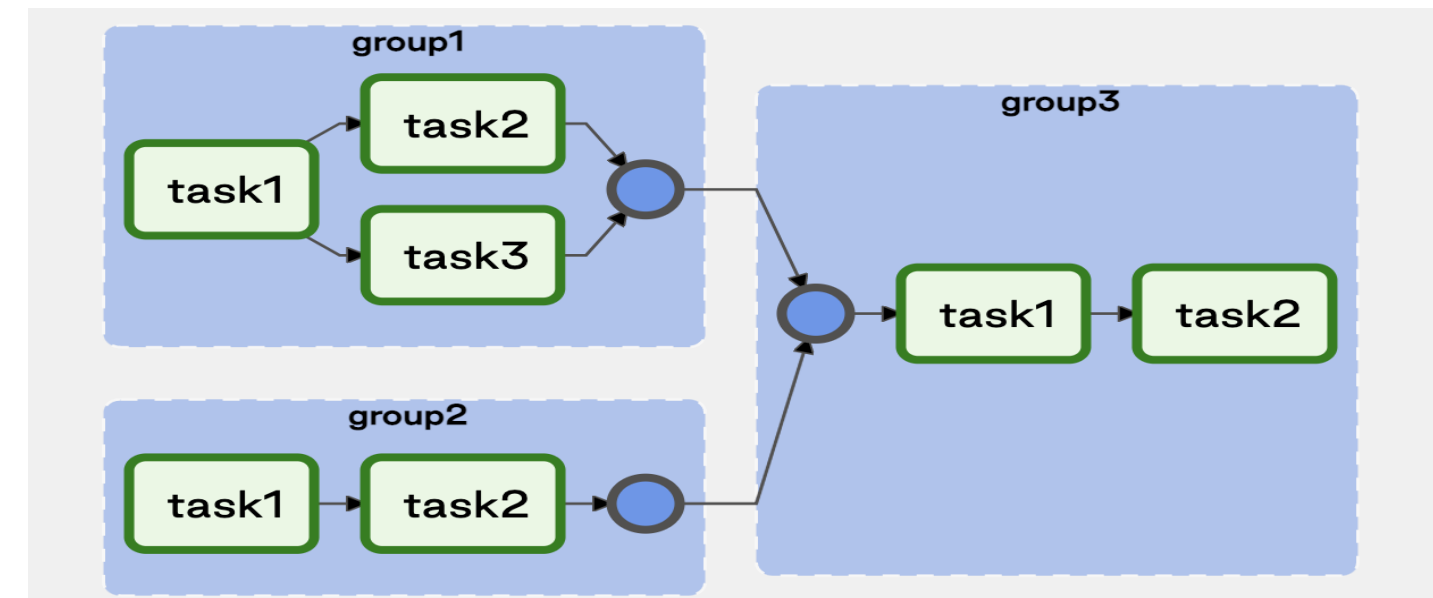
## SubDAG

они позволяют вам создавать «повторно используемые» DAGs, которые вы можете встраивать в другие



