



Objetivo

- Entender o que é um **Data Pipeline**.
- Desenvolver **pipelines e automatizar processos** utilizando o Apache Airflow.
- Ter mais autonomia e independência.
- Saber se **comunicar e trabalhar em equipe**.
- Estar pronto para **enfrentar qualquer desafio**.



Pra quem é esse curso?

- Analistas, Cientistas e Engenheiros de Dados.
- Gerentes e Heads de Tecnologia.



Responsabilidades dos Profissionais

- **Analistas e Cientistas de Dados.**
 - Especificar regras de negócio do projeto.
 - Níveis de acordo de serviço.
 - Disponibilidade dos dados e informações.
 - Desenvolver e manter os Pipelines de Dados.
 - Definir fontes de dados.
 - Especificar e desenvolver tarefas.
 - Monitoramento e análise.
 - Estados das tarefas.
 - Tempo de execução.
 - Agendamento.



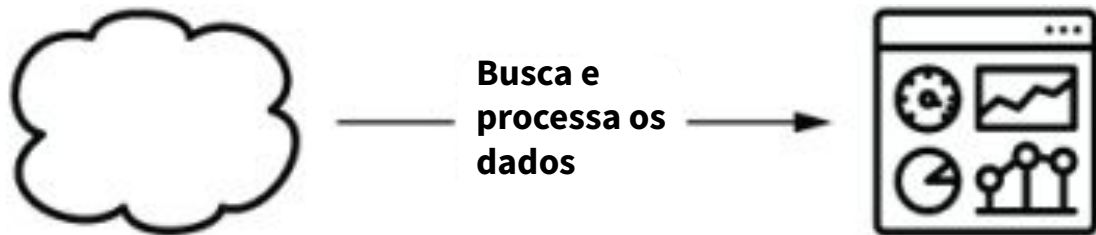
Responsabilidades dos Profissionais

- **Engenheiro de Dados**
 - Desenvolver e manter os Pipelines de Dados
 - Definir ambiente e arquitetura.
 - On Premisse, On Cloud
 - Especificar tecnologias para armazenamento.
 - Data Lakes, Data Warehouse.
 - Definir formas de processamento.
 - Streaming, Clusters.
 - Definir e implementar padrões e segurança e integração com outros sistemas.
 - Monitoramento.



Introdução aos Data Pipelines

Data Pipelines basicamente consiste de algumas tarefas ou ações que devem ser executadas para atingir um resultado.

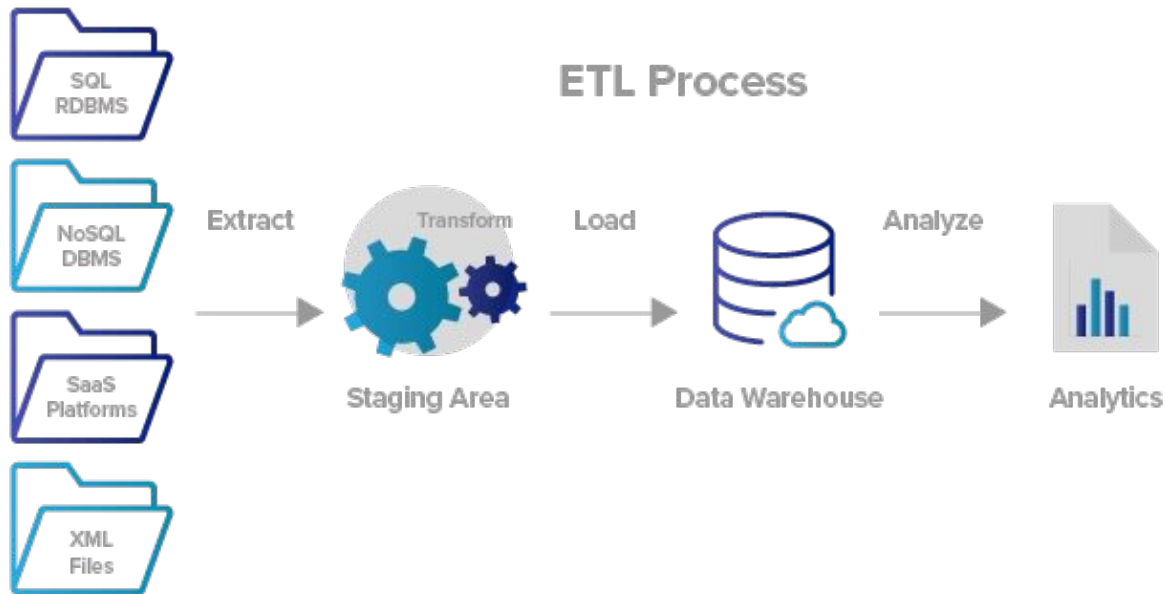


Introdução aos Data Pipelines

1. Busca os dados de temperatura de uma API.
2. Limpa e transforma os dados.
3. Insere os dados transformados em uma base de dados para ser exibido no Dashboard.

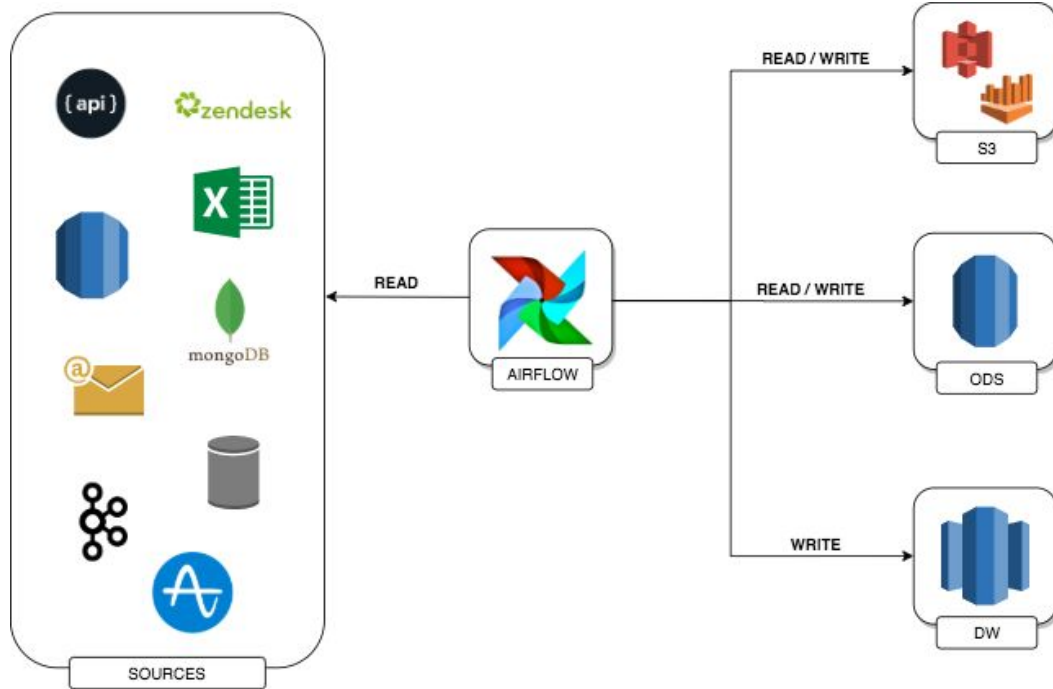
Alguns Pipelines

ETL (Extração, Transformação e Carga).



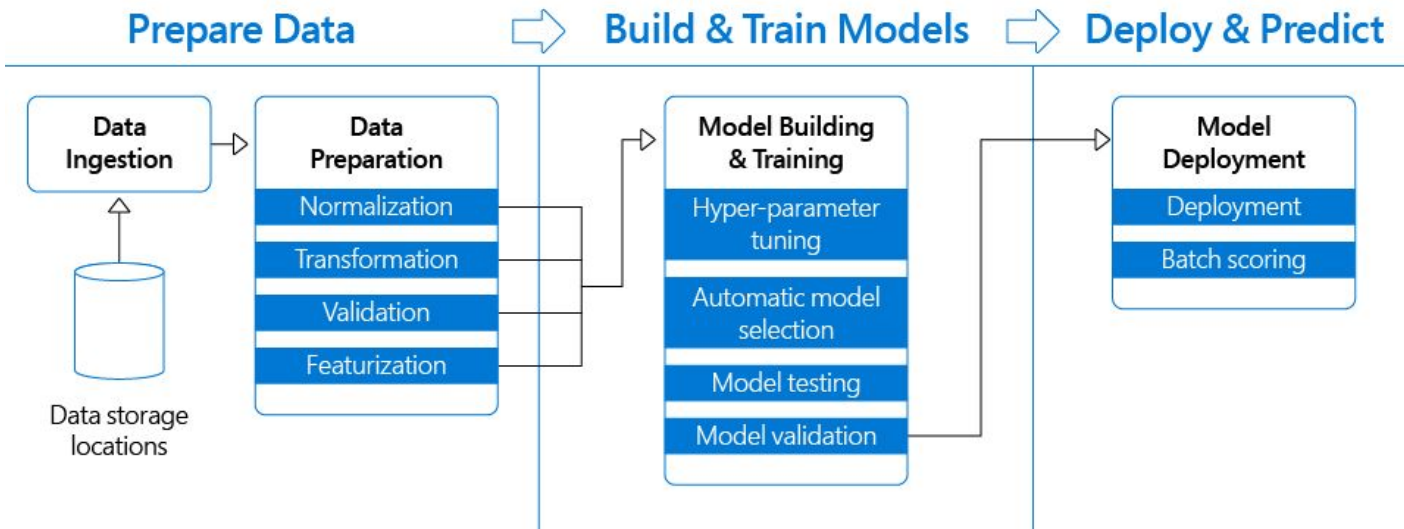
Alguns Pipelines

ETL com Airflow.



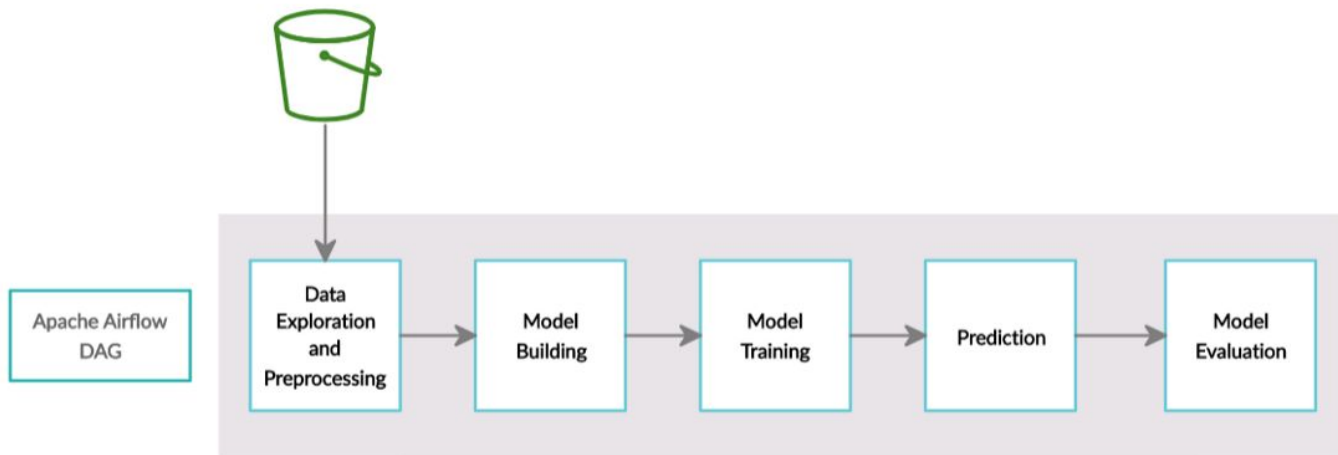
Alguns Pipelines

Machine Learning Pipeline.



