

# Rapport de Compétences

SAE : Collecte Automatisée de Données Web



Aimad Hamdaoui

BUT Sciences des Données

22 juin 2025

## Table des matières

<b>1</b>	<b>Contexte du Projet</b>	<b>2</b>
<b>2</b>	<b>Niveau 1 : Traiter des données structurées</b>	<b>2</b>
2.1	Mesurer l'importance de maîtriser la structure des données . . . . .	2
2.2	Comprendre les structures algorithmiques de base . . . . .	2
<b>3</b>	<b>Niveau 2 : Automatiser le traitement de données multidimensionnelles</b>	<b>2</b>
3.1	Réaliser le rôle central de l'entrepôt de données . . . . .	2
3.2	Identifier et résoudre les problèmes d'intégration de sources . . . . .	3
3.3	Comprendre la nécessité de documenter et tester . . . . .	3
<b>4</b>	<b>Démonstration Visuelle du Script</b>	<b>3</b>
<b>5</b>	<b>Conclusion</b>	<b>5</b>

# 1 Contexte du Projet

Le projet "SAE - Collecte Automatisée de Données Web" a pour objectif de développer un script Python capable d'automatiser le processus de collecte, de stockage et d'interrogation de données scientifiques. Le script interagit avec une API pour récupérer des informations sur des articles, les stocke dans une base de données SQLite et propose une interface en ligne de commande pour effectuer des recherches.

## 2 Niveau 1 : Traiter des données structurées

### 2.1 Mesurer l'importance de maîtriser la structure des données

Le projet met en évidence la transition cruciale entre les données semi-structurées de l'API (JSON) et un schéma relationnel structuré (SQL). Cette transformation a nécessité une conception réfléchie pour garantir la cohérence et l'efficacité des recherches. La création de deux tables normalisées, `articles` et `authors`, en est la preuve.

```
1 CREATE TABLE IF NOT EXISTS articles (  
2     id INTEGER PRIMARY KEY,  
3     title TEXT,  
4     year INTEGER,  
5     n_citation INTEGER,  
6     doc_type TEXT,  
7     publisher TEXT,  
8     ...  
9 );  
10  
11 CREATE TABLE IF NOT EXISTS authors (  
12     author_id INTEGER PRIMARY KEY AUTOINCREMENT,  
13     article_id INTEGER,  
14     name TEXT,  
15     org TEXT,  
16     FOREIGN KEY (article_id) REFERENCES articles(id)  
17 );
```

Listing 1 – Définition SQL de la structure de la base de données.

### 2.2 Comprendre les structures algorithmiques de base

Le script est organisé autour d'une structure algorithmique procédurale claire dans la fonction `main`, qui orchestre chaque étape du traitement : récupération, connexion, création, insertion, et interaction.

## 3 Niveau 2 : Automatiser le traitement de données multidimensionnelles

### 3.1 Réaliser le rôle central de l'entrepôt de données

La base de données `articles.db` agit comme un mini-entrepôt de données. Elle centralise et structure les informations extraites, les rendant disponibles pour l'analyse via le menu interactif. Ce dernier est l'outil qui interroge l'entrepôt, comme le montre la fonction de recherche par auteur qui joint les deux tables.

```

1 def search_by_author(conn, author_name):
2     sql = """
3     SELECT a.id, a.title, a.year
4     FROM articles a
5     JOIN authors au ON a.id = au.article_id
6     WHERE au.name LIKE ?
7     """
8     cursor = conn.cursor()
9     cursor.execute(sql, ('%' + author_name + '%',))
10    # ... Affiche les resultats

```

Listing 2 – Exemple de requête joignant les tables pour la recherche.

### 3.2 Identifier et résoudre les problèmes d'intégration de sources

Le projet est un cas d'école d'un processus ETL (Extract, Transform, Load). La principale difficulté réside dans la transformation des données brutes de l'API (source hétérogène, Figure ??) en un format propre et structuré pour la base de données. Le script gère les champs potentiellement manquants pour assurer la robustesse de l'insertion.

### 3.3 Comprendre la nécessité de documenter et tester

La robustesse du script est assurée par une gestion systématique des erreurs (`try...except`) lors des opérations critiques (connexion, requêtes, saisie). Le code est commenté et accompagné d'une documentation externe, démontrant une compréhension de l'importance de la maintenabilité.

