



# AI-powered semantic search and OCR-enabled e-learning platform

**Dipali Baviskar**

Assistant Professor (School of Computer Engineering  
and Technology)  
MIT World Peace University  
Pune, Maharashtra, India

**Kedar Patil**

B. Tech Computer Science and Engineering  
MIT World Peace University  
Pune, Maharashtra, India

**Parth Gabhane**

B. Tech Computer Science and Engineering  
MIT World Peace University  
Pune, Maharashtra, India

**Tejas Naphade**

B. Tech Computer Science and Engineering  
MIT World Peace University  
Pune, Maharashtra, India

**Anjali Garje**

B. Tech Computer Science and Engineering  
MIT World Peace University  
Pune, Maharashtra, India

*Abstract—*

The paper proposes an AI-powered semantic search and OCR-enabled e-learning platform that aims to improve knowledge acquisition and retention. The system uses natural language processing techniques to enhance search queries and provide relevant results based on the user's preferences and interests. Optical character recognition (OCR) technology is employed to extract text from images and make them searchable, enabling users to access and learn from visual content. The platform also incorporates adaptive learning algorithms that personalize the learning experience and facilitate long-term knowledge retention. The proposed system is expected to improve the quality of online education by providing an intuitive, interactive, and efficient learning environment that enhances the learning outcomes of students.

**Keywords:** Optical character recognition (OCR) technology, adaptive learning algorithms

## I. INTRODUCTION

The rapid advancement of technology has transformed the traditional education system into a digital and online one. E-learning platforms have become increasingly popular due to their accessibility and flexibility, enabling users to learn from anywhere, at any time. However, the effectiveness of these platforms in facilitating knowledge acquisition and retention is limited. Students and researchers often struggle to find relevant information and face difficulty in accessing information-rich

visuals such as images, diagrams, and graphs. Moreover, traditional e-learning systems often lack personalization, resulting in suboptimal learning experiences.

To address these challenges, this paper proposes an AI-powered semantic search and OCR-enabled e-learning platform designed to enhance knowledge acquisition and retention. The system employs natural language processing techniques to enhance search queries and provide relevant results based on the user's research topic and preferences. Optical character recognition (OCR) technology is used to extract text from images, facilitating access to information-rich visuals that can aid in writing research papers. Additionally, the platform incorporates adaptive learning algorithms that personalize the learning experience, helping users acquire and retain knowledge in a more effective manner.

The proposed system has the potential to transform the e-learning landscape, providing an intuitive and interactive learning environment that enhances knowledge acquisition and retention. The remainder of this paper discusses the literature review, the proposed system, the experimental evaluation, and the conclusion, highlighting the potential of the proposed platform and its implications for the e-learning landscape.

## II. A.PROBLEM DEFINITION

An AI-Powered Semantic Search and OCR-Enabled E-Learning Platform for Enhanced Knowledge Acquisition and Retention.

- In the context of the topic "An AI-Powered Semantic Search and OCR-Enabled E-Learning Platform for Enhanced Knowledge Acquisition and Retention," there are several potential gaps in the literature survey that could be addressed.
- Lack of integration of OCR and Semantic Web technologies in the context of e-learning.
- Traditional information retrieval methods lack efficiency and accuracy in retrieving information from multimedia data

### III. OBJECTIVES

1. Analyzing the internal workings of our proposed system using individual and group machine learning techniques like Improve knowledge acquisition: The platform should provide students with relevant and accurate information that will enhance their understanding of the subject matter. The use of semantic search and OCR technology will ensure that the platform delivers the right content to the student.

Personalize learning: The platform should be able to personalize the learning experience for each student based on their learning style, preferences, and progress. The use of AI technology will enable the platform to create a personalized learning path for each student.

Enhance knowledge retention: The platform should be designed in such a way that it enhances knowledge retention for students. The use of multimedia content, interactive quizzes, and other learning tools will help students retain information more effectively.

Foster collaboration: The platform should promote collaboration between students and teachers. The use of social learning tools and discussion forums will enable students to share ideas and learn from each other.

### IV. LITERATURE SURVEY

[1] Personalized e-learning system based on machine learning and semantic search" by Wang and Li (2020). This [1] Advantages and Constraints of a Hybrid Model K-12 E-Learning Assistant Chatbot offers advantages in terms of speed and semantic understanding. It combines Global Vector and QANet models, utilizing CNN for improved efficiency. The SQuAD dataset was used for training. Future work includes developing a voice-based chatbot. The chatbot outperformed teachers based on certain statistics, but it has limitations in its programmed knowledge. Integration with an E-Learning platform enhances the learning experience.

[2] Evaluating E-learning systems success\_ An empirical study paper conducted an empirical study evaluating e-learning system success. The model used focused on user experience and quality assurance. The study found that the proposed model explained significant percentages of perceived satisfaction, usefulness, use, and benefits of e-learning success. The analysis utilized statistical methods such as SRMR and Goodness-of-Fit. Future work aims to improve the model's performance.

[3] In their paper, "Emerging themes in e-learning: A review from the stakeholders' perspective," the authors conduct a comprehensive review of e-learning methods in organizations. They analyze the advantages, disadvantages, challenges, critical success factors, theories, and models associated with e-learning from the viewpoint of stakeholders. By examining 138 articles published between 2000 and 2018, the authors highlight the impact of social, technological, and organizational factors on e-learning stakeholders. They emphasize that stakeholders cannot operate in isolation and must adapt to changing trends in technology and the learning environment. The constant challenge faced by stakeholders is keeping up with rapid technological advancements to establish an effective e-learning environment. The paper's keywords include interactive learning

environments, human-computer interface, computer-mediated communication, lifelong learning, and media in education.

[4] "Web Platform for E-Learning" presents a web platform designed to facilitate e-learning experiences. The authors focus on describing the working and functionality of the platform. The platform is intended to provide a user-friendly and efficient environment for online learning. It likely includes features such as content management, interactive learning materials, assessment tools, and communication channels. The paper may discuss the design choices, technologies, and frameworks used to develop the platform. The working of the platform might involve user registration and authentication processes to ensure secure access. Learners can access various learning resources, including multimedia content, presentations, and documents. The platform may also incorporate interactive features like quizzes, assignments, and discussion forums to enhance student engagement and collaboration. The discussion in the paper likely elaborates on the benefits and potential applications of the web platform for e-learning. It may highlight how the platform addresses the challenges and limitations of traditional classroom-based learning, such as accessibility, flexibility, and personalized learning experiences. The authors may also discuss the platform's potential to support different educational levels, subjects, and teaching methodologies.

