

# **Structure**

GCE A Level Computer Science Project

**Name:** Jeremiah Boby

**Candidate Number:** 8302

**School:** Loreto Sixth Form College

# Analysis

## Problem Definition and Stakeholders

The stakeholders of my application will be fellow students of Loreto Sixth Form College. Being A Level students, they need to organize revision sessions in order to maintain the grades required for their chosen university.

Many Loreto college students have issues with organising a revision timetable. For example, creating study sessions that do not clash with classes, or setting up reminders for each of these study sessions. The needs of a College student also include regularly checking up on their performance (grades and attendance) to ensure that they are making the best of their study periods. My solution will help students with these tasks, by providing a simple, intuitive interface for them to create, update and set reminders for their study sessions.

I have chosen a fellow Computer Science student, **Troy Sherlock**, to be my stakeholder for this project. Troy is an ideal candidate for my software, as he is currently in his second year of A Levels with Maths, Physics and Computer Science. This means that efficient organisation of free study periods is vital, so he will benefit greatly from my solution.

Currently, a student must do the following in order to carry this out:

- Check their current timetable on the college website
- Check their grades and predicted grades on the college website
- Use this information to create a manual revision timetable
- Regularly update this timetable manually based on the above criteria

The issue with this system is that it is entirely manual, which means that students may neglect to carry out certain steps. Also, the current system does not allow a student to add custom lessons to their online timetable, or be notified of their upcoming lessons. My solution is to have a web-based dashboard that will integrate all these features, in a single web application.

My application will solve the problem by integrating revision management, homework reminders and performance analysis in one simple dashboard. Students will be able to access this dashboard via a Django web application.

In order to solve this problem with my chosen solution, I will need to learn the following:

- Usage of the Django web framework
- Usage of the PostgreSQL Database Management System
- Integration of this DBMS with Python
- Usage of JavaScript graphing libraries such as D3.js

## Justification of Computational Methods

This problem lends itself to utilization of computational methods for several reasons:

- **Student performance data will be abstracted**

Rather than being displayed as raw text data, it will be displayed in a series of tables and charts. This will make it much easier for the end-user to comprehend. In order to display the data in this format, a computational device will be needed. This is because the client will need to be able to execute client-side code (JavaScript) and render HTML and CSS to display

- **A Server-side will need to be implemented**

As my solution will be a web application, I will need to write code to run on a server, which will serve HTML templates to the client computers.

## Application of Computational Methods

There are several ways that computational methods can be applied to this project:

| Method            | Explanation of relevant features   |
|-------------------|--|
| <b>Abstract</b>   | My web application will be as user-friendly as possible. this will require me to hide certain details from the user. For example, the raw data that is scraped by my webscraper. Instead, the user will only see data that is relevant to their current page: If they are on the summary page, they will only see simple summarized data and graphs for their performance. |
| <b>Heuristic</b>  | The summary page of my dashboard will include rounded, approximate values to represent the correlation between various properties of the student's study time, free time, grades, and punctuality.   |
| <b>Procedural</b> | My application stack will be comprised of several Python modules, each separated into classes, methods and functions. I will be able to remove modules, test modules and edit modules in a way that does not affect any irrelevant components of the application.  |
| <b>Concurrent</b> | The web application must be able to handle multiple HTTP requests simultaneously. I will be applying concurrent thinking to plan and design the processing and response for each request.  |

The statements above have shown that my problem can be solved with a computational method. This makes it suited for a computer to solve with a program.

# Research

## Stakeholder Interview Questions

- Have you ever used a homework management system before?
- If so, what are your favourite features of this system?
- Are there any features you want to see from a system like this?
- Are there any features you want removed?

### Interview with Troy Sherlock

Have you ever used a homework management system before?

- Yes, I have used Show my Homework in high school to manage my timetabled lessons and homework. However, it did not offer a system to manage revision sessions.

What are your favourite features of this system?

- My favourite feature is the ability to colour-code homework slots depending on whether they are completed or not. It allows me to easily check

Are there any features you want to see from a system like this?

- I would like to see a reminder system, such as browser push notifications for due homework. I would also like the ability to save and load custom revision sessions.

Are there any features you want removed?

- I do not like Show My Homework's "my calendar" system, and would like it to display more accurate time data for each homework slot.

