

**ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH**  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN**



**BÁO CÁO LAB**  
**LAB 3**  
**DỮ LIỆU XML**

**Giáo Viên HD:**  
**Đỗ Văn Tiến**

Thực hiện:

- Ngô Phúc Danh – 21521924

*TP. Hồ Chí Minh, ngày 26 tháng 4 năm 2023*

## ***Lời Cảm Ơn***

Bài lab “Dữ liệu xml” được dựa trên yêu cầu của môn học Tính Toán Đa Phương Tiện, thuộc lớp học CS232.N21, được thực hiện bởi sinh viên thuộc trường Đại học Công Nghệ Thông Tin - Đại học Quốc gia Thành phố Hồ Chí Minh. Bài lab “Dữ liệu xml” thực hành việc thu thập dữ liệu và lưu xuống một cơ sở dữ liệu quan hệ về các bài báo khoa học.

Qua bộ môn “Tính Toán Đa Phương Tiện” sinh viên có thể xây dựng chương trình có khả năng thu thập dữ liệu từ các trang web và lưu xuống nhằm phục vụ các nhu cầu về dữ liệu. Cụ thể với yêu cầu của bài lab 3: “Dữ liệu xml”, Sinh viên có thể thực hành được cách viết Chương trình đọc thông tin từ file xml của trang dblp sau đó lưu xuống một CSDL quan hệ về các bài báo khoa học.

# I. Tổng hợp

## 1. Yêu cầu bài lab:

- Viết Chương trình đọc thông tin từ file xml của trang dblp sau đó lưu xuống một CSDL quan hệ về các bài báo khoa học.

## 2. Xác định yêu cầu:

- Dblp: DBLP (Digital Bibliography & Library Project) là một cơ sở dữ liệu trực tuyến miễn phí chứa các bài báo khoa học trong lĩnh vực khoa học máy tính và khoa học thông tin. DBLP được quản lý bởi một nhóm các nhà khoa học máy tính ở Đức và cung cấp thông tin về các bài báo, tạp chí, cuốn sách và các tài liệu liên quan khác được xuất bản trong lĩnh vực này. Các bài báo trong DBLP được tổ chức theo tên tác giả, tên tài liệu và các từ khóa liên quan. DBLP là một nguồn tài liệu quan trọng cho các nhà nghiên cứu, giảng viên và sinh viên trong lĩnh vực khoa học máy tính và khoa học thông tin.
- File xml: File có định dạng XML (eXtensible Markup Language) là một loại định dạng tệp tin được sử dụng để lưu trữ và truyền tải dữ liệu giữa các hệ thống khác nhau. XML được sử dụng rộng rãi trong các ứng dụng web, như là một cách để trao đổi dữ liệu giữa các ứng dụng khác nhau, hoặc trong các ứng dụng cơ sở dữ liệu, để lưu trữ và truy vấn dữ liệu. XML được thiết kế để dễ đọc và dễ hiểu cho con người và các máy tính. Nó sử dụng các thẻ để đánh dấu các phần khác nhau của dữ liệu, giúp cho dữ liệu được tổ chức và truy xuất một cách dễ dàng.
- CSDL quan hệ: CSDL quan hệ (Relational Database) là một loại cơ sở dữ liệu được tổ chức và lưu trữ dưới dạng bảng, trong đó mỗi bảng đại diện cho một thực thể và mỗi hàng trong bảng đại diện cho một bản ghi của thực thể đó. Mỗi cột trong bảng tương ứng với một thuộc tính của thực thể. Các bảng được liên kết với nhau thông qua các quan hệ (relationship), trong đó một quan hệ được xác định bởi một cặp khóa (key) giữa hai bảng.

➔ Sinh viên cần xây dựng crawler nhằm thu thập dữ liệu trong file xml thuộc trang web dblp về danh sách các bài báo khoa học sau đó lưu xuống database.

## 3. Công nghệ sử dụng:

- Ngôn ngữ lập trình(programming language): Python.
- Tool: selenium.
- Libraries: numpy, pypyodbc, xml.etree.ElementTree.