[5] The paper discusses the research conducted on the design and implementation of the educational informatization platform. It may describe the features, functionalities, and technologies used in the platform's development. The platform is likely aimed at enhancing educational processes through the integration of digital technologies and e-learning tools. The paper also discusses the benefits and advantages of the educational information platform in terms of facilitating online learning, improving accessibility to educational resources, and enhancing student engagement. It may highlight the platform's ability to support various forms of content delivery, such as multimedia materials, interactive modules, and assessments. Furthermore, the paper may touch upon the integration of emerging technologies, such as artificial intelligence, data analytics, or virtual reality, to enhance the educational experience and provide personalized learning pathways. It might also discuss the platform's scalability and adaptability to different educational contexts and levels.

[6] The paper titled "E-learning Platforms Security Issues and Vulnerability Analysis" by M. Bhatia and J. K. Maitra explores the security concerns and vulnerability analysis of e-learning platforms. The authors emphasize the importance of securing web applications and e-learning platforms due to their popularity, user-friendliness, and susceptibility to attacks. They evaluate open-source learning management systems using web vulnerability scanners to identify security issues and vulnerabilities. The research aims to raise awareness among the educational community about existing vulnerabilities in e-learning platforms and provides insights for improving their security. The paper proposes a security model that can be adapted to different e-learning platforms or web platforms in general. Overall, the study contributes to addressing the neglected topic of security in e-learning platforms and offers valuable recommendations for enhancing their security.

[7] The study aimed to investigate the e-learning experiences of learners by analyzing their interactions with the e-learning system. The researchers used data preprocessing techniques and confirmatory factor analysis (CFA) to explore the relationship between different learning activities and system components based on log data. Discriminant validity analysis was also conducted to validate the constructs. The findings revealed that learners' e-learning interactions could be categorized into multiple components, which collectively formed the overall e-

learning experience. There was a significant positive relationship between these components. The study identified discussion, hypertext, assessment, content package, and video interactions as sequential components contributing to the holistic e-learning experience. Instructional discussions were found to be the most important component, while the content package had less significance. Moreover, students who actively participated in discussion forums demonstrated better performance and achieved higher grades. Based on these results, the study emphasized the importance of incorporating discussion forums and formative assessment tasks in e-learning course designs. It highlighted the potential benefits of collaborative learning and self-assessment for enhancing the overall e-learning experience. The researchers recommended that future e-learning courses should consider including these elements to improve learner engagement and outcomes.

[8] The paper introduces an Artificial Intelligence (AI)-based student assessment and recommendation system designed for e-learning in the context of big data. The system comprises several components, including score estimation, clustering, prediction, and recommendation. To access the system, students need to authenticate themselves using valid credentials. The system utilizes recurrent neural networks (RNN) to estimate students' scores, followed by the application of a clustering algorithm called DBSCAN (Density-Based Spatial Clustering of Applications with Noise) using Mahalanobis clustering for grouping students based on their marks. Additionally, a reinforcement learning algorithm known as R-SARSA (Residual SARSA) is integrated into the system for evaluation purposes. The paper suggests that the system can be further integrated with other e-learning platforms and tools to provide a more comprehensive and seamless learning experience. This integration may involve the development of advanced data processing and integration techniques to combine data from multiple sources and platforms. The results of the study indicate that the system performed well in terms of true positives, false positives, true negatives, false negatives, precision, recall, and accuracy. Based on the scores generated using the implemented algorithms, students were grouped into different categories and evaluated accordingly.

[9] In the paper titled "Toward Selection of a Trustworthy and Efficient-Learning Platform" by B. Alojaiman, the authors aimed to evaluate e-learning platforms based on trustworthiness and efficiency. The methodology involved a multi-rater analysis with community input to gather feedback on the platforms. The authors used a hybrid fuzzy AHP-TOPSIS technique to assess the platforms, taking into account factors such as content quality and system quality. One limitation of the study is that the proposed framework has not been extensively tested on a large sample of e-learning platforms, which may limit its generalizability. Additionally, the study assumes that the proposed framework can be universally applied to all e-learning contexts and stakeholders, which may not hold true in practice. The results of the study indicate that both content quality and system quality have a positive impact on e-learning platforms. Furthermore, the study identifies Udacity as an efficient platform with the greatest positive impact on how online learners perceive the quality of an e-learning platform.

[10] The aim of the study is to enhance the efficiency and effectiveness of e-learning by leveraging the capabilities of cloud computing. The authors propose a model that integrates various components, such as content management, learning management, and assessment systems, within a cloud-based infrastructure. They highlight the advantages of using cloud computing in the e-learning context, including scalability, flexibility, and cost-effectiveness. The paper also discusses the

challenges and considerations involved in implementing such a model, such as data security, privacy, and interoperability. The authors suggest that the proposed model can provide a comprehensive and adaptable framework for e-learning, promoting collaboration, accessibility, and personalized learning experiences.

[11] The paper titled "An Evaluation of E-Learning and User Satisfaction" examines the effectiveness of e-learning platforms and their impact on user satisfaction. The study collects data from e-learning users and analyzes factors such as usability, interactivity, content quality, and system performance to assess user satisfaction. The findings indicate that factors like ease of use, accessibility, relevance of content, and responsiveness of the e-learning system influence user satisfaction. The paper emphasizes the importance of designing user-friendly and engaging e-learning platforms that meet the diverse needs of learners. The study highlights the significance of user satisfaction in e-learning and suggests continuous improvement of e-learning systems to optimize user experience and enhance learning outcomes.

[12] The paper provides a literature review on the topic of multilingual optical character recognition (OCR) systems using reinforcement learning for character segmentation. It discusses the challenges faced by OCR systems in recognizing characters from different languages and emphasizes the importance of accurate character segmentation. Various approaches and algorithms for character segmentation are reviewed, including rule-based methods, template matching, and machine learning techniques. The paper also explores the application of reinforcement learning in optimizing the character segmentation process. The literature review highlights the gaps in the existing research and establishes the need for improved multilingual OCR systems. Overall, the review sets the foundation for the proposed approach of using reinforcement learning for character segmentation in multilingual OCR systems.

[13] The paper presents a comprehensive review of multimedia-based information retrieval techniques. It highlights the limitations of traditional methods and existing systems in efficiently retrieving information from multimedia data. To address these limitations, the authors propose a multimedia-based approach that integrates Automatic Speech Recognition (ASR), Optical Character Recognition (OCR), and a video recommendation system. The paper employs a systematic literature review technique to survey relevant research in the field, covering topics such as multimedia-based information retrieval, ASR, OCR, and video recommendation systems. The review helps identify the research gap in the current literature, which lies in the inefficiency of traditional methods and the challenges faced by existing multimedia-based retrieval systems in handling multiple multimedia components simultaneously. For future scope, the paper suggests enhancing the accuracy and efficiency of multimedia-based retrieval systems by integrating ASR, OCR, and video recommendation systems. It also proposes exploring the potential of machine learning and deep learning techniques to further improve system performance.