Alguns Pipelines

Machine Learning Pipeline.



Representação em grafos do Pipeline de exemplo



Legenda

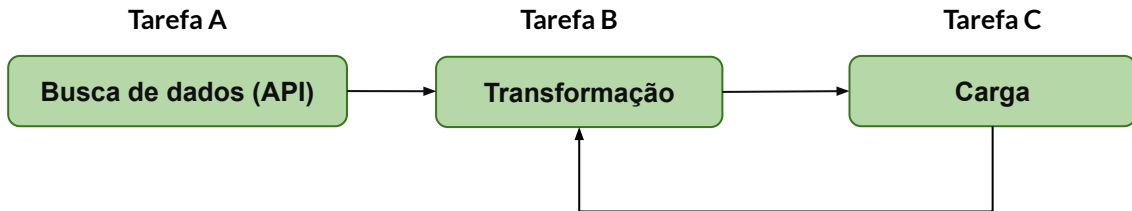


Representação em grafos do Pipeline de exemplo

Grafo direcionado acíclico.

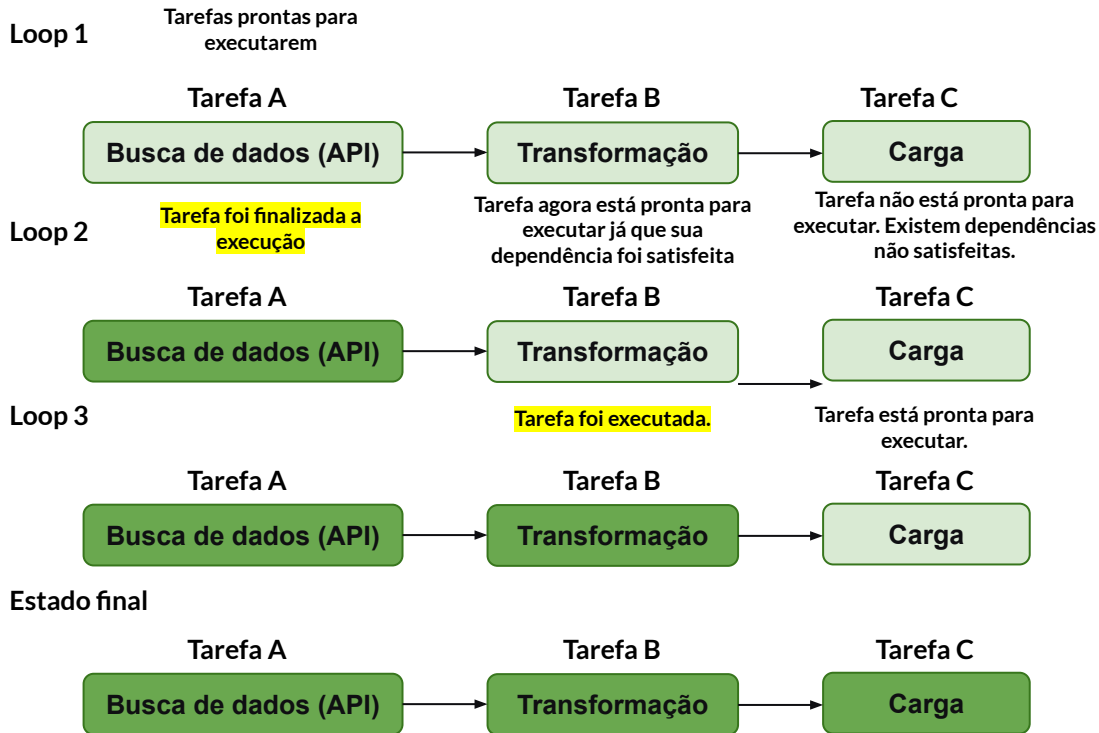


Grafo direcionado cíclico.

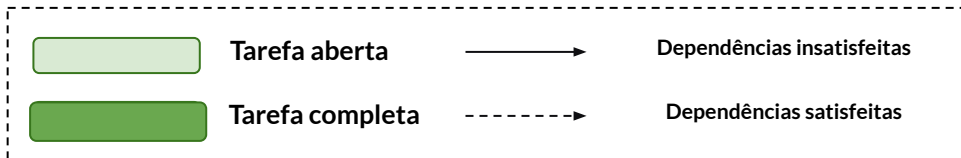


A Tarefa B não será executada pois precisa da execução da tarefa C. A tarefa C também não pode executar pois é dependente da tarefa B (Deadlock)

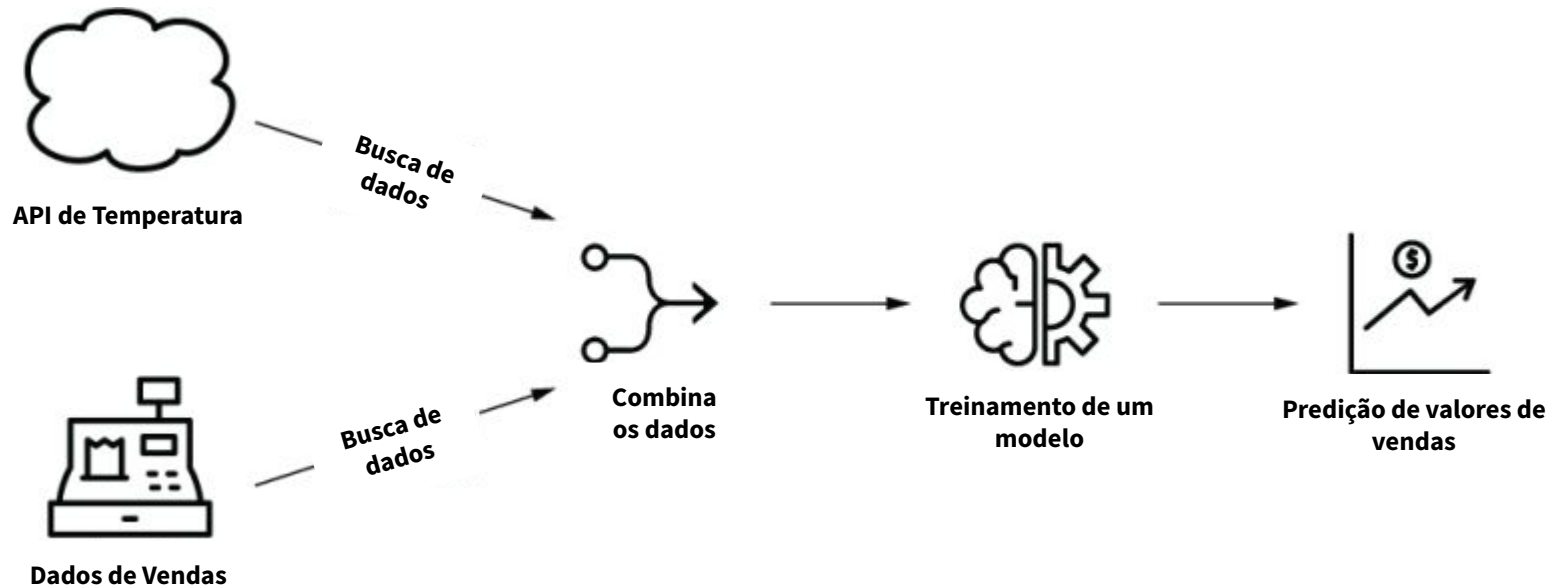
Loops e Estados das Tarefas



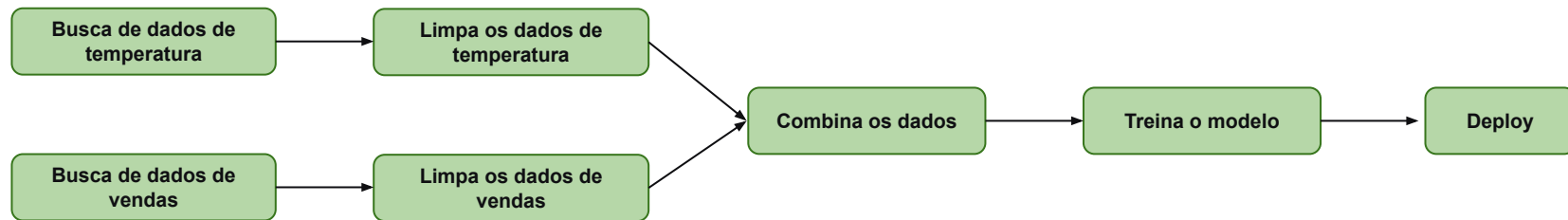
Legenda



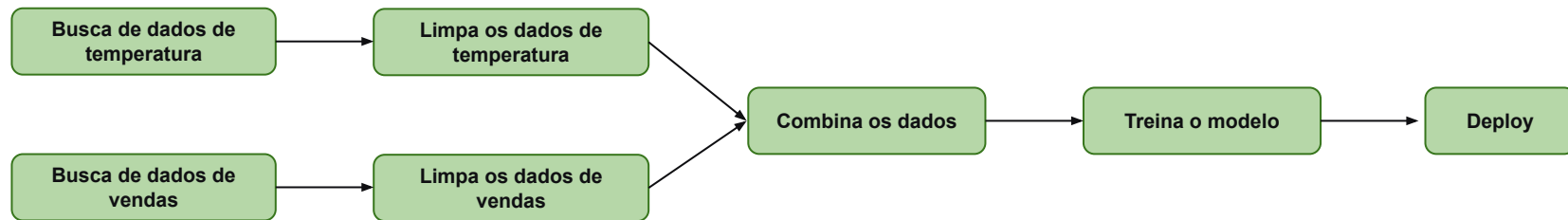
Pipelines em grafos vs scripts sequenciais



Pipelines em grafos vs scripts sequenciais



Pipelines em grafos vs scripts sequenciais





Apache
Airflow

O que é o Apache Airflow ?

- Tecnologia open source para o desenvolvimento e monitoramento de Workflows.

Porque escolher o Airflow?

- Implementação de completos e complexos Pipelines utilizando código Python.
 - Versionamento.
- Fácil integração com outros sistemas por conta da comunidade.
- Ricas opções para trabalhar com agendamento.
 - Processamento incremental.
 - Redução de custos.
- Backfilling para reprocessar dados históricos.
- Interface Web para Monitoramento e Debugging.
- Open Source e API's



Airflow vs Scripts Bash

- Operadores
 - Fácil e Rápido desenvolvimento.
 - Integração com outros sistemas.
 - Sensores.
- Gerenciamento de erro, notificações e log.
- Controle de fluxo de execução.
- Sistema de agendamento avançado.
 - Backfilling.
- Executores para processamento distribuído.
- Interface Web amigável.
- Implementação em ambientes Cloud.



Airflow vs outros Workloads Managers

Table 1.1 Overview of several well-known workflow managers and their key characteristics.