## II. Nội dung

### 1. Thử thách (challenge):

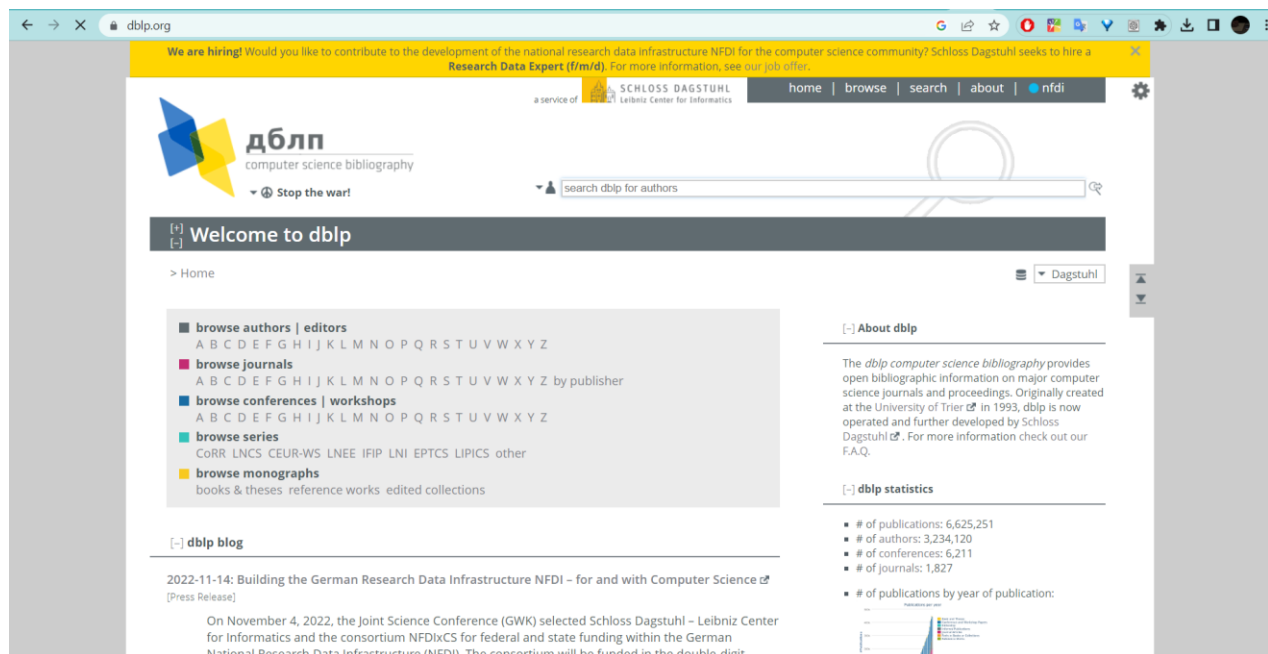
- Crawl dữ liệu trên dblp: Để crawl được dữ liệu về các bài báo trên dblp chúng ta cần phải truy cập vào trang web, chọn khung author và nhập vào tên tác giả sau đó nhấn vào tác giả vừa tìm được chúng ta sẽ được chuyển sang trang mới chứa danh sách các bài báo có sự tham gia của tên tác giả mà chúng ta vừa nhập, cuối cùng là crawl về máy tính cá nhân file xml.
- SQL server: Sau khi chúng ta giải quyết vấn đề crawl data thì chúng ta cần lưu xuống một cơ sở dữ liệu quan hệ.
- Vấn đề về chặn ip: trang web không có thao tác chặn bot như captcha

### 2. Hướng tiếp cận (approach):

- Crawl dữ liệu: ta thấy rằng để đến được danh sách các bài báo cần thông qua nhiều quá trình tương tác, chính vì thế việc sử dụng selenium để tự động hóa quá trình crawl dữ liệu xml về máy tính cá nhân.
- SQL server: Mở SQL server để connect thông qua server name và tạo database. Để tiến hành phân chia các thẻ trong file xml thành các cột trong một bảng chúng ta dùng thư viện xml.etree.ElementTree.
- Sử dụng openvpn để đổi ip kết hợp với tùy chỉnh thời gian truy cập random.

### 3. Minh họa bằng hình ảnh (images illustration):

- Tự động hóa crawl dữ liệu bằng selenium:



Bước 1: Truy cập trang web.