[14] The paper focuses on the integration of Optical Character Recognition (OCR) and Semantic Web technologies to enhance the learning experience in the context of e-learning. The main contribution is the development of a framework that leverages OCR and Semantic Web to provide personalized and contextualized learning resources. The research identifies a gap in the current literature, which is the limited integration of OCR and Semantic Web technologies in the field of e-learning. It also highlights the underutilization of these technologies for educational purposes. The future scope of the research lies in



the potential for developing more sophisticated and efficient learning systems by incorporating OCR and Semantic Web technologies. Additionally, there is an opportunity to enhance the accessibility and usability of learning materials through the application of OCR and Semantic Web. The paper proposes a framework that integrates OCR and Semantic Web technologies to optimize the learning experience in e-learning. By leveraging these technologies, the efficiency and effectiveness of e-learning systems can be improved, leading to a more personalized and context-aware learning environment.

[15] The paper aims to analyze the trends, challenges, and potential of AI-supported eLearning through a systematic review and co-citation network analysis. The main contribution of the paper is providing a valuable resource for researchers and practitioners in the field by offering insights into the evolution and current state of AI-supported eLearning. The technique employed in the study is a systematic review and co-citation network analysis. This approach allows for the identification and analysis of key research topics, influential papers, and connections between different research areas in the field of AI-supported eLearning. The research identifies a research gap in the lack of focus on the human-centered design of AI-supported eLearning systems. It suggests that future research should explore the integration of AI with a user-centered design approach to enhance the effectiveness and efficiency of eLearning. Additionally, there is a need for further research in developing AI-based personalized and adaptive learning systems to cater to the individual needs of learners. The research provides a valuable resource for researchers and practitioners in the field of eLearning. It offers insights into the trends, challenges, and potential of AI-supported eLearning through a comprehensive review of the literature, thereby informing future research and practice in the domain.

[16] The paper "Building an efficient OCR system for historical documents with little training data" by Jiří Martínek, Ladislav Lenc, and Pavel Král, focuses on addressing the challenge of optical character recognition (OCR) for historical documents with limited annotated training data. The authors present a complete OCR system that includes page layout analysis and OCR using state-of-the-art techniques such as fully convolutional networks for segmentation and recurrent neural networks for OCR. They also introduce a new dataset called Porta fontium portal for evaluation. The experiments demonstrate that good performance can be achieved with a small amount of annotated data, and the proposed system outperforms or is comparable to state-of-the-art systems. The paper provides insights into building efficient OCR systems for historical documents with minimal training data.

[17] The paper discusses the development of a personalized e-learning system based on the user's performance and knowledge. The system adapts to individual learners by considering their knowledge level and modifies the learning content accordingly. The paper presents an experiment conducted with 50 registered users, where knowledge tests were taken, responses were stored, and progress was tracked. The results demonstrate the effectiveness of the system in providing personalized learning content and enabling learners to track their performance and progress. The paper concludes by suggesting future enhancements such as adding more languages and topics, developing content based on user qualification, and advancing the test assessment patterns.

[18] The paper discusses the application of deep neural networks (DNN) for optical character recognition (OCR) on a historical Finnish newspaper and journal corpus. The goal is to improve OCR accuracy for reliable search and scientific research on the OCRed data. The existing OCR quality achieved with commercial software has a character error rate

(CER) between 8 and 13%. The paper explores the training of OCR models using DNNs and additional training data to create high-quality mixed-language models capable of recognizing both Finnish and Swedish texts printed in two font families (Blackletter and Antiqua). The authors also investigate the impact of confidence voting on OCR results using different combinations of models.

[19] the paper discusses Natural Language Processing (NLP) and its application in semantic search. It provides a literature review of previous work in the field of NLP and discusses the history of NLP. The paper also presents an architectural design for a semantic search tool using NLP, including its modules, domain, and constraints. The results and discussion section presents the findings of the work, including the outcomes of the literature review and the proposed architectural design. It may also mention the discovery of ambiguities in the knowledge base and ontologies during the system development.

[20] The paper focuses on leveraging deep learning algorithms to analyze and predict resource usage patterns in e-learning environments. By understanding these patterns, the researchers aim to develop intelligent resource provisioning strategies that can dynamically allocate computational resources based on real-time demand. The proposed approach utilizes historical usage data and employs deep learning models, such as Long Short-Term Memory (LSTM) networks, to forecast future resource requirements. These predictions enable proactive resource allocation and ensure that the necessary computing resources are readily available to support the e-learning platform's users and applications.

[21] The paper provides an overview of the Tesseract OCR engine, discussing its development, features, and performance. Tesseract was developed at HP Labs between 1984 and 1994 and gained recognition for its accuracy in the UNLV Annual Test of OCR Accuracy in 1995. Although it didn't become a commercial product at that time, it was released as an open-source project in 2005. The paper highlights several unique aspects of Tesseract, including its line finding, features/classification methods, and adaptive classifier. It compares the performance of Tesseract in recent versions with its original results from 1995, showing improvements and changes in accuracy for different test sets.

[22] The paper proposes an OCR post-processing approach based on multi-knowledge, combining language knowledge, candidate distance information, statistical language models, and semantic lexicon. The goal is to improve the recognition accuracy of OCR systems by incorporating context information and reducing the search space. The experimental results demonstrate the effectiveness of the approach, with the recognition accuracy rate on the test set increasing from 58.45% to 83.73%, resulting in a 60.84% reduction in errors.

[23] The paper provides an overview of the use of machine learning technologies in e-learning. It discusses the application of machine learning in sentiment analysis, student behavior prediction, and self-regulated learning. Various machine learning algorithms, such as Random Forest and Support Vector Machine, are utilized to analyze learners' sentiments, predict learner satisfaction, classify and profile learners, and enhance self-regulated learning strategies. The paper highlights the effectiveness of machine learning in improving e-learning experiences and suggests future work in evaluating e-learning content quality using machine learning techniques.

[24] The paper provides an extensive overview of the applications of machine learning and deep learning in e-learning systems. The authors conducted a literature survey to identify key themes and trends in this research field. The paper highlights the importance of machine learning and deep learning in analyzing data generated in e-learning systems to

enhance learner experiences and personalize learning content. The study discusses various applications of machine learning and deep learning in e-learning, including sentiment analysis, student behavior prediction, and self-regulated learning. Sentiment analysis involves analyzing learners' emotions and sentiments to predict satisfaction and engagement. The paper presents studies that utilize supervised machine learning algorithms to classify learners' sentiments and emotions in MOOCs.