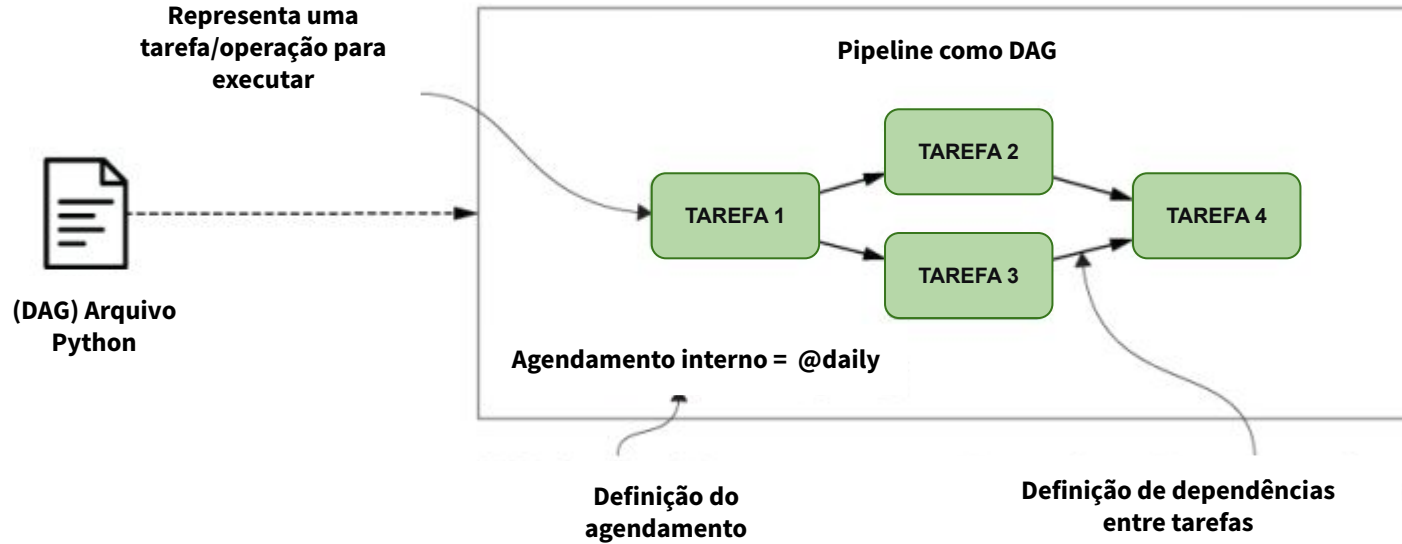
Name	Originated at ^a	Workflows defined in	Written in	Scheduling	Backfilling	User interface ^b	Installation platform	Horizontally scalable
Airflow	Airbnb	Python	Python	Yes	Yes	Yes	Anywhere	Yes
Argo	Applatix	YAML	Go	Third party ^c		Yes	Kubernetes	Yes
Azkaban	LinkedIn	YAML	Java	Yes	No	Yes	Anywhere	
Conductor	Netflix	JSON	Java	No		Yes	Anywhere	Yes
Luigi	Spotify	Python	Python	No	Yes	Yes	Anywhere	Yes
Make		Custom DSL	C	No	No	No	Anywhere	No
Metaflow	Netflix	Python	Python	No		No	Anywhere	Yes
Nifi	NSA	UI	Java	Yes	No	Yes	Anywhere	Yes
Oozie		XML	Java	Yes	Yes	Yes	Hadoop	Yes



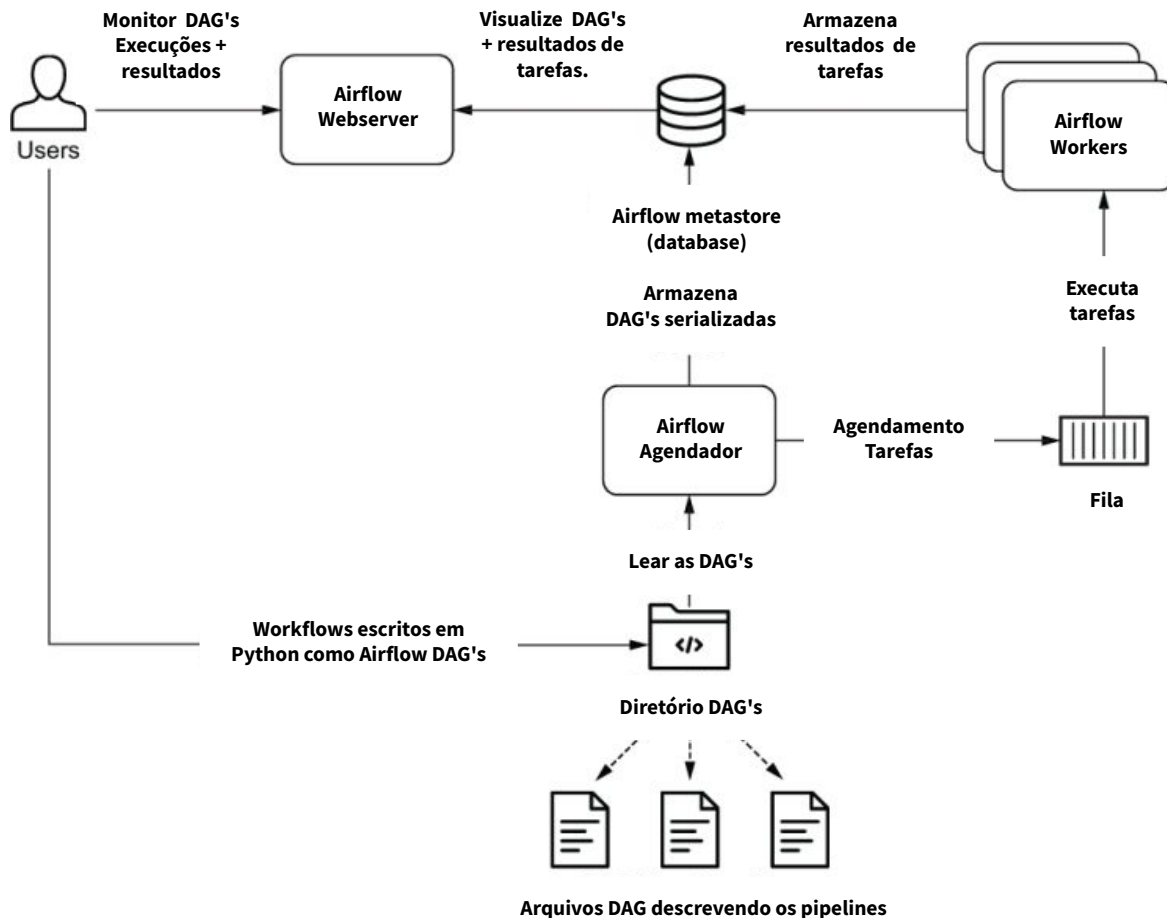
Quando NÃO utilizar Airflow?

- Necessidade de manipular pipelines em Streaming.
- Equipe com pouco ou nenhum conhecimento em Python.
 - Alternativas Azure Data Factory, NIFI, Pentaho.
- Pipelines grandes e complexos podem exigir uma grande organização.

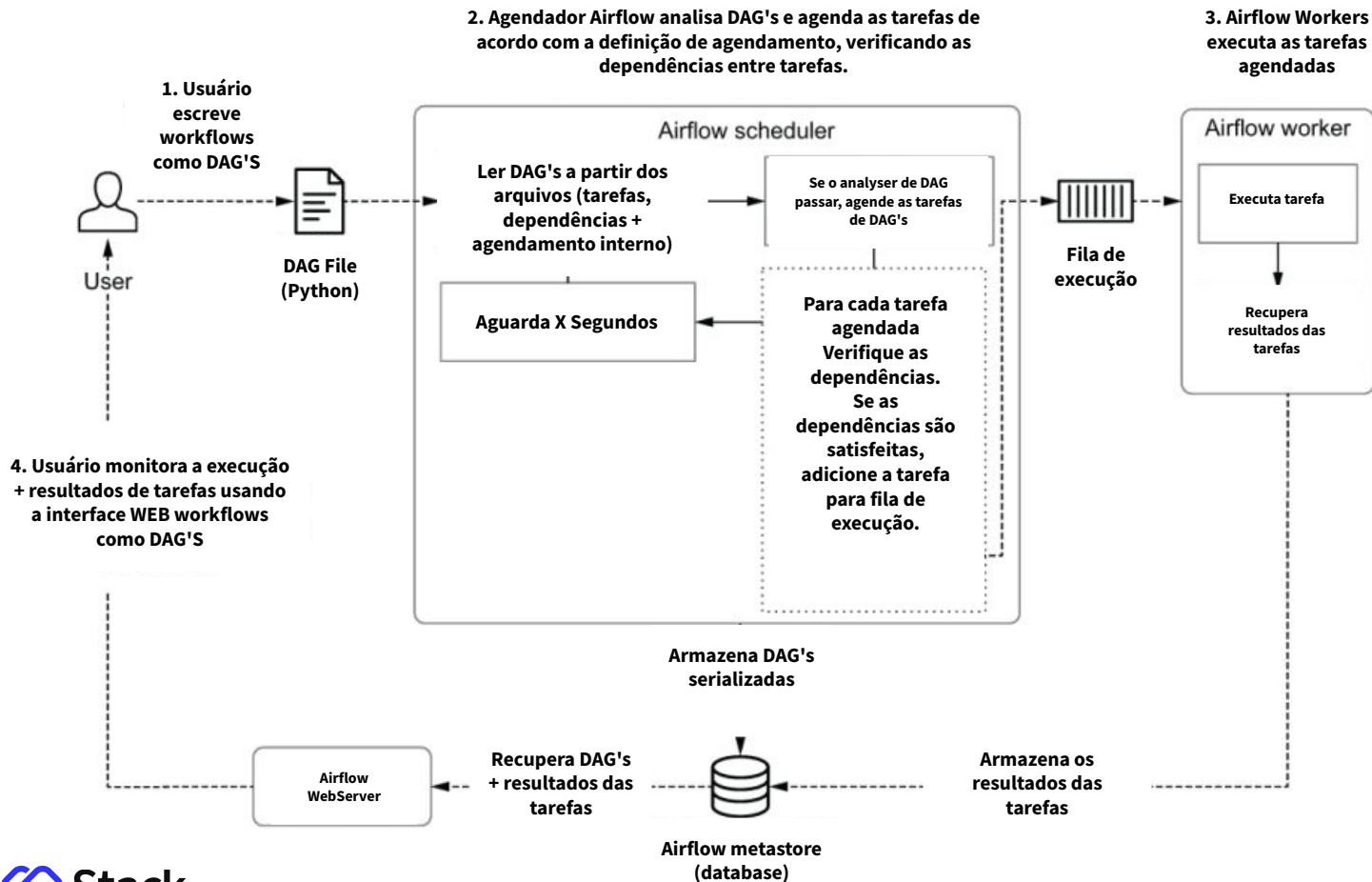
Definição de Pipelines com a Flexibilidade do Código Python



Componentes/Arquitetura do Airflow

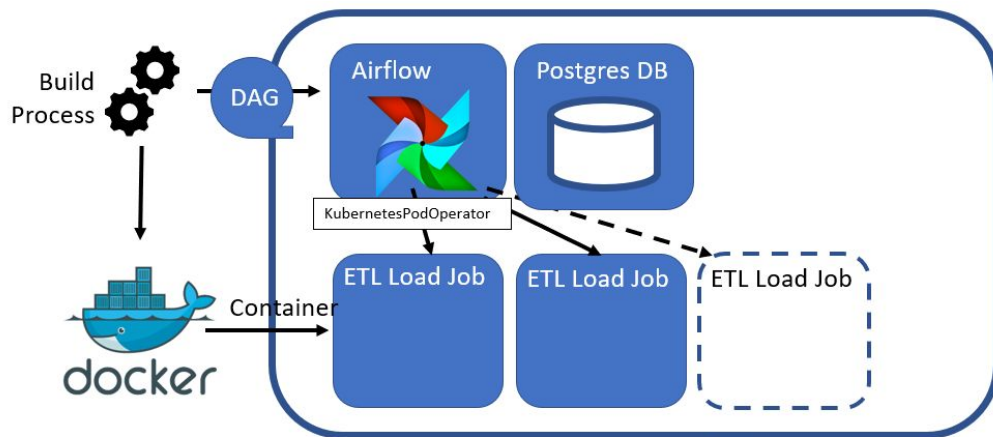


Componentes/Arquitetura do Airflow



**Vamos colocar a mão na
massa.**

Instalando o Airflow utilizando o Docker



Instalando o Airflow utilizando o Docker

```
docker run -d -p 8080:8080 -v "$PWD/dags:/opt/airflow/dags/" \
--entrypoint=/bin/bash \
--name airflow apache/airflow:2.1.1-python3.8 \
-c '(airflow db init && \
    airflow users create --username admin --password abc123 --firstname Felipe --lastname
    Lastname --role Admin --email admin@example.org
); \
airflow webserver & \
airflow scheduler \
'
```



DAG Hello World

```
1 from airflow import DAG
2 import airflow.utils.dates
3 from airflow.operators.bash import BashOperator
4
5 dag = DAG(
6     dag_id="hello_world",
7     description="A primeira DAG de teste do airflow",
8     start_date=airflow.utils.dates.days_ago(14),
9     schedule_interval="@daily",
10 )
11
12 task_echo_message = BashOperator(
13     task_id="echo_message",
14     bash_command="echo Hello World!",
15     dag=dag,
16 )
17
18
19 task_echo_message
```

Airflow Command Line Interface

Conectando ao container.

docker container exec -it airflow bash

Listando os comandos disponíveis.

airflow --help

Verificando as informações do ambiente.

airflow info

Verificando as configurações.

airflow config list

Airflow Command Line Interface

Verificando as conexões do ambiente.

airflow connections list

Listando as dags disponíveis.

airflow dags list

Listando as tasks de uma determinada dag.

airflow tasks list <dag-id>

airflow tasks list hello_world

Testando a execução de uma determinada task.

airflow <tasks> <test> <dag_id> <task_id> <execution-date>

airflow tasks test hello_world echo_message 2021-07-01



Conceitos

Workloads

Uma DAG é composta por tarefas que podem ser de três tipos.

