

---

## Atelier 01 : La Pratique de l'ETL avec Apache Airflow

---



### 1. Objectif :

L'objectif de cet atelier est de vous familiariser avec les concepts de l'ETL (Extract, Transform, Load) en utilisant Apache Airflow, une plateforme open-source pour développer, planifier et surveiller des workflows par lots.

### 2. Introduction :

Les données sont essentielles à toute application et sont utilisées dans la conception d'un pipeline efficace pour la livraison et la gestion de l'information au sein d'une organisation. En général, on définit un pipeline de données lorsqu'il est nécessaire de traiter les données tout au long de leur cycle de vie. Le pipeline peut commencer là où les données sont générées et stockées dans n'importe quel format. Le pipeline peut se terminer par l'analyse des données, leur utilisation en tant qu'informations commerciales, leur stockage dans un entrepôt de données ou leur traitement dans un modèle d'apprentissage automatique (figure 1).

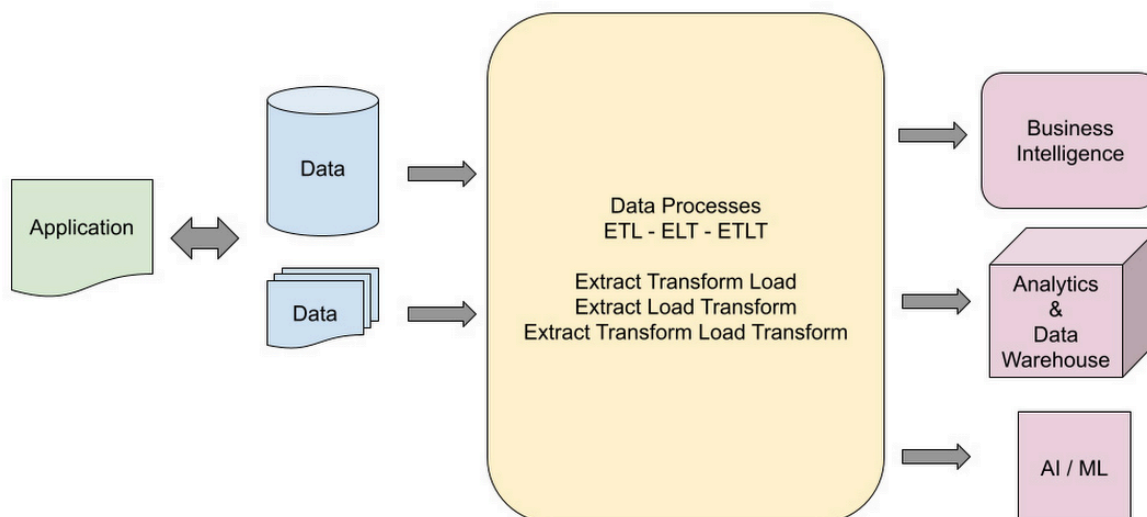


Figure 1 : le processus ETL

Les données sont extraites, traitées et transformées en plusieurs étapes en fonction des besoins des systèmes en aval. Les étapes de traitement et de transformation sont définies dans un pipeline de données. Selon les exigences, les pipelines peuvent être aussi simples qu'une seule étape ou aussi complexes que plusieurs étapes de transformation et de traitement.

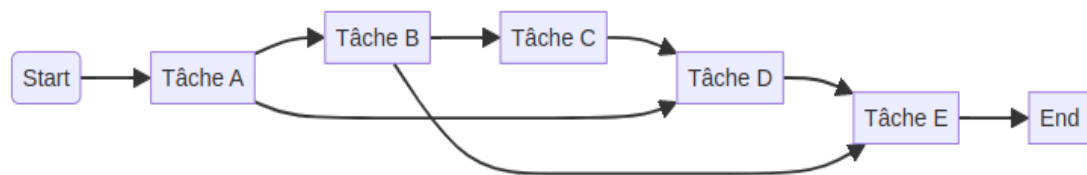
Les processus ETL garantissent que les données sont préparées et disponibles pour les workflows de traitement **en temps réel** et **par lots**. Le traitement en temps réel gère les flux de données au fur et à mesure qu'ils se produisent, tandis que le traitement par lots gère les données à des intervalles prédéterminés et de grande taille.

