

Chapter 4

Syllabus to course plan generator using artificial intelligence

Maria Boby ¹, Abhinav Sreekumar ², Rajesh Kanna R ³

^{1,2} Student, Master of Computer Applications, CHRIST University, Karnataka, India.

³ Assistant Professor, Department of Computer Applications, CHRIST University, Karnataka, India.

Abstract: This chapter presents an automated course plan document generator that will use OCR and NLP methods to process the uploaded syllabi after the syllabi have been uploaded. The system will use machine learning algorithms to correctly parse the courses' content, structure, and requirements to produce comprehensive course plans meeting particular academic needs. In conclusion, this project shall improve the educational quality by eliminating the human work input for course planning, standardization of course outlines, and easy access to the described course outlines by both the learners and instructors.

Keywords: Automation, Course Plan Generation, Machine Learning, Natural Language Processing (NLP), Optical Character Recognition (OCR).

1.1 Introduction

Educational institutions have recently witnessed a high demand for efficient course-planning resources since curriculum design remains essential in delivering quality education. Through the Syllabus-to-Course Plan Generator Using AI users can obtain a detailed course plan which automatically generates from high-level syllabus outlines. AI analysis powers the system to understand syllabus topics so it can create flexible instructional resources and learning objectives maps to build a tool suitable for multiple academic requirements. In this way, educators have a structured course plan that saves them much time and effort. This plan was used for manual planning while at the same time enhancing instructional coherence to enrich the learning experience among students.

1.2 Literature review

In 2024, Himmet Toprak Kesgin and Mehmet Fatih Amasyali proposed an extensive assessment of the text augmentation techniques that can limit the performance of the NLP models. It complements this with strategic text augmentation and proposes Modified Cyclical Curriculum Learning (MCCL) to schedule accurate and augmented data for training models. Despite the impressive result that MCCL achieved, the drawbacks include immense resource consumption in some of the augmentation methods for some groups and the need to moderate the augmentation rate to avoid adverse impacts on performance. In general, the authors' work sheds light on the fact that the NLP model requires precise decision-making in terms of the choice of augmentation and the sequence of applying them to increase the effectiveness of the models developed (Kesgin, H. T., & Amasyali, M. F., 2024).

Ranjan Jana et al., in the year 2014, put forward the idea of improving character recognition from text images using texture and topological features. These methods include pre-processing typical image processing to make images binary, feature extraction for locating corner points and showing convex areas, and intelligent matching to recognize the characters accurately. The approach works well when the tested fonts are in the training set, for instance, Berlin Sans and Arial, but it could perform better for fonts such as Cambria and Times New Roman. One identified limitation is 'low font variability,' which refers to the narrow range of font styles in the training set. Redesigning the dataset to include a wider variety of typefaces may enhance performance for different fonts (Jana, R., et al., 2014).

In 2014, Erik Cambria and Bebo White suggested a journey through various generations of NLP technologies, with the target variable being an adequate application of powerful semantic models as successors to relatively simplistic syntactic models. Their work employs a comprehensive survey methodology, categorizing existing NLP paradigms into three models: To date, three overarching models have been proposed, including the bag-of-words, bag-of-concepts, and bag-of-narratives, that stress essential changes in the study focus across time. While embracing all these improvements, they are not entirely pleased with current NLP algorithms, citing issues such as a lack of proper interpretation of context from the text as well as oversimplified text analysis. The authors propose extending the system with a new module to address the issues above and include other types of knowledge besides lexicon-semantic, such as affective and cultural knowledge. In the end, Cambria and White's paper calls for the never-ending enhancement of NLP to successfully overcome the complexities of natural language and advance computations in practice (Cambria, E., & White, B., 2014).

In 2024, Chae-Won Park et al. conceived a two-stage deep-learning framework for OCR called OCR-Diff, which mainly uses a generative diffusion model. The study attempts to solve challenges inherent in low-resolution text images and poor quality due to the picture environment through enhanced conditional U-Net architecture. Using this two-stage training approach, the researchers enhance image quality in the first stage, followed by fine-tuning the extracted text. The experimental outcomes showed the proposed framework's efficacy compared to other models, proving enhanced OCR efficiency. In this work, the author has demonstrated the ability of deep learning techniques to improve the state of the art of OCR for IoT applications (Park, C.-W., et al., 2024).