- **Operadores** que são tarefas predefinidas que agrupadas formam a execução da DAG.
- **Sensores** que é uma subclasse especial de Operadores que trabalham aguardando um evento externo acontecer.
- **Taskflow** que é uma função Python personalizada empacotada como uma Tarefa.

Conceitos

DAGs - Core do Airflow.

- Reuni tarefas juntas.
- Especifica suas dependências e relacionamentos.
- Define como e quando serão executadas.

Declarando uma DAG.

Através de um Context manager.

- **with DAG("etl-db-producao") as dag:**
 op = DummyOperator(task_id="task")

Através de um construtor padrão.

- **my_dag = DAG("etl-db-producao")**
 op = DummyOperator(task_id="task", dag=my_dag)

Conceitos

Através de um decorator - transforma uma função em uma dag.

- `@dag(start_date=days_ago(2))`
 `def generate_dag():`
 `op = DummyOperator(task_id="task")`
 `dag = generate_dag()`

Conceitos

Executando DAG's

DAG's são executadas de duas formas.

- Acionamento manual ou por API.
- Intervalo de agendamento.
 - `with DAG("etl-db-producao", schedule_interval="@daily"):`
 - `with DAG("etl-db-producao", schedule_interval="0 * * * *"):`

Argumentos padrão

```
default_args = {  
    'start_date': datetime(2016, 1, 1),  
    'owner': 'airflow'  
}
```

```
with DAG('etl-db-producao', default_args = default_args) as dag:  
    op = DummyOperator(task_id='dummy')  
    print(op.owner)
```

Conceitos

Control Flow

Tarefas têm dependências entre elas

- `extracao >> [transformacao, carga]`
- `carga << notificacao`

Podemos definir as dependências através dos métodos.

- `extracao.set_downstream([transformacao, carga])`
- `carga.set_upstream(notificacao)`

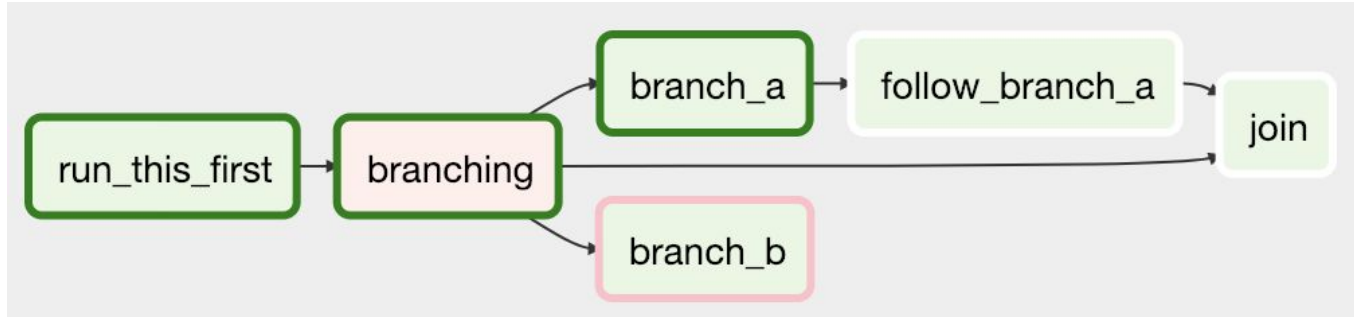
Control Flow

Formas de controlar a execução de tarefas.

- **Branching**
 - Determina qual tarefa mover a partir de uma condição.
- **Latest Only**
 - Só é executada em DAGs em execução no presente.
- **Depends on Past**
 - Tarefas podem depender de si mesmas de uma execução anterior.
- **Trigger Rules**
 - Permite definir as condições para um DAG executar uma tarefa.

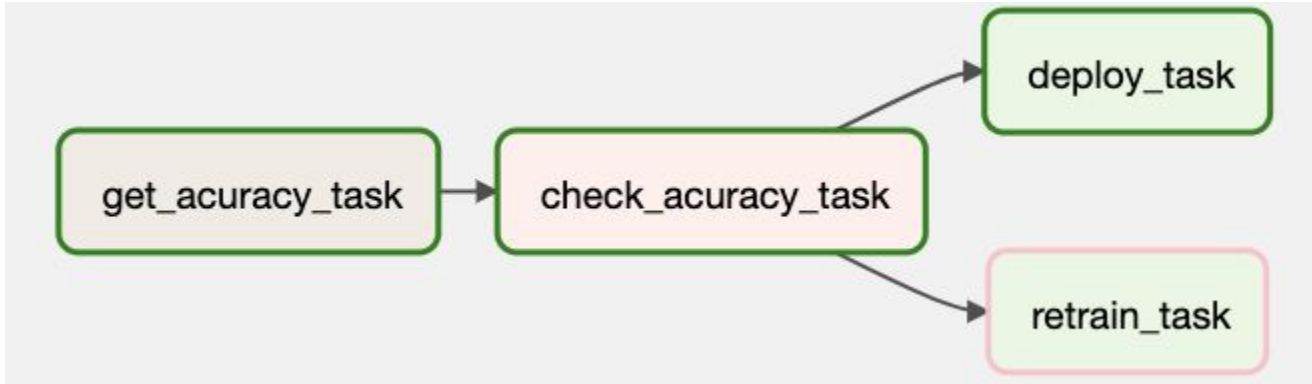
Control Flow

Branching



Control Flow

Branching



Control Flow

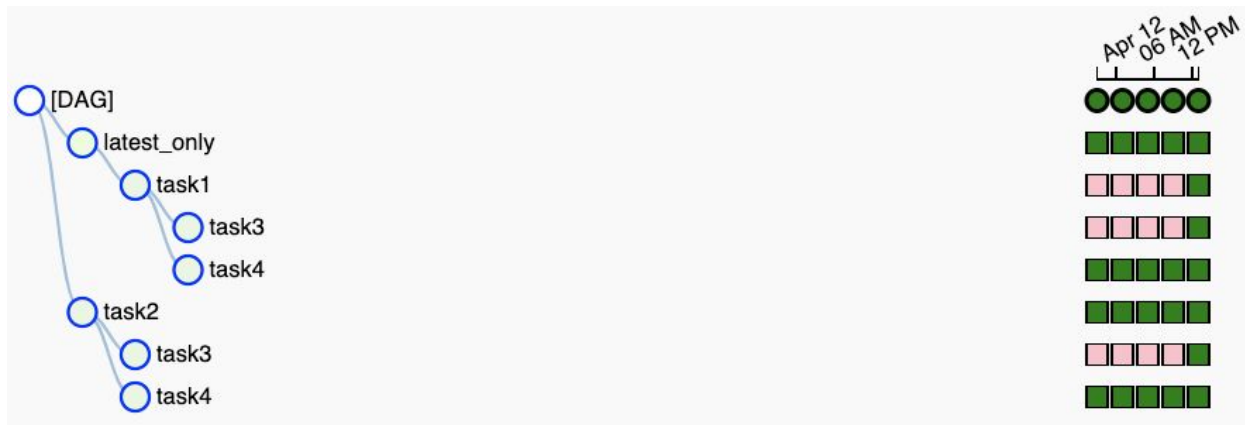
DAG - Branching

```
1 from airflow import DAG
2 import airflow.utils.dates
3 from airflow.operators.bash import BashOperator
4 from airflow.operators.python import BranchPythonOperator
5 from airflow.operators.dummy import DummyOperator
6
7 dag = DAG(
8     dag_id="branching-teste",
9     description="DAG de exemplo para o conceito de Branching",
10    start_date=airflow.utils.dates.days_ago(14),
11    schedule_interval=None,
12 )
13
14 def check_acuracy(ti):
15     acuracy_value = int(ti.xcom_pull(key="model_acuracy"))
16     if acuracy_value >= 90:
17         return 'deploy_task'
18     else:
19         return 'retrain_task'
20
21 get_acuracy_op = BashOperator(
22     task_id='get_acuracy_task',
23     bash_command='echo "{ ti.xcom_push(key="model_acuracy", value=90) }"',
24     dag=dag,
25 )
26
27 check_acuracy_op = BranchPythonOperator(
28     task_id='check_acuracy_task',
29     python_callable=check_acuracy,
30     dag=dag,
31 )
32
33 deploy_op = DummyOperator(task_id='deploy_task', dag=dag)
34 retrain_op = DummyOperator(task_id='retrain_task', dag=dag)
35
36 get_acuracy_op >> check_acuracy_op >> [deploy_op, retrain_op]
```

Control Flow

Latest Only

Este Operador especial ignora todas as tarefas posteriores de si mesmo se você não estiver na "última" execução do DAG.



Control Flow

Depends On Past

Uma tarefa só pode ser executada se a execução anterior da tarefa no DAG Run anterior foi bem-sucedida.

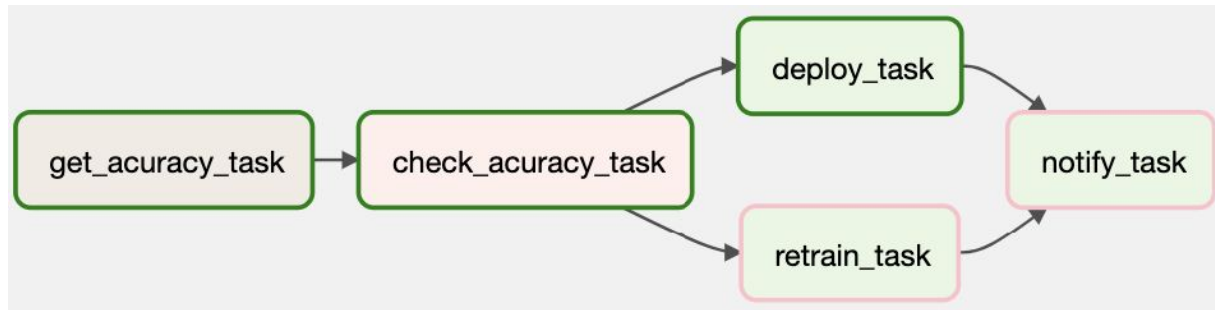
Trigger Rules

Nos permite controlar o comportamento de acionamento das tarefas.