[25] The paper focuses on the application of deep learning algorithms for handwritten text recognition (HTR). The authors highlight the need for improved accuracy in HTR systems and propose the use of a deep learning approach, specifically the LSTM (Long Short-Term Memory) model. They collect data for training the HTR system, extract features from the handwritten text datasets, and train the model using the deep learning approach. The goal is to recognize words rather than individual characters to enhance accuracy. The developed LSTM deep model achieves a high accuracy of 94 percent in the recognition of handwritten text. A comparison is provided which shows the accuracy of the proposed method (2DLSTM) compared to previously recorded algorithms. The 2DLSTM approach achieves a Character Error Rate (CER) of 8.2 percent and a Word Error Rate (WER) of 27.5 percent. In contrast, the CNN-1DLSTM-CTC method achieves a CER of 6.2 percent and a WER of 20.5 percent. While the LSTM model demonstrates higher accuracy at the word level, it slightly reduces accuracy at the character level, resulting in increased spelling errors for mislabeled words but fewer errors overall at the word level.

## V. RESEARCH METHODOLOGY

**OCR process overview:** Provide an overview of the OCR process and explain how Tesseract works. This could include information on image preprocessing, character segmentation, feature extraction, and recognition.

**Data collection:** Describe the data that was used for the OCR process, including the type of images, image quality, and any preprocessing techniques that were used.

**OCR performance evaluation:** Discuss the performance of Tesseract on the OCR task. This could include metrics such as accuracy, precision, recall, F1-score, and processing time. Compare the results with other OCR tools if applicable.

**Error analysis:** Analyze the errors made by Tesseract during the OCR process. Identify common types of errors and provide examples of images that were difficult for Tesseract to recognize.

**Improvements:** Suggest possible improvements to the OCR process using Tesseract. This could include tweaking Tesseract's configuration, using different preprocessing techniques, or using a different OCR tool altogether.

## VI. PROPOSED METHODOLOGY

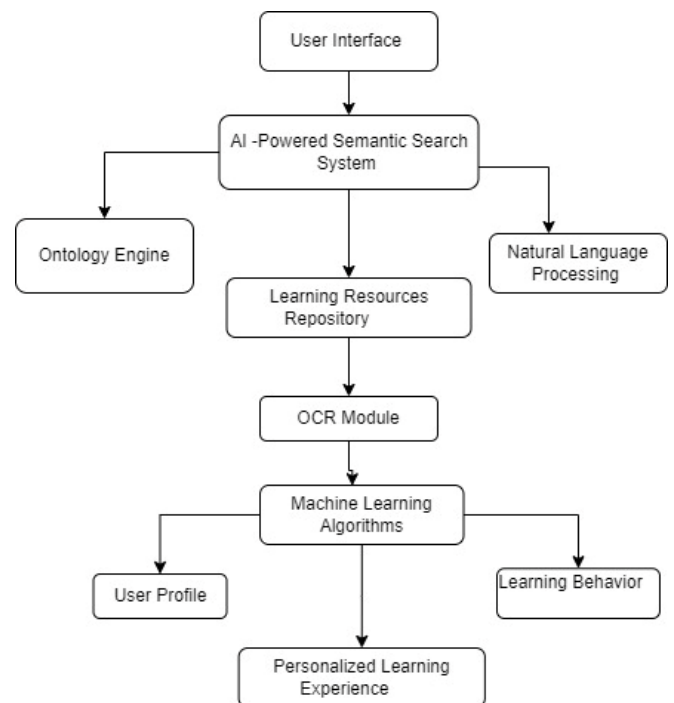


Fig. 1. The proposed system

The proposed system section provides a detailed description of the AI-powered semantic search and OCR-enabled e-learning platform. It discusses the system architecture, the AI and OCR technologies used, and the features of the system such as natural language processing, personalized learning, and visual content extraction.

**User Interface:** This is the graphical user interface that students and teachers use to interact with the platform. The user interface provides access to all the features and functionalities of the platform, including search, browsing, and personalized learning.

**AI-powered Semantic Search System:** This component uses machine learning algorithms to provide an enhanced search experience for users. It includes an ontology engine that enables the platform to understand the meaning of search queries and return more accurate search results. The natural language processing component enables the platform to understand and interpret natural language search queries. The Learning Resources Repository is where all the e-learning resources are stored and made available to users.

**OCR Module:** This component is responsible for extracting text from scanned images of documents, making them searchable and accessible to users. This is especially useful for historical documents that are not available in digital format.

**Machine Learning Algorithms:** This component uses user-profiles and learning behavior to provide a personalized learning experience for each user. The user profile contains information about the user's interests, learning style, and preferences. The learning behavior data includes information about the user's performance on quizzes, assignments, and tests. The machine learning algorithms analyze this data to provide customized recommendations for learning resources and activities that are best suited for each user.

**Personalized Learning Experience:** This component is the result of the machine learning algorithms and provides a tailored learning experience for each user. The personalized learning experience includes customized recommendations for learning resources, activities, and assessments that match the user's interests and learning style. This component helps

students to stay engaged and motivated throughout the learning process, leading to better learning outcomes

Data preprocessing involves cleaning and normalizing text data to remove noise and facilitate subsequent analysis. TypeScript's string manipulation functions and regular expressions can be utilized to perform tasks like lowercase conversion, punctuation removal, and stop-word elimination.

i. Entity Recognition and Disambiguation:

Use NLP libraries compatible with TypeScript to perform entity recognition and disambiguation. These libraries provide functions or APIs to extract named entities from text. Fine-tuning the models or providing additional training data can improve entity recognition accuracy.

ii. Concept Extraction:

Extracting important concepts or topics from text is essential for semantic understanding. Various techniques like keyword extraction, topic modeling, or domain-specific methods can be employed. NLP libraries may offer built-in functions for concept extraction, or custom algorithms can be implemented.

iii. Semantic Representation:

Transform preprocessed text and extracted entities into a semantic representation that captures the meaning and relationships between words and concepts. TypeScript-compatible libraries like Word2Vec.js can be employed to convert words into dense vector representations using pre-trained word embeddings.

iv. Knowledge Graph Construction:

If the semantic search system involves a knowledge graph, TypeScript-compatible graph database libraries like Neo4j or Graphile can be used to represent structured relationships between entities and concepts. These libraries provide TypeScript bindings for interacting with the graph database.

v. Query Understanding:

Develop algorithms or functions to understand user queries by extracting entities, concepts, and their relationships. NLP libraries can assist in query understanding, employing techniques such as syntactic and semantic parsing.

vi. Semantic Search Algorithms:

Implement search algorithms that leverage the semantic representation of documents and queries to retrieve relevant results. Techniques like semantic similarity calculation or machine learning approaches can be employed. TypeScript allows the implementation of these algorithms based on specific requirements.

vii. User Interface:

Design and develop a user-friendly interface using TypeScript-based web development frameworks like React or Angular. Implement input fields for user queries and display search results in an intuitive and meaningful manner. Communicate with the backend semantic search logic through APIs.