The screenshot shows the dblp (computer science bibliography) search results for the author 'Thanh Duc Ngo'. The page header includes the dblp logo, a navigation bar with links like 'home', 'browse', 'search', 'about', and 'nfdi', and a search bar. The search results are displayed under the heading 'Search dblp for Authors'. It shows 'Exact matches' and 'All 2 matches' for the author. The first match is 'Thanh Duc Ngo' with the affiliation 'Vietnam National University Ho Chi Minh City, Vietnam'. The second match is 'Ngô Đức Thanh'.

Bước 2: Nhập tên tác giả.

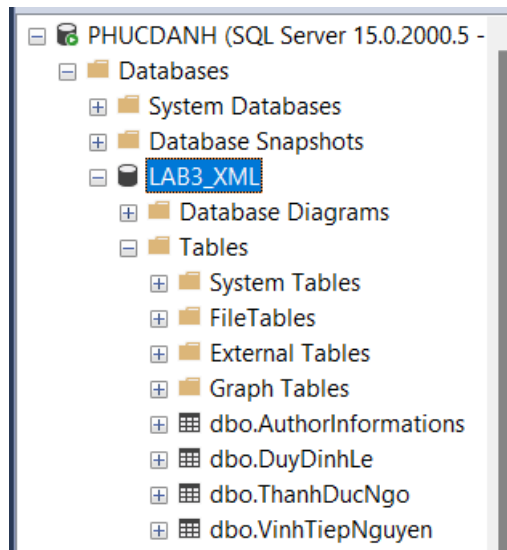
The screenshot shows the author profile page for 'Thanh Duc Ngo' on the dblp website. The page includes a header with the dblp logo and navigation links. The main content area shows the author's name and a list of publications. A dropdown menu is open, showing options to 'export bibliography' in various formats: BibTeX, RIS, RDF N-Triples, RDF Turtle, RDF/XML, XML, and RSS. The 'XML' option is highlighted.

Bước 3: Truy cập vào đường dẫn XML và crawl về máy tính.

- Tự động hóa việc lưu xuống cơ sở dữ liệu:

The screenshot shows a 'Connect to Server' dialog box. The title is 'SQL Server'. The 'Server type' is set to 'Database Engine'. The 'Server name' is 'PHUCDANH'. The 'Authentication' is set to 'Windows Authentication'. The 'User name' is 'PHUCDANH\Admin'. The 'Password' field is empty. There is a checkbox for 'Remember password' which is unchecked. At the bottom, there are buttons for 'Connect', 'Cancel', 'Help', and 'Options >>'.

Bước 1: Mở ứng dụng SQL server và tạo database



Bước 2: Sử dụng thư viện xml để tiến hành tạo bảng và lưu vào các bảng.