- **all_success (padrão):** Todas as tarefas anteriores foram bem-sucedidas.
- **all_failed:** Todas as tarefas anteriores estão em um estado de falha ou **upstream_failed**.
- **all_done:** todas as tarefas upstream são feitas com sua execução.
- **one_failed:** pelo menos uma tarefa anterior falhou (não espera que todas as tarefas anteriores já tenham executado).
- **one_success:** pelo menos uma tarefa anterior foi bem-sucedida (não espera que todas as tarefas anteriores tenham executado).
- **none_failed:** Todas as tarefas anteriores não falharam ou **upstream_failed** - ou seja, todas as tarefas anteriores foram bem-sucedidas ou foram ignoradas.
- **none_failed_or_skipped:** Todas as tarefas anteriores não falharam ou **upstream_failed**, e pelo menos uma tarefa anterior foi bem-sucedida.
- **none_skipped:** nenhuma tarefa anterior está em um estado ignorado - ou seja, todas as tarefas anteriores estão em um estado de sucesso, falha ou **upstream_failed**.
- **dummy:** sem dependências, execute esta tarefa a qualquer momento.

Control Flow

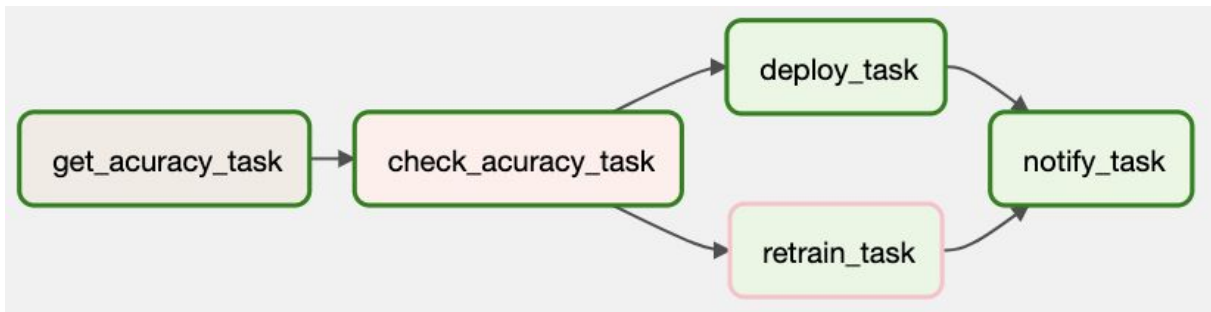
Exemplo da aplicação de trigger rules.



Control Flow

Exemplo da aplicação de trigger rules.

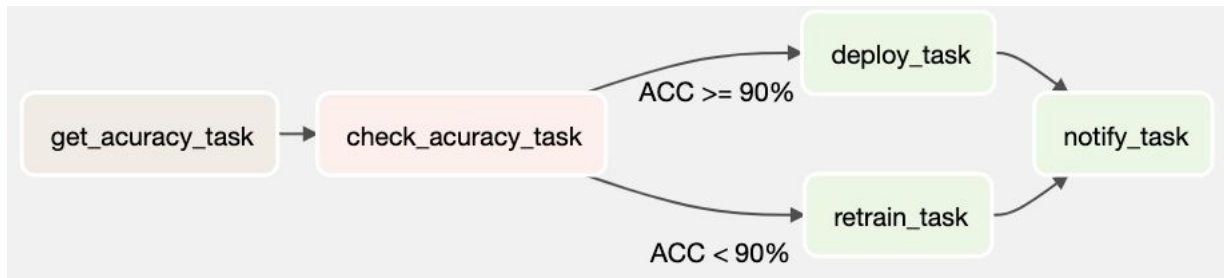
Alterando a sua trigger rule para **NONE_FAILED**.



Edges Labels

Podemos definir labels para documentar as relações e dependências entre tarefas.

- `get_acuracy_op >> check_acuracy_op >> Label("Limit 90%") >> [deploy_op, retrain_op] >> notify_op`
- `get_acuracy_op.set_downstream(check_acuracy, Label("Métrica ACC"))`



Tasks

A unidade de execução mais básica do Airflow.

Podem ser de três tipos.

- **Operadores , Sensores e Taskflow.**

O relacionamento é definido através das dependências (upstream e downstreams).

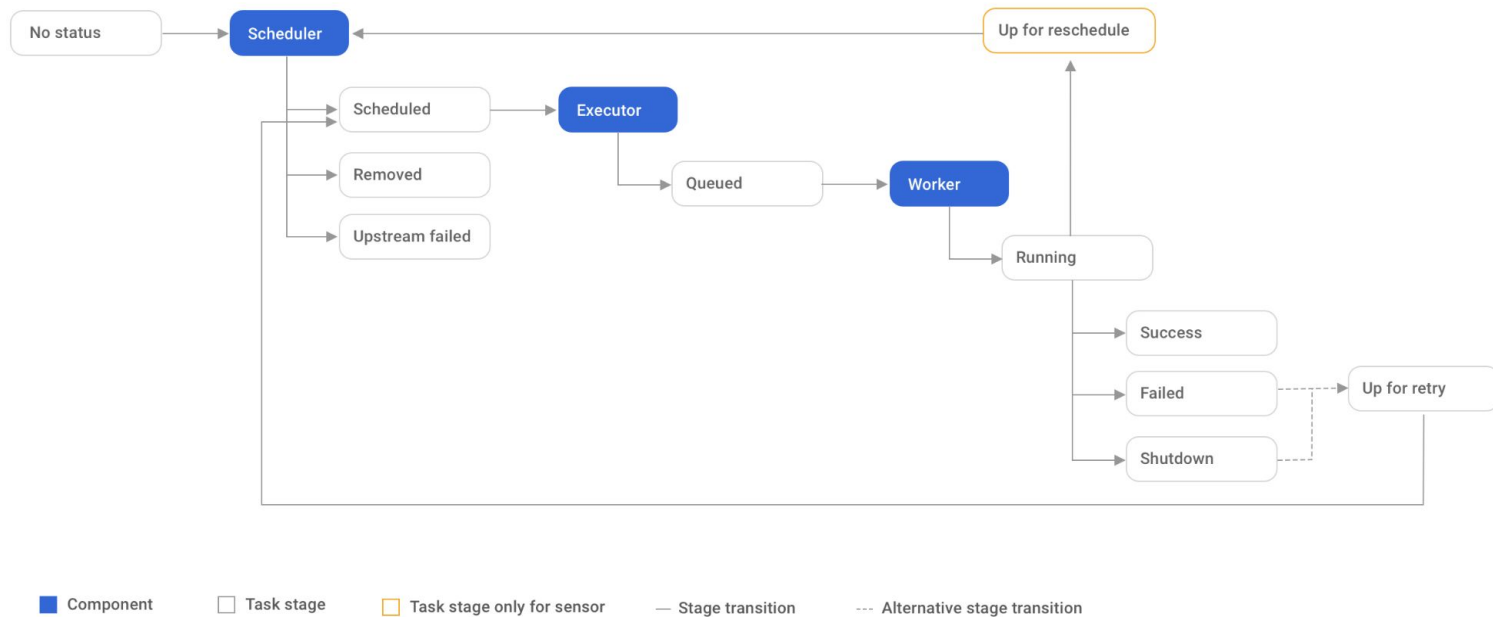
Estados de cada tarefa.

Tasks

Estados de cada tarefa.

- **none:** a tarefa ainda não foi enfileirada para execução (suas dependências ainda não foram atendidas).
- **scheduled:** O agendador determinou que as dependências da Tarefa são atendidas e deve ser executado.
- **queued:** A tarefa foi atribuída a um executor e está aguardando um trabalhador.
- **running:** a tarefa está sendo executada em um trabalhador (ou em um executor local / síncrono).
- **success:** a tarefa terminou em execução sem erro.
- **failed:** a tarefa teve um erro durante a execução e falhou ao executar.
- **skipped:** a tarefa foi ignorada devido a ramificação, LatestOnly ou semelhante.
- **upstream_failed:** uma tarefa upstream falhou e a regra de acionamento diz que precisávamos dela.
- **up_for_retry:** A tarefa falhou, mas ainda restam novas tentativas e será reprogramada.
- **up_for_reschedule:** A tarefa é um Sensor que está em modo de reprogramação.
- **sensing:** a tarefa é um sensor inteligente
- **removed:** a tarefa desapareceu do DAG desde o início da execução.

Tasks



Conceitos

Transferência de dados entre tarefas.

- **X-coms (Cross-communications) - push e pull metadados.**
- **Uploading e Downloading em serviços de armazenamento (Bancos de dados, serviços, arquivos)**

Operators

Um operador é um template para uma tarefa pré-definida.

Exemplos:

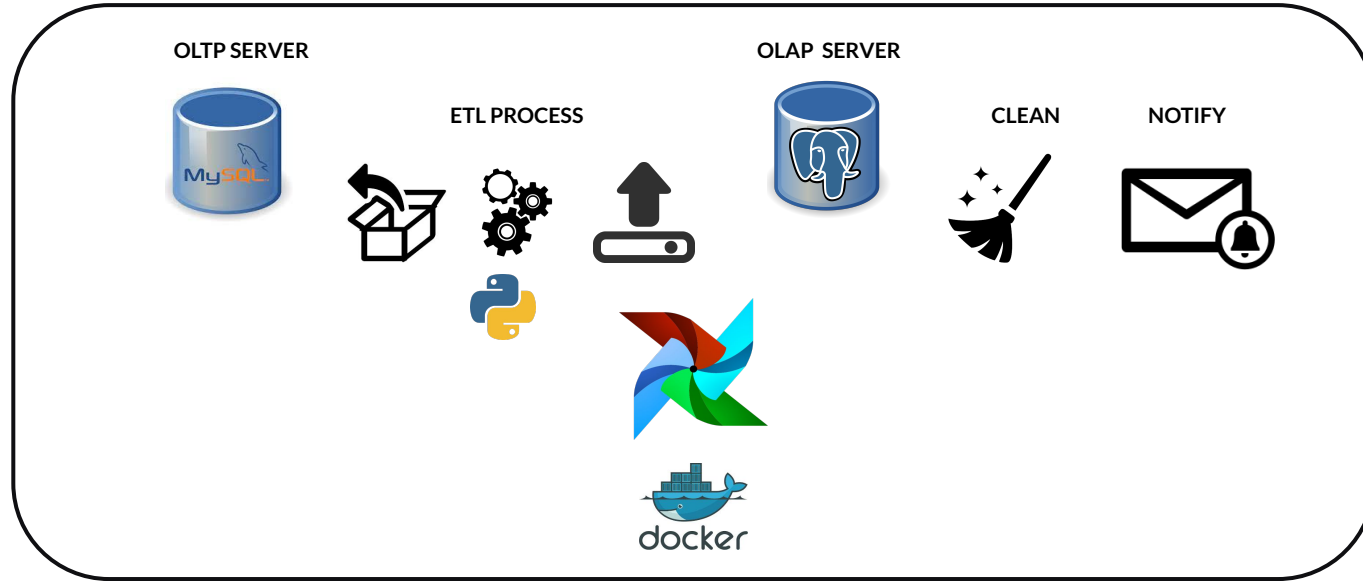
- BashOperator - executa um comando bash.
- PythonOperator - chama uma função Python.
- EmailOperator - envia um e-mail.
- SimpleHttpOperator - envia uma requisição http.
- SqliteOperator - operador para trabalhar com sqlite.

Outros operadores disponíveis...

- MySQLOperator
- PostgresOperator
- MsSqlOperator
- OracleOperator
- JdbcOperator
- DockerOperator

**Vamos colocar a
mão na massa.**

Data Pipeline 01



Data Pipeline 01

Requisitos do Pipeline.