Aspect	Traitement par Lots	Traitement en Temps Réel
Latence	Haute (retardée)	Basse (instantanée)
Complexité	Moins complexe	Plus complexe
Coût	Moins coûteux	Plus coûteux
Volume de Données	Large volume traité à intervalles	Données continues et flux de données
Réactivité	Retardée	Immédiate
Cas d'Utilisation	Rapports de fin de journée, agrégations	Surveillance en temps réel, alertes

### 3. Concepts Clés d'Apache Airflow :

Apache Airflow repose sur plusieurs concepts fondamentaux qui définissent son fonctionnement et son approche en matière de gestion de workflows. Comprendre ces concepts est essentiel pour tirer le meilleur parti d'Airflow.

- **Directed Acyclic Graphs (DAGs)** : Le cœur d'Airflow est le DAG, un graphe acyclique dirigé. Un DAG dans Airflow représente l'ensemble des tâches à exécuter, organisées de manière à refléter leurs dépendances et leur ordre d'exécution. Chaque nœud du DAG représente une tâche et les arêtes indiquent les dépendances entre ces tâches. Cette structure assure qu'une tâche ne s'exécute que lorsque ses dépendances sont satisfaites, apportant clarté et efficacité dans la gestion des workflows.



- Les Operators sont les blocs de construction d'un DAG. Ils définissent ce qu'une tâche doit faire. Airflow propose une variété d'Operators intégrés pour des tâches courantes comme exécuter une commande Python, un script Bash, ou interagir avec des bases de données et des systèmes de fichiers. La personnalisation et la création d'Operators spécifiques permettent d'étendre les fonctionnalités d'Airflow pour répondre aux besoins spécifiques d'un workflow.
- Une **tâche** (Task) dans **Airflow** est une instance d'un Operator. Lorsqu'un DAG est exécuté, **Airflow** crée des instances de tâches (Task Instances) pour chaque tâche définie dans le DAG. Ces instances de tâches sont les éléments qui sont réellement exécutés par les **workers** d'**Airflow**. Le suivi de l'état de chaque instance de tâche (comme réussi, échoué, en cours etc.) est essentiel pour la gestion et le dépannage des workflows.
- **Workflow Scheduling** : Airflow n'est pas seulement un système d'exécution de tâches, mais aussi un planificateur de tâches. Il permet de définir des plannings complexes pour l'exécution des workflows, en utilisant des expressions cron ou des intervalles personnalisés. Cette fonctionnalité rend

Airflow particulièrement puissant pour les ETL, où les données doivent être traitées et déplacées à des intervalles réguliers.

#### 4. Les « workflows en tant que code » de Airflow :

- Dynamique : les pipelines Airflow sont configurés sous forme de code Python, ce qui permet la génération dynamique de pipelines.
- Extensible : le framework Airflow® contient des opérateurs permettant de se connecter à de nombreuses technologies. Tous les composants Airflow sont extensibles pour s'adapter facilement à votre environnement.
- Flexible : le paramétrage du workflow est intégré en exploitant le moteur de création de modèles Jinja.

#### 5. Installation d'Airflow :

- Pour installer Apache Airflow en utilisant pip, suivez les étapes ci-dessous. Notez qu'Apache Airflow nécessite un environnement Python et quelques dépendances spécifiques. Nous allons installer Airflow dans un environnement virtuel pour éviter les conflits avec d'autres packages Python. Assurez-vous d'avoir Python 3.6+ et pip installés sur votre système:

```
msakkari@msakkari-ThinkPad-E580:~$ python3 -V
Python 3.10.12
pip -V
pip 22.0.2 from /usr/lib/python3/dist-packages/pip (python 3.10)
```

- Il est recommandé de créer un environnement virtuel pour l'installation d'Airflow. Cela permet de gérer les dépendances de manière isolée :

(virtualenv airflow\_env, source airflow\_env/bin/activate)

- Install : pip install apache-airflow
- définir la variable d'environnement AIRFLOW\_HOME :

```
export AIRFLOW_HOME=${PWD}
```

- créer le fichier de configuration :

```
airflow config list --defaults > "${AIRFLOW_HOME}/airflow.cfg"
```