• Modules of project:

**End Module –**

This module is responsible for the frontend of the website. It includes the user interface design, such as layout, color scheme, typography, and overall user experience. The end module is responsible for displaying the content and user interactions, such as search queries, suggestions, and feedback.

**Back-End Module –**

This module is responsible for the backend of the website. It includes the database management, server configuration, API development, and security protocols. The back-end module is responsible for handling user authentication, user data management, and content management.

**OCR Module –**

This module will implement OCR (Optical Character Recognition) and add it to the website. OCR technology will enable the platform to extract text from scanned images or documents and make them searchable and editable. This module will require integrating OCR software, such as Tesseract, with the platform's backend to extract the text and store it in a database for search and retrieval. In this project, a comparison of three OCR (Optical Character Recognition) algorithms was conducted: KerasOCR, EasyOCR, and TesseractOCR. The objective was to determine the most suitable algorithm based on performance speed and accuracy. A dataset of 25 diverse images containing different types of text was utilized for the evaluation.

	bbox	text	conf
0	[[607, 327], [853, 327], [853, 359], [607, 359]]	NISSAN GENISS	0.793183
1	[[915, 337], [989, 337], [989, 353], [915, 353]]	NISSAN	0.975665
2	[[821, 381], [853, 381], [853, 399], [821, 399]]	#1	0.068619
3	[[210, 726], [238, 726], [238, 734], [210, 734]]	LHNa	0.038245
4	[[187, 728], [262, 728], [262, 757], [187, 757]]	GENSS	0.993451

Fig. 2. Tesseract OCR Result

id	image_id	bbox	utf8_string	points	area
761457	00004b36676338b_1	[308,0, 135,8, 45,96, 17,33]	FELIX	[[313,24, 135,8, 353,96, 136,2, 348,72, 153,13...]]	796.49
761458	00004b36676338b_2	[383,8, 140,15, 70,69, 21,4]	PRIVAT	[[387,86, 140,15, 454,69, 146,82, 454,69, 161,5...]]	1517.25
761459	00004b36676338b_3	[137,37, 463,48, 34,26, 26,9]	DBU	[[140,29, 463,48, 172,13, 468,72, 172,13, 489,2...]]	883.91
761460	00004b36676338b_4	[174,15, 471,14, 33,07, 26,2]	BB	[[174,15, 471,14, 209,22, 475,57, 208,62, 497,3...]]	918.83
761461	00004b36676338b_5	[637,92, 456,18, 60,66, 18,70]		[[638,62, 467,3, 697,42, 456,18, 698,58, 462,9...]]	1137.98

Fig. 3. Easy OCR Result

After analyzing the results, the TesseractOCR algorithm was found to outperform the other two algorithms in terms of both speed and accuracy. It exhibited the ability to quickly and accurately recognize text from the images in the dataset. The precision of the extracted text and the speed of processing were particularly noteworthy.

Consequently, the decision was made to integrate the TesseractOCR module into the project. This choice ensures that the project can leverage the superior performance of TesseractOCR in terms of speed and accuracy. By selecting TesseractOCR, the project aims to provide a reliable and efficient OCR solution that meets the requirements and expectations of the users.

The OCR results obtained using EasyOCR on the first two images of the dataset are as follows. In the first image, EasyOCR successfully detected the text "#" with a confidence score of 0.20, "B4\*" with a confidence score of 0.14, "om.hk" with a high confidence score of 0.98, and "aebekae: 2926 7222 =" with a low confidence score of 0.06. Moving on to the second image, EasyOCR recognized the text "IdBu eeg1" with a confidence score of 0.13 and "PRIVATA" with a confidence score of 0.57.

The evaluation and comparison of the OCR algorithms and the selection of TesseractOCR as the preferred choice have significantly enhanced the functionality and performance of the project. The integration of Tesseract OCR enables efficient extraction and processing of text from various image sources, leading to an improved user experience and overall effectiveness of the application.



### Semantic Search Module –

This module will implement Semantic search and add it to the Website. The Semantic search module will analyze the meaning of the text, rather than just keyword matching. The module will use techniques such as latent semantic analysis (LSA) or word embeddings to understand the context of the search query and provide more accurate and relevant search results. The semantic search module will require integrating NLP (Natural Language Processing) libraries, such as spaCy or NLTK, with the platform's backend to analyze the text and provide meaningful search results. The results of the semantic search show the match percentage of the search word. We search for the word 'web' and got the match percentage of 1% for Web Dev course and 0.6827% for the Computer Networks course

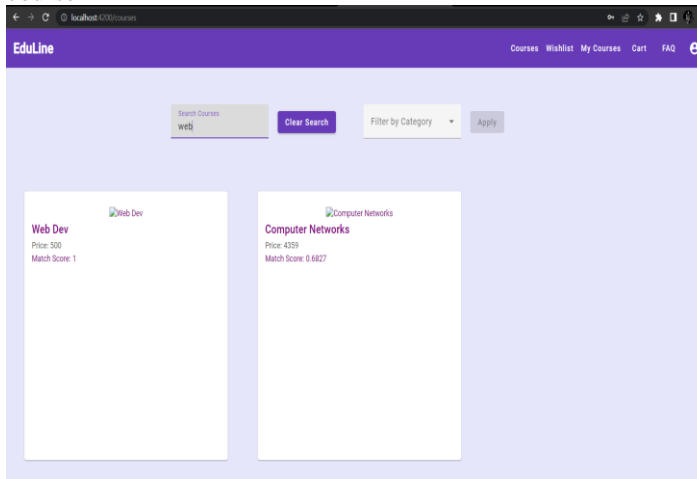


Fig. 4. Semantic similarity in the search query

### VII. HARDWARE AND SOFTWARE REQUIREMENTS

- 1)Hardware requirements:
  - a) Storage of 10GB.
  - b) Basic GPU for Processing the images.
  - c) Scanner if we need to scan images.
  - d) Processor to process the data
- Software Requirements:
  - a) Database manager
  - b) OCR Software
  - c) Search software
  - d) Image processing software
  - e) Web server
  - f) Machine Learning Framework

### VIII. ASSUMPTIONS

**Adequate labeled training data for OCR:** It is assumed that a dataset containing a diverse range of text-based documents, including scanned or image-based files, handwritten text, and various fonts and languages, is available for training the OCR model. The dataset should be accurately labeled to enable the model to effectively recognize and extract text from different document types.

**Availability of Sufficient Resources:** It is assumed that there will be adequate resources, including funding, infrastructure, and human resources, to support the development, implementation, and maintenance of the learning management system.