- Necessidade de separar os dados do ambiente OLTP para o ambiente OLAP.
- Todas as tarefas devem ser executadas diariamente.
- Todo o processo deve ser notificado.
- O ambiente de stage deve ser limpo.
- Todos os logs devem ser mantidos.



Montando o Ambiente OLTP

MySQL (OLTP)

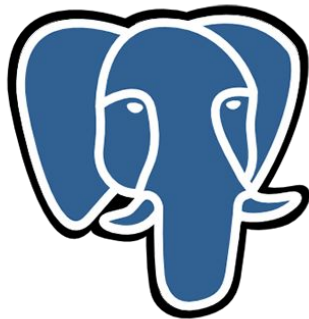
- Criar o container.
- Mapear volumes.
- Instalar e configurar os softwares (apt-get, vim)
- Carregar os dados.
- Realizar Queries.



Montando o Ambiente OLAP

Postgresql (OLAP)

- Criar o container.
- Mapear volumes.
- Instalar e configurar os softwares (apt-get, vim)
- Criar o banco de dados.



Configurando o Airflow para enviar E-mails

Servidor SMTP

- Criar uma conta no Mailtrap.
- Configurar o arquivo de configurações do Airflow.



Agendamento no Airflow

Agendamento em intervalos

- `scheduler_interval= None` - significa que a dag não terá um agendamento.
- Podemos definir execuções uma vez em uma hora, em um dia, um mês com o parâmetro `scheduler_interval`.

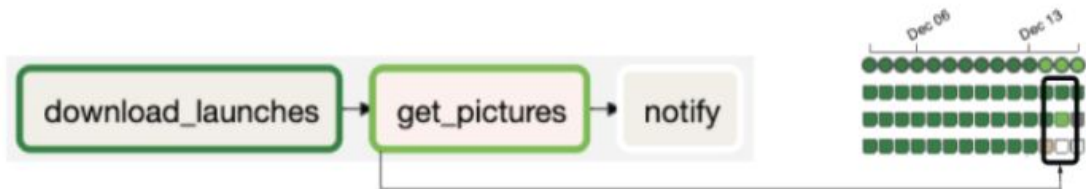


Agendamento no Airflow

Agendamento em intervalos

- `scheduler_interval= None` - significa que a dag não terá um agendamento.
- Podemos definir execuções uma vez em uma hora, em um dia, um mês com o parâmetro `scheduler_interval`.

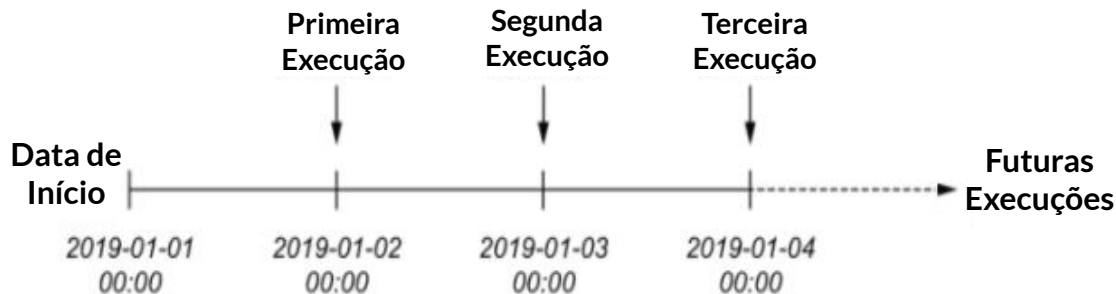
Figure 2.14 Relationship between graph view and tree view



Agendamento no Airflow

Esquema de execução no agendamento diário do Airflow.

- **"@daily"** - Especifica a execução uma vez no dia.
- **end_date** - Especifica a data fim para a execução das dags.



Agendamento no Airflow

Exemplo.

```
dag = DAG (  
    dag_id="teste_date",  
    schedule_interval="@daily",  
    start_date=dt.datetime(year=2019, month=1, day=1),  
    end_date=dt.datetime(year=2019, month=1, day=5),  
)
```

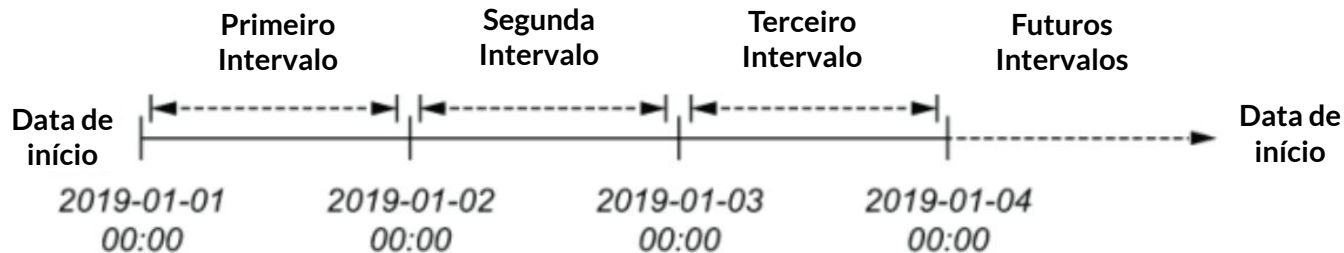


Dag Schedule Daily

```
1 from datetime import datetime
2
3 from airflow import DAG
4 from airflow.operators.bash import BashOperator
5
6 dag = DAG(
7     dag_id="scheduler-daily",
8     schedule_interval="@daily",
9     start_date=datetime(2021, 7, 25),
10 )
11
12 print_date_task = BashOperator(
13     task_id="print_date",
14     bash_command=(
15         "echo Iniciando a tarefa: &&"
16         "date"
17     ),
18     dag=dag,
19 )
20
21 processing_task = BashOperator(
22     task_id="processing_data",
23     bash_command=(
24         "echo Processando os dados... &&"
25         "sleep 50"
26     ),
27     dag=dag,
28 )
29
30 print_date_task >> processing_task
```

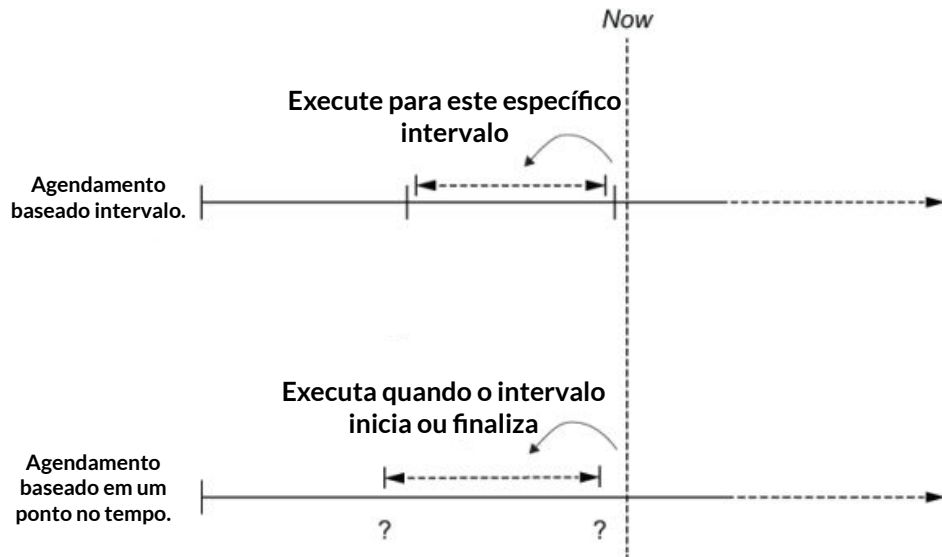
Execution Date

Entendo a execução baseada em intervalos.



Execution Date

Entendo a execução baseada em intervalos.



Cron-based Scheduler

Agendamento em intervalos de tempo mais complexos.

- Que tal todos os dias às 23:30 ? Todos os sábados às 23:45?

```
1 # minute (0 - 59)
2 # hour (0 - 23)
3 # day of the month (1 - 31)
4 # month (1 - 12)
5 # day of the week (0 - 6) (Sunday to Saturday;
6 # 7 is also Sunday on some systems)
7 # * * * * *
```

Cron-based Scheduler

Alguns exemplos de expressões cron.

- **0 * * * *** = hourly (executando em uma hora)
- **00 * * * *** = daily (executando todos os dias à meia noite)
- **00 * * 0** = weekly (executando a meia noite do domingo)

Um pouco mais complicado.

- **001 * * *** = meia noite no primeiro dia de cada mês.
- **45 23 * * SAT** = às 23:45 de cada sábado.
- **00 * * MON-FRI** = executa de segunda à sexta à meia noite.
- **00,12 * * *** = executa todos os dias a meia noite e ao meio dia.

Macros para Agendamento

Algumas macros fornecidas pelo Airflow.

Macro	Significado.
@once	Executa uma vez e apenas uma.
@hourly	Executa uma vez a cada hora e no início da hora.
@daily	Executa uma vez no dia à meia noite.
@weekly	Executa uma vez por semana à meia noite.
@monthly	Executa uma vez por mês à meia noite no primeiro dia do mês.
@yearly	Executa uma vez por ano à meia noite iniciando em primeiro de Janeiro.

Cron-based Scheduler

Crontab Guru

<https://crontab.guru/>

Frequency-based Scheduler

Através do agendamento baseado no cron é difícil definir intervalos baseados em frequência.

Como agendar a execução da Dag uma vez a cada três dias?

Podemos utilizar o objeto `timedelta`.

```
dag = DAG(  
    dag_id="teste",  
    schedule_interval=dt.timedelta(days=3),  
    start_date=dt.datetime(year=2019, month=1, day=1),  
    end_date=dt.datetime(year=2019, month=1, day=5),  
)
```

Para executar a cada 10 minutos

`timedelta(minutes=10)`

A cada duas horas.

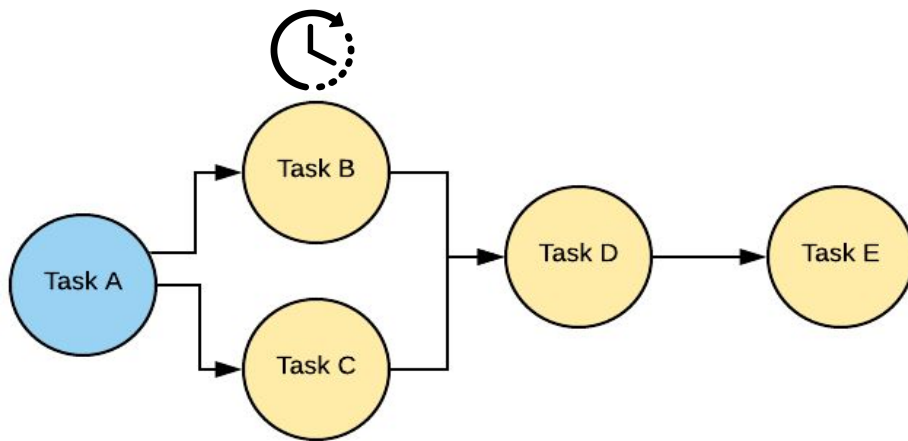
`timedelta(hours=2)`

Frequency-based Scheduler

Exemplo

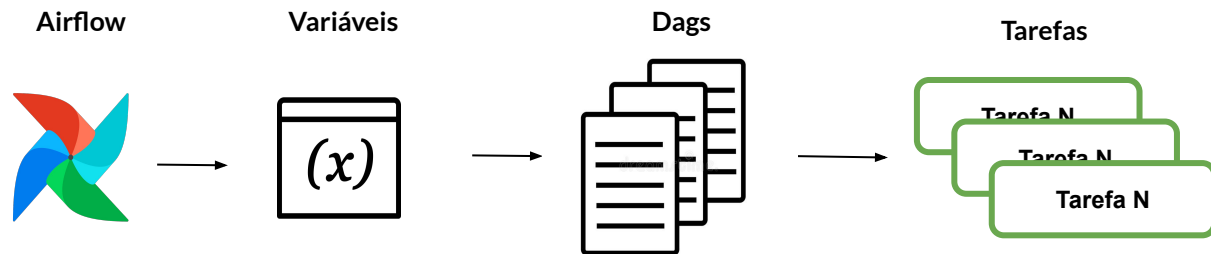
Re-execução de Tarefas

- Existem fatores externos que influenciam na performance da tarefas.
- Continuar o Pipeline de onde parou.