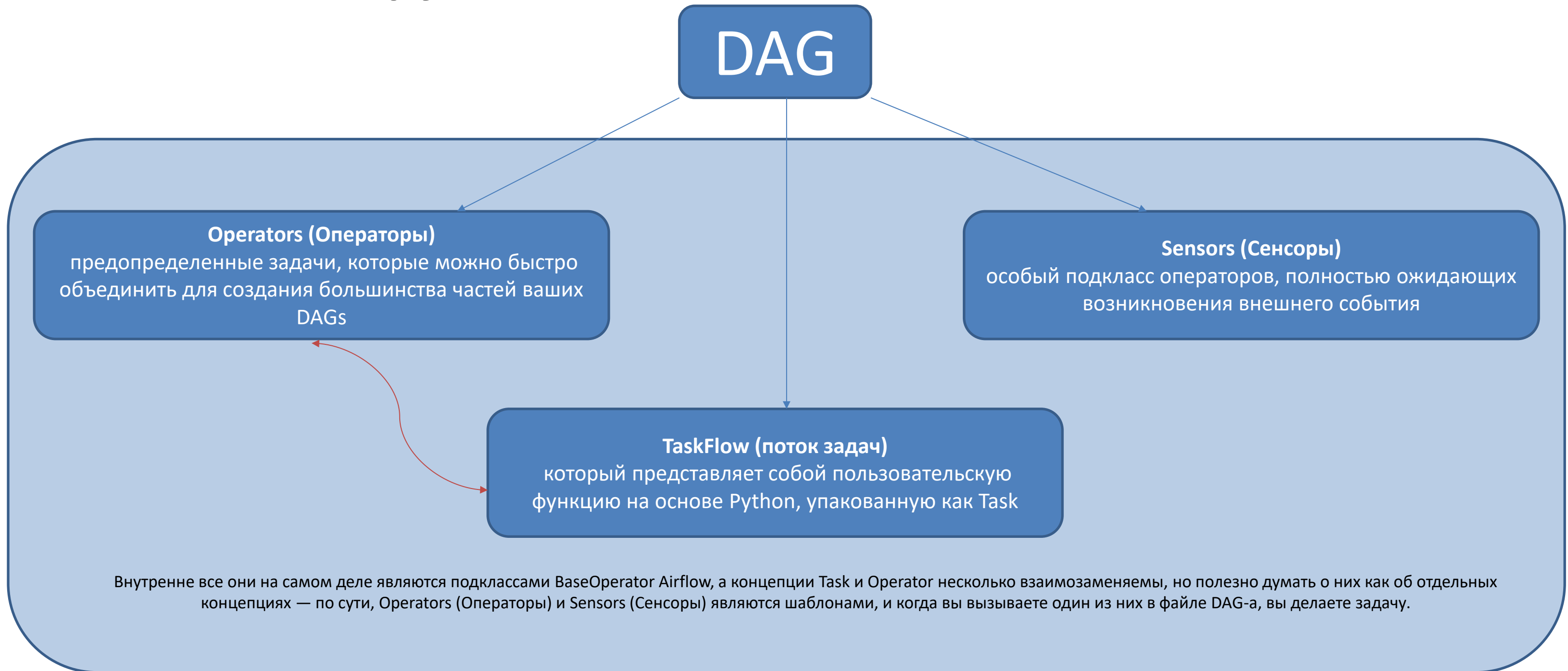
## TaskGroups

позволяют визуально группировать задачи в пользовательский интерфейс



Существуют также функции, позволяющие легко предварительно настроить доступ к центральному ресурсу, такому как хранилище данных, в форме Connections & Hooks, а также для ограничения параллелизма через Pools (пулы)