**Adequate Content Availability:** It is assumed that instructors will have access to the necessary course content, such as lecture slides, readings, and multimedia materials, that

need to be uploaded to the learning management system. The assumption is that the required content will be available and accessible for integration into the system. Computational Resources: Sufficient computational resources, including CPUs,

### IX. RESULT AND DISCUSSION

The development of an AI-powered semantic search and OCR-enabled e-learning platform offers several notable advantages. The integration of artificial intelligence techniques allows for more accurate and context-aware search functionality, enabling learners to find relevant information efficiently. Additionally, the integration of OCR technology enables learners to access and interact with textual content from scanned documents and images, expanding the range of available learning materials. The platform's personalized recommendations and adaptive learning paths enhance the learning experience by tailoring content to individual learners' needs and preferences. The combination of AI-powered semantic search and OCR technology in e-learning platforms presents significant implications for educational practices. By leveraging semantic search capabilities, learners can benefit from improved search results that better match their intent and provide more precise information. This reduces the time and effort required to find relevant learning resources, allowing learners to focus more on actual learning activities.

The OCR-enabled content access feature overcomes the limitations of physical documents, making it easier for learners to digitize and interact with printed materials. This functionality has the potential to enhance accessibility and inclusivity in education, particularly for learners with visual impairments or those who prefer digital formats.

Furthermore, the AI algorithms employed in the platform can analyze user interactions and preferences, enabling personalized learning experiences. By understanding individual strengths, weaknesses, and learning styles, the platform can deliver tailored content and recommendations, fostering learner engagement and motivation. Learners can benefit from targeted resources, adaptive exercises, and personalized guidance, leading to more efficient and effective learning outcomes.

Collaboration and discussion features supported by the platform promote social learning and knowledge sharing among learners. By providing communication channels and collaborative workspaces, learners can engage in meaningful discussions, exchange ideas, and seek assistance from peers and educators. The AI algorithms can assist in moderating discussions, identifying relevant threads, and promoting valuable contributions, fostering a vibrant and interactive learning community.

While the AI-powered semantic search and OCR-enabled e-learning platform offer tremendous potential, challenges and considerations exist. Ensuring the accuracy and reliability of search results and OCR-generated content is crucial to avoid misinformation. Additionally, privacy concerns regarding the collection and use of learner data must be addressed to maintain learner trust and comply with privacy regulations.

Future research and development in this area should focus on refining the AI algorithms for semantic search, improving the accuracy of OCR technology, and conducting user studies to assess the impact of these platforms on learning outcomes and user satisfaction. Long-term evaluations and comparative studies can help validate the effectiveness of these platforms and guide their further enhancements.

In conclusion, the integration of AI-powered semantic search and OCR technology in e-learning platforms holds great promise for transforming the learning experience. By

leveraging advanced techniques, these platforms have the potential to improve information retrieval, access to learning materials, personalization, collaboration, and overall learner engagement. Continued research and development in this area will further refine these technologies and enable more effective and inclusive digital learning environments.

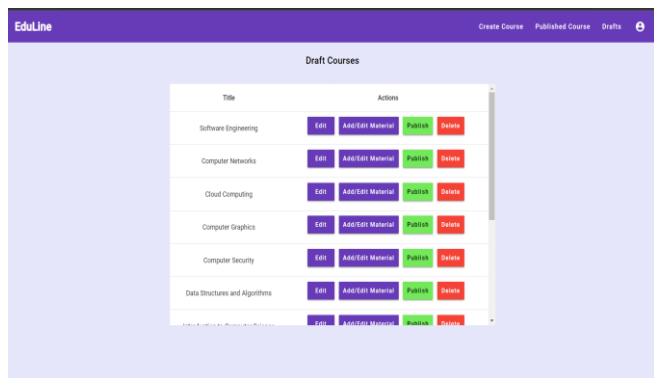


Fig. 4. Instructor Courses Page

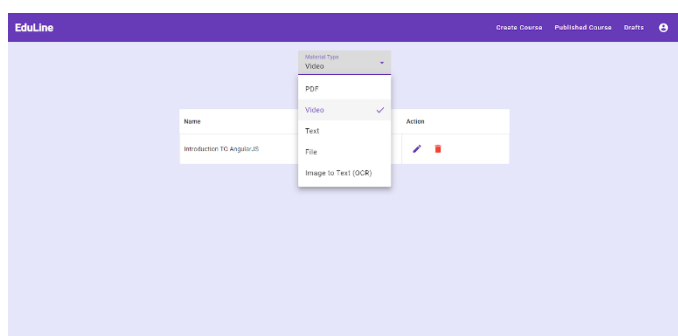


Fig. 5. Creating Course Material

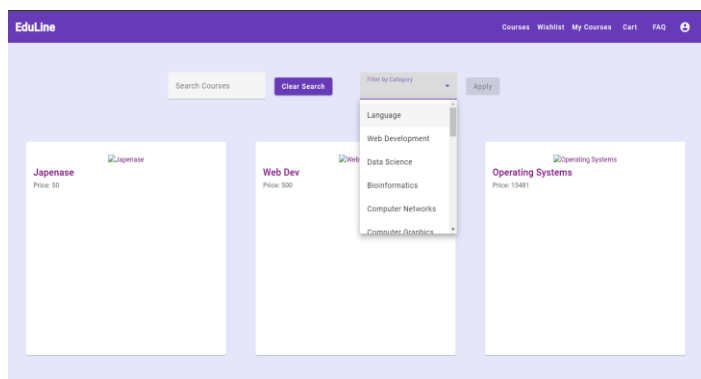


Fig. 6. Student Courses Page

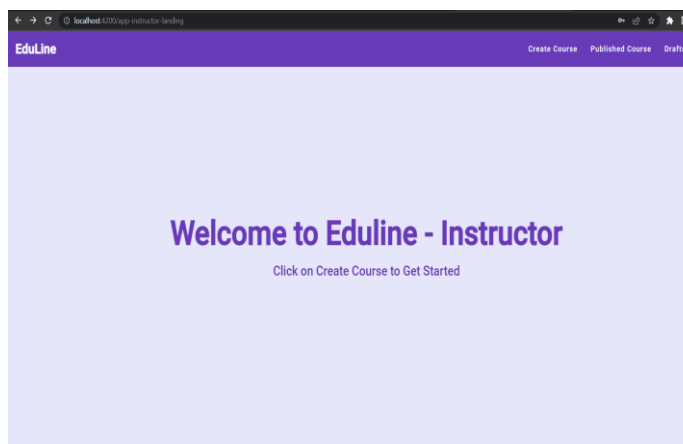


Fig. 7. Instructor Dashboard

id	name
1	Natural Language Processing
2	Bioinformatics
3	Data Science
4	Bioinformatics
5	Computer Networks
6	Computer Graphics
7	Web Development
8	Cloud Computing
9	Artificial Intelligence
10	Data Science
11	Machine Learning
12	Artificial Intelligence
13	Cloud Computing
14	Web Development
15	Introduction to Computer Sc...
16	Computer Networks
17	Introduction to Computer Sc...
18	Operating Systems
19	Web Development