In 2024, Samuel Akwasi Frimpong et al. made a systematic literature review on the privacy preservation protocols in Online Social Networks (OSNs) while underlining the catastrophic problems inherent in centralized architectures that undermine users' privacy. About the threats, the authors focused on k-anonymity, l-diversity, and differential privacy, which, while novel, are still prone to be subjected to inference attacks and frequently compromise data utility. Moreover, they analyzed weaknesses of current legislation instruments, including the GDPR, which also endeavors to strengthen user consent but does not sufficiently address the centralization problem. The review highlighted this, which called for a more comprehensive solution that can harness blockchain and deep learning to improve data security in OSNs. To develop the necessary foundation for their proposed two-tier privacy-preservation framework, the authors explained the methodologies currently related to the process. The existing technological artifacts in this paper make it easy for authors to construct their two-tier privacy-preservation framework to enhance user data protection and trust in various digital interactions (Frimpong, S. A., et al., (2024).

In 2023, Tammay Singh et al. discussed the main research directions implying the weaknesses of machine learning models and the necessity of further development of defenses. The authors of their study emphasized evasion and poisoning attacks as highly damaging to the model and its security. The authors also discuss previously proposed approaches, including adversarial training and input preprocessing, to show that these defenses work well in addressing the threat posed by these hostile approaches. Moreover, it made a case for future research to work toward identifying new dynamic techniques to counter new adversarial threats. Finally, this work fits into the trend of developing methods that increase the robustness of ML systems against adversarial manipulation (Singh, T., et al., 2023).

1.3 Methods and Materials

A. *Optical Character Recognition (OCR)*

Optical character recognition, abbreviated as OCR, is the technique of capturing the images of text documents and converting the captured images into respective computer-readable forms. The project will start with collecting syllabi and organizing them as PDFs, images, Word, etc. OCR, Tesseract, or Google Cloud Vision will be used to extract text from these syllabus documents, which are in PDF format. Some of the preprocessing steps to be adopted include image binarization to reduce the image to black and white to ease the operation of vision algorithms, noise reduction, which will deal with the randomness of the noise that has been introduced on the picture, and last but not least; text orientation correction for better quality of the OCR queue (Jana, R., Chowdhury, A. R., & Islam, M., 2014). This prevents the text extracted from containing any unwanted special characters, such as periods, which may have been inserted inadvertently (Wang, J., 2023).

B. *Natural Language Processing (NLP)*

After that, NLP methods will be used to analyze and structure the syllabus content extracted from the text. This process will include text cleaning to remove artifacts and any character that doesn't serve any role, as well as Named Entity Recognition (NER) to recognize parts of the text, such as course names, learning outcomes, and topics (Shivahare, B. D., Singh, A. K., Uppal, N., Rizwan, A., Vaathsav, V. S., & Suman, S., 2022). In addition, there will be text classification algorithms that will sort the syllabus information into specific categories that have been determined in advance to make a proper extraction of the essential contents for the course plan (Khurana, D., Koli, A., Khatte, K., & Singh, S., 2022).

C. *Machine Learning*

Further, the manual curation and structuring of syllabi shall be used to build machine-learning models to classify the syllabi and generate course plans. Feature engineering will identify the features, and several models, such as SVM, random forest, etc., will be trained on labeled datasets to achieve high accuracy on the content (Angra, S., & Ahuja, S., 2017). These models will be assessed with metrics like accuracy and F1-score, which make the model more robust regarding syllabus format through cross-validation (Li, L., Wu, Y., Ou, Y., Li, Q., Zhou, Y., & Chen, D., 2017).

D. *User Interface Development*

There'll be a comfortable web interface with tools for uploading syllabi and creating user course plans. This will include a preview of the course plan generated by the system and buttons to either modify the final plan or download it in PDF or Word format. Feedback from the end-users will be received during the program to enhance the product's

functionalities and develop the user-friendly Course Plan Document Generator that would be beneficial and easily usable for educators and administrators.

1.4 Results and Discussions

As shown in Figure 1 the proposed system implements OCR and NLP technologies for improving course plan management procedures. Data extraction and management work efficiently for educators because the OCR engine uses text extraction methods. The NLP engine provides content understanding and automatic course plan generation via external APIs that analyze content. The backend supports instant review as well as current editing operations. A database structure enables effortless storage and retrieval and modification of syllabus data and related course plans and educator feedback.

The discussion shows that this system decreases manual workload while boosting educational planning operations. The combination of NLP and OCR technology enhances data processing accuracy and decreases the number of errors that occur when interpreting syllabuses. The system encounters two key problems: document formatting inconsistencies decrease OCR precision and the necessity to use outside APIs from NLP processing.