#### 4. Source code:

- Ý tưởng:
  - Crawler: sử dụng selenium để tự động hóa việc lấy dữ liệu.
    - Bước 1: Truy cập vào đường link: <https://dblp.org/>.
    - Bước 2: Chọn chế độ tìm kiếm bằng tên tác giả.
    - Bước 3: Thực hiện việc nhập tên tác giả vào ô tìm kiếm.
    - Bước 4: Truy cập vào mục xml và tiến hành crawl file.xml.
  - Sql: sử dụng pypyodbc kết hợp với xml.etree.ElementTree.
    - Bước 1: Tạo bảng (table) cho từng file xml.
    - Bước 2: Tạo từng cột với tên ứng với từng thẻ tag.
    - Bước 3: Lưu dữ liệu xuống cơ sở dữ liệu.
  - Run code:
    - Bước 1: Import 2 module crawler và sql đã thực hiện ở trên.
    - Bước 2: Nhập tên Server và database.
    - Bước 3: Nhập tên các tác giả umôn crawl bài báo.
    - Bước 5: Tiến hành chạy chương trình.
- Code:
  - Crawler:

```

Crawler

# import libraries
from selenium import webdriver
from selenium.webdriver.common.keys import Keys
from selenium.webdriver.common.by import By
from time import sleep
import numpy as np
import requests

# implement crawler for crawling xml file
class crawling:
    def __init__(self, browser):
        self.browser = browser

    def navigate_to_dblp(self):
        sleep(np.random.randint(2, 6))
        dblp_url = "https://dblp.org/"
        self.browser.get(dblp_url)
        sleep(np.random.randint(2, 6))
        return self.browser

    def author_searching_function(self, authorname):
        searching_options_selector = "#search-mode-selector > div.head > img"
        searching_options = self.browser.find_element(
            By.CSS_SELECTOR,
            searching_options_selector,
        )
        searching_options.click()
        searching_author_selector = "#search-mode-author"
        searching_authors = self.browser.find_element(
            By.CSS_SELECTOR,
            searching_author_selector,
        )
        searching_authors.click()
        sleep(np.random.randint(2, 6))
        author_name_xpath = "/html/body/div[2]/div[2]/form/input"
        author_name_searching = self.browser.find_element(
            By.XPATH,
            author_name_xpath,
        )
        author_name_searching.send_keys(authorname)
        author_name_searching.send_keys(Keys.ENTER)
        sleep(np.random.randint(2, 6))
        author_chosen_selector = (
            "#completeSearch-authors > div > ul:nth-child(2) > li > a"
        )
        author_chosen = self.browser.find_element(
            By.CSS_SELECTOR,
            author_chosen_selector,
        )
        author_chosen.click()
        return self.browser

    def access_author_link(self):
        sleep(np.random.randint(2, 6))
        xml_selector = "#headline > nav > ul > li.export.drop-down > div.body > ul:nth-child(2) > li:nth-child(6) > a"
        xml_option = self.browser.find_element(By.CSS_SELECTOR, xml_selector)
        xml_link = xml_option.get_attribute("href")
        return xml_link

    def download_xml(self, author_name, xml_link):
        author_name = author_name.replace(" ", "")
        filename = f"{author_name}.xml"
        response = requests.get(xml_link)
        if response.status_code == 200:
            with open(filename, "wb") as f:
                f.write(response.content)
            print(f"Downloaded {filename} successfully.")
        else:
            print(f"Failed to download {filename}. Status code: {response.status_code}")
            sleep(5)
        self.browser.close()
        return filename

```

Hình 2.4.1: Source code crawler.

## ➤ Sql:

```

# import Libraries
import pypodbc
import xml.etree.ElementTree as ET

# implement the connection to SQL SERVER
class mysql_server:
    def __init__(self, server, database):
        self.server = server
        self.database = database

    def connect_server(self):
        conn = pypodbc.connect(
            "Driver={SQL SERVER};"
            f"Server={self.server};"
            f"Database={self.database};"
            "Trusted_Connection=yes;"
        )
        return conn

    def create_author_table(self, cursor):
        cursor.execute(
            """
            CREATE TABLE AuthorInformations
            (ID INT PRIMARY KEY, Author VARCHAR(255), Title VARCHAR(355))
            """
        )
        return cursor

    def insert_author_table(self, cursor, name, articles, conn, id):
        id += 1
        for ar in articles:
            try:
                title = ar.find("./title").text
            except:
                title = ar.find("./i").text
            cursor.execute(
                "INSERT INTO AuthorInformations (ID, Author, Title) VALUES (?, ?, ?)",
                (
                    id,
                    name,
                    title,
                ),
            )
            conn.commit()
            id += 1
        return conn

    def create_table(self, cursor, file_name):
        name = file_name.replace(".xml", "")
        cursor.execute(
            """
            CREATE TABLE {}
            (ID INT PRIMARY KEY, IDPaper INT, AuthorAndCoauthor VARCHAR(955), Title
            VARCHAR(355), Pages VARCHAR(255), Year INT, Volume VARCHAR(255), Booktitle VARCHAR(255),
            Journal VARCHAR(255), ee VARCHAR(255), crossref VARCHAR(255), url_link VARCHAR(255))
            """
            .format(
                name
            )
        )
        return cursor

    def parse_xml(self, file_xml):
        tree = ET.parse(file_xml)
        articles = tree.findall("r")
        return articles

def insert_information(self, articles, cursor, conn, file_name, id_paper):
    id = 1
    id_paper += 1
    name = file_name.replace(".xml", "")
    for ar in articles:
        try:
            authors = [author.text for author in ar.findall("./author")]
            author = ", ".join(authors)
        except:
            author = None
        title = ar.find("./title").text
        if title == None:
            title = f"Not finding {id}"
        try:
            page = ar.find("./pages").text
        except:
            page = None
        try:
            year = ar.find("./year").text
        except:
            year = None
        try:
            volume = ar.find("./volume").text
        except:
            volume = None
        try:
            booktitle = ar.find("./booktitle").text
        except:
            booktitle = None
        try:
            journal = ar.find("./journal").text
        except:
            journal = None
        try:
            ee = ar.find("./ee").text
        except:
            ee = None
        try:
            crossref = ar.find("./crossref").text
        except:
            crossref = None
        try:
            url_link = ar.find("./url").text
        except:
            url_link = None
        cursor.execute(
            "INSERT INTO {} (ID, IDPaper, AuthorAndCoauthor, Title, Pages, Year,
            Volume, Booktitle, Journal, ee, crossref, url_link) VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?,
            ?, ?)".format(
                name
            ),
            (
                id,
                id_paper,
                author,
                title,
                page,
                year,
                volume,
                booktitle,
                journal,
                ee,
                crossref,
                url_link,
            ),
        )
        conn.commit()
    except:
        continue
    id_paper += 1
    id += 1