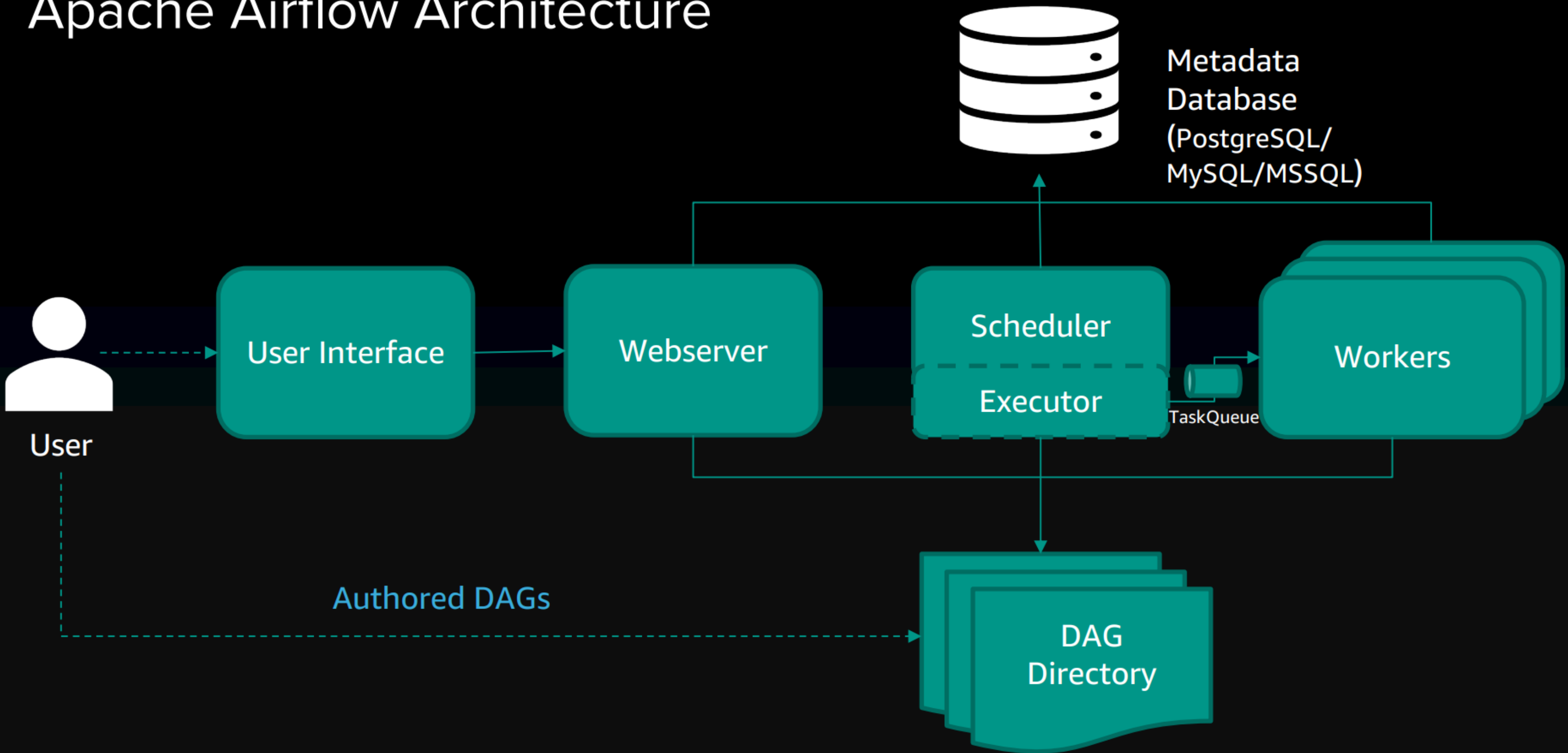
# Рабочие нагрузки



# Операторы

- BashOperator
- PythonOperator
- PostgresOperator
- DockerOperator
- **SparkSubmitOperator**
- etc (>100)

# Apache Airflow Architecture



# Типы Executors



## LocalExecutor

Локальное исполнение  
DAG`ов внутри scheduler



## KubernetesExecutor

Интеграция с Kubernetes  
позволяет использовать  
для каждой Task`и отдельный POD



## CeleryExecutor

Распределённое исполнение с  
помощью механизмов celery.  
Легко масштабируется



## LocalKubernetesExecutor CeleryKubernetesExecutor

Гибридные экзекьюторы, совмещающие  
в себе разные технологии

# DAGs

All 26

Active 10

Paused 16

Filter DAGs by tag

Search DAGs

<div><div>i</div></div> DAG	Owner	Runs <div><div>i</div></div>	Schedule	Last Run <div><div>i</div></div>	Recent Tasks <div><div>i</div></div>	Actions	Links
<div><div><div></div></div><div>example_bash_operator</div><div><div>example</div><div>example2</div></div></div>	airflow	<div><div>2</div><div></div><div></div></div>	0 0 ***	2020-10-26, 21:08:11 <div><div>i</div></div>	<div><div>6</div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_branch_dop_operator_v3</div><div><div>example</div></div></div>	airflow	<div><div></div><div></div><div></div></div>	* / 1 * * * *		<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_branch_operator</div><div><div>example</div><div>example2</div></div></div>	airflow	<div><div></div><div>1</div><div></div></div>	@daily	2020-10-23, 14:09:17 <div><div>i</div></div>	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div>11</div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_complex</div><div><div>example</div><div>example2</div><div>example3</div></div></div>	airflow	<div><div>1</div><div>1</div><div></div></div>	None	2020-10-26, 21:08:04 <div><div>i</div></div>	<div><div>37</div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_external_task_marker_child</div><div></div></div>	airflow	<div><div></div><div>1</div><div></div></div>	None	2020-10-26, 21:07:33 <div><div>i</div></div>	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div>2</div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_external_task_marker_parent</div><div></div></div>	airflow	<div><div></div><div>1</div><div></div></div>	None	2020-10-26, 21:08:34 <div><div>i</div></div>	<div><div>1</div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_kubernetes_executor</div><div><div>example</div><div>example2</div></div></div>	airflow	<div><div></div><div></div><div></div></div>	None		<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_kubernetes_executor_config</div><div><div>example3</div></div></div>	airflow	<div><div></div><div>1</div><div></div></div>	None	2020-10-26, 21:07:40 <div><div>i</div></div>	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div>5</div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_nested_branch_dag</div><div><div>example</div></div></div>	airflow	<div><div></div><div>1</div><div></div></div>	@daily	2020-10-26, 21:07:37 <div><div>i</div></div>	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div>9</div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...
<div><div><div></div></div><div>example_passing_params_via_test_command</div><div><div>example</div></div></div>	airflow	<div><div></div><div></div><div></div></div>	* / 1 * * * *		<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>	...