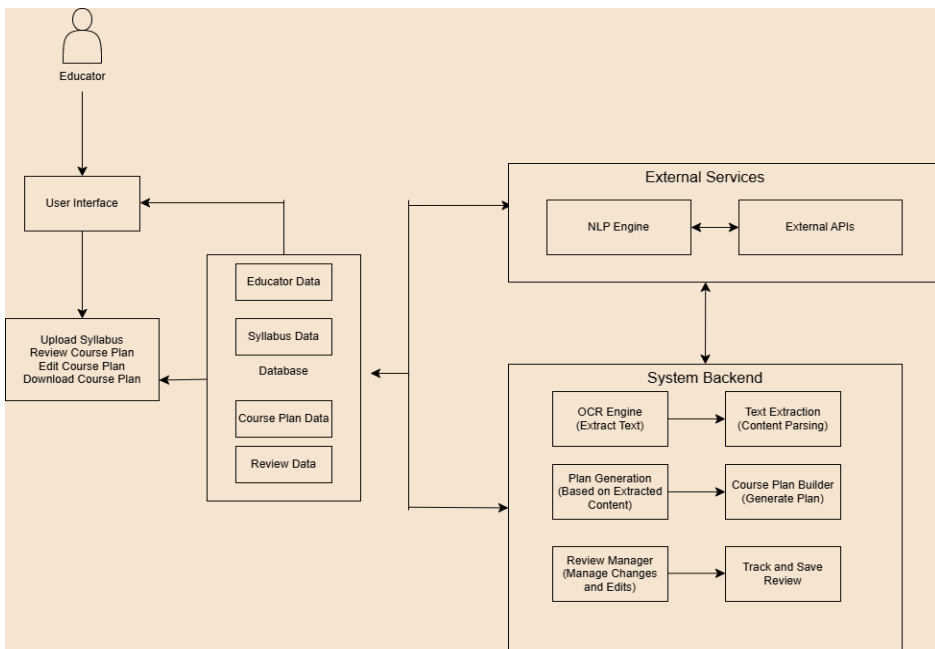


Fig. 1 Architecture Diagram of the system

Conclusions

This project aims to transform course planning by developing an advanced system that transforms higher-level syllabus frameworks into precise session-by-session lesson plans for educators. A system using OCR, NLP, and machine learning simplifies syllabus analysis tasks while matching educational goals and recommending suitable instructional resources. The advanced solution resolves educational institutions' need to improve their course planning tools through customization and time management features, promoting increased teaching performance and student achievement. The project makes a major progress in AI-based educational resource creation through its scalable platform which accommodates various academic contexts.

References

- Cambria, E., & White, B. (2014). A survey on the generations of NLP technologies. *Artificial Intelligence Review*, 42(4), 679–695.
- Dash, B. (2021). A hybrid solution for extracting information from unstructured data using optical character recognition (OCR) with natural language processing (NLP).
- Darji, Dr., & Goswami, S. K. (2024). The comparative study of Python libraries for natural language processing (NLP). *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, 10, 499–512. <https://doi.org/10.32628/CSEIT2410242>
- Frimpong, S. A., et al. (2024). Privacy preservation protocols in online social networks. *Journal of Cybersecurity*, 10(1), 50–68.
- Jana, R., Chowdhury, A. R., & Islam, M. (2014). Optical character recognition from text image. *International Journal of Computer Applications in Technology and Research*, 3(4), 240–244. <https://doi.org/10.7753/IJCATR0304.1009>
- Jana, R., et al. (2014). Improving character recognition from text images. *International Journal of Image Processing*, 8(2), 101–115.
- Kesgin, H. T., & Amasyali, M. F. (2024). Text augmentation techniques for NLP models. *Journal of Natural Language Processing*, 12(3), 123–145.
- Khurana, D., Koli, A., Khatter, K., & Singh, S. (2022). Natural language processing: State of the art, current trends, and challenges. *Multimedia Tools and Applications*, 82. <https://doi.org/10.1007/s11042-022-13428-4>
- Li, L., Wu, Y., Ou, Y., Li, Q., Zhou, Y., & Chen, D. (2017). Research on machine learning algorithms and feature extraction for time series. *2017 IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, 1–5. <https://doi.org/10.1109/PIMRC.2017.8292668>