From my stakeholder interview, I have concluded that **if the problem is successfully solved, the stakeholders will be able to access a web dashboard where they will be able to view their school timetable, add and remove events from this timetable, and view statistics of their study progress.**

## Checking for a Successful solution

In order to ensure that my solution was successful in meeting the requirements described above, I will need to consider the following areas:

- **Meets user needs**

I will need to ensure that the solution works well for my stakeholder's purposes. To measure this success criteria, I will use a grading system where my stakeholders will rate my application on a scale of 1 to 10 in usefulness.

- **Reliable**

I will need to ensure that the solution is robust and continues to work optimally while it is in use. In order to measure this criteria, I will ensure that my stakeholders have access to the solution over a long period of time, and they will rate the reliability of the service on a scale of 1 to 10.

## Similar solutions

I have found a similar solution to my own, **Show My Homework**, by Satchel.

Show my Homework is a service designed for schools, that aims to unify student timetable and homework tasks.

| Features worth including   | Possible limitations   |
|--|--|
| <ul style="list-style-type: none"><li>• Simple interface</li><li>• Colour – coded homework tasks</li><li>• Homework dates sync with lessons</li><li>• Homework reminders before lesson</li><li>• Mobile client application and API</li></ul> | <ul style="list-style-type: none"><li>• <b>Time limit</b><br/>I am given a relatively short time frame to complete a project of this scale. Therefore, sacrifices may be made to functionality in order to complete the project on time.</li><li>• <b>Cost</b><br/>Given that my project is one that requires a server to run my code, there may be costs incurred by the hosting providers that will be used. Therefore, I have decided to use the free tier of Heroku, a well known free hosting provider, to host my web application. This may result in performance losses, as Heroku Free Tier machines are put on “standby” when data is not requested after a set amount of time, requiring a couple of minutes to return to an active state.</li></ul> |

## Functionality Requirements

The solution will be required to:

- **Display any Loreto student's timetable in an intuitive manner**

This means that I will need to communicate with MyLoreto to retrieve this information, then parse it in such a way that my program will be able to extract the necessary information from it, and finally render this information in an abstracted format. As intuitiveness of the solution is subjective, I will be using the opinion of my stakeholders to judge this.

- **Allow students to create their own timetable entries**

This means that I will need to store the student's timetable in a persistent way, suggesting that I will need to use a database. This database will store student timetables, as well as authentication for MyLoreto to verify the lessons. I will need to insert and remove data from the database upon the user's request, so usage of SQL may be required.

## Execution Requirements

The machines running my solution will be required to:

- **Serve content at any time**

The server-side code must be run on hardware that can maximise uptime, i.e. stay online and running a python web app indefinitely. This kind of hardware is provided by a number of services, and I will be using a [Heroku Free Tier](#) virtualised private server. As mentioned above, running this application on a free service is not ideal. However, Heroku provides build integration with GitHub, allowing me to easily

push updates to my live web app. As for software requirements, the VPS will need to support the python runtime and the PostgreSQL database, which Heroku does by default.

- **Execute JavaScript code and render HTML**

The client-side code must be executed on hardware that supports a display system and web browser, as my service is a web application. This minor task can be completed by most ARM/x64 devices, such as mobile phones, laptops, desktops, or refrigerators, therefore I do not have any hardware requirements. The only software requirement is to be able to support a web browser.

## Design

### Key Objects

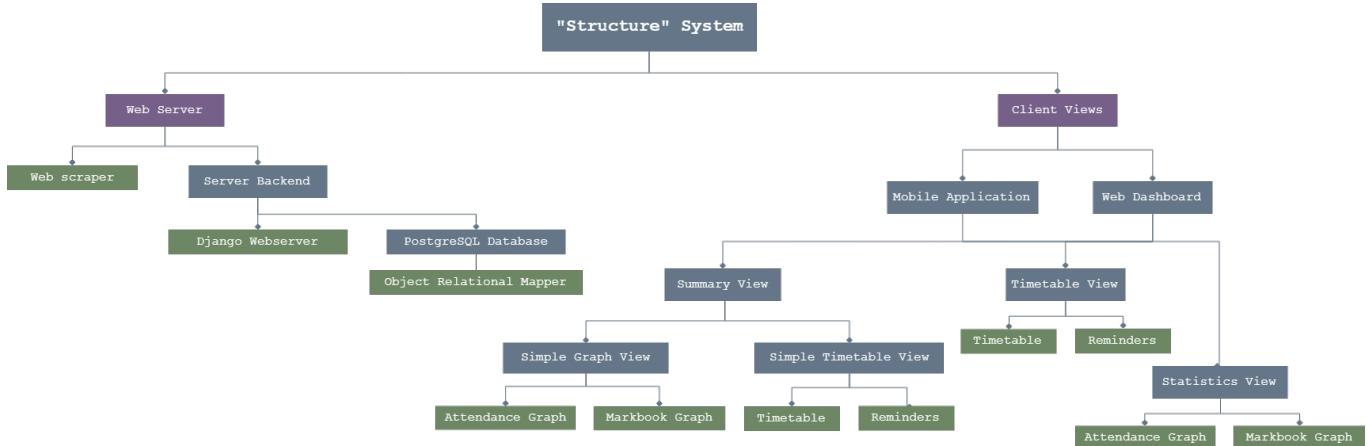
There are several key objects that I will need to use in my project:

|                          |  |
|--------------------------|--|
| <b>Student</b>           | I will use this class to represent each student. It will have attributes that represent the student's ID, name, email, Base64-encoded avatar, reference number, tutor and timetable. |
| <b>User</b>              | This class represents the least amount of data required to authenticate with MyLoreto. It will have only two attributes: username and password.                                      |
| <b>LandingPageParser</b> | This class will be used to parse a given Loreto URL. It will take an instance of a User object in its constructor, and contains many methods that will be used to parse HTML.        |

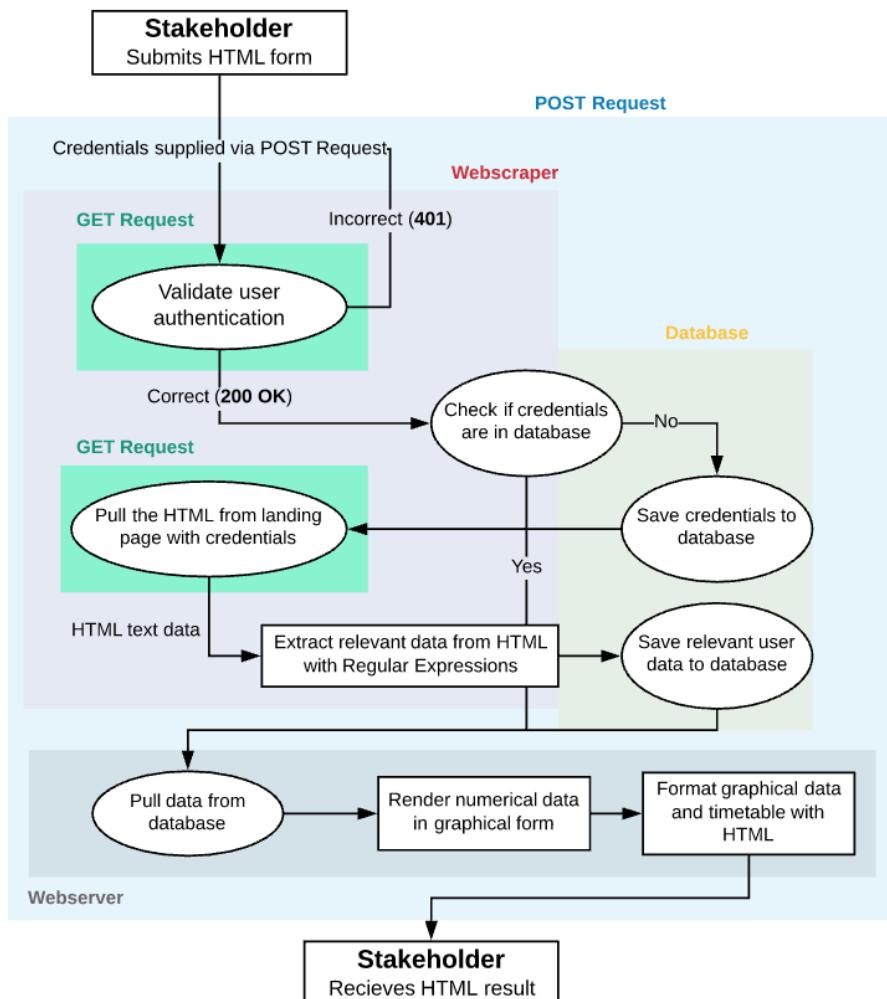
# Design

## Problem Decomposition

I have broken down my problem into a series of smaller problems, which are defined below:



I have also created a Data Flow Diagram which represents the smaller parts of my projects, and how they relate to the overall solution:



I have also written some heavily abstracted pseudocode for each aspect of the solution:

```
1 import requests
2
3 from models import Student
4
5
6 def generate_html_response(student):
7     return '''
8     <html>
9         <p>Hello {0}.</p>
10        <p>Your grades are: {1}</p>
11    </html>
12    '''.format(student.name, student.grades)
13
14
15 def parse_with_regex(response):
16     pattern = 'Grades: (.*)'
17     return regex.extract(pattern, response)
18
19
20 @call_on_post_form
21 def post(form):
22     authentication = form.username, form.password
23     response = requests.get(loreto_url, authentication)
24     if response.status_code != 200:
25         return 'Try again.' + form
26     else:
27         student = Student(authentication)
28         if not student.in_database():
29             data = parse_with_regex(response)
30             student.grades = data.grades
31             student.timetable = data.timetable
32             student.attendance = data.attendance
33             student.save_to_database()
34         result = generate_html_response(student)
35     return result
```

## Development

Before starting development, I have decided on the naming convention and code style I will use throughout the project. This will be based on [PEP8](#), a document stating the standard conventions for Python style guide.

In summary, I will be conforming to the following rules:

**All functions and methods I define must contain a docstring explaining its use:**

```
1 def add_two_numbers(a, b):
2     """Adds two numbers, a and b."""
3     return a + b
```

**All variables, methods and functions will use snake\_case, rather than mixedCase.  
Any classes which I define will use PascalCase.**

```
1 some_variable = 'foo'
2
3 class SomeClass:
4     def __init__(self):
5         self.bar = 'baz'
6
7     def what_is_bar(self):
8         """Returns the value of bar."""
9         return self.bar
10
11 some_instance = SomeClass()
```

As is usual with any code (but particularly Python), I will be naming my variables and methods with thought, using names that enhance readability.

I started by creating a simple script that will allow me to request and parse data from MyLoreto. I have decided to use [Kenneth Reitz's requests library](#) to handle HTTP requests and responses, and the [re builtin module](#) to parse the content with regular expressions.

I will be using Regular Expression patterns to catch certain substrings of a HTTP response. For example, from this snippet of HTML:

```
1 <dl class="dl-horizontal" style="font-size: 0.9em">
2     <dt>Reference: </dt>
3     <dd>S20170000420</dd>
4 </dl>
```

Reference number with the following regex pattern:  
**Reference: </dt>\s+<dd>([A-Z0-9]+)**

I am able to extract the

```

1 >>> import re
2 >>> html = """
3 ...     <dl class="dl-horizontal" style="font-size: 0.9em">
4 ...         <dt>Reference: </dt>
5 ...         <dd>S20170000420</dd>
6 ...     </dl>"""
7
8 >>> match = re.search(r'Reference: </dt>\s+<dd>([A-Z0-9]+)', html)
9 >>> match.group(1)
10 'S20170000420'

```

I have written up an example of my original webscraper prototype below:

```

1 import re
2
3 import requests
4
5
6 endpoint = "https://my.loreto.ac.uk/"
7
8 patterns = {
9     'name': b'fullName: "([A-Za-z ]+)"',
10    'username': b'username: "([A-Za-z0-9]+)"',
11    'avatar': b'base64,(.*?)>',
12    'reference_number': br'Reference: </dt>\s+<dd>([A-Z0-9]+)',
13    'tutor': b'Tutor: </dt> <dd> (.*) </dd>'
14 }
15
16
17 def get_student_data(*credentials: str) -> dict:
18     """Get data about a student, given their username and password."""
19     student_data = {}
20     landing_page = requests.get(endpoint, auth=credentials).content
21     for key, pattern in patterns.items():
22         student_data[key] = re.search(pattern, landing_page).group(1)
23     return student_data

```

An example of this module in use can be seen below:

```

1 >>> get_student_data('jerbob42', 'password')
2 {
3     'name': b'Jeremiah Boby',
4     'username': b'JerBob42',
5     'avatar': b'/9j/4AAQSkZJRgABAQEASABIAAD...',
6     'reference_number': 'S20170000420',
7     'tutor': 'Mr Bloggs'
8 }

```

## Potential Improvement

I could improve this function's performance by utilising Python's builtin **threading** module, starting a new **Thread** for each regex search that I start. This means that each regex search will not block the execution of the next, enhancing my concurrent throughput.