Press **shift** + **/** for Shortcuts

deferred

failed

queued

removed

restarting

running

scheduled

skipped

success

up\_for\_reschedule

up\_for\_retry

upstream\_failed

no\_status

« »

DAG

Run

dbt\_workflow / 2025-07-01, 07:00:00 MSK

Clear

Mark state as...

Details

Graph

Gantt

Code

Duration

Jun 29, 13:00

03:02:16

01:31:08

00:00:00

stage\_sources

dbt\_run ^

source\_access\_groups

source\_access\_groups\_postgres

source\_accounts

source\_accounts\_postgres

source\_calculator\_intervalpayments

source\_calculator\_intervalpayments\_postgres

source\_calculator\_orderprices\_fixed

source\_calculator\_orderprices\_fixed\_postgres

source\_calculator\_orders

source\_calculator\_orders\_postgres

source\_calculator\_resources

source\_calculator\_resources\_orderstatus

source\_calculator\_resources\_postgres

source\_data\_centers

source\_data\_centers\_postgres

source\_environment\_prefixes

source\_environment\_prefixes\_postgres

source\_events

source\_events\_lastitem

source\_events\_lastitem\_postgres

source\_flavors

source\_flavors\_postgres

source\_folder\_account\_link

source\_folder\_account\_link\_postgres

source\_folders

source\_folders\_postgres

source\_information\_systems

source\_information\_systems\_postgres

source\_order\_actions

source\_order\_actions\_postgres

source\_orders

source\_orders\_postgres

source\_organizations

source\_organizations\_postgres

source\_portal\_options

source\_portal\_options\_postgres

source\_product\_graphs

source\_product\_graphs\_postgres

source\_product\_graphversions

source\_product\_graphversions\_postgres

source\_product\_versions

source\_product\_versions\_postgres

source\_tags\_inventory

source\_information\_systems

source\_information\_systems

source\_events

source\_orders

dict\_projects

source\_folders

source\_service\_manager\_p...

source\_project\_environments

source\_projects

source\_access\_groups

source\_calculator\_intervalp...

source\_information\_system...

dict\_items

source\_folders\_postgres

source\_service\_manager\_p...

source\_projects\_postgres

source\_data\_centers

source\_products

source\_access\_groups\_po...

planner\_report

source\_calculator\_intervalp...

dict\_information\_system

facts\_spends

dict\_orders

dict\_projects\_postgres

dict\_folders

source\_folder\_account\_link

source\_projects\_postgres

source\_data\_centers\_postg...

dict\_products

source\_products\_postgres

dict\_information\_system\_p...

facts\_spends\_postgres

project\_product\_context

facts\_negative\_balance\_t...

cte\_4160\_ris\_spends\_by...

cte\_4160\_ris\_spends\_by...

cte\_4515\_ips

cte\_4515\_ips\_postgres

dict\_folders\_postgres

source\_folder\_account\_link...

source\_project\_environman...

resource\_limits

copy\_constraints

spent\_validation

create\_parquet\_tables

dbt\_clickhouse\_infportal\_run

resource\_pools

planner\_history

resource\_limits\_history

Layout:

Left -> Right

+

-

↺

React Flow



# Apache Airflow Code Example

```
1 from airflow.contrib.sensors.file_sensor import FileSensor
2 from airflow.operators.dummy_operator import DummyOperator
3
4 import datetime
5 from datetime import date, timedelta
6 import airflow
7
8 default_args = {
9     "depends_on_past": False,
10    "start_date": airflow.utils.dates.days_ago(1),
11    "retries": 1,
12    "retry_delay": datetime.timedelta(hours=5),
13 }
14 today = datetime.datetime.today()
15 yesterday = date.today() - timedelta(days=1)
16
17 with airflow.DAG("file_sensor_example", default_args=default_args,
18                 schedule_interval= "*/* * * * *") as dag:
19
20     start_task = DummyOperator(task_id="start")
21     stop_task = DummyOperator(task_id="stop")
22     sensor_task = FileSensor(task_id="file_sensor_task",
23                             poke_interval=30,
24                             fs_conn_id= "<path>",
25                             filepath= "<file or directory name>")
26
27 start_task >> sensor_task >> stop_task
```

Imported Libraries

Default Arguments/dictionary of default parameters

Unique identifier + schedule interval

Task definitions

Dependency



# Demo

<https://tinyurl.com/itmo-airflow>