## 4 Démonstration Visuelle du Script

L'enchaînement des figures ci-dessous illustre une session utilisateur complète, de l'exécution du script à l'obtention d'informations détaillées sur un article.

```

1. Hrachya Astsatryan (Org: Institute for Informatics and Automation Problems of the National Academy, of Sciences
of the Republic of Armenia, Yerevan, Armenia)
2. Vladimir Sahakyan (Org: Institute for Informatics and Automation Problems of the National Academy, of Sciences
of the Republic of Armenia, Yerevan, Armenia)
3. Yuri Shoukouryan (Org: Institute for Informatics and Automation Problems of the National Academy, of Sciences o
f the Republic of Armenia, Yerevan, Armenia)
4. Michel Daydé (Org: Institut de Recherche en Informatique de Toulouse, École Nationale Supérieure d'Électrotechni
que, d'Électronique, d'Informatique, d'Hydraulique et des Télécommunicati ...#TAB#)
5. Aurelie Hurault (Org: Institut de Recherche en Informatique de Toulouse, École Nationale Supérieure d'Électrote
chnique, d'Électronique, d'Informatique, d'Hydraulique et des Télécommunicati ...#TAB#)
6. Marc Pantel (Org: Institut de Recherche en Informatique de Toulouse, École Nationale Supérieure d'Électrotechni
que, d'Électronique, d'Informatique, d'Hydraulique et des Télécommunicati ...#TAB#)
7. Eddy Caron (Org: Laboratoire de l'Informatique du Parallélisme, Université de Lyon - UMR ENS-Lyon, INRIA, CNRS,
UCBL1 5668, Lyon, France 69364#TAB#)

```

FIGURE 1 – Lancement du script : initialisation de la collecte de données depuis l'API et connexion à la base de données locale `articles.db`.

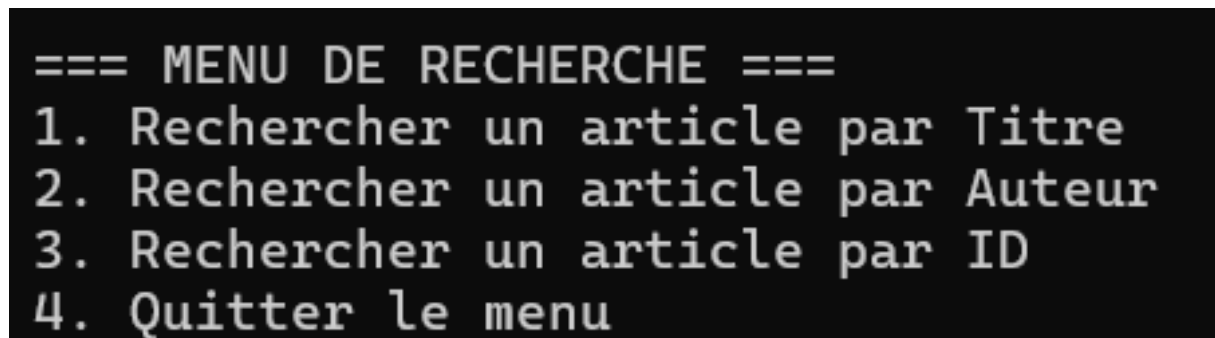


FIGURE 2 – Interaction avec le menu : l'utilisateur choisit la recherche par auteur (choix 2) et saisit le mot-clé "Caron".

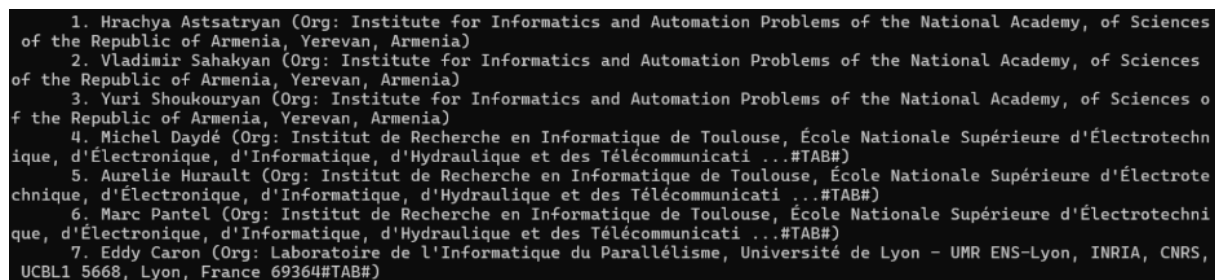


FIGURE 3 – Affichage des résultats : le script liste les articles correspondants, puis, après une nouvelle recherche par ID, affiche les détails complets d'un article, y compris la liste de ses auteurs.

## 5 Conclusion

Le projet de collecte automatisée a permis de mettre en pratique l'ensemble du cycle de vie de la donnée. Il valide solidement les compétences du **Niveau 1** (manipulation de structures, syntaxe) et démontre une maîtrise approfondie des concepts du **Niveau 2**.

Comme l'illustrent les captures d'écran, le projet a abouti à un outil fonctionnel qui automatise un flux de données complexe, depuis une source web hétérogène jusqu'à un entrepôt de données local et interrogeable. Cela prouve la capacité à concevoir et réaliser une solution logicielle complète, robuste et documentée pour un problème concret de traitement de données.