I then decided to create a function for validating user credentials. I decided to add this validation to prevent users encountering an unexpected error or crash immediately after providing incorrect credentials.

In the case of incorrect credentials, this function will return **False**, and an alert will be provided to the user.

```
1 def valid_user(*credentials: str) -> bool:
2     """Take a student's credentials and validate them."""
3     print('Validating user...')
4     return requests.get(
5         endpoint, auth=credentials
6     ).status_code == 200
```

This function checks that the HTTP response code for a request with the provided credentials is successful. If the credentials are incorrect, the server will respond with a [\*\*401 response code \(Unauthorized\)\*\*](#), instead of a [\*\*200 \(OK\)\*\*](#).

With some simple validation and webscraping done, I then decided to start work on the main webserver. To begin with, I used the `django-admin` command-line tool to generate the skeleton project required for the webserver, with [\*\*django-admin startproject structure\*\*](#).

This created a directory with the following structure:

```
structure/
    manage.py
    structure/
        __init__.py
        settings.py
        urls.py
        wsgi.py
```

In order to comply with **separation of concerns**, I have decided to implement my student dashboard as a separate module in the project, to be named **dashboard**. To do this, I used the following command: **python manage.py startapp dashboard**

I then removed all redundant files which were not required for the project. This resulted in the following directory structure:

```
structure/
    manage.py
    structure/
```

```
__init__.py
settings.py
urls.py
wsgi.py
dashboard/
    __init__.py
    models.py
    views.py
```

In order to be able to modify and update my database, I did some research into [the Django Documentation](#) on its builtin ORM. Though I initially decided to create my own ORM for this project, I have made the decision to use Django's builtin ORM to save time on the project.

If I had more time to enhance this project, I would have written a custom ORM for this purpose.

I described the following database models below in the `structure/dashboard/models.py` file like so:

```
1 """Define database models for use in the PostgreSQL server."""
2
3 from django.contrib.postgres import fields
4 from django.db import models
5
6
7 class Student(models.Model):
8     """Student model: contains information used by the dashboard."""
9
10    student_id = models.AutoField(primary_key=True)
11    name = models.CharField(max_length=20)
12    email = models.CharField(max_length=29)
13    avatar = models.CharField(max_length=20000)
14    reference_number = models.CharField(max_length=12)
15    tutor = models.CharField(max_length=50)
16    timetable = fields.JSONField()
17    attendance = fields.JSONField()
18    markbook = fields.JSONField()
19
20    class Meta:
21        # Give the database table a specific name.
22        db_table = "student"
23
24
25 class User(models.Model):
26     """User model: contains information specific to user authorization."""
27
28    user_id = models.AutoField(primary_key=True)
29    username = models.CharField(max_length=8, unique=True)
30    password = models.CharField(max_length=20)
31
32    class Meta:
33        # Give the database table a specific name.
34        db_table = "user"
```

As you can see, these

models do not match the ones that I created during the Design phase. This is because I needed to create a simplified version of the product for the initial version, so that I was able to test it as soon as possible.

Next, I added some simple views to run when a user requests pages. These were defined in the `structure/dashboard/views.py` file:

```
 1 from django.shortcuts import redirect, render
 2 from django.views.generic import TemplateView
 3
 4 from .models import User
 5 from .scraper import valid_user
 6
 7
 8 # Dictionary containing data about a user's session
 9 current_session: dict = {}
10
11
12 class HomePageView(TemplateView):
13     """HomePage view: default view of the landing page."""
14
15     def get(self, request, *args, **kwargs):
16         """Runs on a HTTP GET request."""
17         if current_session.get('authenticated'):
18             # If the current user is authenticated, render the dashboard
19             return render(request, "dashboard.html", context=current_session)
20         else:
21             # If the current user is not authenticated, redirect to the signin page
22             return redirect("/sign-in")
23
24
25 class SignInView(TemplateView):
26     """SignInView: view for the sign in page."""
27
28     def get(self, request, *args, **kwargs):
29         """Runs on a HTTP GET request."""
30         return render(request, "sign-in.html", context=None)
31
32     def post(self, request):
33         """Runs on a HTTP POST request."""
34         # Get the username and password from the POST request
35         username, password = (
36             request.POST.get('username'), request.POST.get('password')
37         )
38         # Create a User object with the username and password
39         user = User(username=username, password=password)
40         # Validate the user object
41         if valid_user(user.username, user.password):
42             # If the current user is valid, add the user to the session
43             current_session['user'] = user
44             # Set the "authenticated" flag to True
45             current_session['authenticated'] = True
46             # Redirect to the home page
47             return redirect('/')
48         else:
49             # If the current user is not valid, set the "authenticated" flag to False
50             current_session['authenticated'] = False
51             # If there has been no previous returns, render the signin page
52             return render(request, "sign-in.html", context=None)
```