- Dans le fichier décommenter les lignes suivantes :

```
7 dags_folder = /home/msakkari/dags
```

```
124 load_examples = True
```

```
638 base_log_folder = /home/msakkari/logs
```

- mkdir dags/logs

## 6. Création de la base de données

- migrer la base de données : `airflow db migrate`
- créer un utilisateur pour accéder à l'interface Web d'Airflow

```
airflow users create --username admin --password admin --firstname
FIRST_NAME --lastname LAST_NAME --role Admin --email
EMAIL_ADDRESS
```

```
example : airflow users create --username admin --password admin
--firstname med --lastname sakkari --role Admin --email
sakkari.tn@gmail.com
```

## 7. Interface Utilisateur d'Airflow

- démarrer le serveur Web et le scheduler avec : `airflow webserver`
- L'interface utilisateur web d'Apache Airflow offre une visibilité et un contrôle complets sur les workflows. Dès que vous avez configuré et lancé le serveur web Airflow, vous pouvez accéder à cette interface via votre navigateur à l'adresse `http://localhost:8080`. Connectez-vous avec le user admin créé précédemment.
- Le tableau de bord principal affiche une liste de tous les DAGs disponibles. Pour chaque DAG, vous pouvez voir son état (actif ou non), la fréquence d'exécution, la dernière exécution et d'autres détails utiles. C'est le point de départ pour surveiller et gérer vos workflows.

The scheduler does not appear to be running.  
The DAGs list may not update, and new tasks will not be scheduled.

Do not use **SQLite** as metadata DB in production – it should only be used for dev/testing. We recommend using Postgres or MySQL. [Click here](#) for more information.

Do not use the **SequentialExecutor** in production. [Click here](#) for more information.

## DAGs

All 0 Active 0 Paused 0 Running 0 Failed 0
 

☒ Auto-refresh

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions	Links
No results								

Showing 0-0 of 0 DAGs

## 8. Création d'un Workflow Simple :

Supposons que nous voulions créer un workflow quotidien qui exécute une tâche de script Python. Ce workflow peut être défini comme suit :

- Définition du DAG : Tout d'abord, nous devons définir le DAG dans Airflow. Créez un fichier Python dans le répertoire dags d'Airflow. Nommons ce fichier `data5_atelier01_exemple1_dag.py`
- Importation des Modules
- Définition de la Fonction de la Tâche
- Paramètres du DAG
- Création de l'Operator

```

instance connection      airflow.cfg      *data5_atelier01_exemple1_dag.py
1 # B. Importation des Modules
2 from datetime import timedelta
3 from airflow import DAG
4 from airflow.operators.python_operator import PythonOperator
5 from airflow.utils.dates import days_ago
6
7 # C. Définition de la Fonction de la Tâche : Définissez une fonction Python pour la tâche que vous voulez exécuter. Par exemple :
8 def my_simple_task():
9     print("Exécution de ma tâche simple")
10
11 # D. Paramètres du DAG : Définissez les paramètres de votre DAG, comme l'identifiant du DAG, le calendrier et les paramètres par défaut :
12 default_args = {
13     'owner': 'airflow',
14     'depends_on_past': False,
15     'start_date': days_ago(1),
16     'email_on_failure': False,
17     'email_on_retry': False,
18     'retries': 1,
19     'retry_delay': timedelta(minutes=5),
20 }
21
22 dag = DAG(
23     'data5_atelier01_exemple1_dag',
24     default_args=default_args,
25     description='Un exemple de DAG simple',
26     schedule_interval=timedelta(days=1),
27 )
28 # E. Création de l'Operator : Créez un operator pour exécuter la fonction de tâche. Dans cet exemple, nous utilisons un PythonOperator :
29 t1 = PythonOperator(
30     task_id='simple_task',
31     python_callable=my_simple_task,
32     dag=dag,
33 )

```

F. **Démarrage du scheduler** : Lancez le scheduler d'**Airflow**, qui surveille les DAGs et exécute les tâches : Airflow scheduler,

The screenshot shows the Airflow web interface at localhost:8080. The top navigation bar includes links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The main section is titled "DAGs" and displays a table of DAGs. The table has columns for DAG name, Owner, Runs, Schedule, Last Run, Next Run, and Recent Tasks. The first five DAGs are listed, including 'dataset\_consumes\_1', 'dataset\_consumes\_1\_and\_2', 'dataset\_consumes\_1\_never\_scheduled', 'read\_dataset\_event', and 'read\_dataset\_event\_from\_classic'. The interface also shows filters for DAG status (All, Active, Paused, Running, Failed) and a search bar.