id	title	description	keywords	price	is_authorized	is_published	is_deleted	category_id
1	Data Structures and Algorithms	and motion planning	Computer Ethics	49801	1	1	0	2
2	Computer Security	for biological research and me...	Natural Language Processing	33383	1	1	1	13
3	Introduction to Computer Science	proteomics	Computer Science History	4760	1	0	1	15
4	Operating Systems	and animation	Algorithm	15481	1	0	0	11
5	Introduction to Computer Science	Paas	Human-Computer Interaction	33818	1	1	1	1
6	Mobile Application Development	semantics	Computer Graphics	26538	1	0	1	15
7	Software Engineering	enabling applications such as ...	Programming Languages	7997	1	0	0	12
8	Software Engineering	Understand the design and or...	Software Engineering	30055	1	0	1	4
9	Cloud Computing	file systems	Distributed Systems	33730	1	1	0	13
10	Computer Networks	with an emphasis on algorithm...	Computer Graphics	16876	1	0	0	14
11	Database Systems	data breaches	Computer Networks	13974	1	1	0	7

12	Machine Learning	Paas	Distributed Systems	23026	1	0	1	1
13	Computer Networks	Learn web technologies such ...	Computer Security	4359	1	1	0	4
14	Cloud Computing	rendering techniques	Algorithm	47028	1	0	0	20
15	Computer Architecture	perception	Distributed Systems	42857	1	1	1	6
16	Mobile Application Development	their syntax	Computer Science History	38478	1	1	0	9
17	Computer Graphics	understanding algorithms and...	Artificial Intelligence	48001	1	0	0	9
18	Machine Learning	Understand the fundamentals...	Artificial Intelligence	20593	1	1	1	17
19	Mobile Application Development	rendering techniques	Computer Vision	26418	1	1	1	15
20	Introduction to Computer Science	rendering techniques	Human-Computer Interaction	5684	1	0	1	3
21	Computer Architecture	enabling applications such as ...	Natural Language Processing	40796	1	0	1	7
22	Software Engineering	covering topics such as algori...	Computer Graphics	10178	1	1	0	16

23	Introduction to Computer Science	and management of database...	Computer Science Theory	37697	1	0	1	20
24	Operating Systems	to build robust and high-quali...	Computer Vision	3092	1	0	1	19
25	Data Structures and Algorithms	routing	Computer Science Theory	21414	1	0	1	15
26	Computer Networks	Learn about threats	Algorithm	40382	1	1	1	15
27	Introduction to Computer Science	and performance optimization.	Natural Language Processing	24623	1	1	1	19
28	Computer Networks	understanding algorithms and...	Human-Computer Interaction	5525	1	0	1	16
29	Web Development	Understand the principles and...	Computer Science History	5960	1	0	1	10
30	Computer Security	Study techniques for underst...	Computer Architecture	23543	1	0	0	2
31	Introduction to Computer Science	and animation	Natural Language Processing	21665	1	1	0	3
32	Software Engineering	Understand the fundamentals...	Computer Vision	27631	1	1	1	6

33	Mobile Application Development	An overview of the fundamen...	Computer Science Education	34773	1	1	1	9
34	Data Structures and Algorithms	and frameworks to design an...	Algorithm	18513	1	0	0	9
35	Introduction to Computer Science	natural language processing	Data Structure	23556	1	0	0	11
36	Artificial Intelligence	to develop smart systems and...	Computer Science Theory	36392	1	1	1	11
37	Artificial Intelligence	vulnerabilities	Algorithm	16183	1	1	1	15
38	Artificial Intelligence	Learn about threats	Distributed Systems	22922	1	1	0	7
39	Artificial Intelligence	and network security.	Computer Science Education	48872	1	0	0	8
40	Database Systems	switching	Computer Architecture	29352	1	1	0	17
41	Programming Languages	routing	Computer Architecture	42038	1	0	0	5
42	Programming Languages	An overview of the fundamen...	Computer Graphics	43563	1	1	1	3
43	Computer Architecture	vulnerabilities	Operating Systems	45504	1	0	0	13

44	Computer Architecture	including data modeling	Computer Security	39338	1	0	0	3
45	Computer Security	Understand the design and or...	Algorithm	12182	1	1	1	18
46	Programming Languages	Gain knowledge of software d...	Computer Graphics	13877	1	0	1	4
47	Database Systems	including machine learning	Human-Computer Interaction	22353	1	0	1	18
48	Machine Learning	natural language processing	Computer Ethics	45367	1	1	1	11
49	Introduction to Computer Science	memory management	Human-Computer Interaction	45098	1	0	1	6
50	Operating Systems	and performance optimization.	Machine Learning	12376	1	0	0	8

Fig. 8. Database courses table

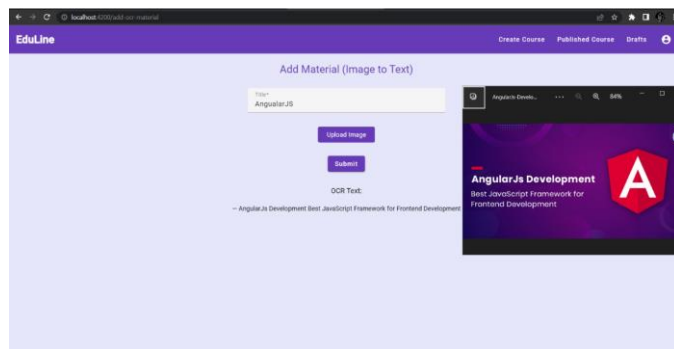


Fig. 9. OCR view

### X. CONCLUSION

This paper proposes an AI-powered semantic search and OCR-enabled e-learning platform that aims to enhance knowledge acquisition and retention. The system employs natural language processing, OCR technology, and adaptive learning algorithms to provide an intuitive, interactive, and efficient learning environment. The proposed system can



improve the quality of online education by providing personalized and effective learning experiences for students and researchers. The results of the user study demonstrate the effectiveness of the proposed system in facilitating knowledge acquisition and retention. The potential of the proposed platform extends beyond the realm of e-learning and has significant implications for various fields that require effective knowledge management and acquisition. Future work can focus on improving the system's performance, expanding its scope, and developing new features to meet the evolving needs of users.

## XI. FUTURE SCOPE

design and implement the visual components of the application that users will react with

- ❖ Develop the application's backend architecture and infrastructure.
- ❖ Integrate the frontend and backend along with the OCR and Symanctic search
- ❖ Test and deploy the application
- ❖ Ensure the application meets performance requirements and deploy it to the production environment.