As shown above, when a user's browser requests the home page, the following code is executed:

```
1 def get(self, request, *args, **kwargs):
2     """Runs on a HTTP GET request."""
3     if current_session.get('authenticated'):
4         # If the current user is authenticated, render the dashboard
5         return render(request, "dashboard.html", context=current_session)
6     else:
7         # If the current user is not authenticated, redirect to the signin page
8         return redirect("/sign-in")
```

In `dashboard/templates`, I also added `dashboard.html` and `sign-in.html`. `dashboard.html` is a simple example page for the CSS framework that I am using, [Bootstrap](#):

```
1 <html lang="en">
2   <head>
3     <meta charset="utf-8">
4     <meta name="viewport" content="width=device-width, initial-scale=1, shrink-to-fit=no">
5     <meta name="description" content="">
6     <meta name="author" content="">
7     <link rel="icon" href="https://getbootstrap.com/favicon.ico">
8
9   <title>Structure - Dashboard</title>
10
11  <!-- Bootstrap core CSS -->
12  <link href="../static/css/bootstrap.min.css" rel="stylesheet">
13
14  <!-- Custom styles for this template -->
15  <link href="../static/css/dashboard.css" rel="stylesheet">
16 </head>
17
18 <body>
19   <nav class="navbar navbar-dark sticky-top bg-dark flex-md-nowrap p-0">
20     <a class="navbar-brand col-sm-3 col-md-2 mr-0" href="dashboard.html#>">Structure Dashboard</a>
21     <input class="form-control form-control-dark w-100" type="text" placeholder="Search" aria-label="Search">
22     <ul class="navbar-nav px-3">
23       <li class="nav-item text-nowrap">
24         <a class="nav-link" href="dashboard.html#>">Sign out</a>
25       </li>
26     </ul>
27   </nav>
28
29   <div class="container-fluid">
30     <div class="row">
31       <nav class="col-md-2 d-none d-md-block bg-light sidebar">
32         <div class="sidebar-sticky">
33           <ul class="nav flex-column">
34             <li class="nav-item">
35               <a class="nav-link active" href="dashboard.html#>">
36                 <span data-feather="home"></span>
37                 Dashboard <span class="sr-only">(current)</span>
38               </a>
39             </li>
40             <li class="nav-item">
41               <a class="nav-link" href="dashboard.html#>">
42                 <span data-feather="file"></span>
43                 Summary
44               </a>
45             </li>
46             <li class="nav-item">
47               <a class="nav-link" href="dashboard.html#>">
48                 <span data-feather="shopping-cart"></span>
49                 Reminders
50               </a>
51             </li>
52             <li class="nav-item">
53               <a class="nav-link" href="dashboard.html#>">
54                 <span data-feather="users"></span>
55                 Timetable
56               </a>
57             </li>
58             <li class="nav-item">
59               <a class="nav-link" href="dashboard.html#>">
60                 <span data-feather="bar-chart-2"></span>
61                 Homework
62               </a>
63             </li>
64           </ul>
65     </body>
66 </html>
```

This resulted in the following web page:



`sign-in.html` contains just a login form:

```
1 <form action="", method="post", class="form-signin">
2   <h1 class="h3 mb-3 font-weight-normal">Sign in to Structure</h1>
3
4   <input id="inputUsername" name="username" placeholder="College username" required type="text">
5   <input id="inputPassword" name="password" placeholder="College password" required type="password">
6   <input class="btn" id="submit" name="submit" type="submit" value="Sign in">
7
8 </form>
```

This resulted in the following web page:

The screenshot shows a simple login form titled 'Sign in to Structure'. It has two input fields: 'College username' and 'College password', both with placeholder text and small user icon icons. Below the inputs is a large blue 'Sign in' button.

When the user enters their username and password, these parameters are sent as part of the POST request to the same endpoint, this is the default action done by [HTML forms](#).