```

Hình 2.4.2: Source code Sql connecting.



➤ Run code:

```
Run code

# import module
import crawler
import SQLconnection
from crawler import crawling
from SQLconnection import mysql_server

def main():

    # insert server, database, number of authors and their names
    server = input("Enter the server name: ")
    database = input("Enter the database name: ")
    n = int(input("Enter the number of author: "))
    author_name_list = []
    for i in range(n):
        temp_author_name = input(f"Enter the author's name ({i+1}): ")
        author_name_list.append(temp_author_name)

    # create an instance of mysql_server
    SQL_connect = mysql_server(server, database)
    conn = SQL_connect.connect_server()
    cursor = conn.cursor()
    SQL_connect.create_author_table(cursor)
    table_name_list = []
    id = 1000
```

```
Run code

# start crawling and save into database
for name in author_name_list:
    browser = crawler.webdriver.Chrome(executable_path="chromedriver.exe")
    crawler_dbfp = crawling(browser)
    crawler_dbfp.navigate_to_dbfp()
    crawler_dbfp.author_searching_function(name)
    xml_link = crawler_dbfp.access_author_link()
    file_name = crawler_dbfp.download_xml(name, xml_link)
    cursor = SQL_connect.create_table(cursor, file_name)
    articles = SQL_connect.parse_xml(file_name)
    SQL_connect.insert_information(articles, cursor, conn, file_name, id)
    SQL_connect.insert_author_table(cursor, name, articles, conn, id)
    table_name_list.append(file_name)
    id += 1000

# create relational database
for table_name in table_name_list:
    name = table_name.replace(".xml", "")
    try:
        cursor.execute(
            "ALTER TABLE {} ADD CONSTRAINT {} FOREIGN KEY (IDPaper) REFERENCES
            AuthorInformations (ID)".format(
                name, "fk_" + name
            )
        )
    except Exception as e:
        print(
            "Error creating foreign key constraint for table {}: {}".format(name, e)
        )

# close conn and cursor
conn.close()
cursor.close()

# run code
if __name__ == "__main__":
    main()
```

Hình 2.4.3: Source code main.py.

### III. Cài đặt và kiểm thử kết quả

#### 1. Cài đặt:

- Folder chứa source code gồm 3 file .py bao gồm (crawler.py, SQLconnection, main.py).
- Người dùng thực hiện việc download chromedriver trước khi thực hiện run code.
- Người dùng mở ứng dụng sql server management chọn server sau đó kết nối và tạo database. Ở bài báo cáo này, sinh viên thực hiện sử dụng server: **PHUCDANH** và database: **LAB3\_XML**.
- Sau khi chuẩn bị xong, người dùng chạy file main.py.
- Nhập vào tên server và database để lưu trữ.
- Nhập vào số lượng tác giả và tên của từng tác giả, ở đây sinh viên thực hiện kiểm thử với 3 tác giả.
- Chờ và ghi lại kết quả.