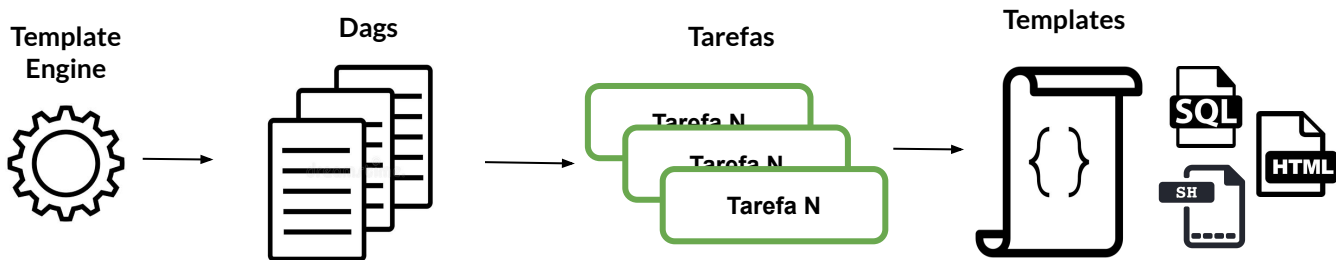
Variáveis, Templates e Macros

- Atribuição de valores e parâmetros comuns para várias dags e tarefas.
- Atualização e manutenção centralizada.
- Dados independente do código.
- Melhor legibilidade e organização.



Variáveis, Templates e Macros

- Atribuição de valores e parâmetros comuns para várias dags e tarefas.
- Atualização e manutenção centralizada.
- Dados independente do código.
- Melhor legibilidade e organização.

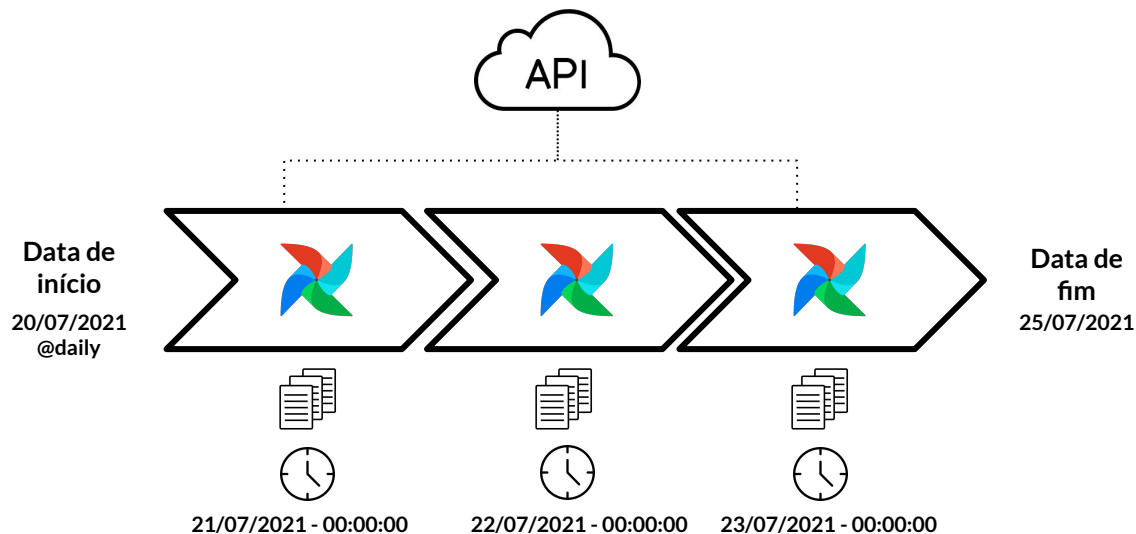


Variáveis, Templates e Macros

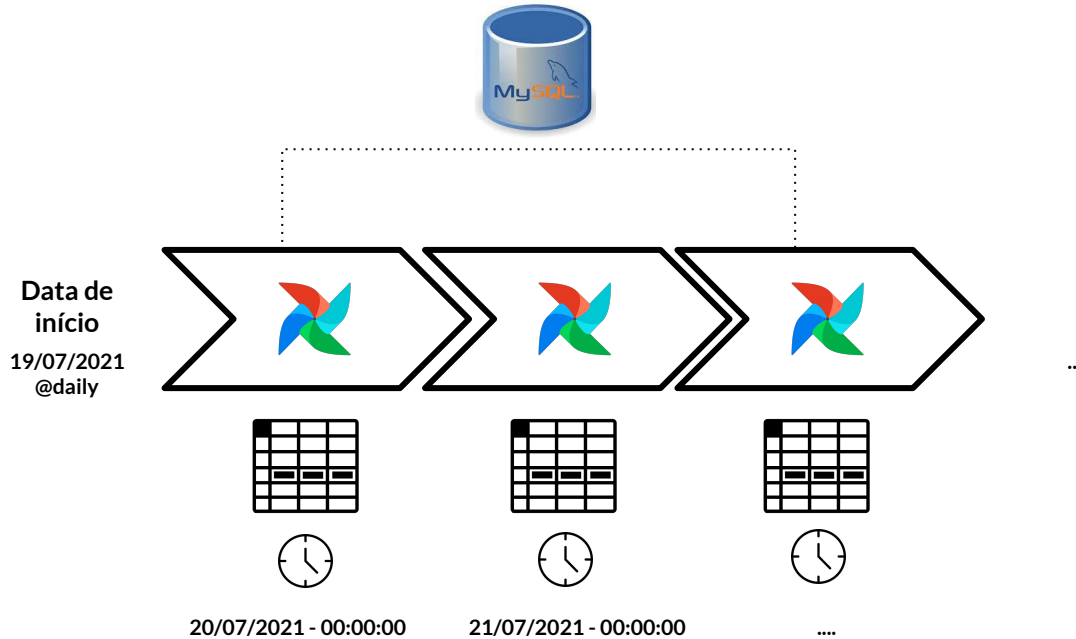
- Macros padrão definidas pelo Airflow e que são utilizáveis em todo template.
- Documentação:
<https://airflow.apache.org/docs/apache-airflow/stable/macros-ref.html>

<code>{{ ds }}</code>	A data de execução da dag no formato YYYY-MM-DD.
<code>{{ ti }}</code>	equivale a macro <code>{{ task_instance }}</code>
<code>{{ params }}</code>	Uma referência para parâmetros definidos pelo usuário.
<code>{{ var.value.my_var }}</code>	Acesso a variáveis globais definidas e representadas como um dicionário.

Exemplo Coleta de Dados de API do Mercado Financeiro



Coletando e Carregando Dados Incrementais

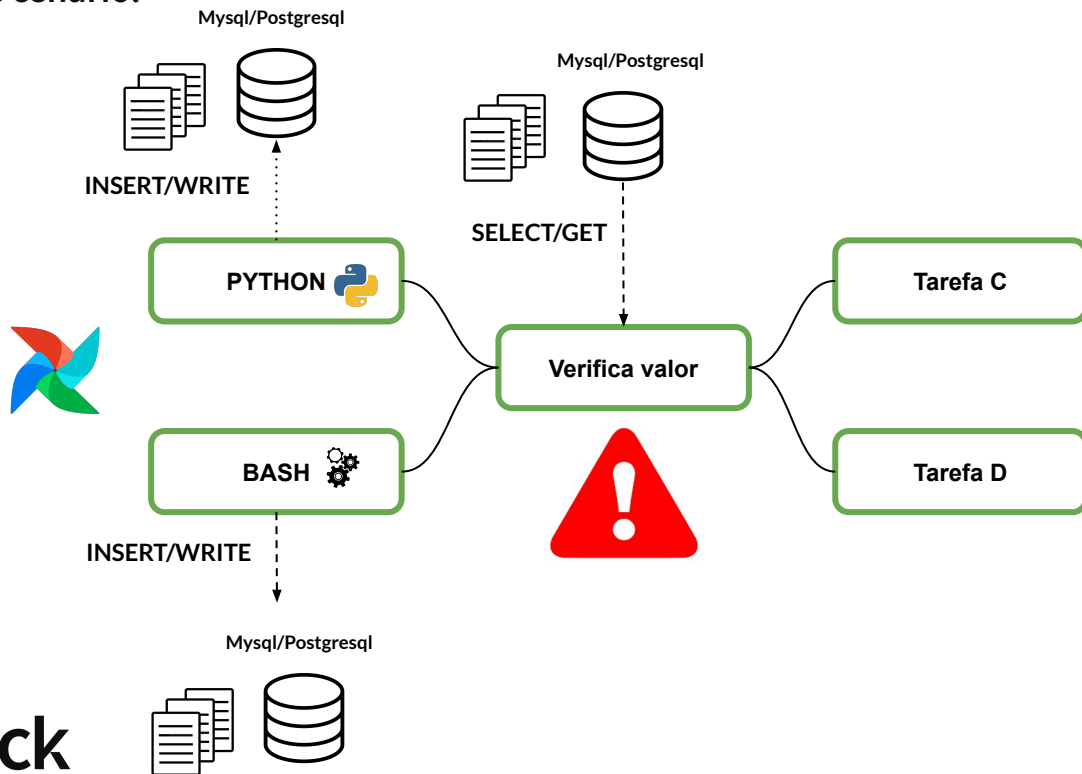


Comunicação entre Tarefas

- Como permitir a comunicação entre tarefas?
- Como compartilhar dados entre cada tarefa?
- Utilizar arquivos ou bancos de dados como armazenamento temporário?
 - Cada tarefa/operador implementar sua persistência?
 - Qual será o formato de armazenamento e recuperação dos dados?

Comunicação entre Tarefas

- Análise de cenário.