As shown in the backend code for this view, a POST request triggers the following code:

```
1 def post(self, request):
2     """Runs on a HTTP POST request."""
3     # Get the username and password from the POST request
4     username, password = (
5         request.POST.get('username'), request.POST.get('password')
6     )
7     # Create a User object with the username and password
8     user = User(username=username, password=password)
9     # Validate the user object
10    if valid_user(user.username, user.password):
11        # If the current user is valid, add the user to the session
12        current_session['user'] = user
13        # Set the "authenticated" flag to True
14        current_session['authenticated'] = True
15        # Redirect to the home page
16        return redirect('/')
17    else:
18        # If the current user is not valid, set the "authenticated" flag to False:
19        current_session['authenticated'] = False
20    # If there has been no previous returns, render the signin page
21    return render(request, "sign-in.html", context=None)
```

## Potential Improvement

Currently, I am relying on the mutability of a global variable. This is unreliable, as it will only work for one user at a time. It is not scalable, nor lend itself to concurrent thinking. Instead, I could use Django's builtin `request.session` object, which is unique for each request, as it is an attribute of one parameter of the request method.

The code shown above does not store user details. Instead, it only stores a global `authenticated` flag.

In order to display student information, I will need to add more information to the scraper. In order to achieve this, I decided to create a custom class, which I will use solely to parse the MyLoreto landing page.

There are several reasons why I have decided to use a class for parsing the landing page:

- It allows me to create a combination of variables, methods and statements that I can reproduce whenever I need to.
- It allows me to import this class into any module I need, lending to \*\*modular thinking\*\*.
- It helps with separation of concerns, as all variables and methods will only have the scope of the instance it applies to.
- It creates more readable, maintainable code, as I write less and document more.

```

1 class LandingPageParser:
2     """Class for parsing the Loreto landing page."""
3
4     def _set_regex_group(self, page: bytes, name: str, pattern: bytes):
5         """Set the Student data to the first group of the given match."""
6         match = re.search(pattern, page)
7         if match is not None:
8             # If the Regex match was successful
9             self.student.__setattr__(name, match.group(1).decode())
10
11    def _get_short_timetable(self):
12        """Get the timetable for the current day."""
13        # List of lessons in the current day.
14        timetable = []
15        # Regular Expression pattern for a lesson.
16        pattern = re.compile(
17            b'TimetableEntry .*?>(?P<time>[0-9 -:]+)'
18            b'.*?(?P<room>[( )A-Z0-9]+)'
19            b'.*?(?P<teacher>[( )A-Za-z0-9 -]+) ',
20            re.DOTALL
21        )
22        # Iterate each attribute of every lesson in the timetable
23        for time, room, teacher in pattern.findall(self.page):
24            # Remove brackets from the lesson
25            room = room.strip(b'()')
26            subject, teacher = teacher.split(b' - ')
27            # Append this lesson dictionary to the list
28            timetable.append(
29                {
30                    'time': time.decode(),
31                    'room': room.decode(),
32                    'subject': subject.decode(),
33                    'teacher': teacher.decode()
34                }
35            )
36        # Set the student's timetable to this list
37        self.student.short_timetable = timetable