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
dataset_consumes_1	airflow	0	Dataset		On s3://dag1/output_1.txt	
dataset_consumes_1_and_2	airflow	0	Dataset		0 of 2 datasets updated	
dataset_consumes_1_never_scheduled	airflow	0	Dataset		0 of 2 datasets updated	
read_dataset_event	airflow	0	@daily		2024-09-29, 00:00:00	
read_dataset_event_from_classic	airflow	0	@daily		2024-09-29, 00:00:00	

Showing 1-5 of 5 DAGs

G. **Exécution du Workflow** : Une fois le fichier enregistré dans le répertoire dags, Airflow détectera automatiquement le nouveau DAG. Vous pouvez alors l'activer et le surveiller via l'interface utilisateur web:

airflow webserver

The screenshot shows the Airflow web interface at localhost:8080/home?tags=consumes. The top navigation bar includes links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The main section is titled 'DAGs' and displays a list of DAGs filtered by the tag 'consumes'. The list includes DAGs like 'dataset\_consumes\_1', 'dataset\_consumes\_1\_and\_2', 'dataset\_consumes\_1\_never\_scheduled', 'read\_dataset\_event', and 'read\_dataset\_event\_from\_classic'. A search bar on the right shows the input 'd', and a dropdown menu displays search results, with 'data5\_atelier01\_exemple1\_dag' highlighted. The interface also shows status counts (All 5, Active 0, Paused 5, Running 0, Failed 0) and an 'Auto-refresh' toggle.

The screenshot shows the Airflow web interface at localhost:8080/dags/data5\_atelier01\_exemple1\_dag/grid?tags=consumes&search=data5\_atelier01\_exemple1\_dag&task\_id=simple\_task&tab=graph. The top navigation bar is the same as the previous screenshot. The main section is titled 'DAG: data5\_atelier01\_exemple1\_dag' and displays the DAG's graph. The task 'simple\_task' is highlighted in the left sidebar. The graph view shows a single task 'simple\_task' of type 'PythonOperator'. The interface includes filters for 'All Run Types' and 'All Run States', a 'Clear Filters' button, and an 'Auto-refresh' toggle. A status bar at the top of the graph view shows various task states like 'deferred', 'failed', 'queued', 'removed', 'restarting', 'running', 'scheduled', 'shutdown', 'skipped', 'success', 'up\_for\_reschedule', 'up\_for\_retry', 'upstream\_failed', and 'no\_status'.



Après avoir exécuté votre premier workflow dans Apache Airflow, explorez l'interface web et répondez aux questions suivantes pour mieux comprendre les différentes fonctionnalités :

1. Tableau de Bord (Dashboard) :
  - Que pouvez-vous voir sur le tableau de bord d'Airflow ? (Décrivez chaque onglet avec une capture d'écran)
  - Comment pouvez-vous accéder à la liste de tous les DAGs disponibles ?
2. Visualisation des DAGs :
  - Lorsque vous sélectionnez un DAG, quelles informations sont affichées concernant ses tâches ?
  - Comment pouvez-vous identifier l'état d'exécution de chaque tâche dans le DAG ?
3. Historique des Exécutions :
  - Où pouvez-vous trouver l'historique des exécutions d'un DAG spécifique ?
  - Quelles informations sont fournies pour chaque exécution dans l'historique ?
4. Accès aux Journaux (Logs) :
  - Comment pouvez-vous accéder aux journaux d'une tâche spécifique ?
  - Pourquoi les journaux sont-ils importants pour le suivi et le débogage des workflows ?
5. Configuration et Paramètres :
  - Quels types de paramètres ou configurations pouvez-vous ajuster via l'interface web ?
  - Comment ces paramètres peuvent-ils affecter l'exécution des workflows ?

*Bon travail*