#### 2. Kiểm thử kết quả:

- Kết quả sau khi crawl tất cả bài báo của 3 tác giả là: thầy Ngô Đức Thành, thầy Lê Đình Duy và thầy Nguyễn Vinh Tiệp.

```
ThanhDucNgo.xml - Notepad
File Edit View

<?xml version="1.0"?>
<dblpperson name="Thanh Duc Ngo" pid="65/3565" n="74">
  <person key="homepages/65/3565" mdate="2022-02-15">
    <author pid="65/3565">Thanh Duc Ngo</author>
    <author pid="65/3565">Thanh-Duc Ngo</author>
    <note type="affiliation">Vietnam National University Ho Chi Minh City, Vietnam</note>
    <note type="affiliation" label="PhD">Graduate University for Advanced Studies, Tokyo, Japan</note>
    <url>https://scholar.google.com/citations?user=I8bNZakAAAAJ</url>
    <url>https://dl.acm.org/profile/81392602249</url>
    <url>https://orcid.org/0000-0001-6882-0070</url>
  </person>
  <?><article pubtype="informal" key="journals/corr/abs-2304-07444" mdate="2023-04-24">
    <author pid="305/2136">Thanh-Danh Nguyen</author>
    <author pid="305/1969">Anh-Khoa Nguyen Vu</author>
    <author pid="305/1575">Nhat-Duy Nguyen</author>
    <author pid="60/11111">Vinh-Tiep Nguyen</author>
    <author pid="65/3565">Thanh Duc Ngo</author>
    <author pid="28/8692">Thanh-Toan Do</author>
    <author pid="44/7448">Minh-Triet Tran</author>
    <author pid="119/1364-2">Tam V. Nguyen 0002</author>
    <title>Few-shot Camouflaged Animal Detection and Segmentation.</title>
    <year>2023</year>
    <volume>abs/2304.07444</volume>
  </journal><CoRR</journal>
</dblp>
Ln 1, Col 1 100% Unix (LF) UTF-8
```

Hình 3.2.1: file xml các bài báo có sự tham gia của thầy Thành.

```
DuyDinhLe.xml - Notepad
File Edit View

<?xml version="1.0"?>
<dblpperson name="Duy-Dinh Le" pid="48/4683" n="114">
  <person key="homepages/48/4683" mdate="2009-06-10">
    <author pid="48/4683">Duy-Dinh Le</author>
  </person>
  <?><article key="journals/sivp/NguyenVNPDDNNL22" mdate="2022-10-18">
    <author orcid="0000-0002-9188-8564" pid="314/0623">Xuan-Duong Nguyen</author>
    <author pid="305/1969">Anh-Khoa Nguyen Vu</author>
    <author orcid="0000-0001-6577-2122" pid="305/2136">Thanh-Danh Nguyen</author>
    <author pid="191/3530">Nguyen Phan</author>
    <author pid="330/6486">Bao-Duy Duyen Dinh</author>
    <author orcid="0000-0001-8566-273X" pid="305/1575">Nhat-Duy Nguyen</author>
    <author orcid="0000-0003-0236-7992" pid="119/1364-2">Tam V. Nguyen 0002</author>
    <author orcid="0000-0003-4260-7874" pid="60/11111">Vinh-Tiep Nguyen</author>
    <author orcid="0000-0003-0356-5501" pid="48/4683">Duy-Dinh Le</author>
    <title>Adaptive multi-vehicle motion counting.</title>
    <pages>2193-2201</pages>
    <year>2022</year>
    <volume>16</volume>
  </journal>Signal Image Video Process.</journal>
  <number>8</number>
  <ee>https://doi.org/10.1007/s11760-022-02184-5</ee>
  <url>db/journals/sivp/sivp16.htm#NguyenVNPDDNNL22</url>
</article>
</dblp>
Ln 1, Col 1 100% Unix (LF) UTF-8
```