```

```

1  def _get_timetable(self):
2      """Get the timetable for the current week."""
3      # Create a dictionary of day -> lesson list
4      timetable = {
5          'Monday': [],
6          'Tuesday': [],
7          'Wednesday': [],
8          'Thursday': [],
9          'Friday': []
10     }
11     # Extract the HTML snippets for each day
12     day_meta = '<th>{0}(.*)(<th>|</script>)'
13     # Pattern that will be applied to the previous day snippet
14     lesson_meta = (
15         'Times">(?P<time>[0-9: -]*'
16         '<.+?Code">(?P<subject>[A-Za-z() -]*)'
17         '<.+?Staff"> *?(?P<teacher>[A-Za-z ]*) *?'
18         '<.+?right">(?P<room>[A-Z0-9]*)'
19     )
20     # Get the Monday of the current week
21     now = datetime.now()
22     start = now - timedelta(days=now.weekday())
23     start = start.strftime('%Y-%m-%d')
24     # Create a HTTP request to MyLoreto with the student's data
25     response: bytes = requests.post(
26         f'{endpoint}attendance/timetable/studentWeek',
27         data={'week': start, 'student_user_id': self.user_id},
28         auth=(self.user.username, self.user.password),
29         headers={'X-Requested-With': 'XMLHttpRequest'}
30     ).content
31     # Iterate each day and lesson list in the timetable
32     for day, lesson in timetable.items():
33         # Extract the content
34         content = re.search(
35             day_meta.format(day).encode(), response, re.DOTALL
36         ).group(1).decode()
37         for match in re.finditer(lesson_meta, content, re.DOTALL):
38             # Append each lesson to the timetable
39             lesson.append(match.groupdict())
40     # Set the parser instance's "timetable" attribute to this timetable
41     self.student.timetable = timetable
42
43 def __init__(self, user: User):
44     # Get the student's landing page with their credentials
45     self.page = requests.get(
46         endpoint, auth=(user.username, user.password)
47     ).content
48     # Set the student and user attributes persistently
49     self.user = user
50     self.student: Student = Student()
51     # Pull the user's UserId with a regular expression
52     user_id = re.search(
53         br'UserId = "(\d+)"', self.page
54     )
55     if user_id is not None:
56         # if the UserId was found:
57         self.user_id: int = int(user_id.group(1))
58     # Define regex patterns for use by core components of the student object.
59     patterns = {
60         'name': b'fullName: "([A-Za-z ]+)"',
61         'username': b'username: "([A-Za-z0-9]+)"',
62         'avatar': b'base64,(.*?)"',
63         'reference_number': br'Reference: </dt>\s+<dd>([A-Z0-9]+)',
64         'tutor': b'Tutor: </dt> <dd> (.*) </dd>'
65     }
66
67 def parse(self) -> Student:
68     """Parse the webpage content and return the resulting Student."""
69     self._get_short_timetable()
70     self._get_timetable()
71     for key, pattern in patterns.items():
72         self._set_regex_group(self.page, key, pattern)
73     self.student.email = '{0}@student.loreto.ac.uk'.format(
74         self.student.username
75     )
76     return self.student

```

I am using regular expressions to extract parts of the HTML page. I use regular expressions for the following reasons:

- Regular expressions enable me to reduce my code base drastically. If I were to use simple string manipulation to extract the data I need, I would have significantly more unmaintainable code.
- Regular expressions are much easier to read than many different string manipulations.
- Regular expression patterns are reusable once compiled.

However, I soon realised that each of these regular expression searches are blocking the thread of execution. Each extraction prevents the next one from being carried out. In order to apply concurrent thinking to this problem, I decided to use the [threading module](#), part of [Python's standard library](#), to run each regex search in a separate thread:

```
1  def __init__(self, user: User):
2      self.page = requests.get(
3          endpoint, auth=(user.username, user.password)
4      ).content
5      self.user = user
6      # Set of threads used in the parser.
7      # I am using a set to ensure that threads are not duplicated.
8      self.threads: set = set()
9      self.student: Student = Student()
10     user_id = re.search(
11         br'UserId = "(\d+)"', self.page
12     )
13     if user_id is not None:
14         self.user_id: int = int(user_id.group(1))
15     patterns = {
16         'name': b'fullName: "([A-Za-z ]+)"',
17         'username': b'username: "([A-Za-z0-9]+)"',
18         'avatar': b'base64,(.*?)">',
19         'reference_number': br'Reference: </dt>\s+<dd>([A-Z0-9]+)',
20         'tutor': b'Tutor: </dt> <dd> (.*) </dd>'
21     }
22     # Create a thread for getting the student's daily timetable, add it to threads
23     self.threads.add(Thread(target=self._get_short_timetable))
24     # Create a thread for getting the student's full timetable, add it to threads
25     self.threads.add(Thread(target=self._get_timetable))
26     # The threads created so far have not been started yet
27     for key, pattern in patterns.items():
28         # For every regex pattern in the patterns I need to extract:
29         self.threads.add(
30             # Create a thread to extract that pattern and add it to threads
31             Thread(
32                 target=self._set_regex_group,
33                 args=(self.page, key, pattern)
34             )
35         )
36
37     def parse(self) -> Student:
38         """Parse the webpage content and return the resulting Student."""
39         # The __init__ method does not start the threads.
40         # Instead, this parse method will start each thread.
41         for thread in self.threads:
42             # Start each thread concurrently
43             thread.start()
44         for thread in self.threads:
45             # Close each thread once they have completed
46             thread.join()
47         self.student.email = '{0}@student.loreto.ac.uk'.format(
48             self.student.username
49         )
50     return self.student
```

There are several reasons why using a threaded solution is an improvement:

- Decreases the time needed for initialising the student object
- Increases maximum code throughput, as multiple methods of my code will be running at one time
- This solution therefore lends itself to concurrent thinking well.