# Project Report: Website Data Scraper

## Objective

The primary goal of this project was to develop a solution that scrapes data from a list of websites and extracts the following information:

- Social Media Links
- Tech Stack (MVC, CMS, JS type, etc.)
- Meta Title
- Meta Description
- Payment Gateways (e.g., PayPal, Stripe, Razorpay)
- Website Language
- Category of Website

The extracted data was to be stored in a MySQL database for further analysis or use.

# Approach

## 1. Requirements Gathering

To begin with, I identified the key pieces of information that needed to be extracted from each website. These included meta tags, social media links, tech stack, payment gateways, and other specific attributes.

## 2. Setting Up the Environment

The project was developed using Python, leveraging its powerful libraries for web scraping and database interaction. The libraries used included:

- `requests` for making HTTP requests to websites.
- `BeautifulSoup` for parsing HTML and XML documents.
- `mysql-connector-python` for connecting and interacting with MySQL databases.
- `re` for regular expressions to identify patterns within the HTML content.

## 3. Database Setup

A MySQL database named `website_data` was created to store the extracted data. The following table schema was defined:

```sql
CREATE DATABASE website_data;

USE website_data;

CREATE TABLE website_info (
    id INT AUTO_INCREMENT PRIMARY KEY,
    url VARCHAR(255) NOT NULL,
    social_media_links TEXT,
    tech_stack TEXT,
    meta_title VARCHAR(255),
    meta_description TEXT,
    payment_gateways TEXT,
    website_language VARCHAR(50),
    category VARCHAR(50)
);
```

## 4. Web Scraping Logic

The core of the project was the Python script which performed the following steps:

- **Import Libraries**

  The required libraries were imported to handle HTTP requests, HTML parsing, and database operations.

- **List of Websites**

  A list of websites was defined from which data needed to be scraped.

- **Connect to MySQL**

  A connection to the MySQL database was established using `mysql-connector-python`.

- **Define Functions to Extract Data**

  Several functions were defined to extract specific pieces of information from the HTML content:

  - `get_meta_data(soup)`: Extracts meta title and description.
  - `get_social_media_links(soup)`: Identifies and extracts social media links.
  - `get_payment_gateways(text)`: Detects payment gateways mentioned in the website text.
  - `get_tech_stack(soup)`: Identifies the tech stack used by the website.
  - `get_website_language(soup)`: Extracts the language attribute from the HTML tag.
  - `categorize_website(url)`: Categorizes the website based on keywords in the URL.

- **Scrape Website Function**

  A function `scrape_website(url)` was created to scrape data from each website and insert the extracted information into the MySQL database.

## 5. Iteration Through Websites

A loop iterated through each website URL in the list, calling the `scrape_website` function for each one.

## 6. Handling Exceptions

Appropriate exception handling was implemented to manage errors during HTTP requests, HTML parsing, and database operations.

## 7. Closing Database Connection

After all websites were scraped, the database connection was closed to ensure proper resource management.

# Challenges Encountered

### 1. Variability in Website Structures

Websites vary significantly in their HTML structure, which made it challenging to extract information consistently. To address this, I implemented robust HTML parsing logic and added multiple checks for each data point.

### 2. Handling Dynamic Content

Some websites load content dynamically using JavaScript, which is not directly accessible through basic HTTP requests. Handling such cases would require using tools like Selenium or Puppeteer, but for this project, I focused on static content.

### 3. Timeouts and Slow Responses

Web scraping can be affected by network issues, server timeouts, and slow responses. I implemented timeouts and retries in the HTTP requests to handle such scenarios.

### 4. Data Quality and Validation

Ensuring the accuracy of extracted data was crucial. I validated the extracted data before inserting it into the database to minimize errors and inconsistencies.

### 5. Legal and Ethical Considerations

Scraping websites must be done responsibly and in compliance with the website's terms of service. I ensured that the scraping logic was respectful of the websites' usage policies and did not overwhelm the servers with requests.

# Conclusion

This project provided a comprehensive solution for scraping and storing website data in a structured manner. The Python script, combined with MySQL for data storage, offered an efficient and scalable way to gather information from multiple websites. Despite the challenges, the project was successful in achieving its objectives and laid a solid foundation for further enhancements, such as handling dynamic content and improving categorization logic.