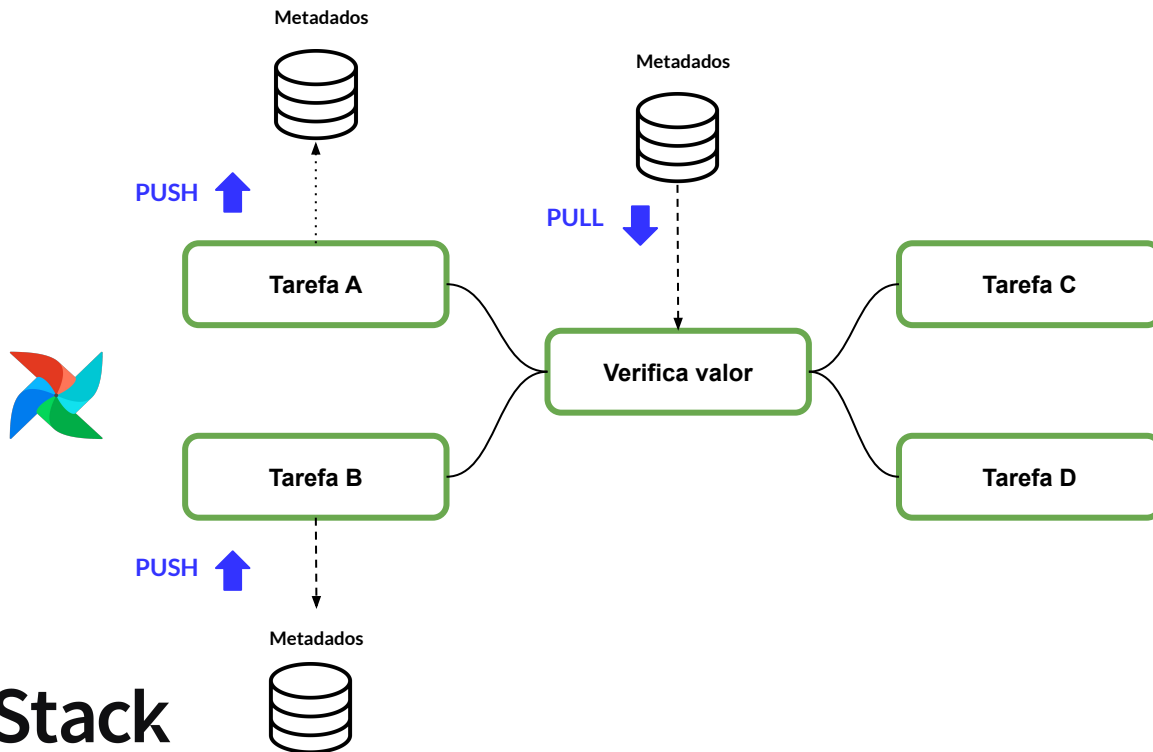
XComs

Cross Communications

- Permite fácil comunicação entre as tarefas.
- Métodos de armazenamento e recuperação dos dados padronizados.
- Formato nativo do Airflow.
- Armazenamento de metadados e logs importantes.

XComs

Cenário de comunicação entre tarefas utilizando XComs.



XComs

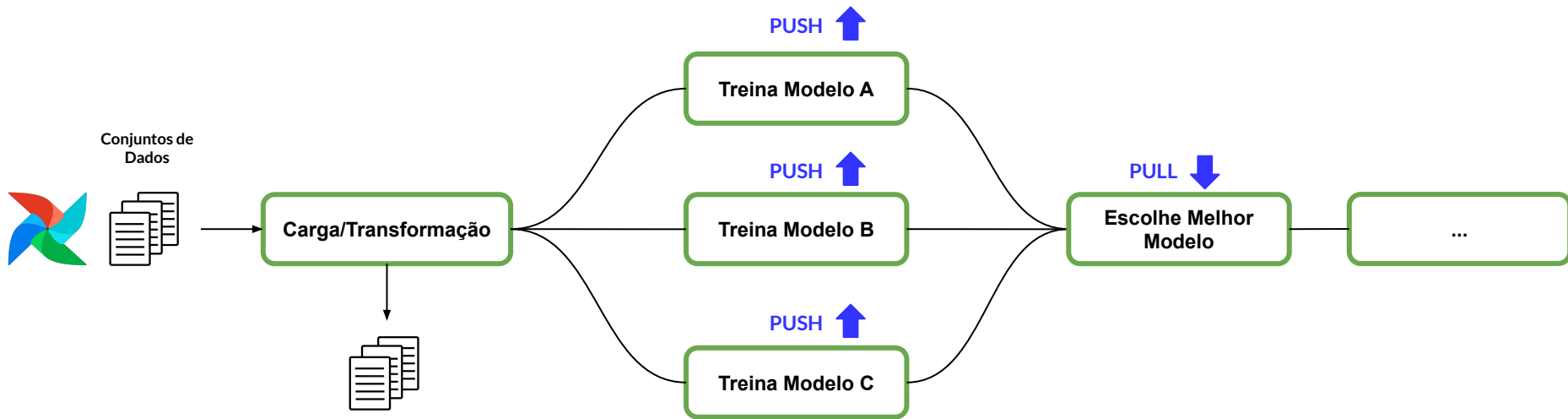
Atributo	Descrição
Key	O identificador do registro. É utilizado para recuperar o valor associado.
Value	O valor que será compartilhado e que será associado a key.
Timestamp	A data e horário que o xcom foi criada.
Execution Date	A data de execução que a DAG gerou a xcom.
Task ID	O ID da tarefa onde a xcom foi criada.
Dag ID	O ID da dag onde a xcom foi criada.

Cross Communications

- Atributos no modelo de dados.
- **Limitação**
 - SQLite: 2 GB
 - Postgres: 1 GB
 - MySQL: 64 KB

XComs Prática

Cenário de comunicação entre tarefas utilizando XCom através de um Pipeline para carregar e transformar os dados para testar três diferentes modelos de Machine Learning.

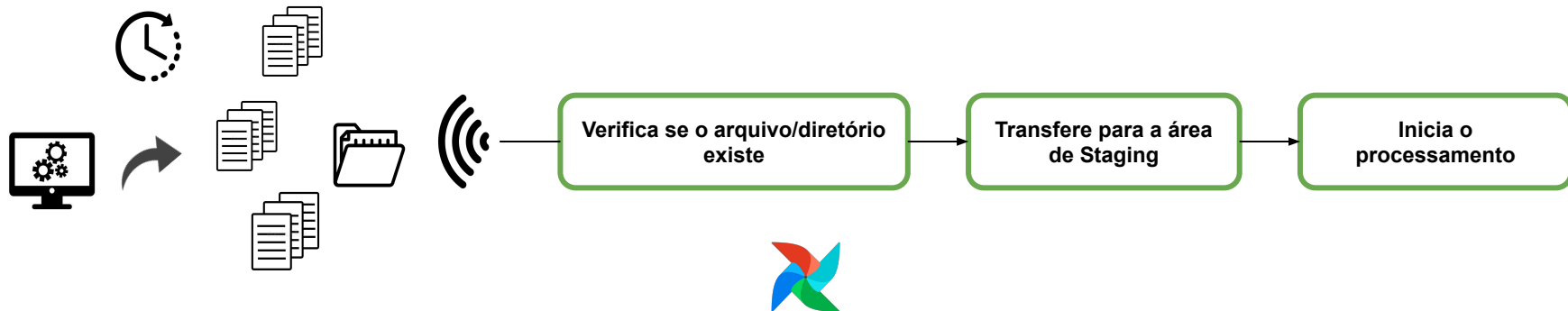


Sensor Operators

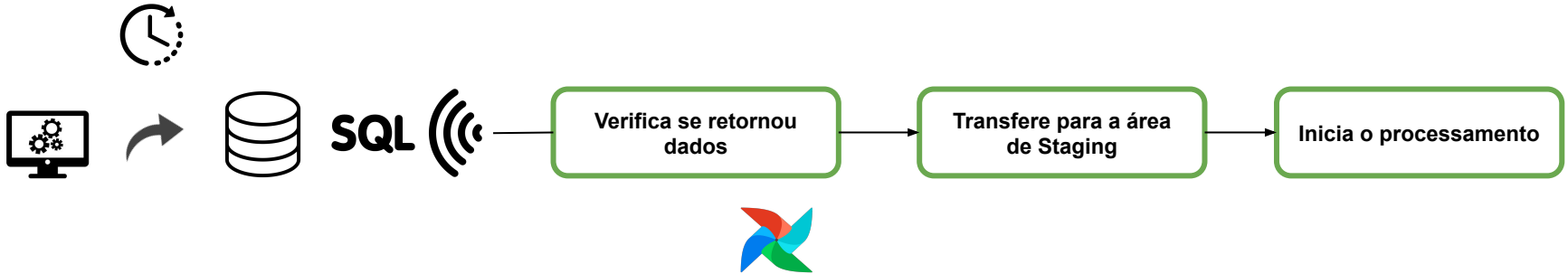
Tipo especial de operadores que nos permitem esperar para um evento acontecer ou uma condição ser satisfeita.

- Aguardar um arquivo estar disponível em um diretório/volume.
- Aguardar dados serem retornados a partir de uma consulta SQL.
- Receber um status de uma requisição http.
- Aguardar uma partição no hive está disponível.
- Aguardar por um arquivo está presente em um bucket S3.

Sensor Operators



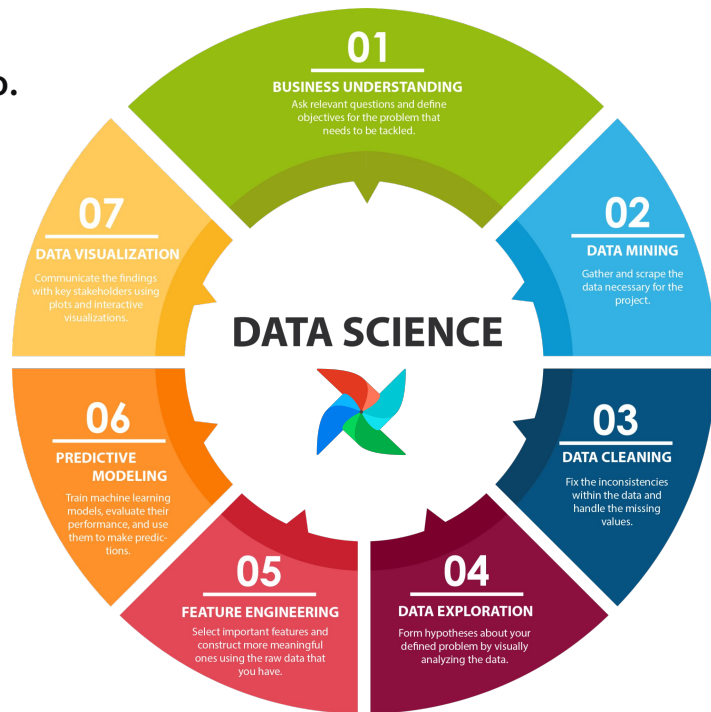
Sensor Operators



Data Science Project Automation

Automação das Etapas do Projeto.

- Carga dos dados.
- Limpeza.
- Análise e Exploração.
- Pré-processamento e transformação.
- Feature Engineering.
- Modelagem.
- Ajustes e Avaliação.
- Deploy.





Ferramenta para **parametrização e execução**
de Jupyter Notebooks.

<https://github.com/nteract/papermill>



Algumas alternativas



Vantagens

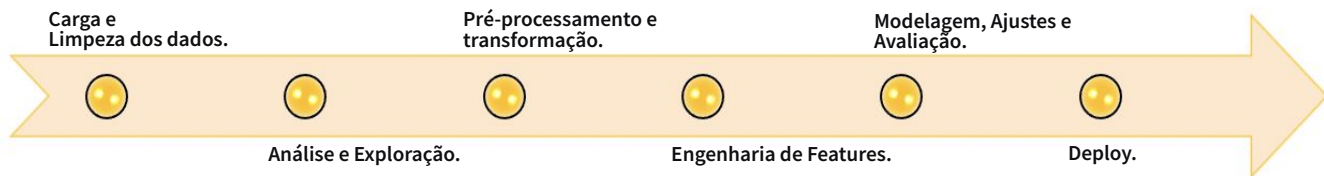
- Alternativa mais rápida e simples.
- Ideal para etapas de prototipação.
- Desenvolvimento simples.
- Sem aumentar a curva de aprendizado.

Desvantagens

- Difícil manutenção e controle de erros.
- Todo o processo é atômico. Difícil dividir as operações.
- Todo o processamento é feito no nível do Kernel.



Data Science Project Automation



Algumas alternativas



Vantagens

- Integração com outros componentes.
- Tolerância a falhas.
- Alertas personalizados.
- Operações independentes.
 - Redução do tempo de processamento e custos.
- Melhor versionamento e log.
- Fácil manutenção e controle de erros.
- Melhor controle de fluxos.
- Fácil integração entre equipes de Engenharia e Ciência de Dados.

Desvantagens

- Maior curva de aprendizado.
- Maior tempo de desenvolvimento.