Hình 3.2.2: file xml các bài báo có sự tham gia của thầy Duy.

```
VinhTiepNguyen.xml - Notepad
File Edit View

<?xml version="1.0"?>
<dblpperson name="Vinh-Tiep Nguyen" pid="60/11111" n="72">
  <person key="homepages/60/11111" mdate="2012-03-24">
    <author pid="60/11111">Vinh-Tiep Nguyen</author>
  </person>
  <?><article pubtype="informal" key="journals/corr/abs-2304-05731" mdate="2023-04-19">
    <author pid="00/11111">Trung-Nghia Le</author>
    <author pid="119/1364-2">Tam V. Nguyen 0002</author>
    <author pid="271/7587">Minh-Quan Le</author>
    <author pid="344/5384">Trong-Thuan Nguyen</author>
    <author pid="335/8998">Viet-Tham Huynh</author>
    <author pid="223/1830">Trong-Le Do</author>
    <author pid="188/2629">Khanh-Duy Le</author>
    <author pid="214/2465">Mai-Khiem Tran</author>
    <author pid="316/4377">Nhat Hoang-Xuan</author>
    <author pid="273/7350">Thang-Long Nguyen-Ho</author>
    <author pid="60/11111">Vinh-Tiep Nguyen</author>
    <author pid="322/6036">Nhat-Quynh Le-Pham</author>
    <author pid="344/6195">Huu-Phuc Pham</author>
    <author pid="344/5600">Trong-Vu Hoang</author>
    <author pid="344/5396">Quang-Binh Nguyen</author>
    <author pid="325/2495">Trong-Hieu Nguyen Mau</author>
    <author pid="324/2487">Tuan-Luc Huynh</author>
    <author pid="326/6523">Thanh-Danh Le</author>
    <author pid="344/5652">Ngoc-Linh Nguyen-Ha</author>
  </article>
</dblp>
Ln 1, Col 1 100% Unix (LF) UTF-8
```

Hình 3.2.3: file xml các bài báo có sự tham gia của thầy Tiệp.



90 %

Results Messages

	ID	Author	TITLE
1	1001	Thanh Duc Ngo	Few-shot Camouflaged Animal Detection and Segmen...
2	1002	Thanh Duc Ngo	Instance-level Few-shot Learning with Class Hierarchy ...
3	1003	Thanh Duc Ngo	Few-shot object detection via baby learning.
4	1004	Thanh Duc Ngo	Antique Photo Restoration and Colorization via Generat...
5	1005	Thanh Duc Ngo	A Crowdsourcing Data Annotation System For Vietnam...
6	1006	Thanh Duc Ngo	UIT at VBS 2022: An Unified and Interactive Video Retri...
7	1007	Thanh Duc Ngo	Dictionary-Guided Scene Text Recognition.
8	1008	Thanh Duc Ngo	Unweighted Bipartite Matching For Robust Vehicle Cou...
9	1009	Thanh Duc Ngo	Multilingual-GAN: A Multilingual GAN-based Approach ...
10	1010	Thanh Duc Ngo	A robust framework for mathematical formula detection.
11	1011	Thanh Duc Ngo	DF-FSOD: A Novel Approach for Few-shot Object Detec...
12	1012	Thanh Duc Ngo	MC-OCR Challenge 2021: An end-to-end recognition fr...
13	1013	Thanh Duc Ngo	Single-image crowd counting: a comparative survey on...
14	1014	Thanh Duc Ngo	An Evaluation of Deep Learning Methods for Small Obj...
15	1015	Thanh Duc Ngo	Interpolation based Anime Face Style Transfer.
16	1016	Thanh Duc Ngo	U15-Logos: Unconstrained Logo Dataset with Evaluati...
17	1017	Thanh Duc Ngo	Searching For Desired Papers Doing Desired Action b...

Query executed successfully.

Hình 3.2.7: Kết quả khi lưu toàn bộ tác và bài báo của cả 3 tác giả.



Hình 3.2.8: Cơ sở dữ liệu quan hệ.

## IV. Đánh giá và kết luận

- **Đánh giá:**

- Kết quả thu thập được sau một quá trình tự động hóa khá tốt.
- Mỗi bài báo lưu trên file xml có cấu trúc thẻ không cố định chính vì thế các thẻ mang thông tin không có trong một bài báo sẽ được gán giá trị None hay NULL.

- **Kết Luận:**

- Ưu điểm:
  - Dễ hiểu: việc source code được chia thành các class theo hướng đối tượng giúp code được sắp xếp rõ ràng, dễ nhìn.
  - Dễ thực thi: để chạy code cần chạy file main.py nên không cần phải chạy nhiều file một lúc.
  - Kết nối nhanh: Việc lưu một lượng lớn dữ liệu xuống database diễn ra khá nhanh.
- Nhược điểm:
  - Tính toán thời gian crawl: Việc sử dụng selenium có tính tương tác cao đòi hỏi người dùng cần canh chỉnh thời gian hợp lý.
  - Yêu cầu kết nối internet ổn định: Việc loading website đòi hỏi đường truyền ổn định tuy nhiên không đòi hỏi đường truyền quá tốt.