

## Assignment Specifications

Aims	<ul style="list-style-type: none"> <li>• Enable students to analyse and employ appropriate Natural Language Processing (NLP) techniques to solve problems and design intelligent systems.</li> <li>• Enable students to use relevant tools and technology, such as Python programming, to develop intelligent computer programs.</li> </ul>
Learning Outcomes Assessed	<ul style="list-style-type: none"> <li>• Discuss the concepts and importance of natural language in real-world scenarios (C2, PLO1).</li> <li>• Practice the NLP algorithm using the features and capabilities of NLP software (P3, PLO3).</li> </ul>
Tutor	Noor Aida binti Husaini Kathleen Tan Swee Neo Nicholas Lim Jun Jie
Outline of Problem	For your assignment, you should critically evaluate current technologies, and then propose a project for the selected NLP topic, and implement an NLP solution to solve the problem in the proposed project by using Python/any suitable tools/ Framework.
Details	<ul style="list-style-type: none"> <li>• The Assignment consists of <b>TWO (2)</b> parts, please refer to <b>Project Details: Part 1</b> and <b>Project Details: Part 2</b>. Both parts are related. For Part 2, you are required to develop the program using Python.</li> <li>• Group yourself into a team of <b>TWO OR THREE (2-3)</b> members. You are given <b>TWO (2) OPTIONS</b> for research areas as listed below. Your group shall select one of them for this assignment purpose.</li> <li>• Please note that the following are the basic requirements of the assignment. Fulfilling the requirements may lead you to an <b>Average</b> or <b>Good</b> grade. Achieving an <b>Excellent</b> grade requires <b>extra effort</b>, such as learning new skills, introducing new ideas, and complex NLP algorithms, demonstrating the ability to process big data, and/or producing excellent reports with working prototypes.</li> <li>• Your group is expected to produce ideas that originated from the members, but not to take the work or an idea of someone else (including the Web) and pass it off as your own. <b>NO GROUP IS ALLOWED to share the same idea</b>, i.e., each group must propose a unique title or solution.</li> </ul>
Type	<p><b>Type A:</b> Develop one approach/system, but each member takes on different phases (pre-processing, main processing, post-processing, user interface (not necessarily GUI)).</p> <p><b>Type B:</b> The group members developed different approaches to achieve the same goal and then compared them.</p>

Title	<p><b>1. Interactive Chat Application for Conversations with PDF or Document Content Using a Language Model (LLM)</b></p> <p>You are required to design and implement a chat application that utilises a Language Model (LLM), such as GPT-3.5, to facilitate conversations between users and the content of a PDF or a document. The application should have the following functionalities:</p> <ol style="list-style-type: none"> <li>a. <b>User Interface:</b> Create a user-friendly interface that allows users to upload a PDF or a document file.</li> <li>b. <b>Text Extraction:</b> Implement a mechanism to extract the text content from the uploaded file.</li> <li>c. <b>Language Model Integration:</b> Integrate the LLM into the chat application to process the extracted text and generate meaningful responses based on user queries or inputs.</li> <li>d. <b>Chat Functionality:</b> Develop the chat functionality that enables users to have interactive conversations with the content of the PDF or document. Users should be able to ask questions, provide inputs, and receive relevant responses based on the extracted text, with the LLM providing natural and contextually appropriate replies.</li> <li>e. <b>User Experience:</b> Ensure that the chat application provides a smooth and intuitive user experience. Consider factors such as response time, error handling, and overall usability, while utilising the capabilities of the LLM to enhance user interactions.</li> <li>f. <b>Testing and Validation:</b> Conduct thorough testing of the application to ensure its functionality and accuracy. Validate the responses generated by the LLM against the original content of the PDF or document.</li> <li>g. <b>Documentation:</b> Provide clear and comprehensive documentation that includes instructions on how to use the application and details about the implementation, including the technologies and libraries used. Explain the integration of the LLM and how it contributes to the chat functionality.</li> </ol> <p><b>Note:</b> You can choose any programming language and framework of your preference for this assignment. Make sure to justify your choices in the documentation, including the selection and configuration of the LLM.</p>
-------	---

## 2. Sentiment Analysis on Out-Of-Vocabulary (OOV) Malaysia Rojak Language

You are required to design and implement a sentiment analysis system that can accurately analyse the sentiment of text written in out-of-vocabulary (OOV) Malaysia Rojak language. The application should have the following functionalities:

- a. **Data Collection:** Collect a sufficient amount of Malaysia Rojak language data from various sources, including social media, online forums, and news articles. This dataset will serve as the training data for your sentiment analysis model.
- b. **Preprocessing:** Develop preprocessing techniques specifically tailored for the Malaysia Rojak language. This may include tokenisation, normalisation, and handling of linguistic nuances specific to the language.
- c. **Model Selection:** Choose an appropriate sentiment analysis model or algorithm that can effectively handle the OOV nature of the Malaysia Rojak language. Consider techniques such as transfer learning, domain adaptation, or multilingual models.
- d. **Training and Evaluation:** Train your selected model using the collected Malaysia Rojak language dataset. Split the dataset into training and evaluation sets to assess the performance of your model accurately. Use appropriate evaluation metrics such as Human Level Performance (HLP) with/without survey.
- e. **Testing:** Evaluate the performance of your trained model on a separate test dataset consisting of Malaysia Rojak language text samples. Measure the model's accuracy in predicting sentiment labels for the given texts.
- f. **Error Analysis and Improvements:** Conduct an error analysis to identify common mistakes or challenges faced by your model in analysing the sentiment of Malaysia Rojak language. Propose improvements or refinements to enhance the model's accuracy and robustness.
- g. **Documentation:** Provide clear and comprehensive documentation that includes details about the data collection process, preprocessing techniques, model selection, training methodology, evaluation results, error analysis, and proposed improvements.

**Note:** You can choose any programming language and existing libraries or frameworks for sentiment analysis and text processing. Justify your choices in the documentation, specifically regarding the handling of OOV Malaysia Rojak language.

Submission Deadlines	<p>Submit Prototype Source Code and Documentation by:  <b>Part 1: Documentation: Week 12, Sunday, before 11.59 pm</b>  <b>Part 2: Source Code: Week 12, Sunday, before 11.59 pm</b></p> <p><i>Please submit it to Google Classroom. Late submissions will be penalised. A demo session to present the prototype is required.</i></p>
Contribution	<p><b>Part 1: Documentation - 40%</b>  <b>Part 2: Source Code - 60%</b></p> <p>Part 1: Each member must submit their documentation through Google Classroom.  Part 2: Only one member from each group should submit the source code on behalf of the group through Google Classroom.</p>
Academic Integrity and Plagiarism	<p>Your work must have originality, i.e., <b>DO NOT</b> copy or refer to other group(s). You may only work with your team member(s) to produce the solution of this assignment. You <b>MUST NOT</b> share with nor refer to any part of the assignment (including the code) of anyone else except your team member(s) and your tutor.</p> <p>Before submitting your assignment, please ensure that you have complied with <b>TAR UMT Plagiarism Policy</b>. Any cheating attempt to cheat, plagiarism, collusion, and any other attempts to gain an unfair advantage in assessment will cause the students concerned to be penalised.</p> <p><b>IMPORTANT:</b> Students found to be dishonest are liable to disciplinary action.</p>
Late Submission	<p>Late submission without a valid reason will <b>NOT</b> be tolerated. For late submissions, there will be a reduction of total marks:</p> <ul style="list-style-type: none"> <li>• Late 1 to 3 days after the deadline of submission: minus 10 marks</li> <li>• Late 4 to 7 days after the deadline of submission: minus 20 marks</li> <li>• Late more than 7 days after the deadline of submission: 0 marks</li> </ul> <p>In certain circumstances, a student may be allowed to submit the assignment late with a valid reason. S/he must contact the tutor at least one week before the assignment is due. The tutor will evaluate whether the circumstance warrants submitting the assignment late, but no guarantee that the students will not be penalised.</p> <p>Failing to submit the reports &amp; code will lead to failure of the coursework.</p> <p><b>Note:</b> For Extenuating Mitigating Circumstances (EMC) cases, students need to apply the Leave of Application from the FOCS office as a record where proof of MCs, Death Certificate. SPM Examination, MUET Examination, etc. is given.</p>

### Project Details: Part 1 - Proposal Writing

Introduction	Your task for this part of the assignment is to <b>identify a problem</b> for the selected NLP topic, perform a <b>literature review</b> , and propose your respective <b>NLP solution(s)</b> that helps in solving the problem in the proposed project.
What to hand in?	Each student is expected to <b>produce individual work</b> regarding the specified items and must not repeat what others have done. The report shall be completed using the given Google Doc template in the Google Classroom.
Format for Deliverable	For your proposal, you are required to complete the report using the Google Docs template available in the Google Classroom. You must <b>turn in the report</b> to Google Classroom before the stipulated deadline. Otherwise, the late submission penalty will be applied.

### Project Details: Part 2 - Prototype Development

Introduction	Your task for this assignment is to implement the NLP solution using Python or any other relevant tools, perform testing and evaluation on the system, and finally present the work to your tutor.
What to hand in?	Submit all your source code to the Google Classroom submission page. Each team member is required to do a presentation of the work and Q&A session on Week 13 or Week 14 based on the arrangement with the tutor.
Format for Deliverable	<b>Compress</b> the entire source code using zip format. Submit the source code individually.