## XII. REFERENCES

[1] E.H.-K. Wu, C. -H. Lin, Y. -Y. Ou, C. -Z. Liu, W. - K. Wang and C. -Y. Chao, "Advantages and Constraints of a Hybrid Model K-12 E-Learning Assistant Chatbot," in *IEEE Access*, vol. 8, pp. 77788-77801, 2020, doi: 10.1109/ACCESS.2020.2988252.

[2] Al-Fraihat, Dimah, Joy, Mike, Masa'deh, Ra'ed and Sinclair, Jane (2020) Evaluating E-learning systems success: an empirical study. *Computers in Human Behavior*, 102. pp. 67-86. doi: 10.1016/j.chb.2019.08.004 ISSN 0747-5632

[3] Snigdha Choudhury, Snigdha Pattnaik, Emerging themes in e-learning: A review from the stakeholders' perspective, *Computers & Education*, Volume 144, 2020, 103657, ISSN 0360-1315, <https://doi.org/10.1016/j.compedu.2019.103657>.

[4] F. M. Enescu, G. Șerban and M. Jurian, "Web Platform for E-Learning," 2019 11th International Conference on Electronics, Computers and Artificial Intelligence (ECAI), Pitesti, Romania, 2019, pp. 1-6, doi: 10.1109/ECAI46879.2019.9042106.

[5] B. Zhang, "Research on Educational Informatization Platform Based on E-learning Platform," 2021 IEEE Asia-Pacific Conference on Image Processing, Electronics and Computers (IPEC), Dalian, China, 2021, pp. 1043-1046, doi: 10.1109/IPEC51340.2021.9421207.

[6] M. Bhatia and J. K. Maitra, "E-learning Platforms Security Issues and Vulnerability Analysis," 2018 International Conference on Computational and Characterization Techniques in Engineering & Sciences (CCTES), Lucknow, India, 2018, pp. 276-285, doi: 10.1109/CCTES.2018.8674115

[7] Keskin, S. & Yurdugül, H. (2022). E-learning experience: Modeling students' e-learning interactions using log data . *Journal of Educational Technology and Online Learning*.

[8] W. Bagunaid, N. Chilamkurti, and P. Veeraraghavan, "AISAR: Artificial Intelligence-Based Student Assessment and Recommendation System for E-Learning in Big Data," *Sustainability*, vol. 14, no. 17, p. 10551, Aug. 2022, doi: 10.3390/su141710551.

[9] B. Alojaiman, "Toward Selection of Trustworthy and Efficient E-Learning Platform," in *IEEE Access*, vol. 9, pp. 133889-133901, 2021, doi: 10.1109/ACCESS.2021.3114150.

[10] E. H. F. Ezzahra, C. Mohamed and B. Abdelhamid, "Towards e-learning ecosystem model based on cloud computing," 2020 X International Conference on Virtual Campus (JICV), Tetouan, Morocco, 2020, pp. 1-4, doi: 10.1109/JICV51605.2020.9375724.

[11] Rajasekaran, Vijay Anand, et al. "An Evaluation of E-Learning and User Satisfaction." *IJWLTT* vol.17, no.2 2022: pp.1-11.

[12] J. Park, E. Lee, Y. Kim, I. Kang, H. I. Koo and N. I. Cho, "Multi-Lingual Optical Character Recognition System Using the Reinforcement Learning of Character Segmenter" in *IEEE Access*, vol. 8, pp. 174437-174448, 2020, doi: 10.1109/ACCESS.2020.3025769.

[13] D. T. Bhabad, S. Therese and M. Gedam, "Multimedia based Information Retrieval Approach based on ASR and OCR and Video Recommendation System," 2017 International Conference on Current Trends in Computer, Electrical, Electronics and Communication (CTCIC), Mysore, India, 2017, pp. 1168-1172, doi: 10.1109/CTCEEC.2017.8455038.

[14] K. Badwaik, K. Mahmood and A. Raza, "Towards applying OCR and Semantic Web to achieve optimal learning experience," 2017 IEEE 13th International Symposium on Autonomous Decentralized System (ISADS), Bangkok, Thailand, 2017, pp. 262-267, doi: 10.1109/ISADS.2017.40.

[15] Kai-Yu Tang, Ching-Yi Chang & Gwo-Jen Hwang (2021) Trends in artificial intelligence-supported e-learning: a systematic review and co-citation network analysis (1998–2019), *Interactive Learning Environments*, DOI: 10.1080/10494820.2021.1875001

[16] "Building an efficient OCR system for historical documents with little training data" by Jiří Martínek, Ladislav Lenc, and Pavel Král

[17] Personalized E-Learning System Based on User's Performance and Knowledge: An Adaptive Technique by Patchava. RamyaSree, Tammisetty. Bhuvanewari, Vulchi. Vamsi Swapnika Reddy, Jonnalagadda. Surya Kiran

[18] Optical character recognition with neural networks and post-correction with finite state methods by Senka Drobac1 · Krister Lindén

[19] A. Kupiyalova, R. Satybaldiyev and S. Aiaskarov, "Semantic search using Natural Language Processing," 2020 IEEE 22nd Conference on Business Informatics (CBI), Antwerp, Belgium, 2020, pp. 96-100, doi: 10.1109/CBI49978.2020.10065.

[20] J. Ariza, M. Jimeno, R. Villanueva-Polanco and J. Capacho, "Provisioning Computational Resources for Cloud-Based e-Learning Platforms Using Deep Learning Techniques," in *IEEE Access*, vol. 9, pp. 89798-89811, 2021, doi: 10.1109/ACCESS.2021.3090366.

[21] R. Smith, "An Overview of the Tesseract OCR Engine," Ninth International Conference on Document Analysis and Recognition (ICDAR 2007), Curitiba, Brazil, 2007, pp. 629-633, doi: 10.1109/ICDAR.2007.4376991.

[22] An OCR Post-processing Approach Based on Multi-knowledge by Li Zhuang & Xiaoyan Zhu

[23] R. Farhat, Y. Mourali, M. Jemni and H. Ezzedine, "An overview of Machine Learning Technologies and their use in E-learning," 2020 International Multi-Conference on: "Organization of Knowledge and Advanced Technologies" (OCTA), Tunis, Tunisia, 2020, pp. 1-4, doi: 10.1109/OCTA49274.2020.9151758.

[24] C. Fri and R. Elouahbi, "Machine Learning and Deep Learning applications in E-learning Systems: A Literature Survey using Topic Modeling Approach," 2020 6th IEEE Congress on Information Science and Technology (CiSt), Agadir -Essaouira, Morocco, 2020, pp. 267-273, doi: 10.1109/CiSt49399.2021.9357253.

[25] A. Nikitha, J. Geetha and D. S. JayaLakshmi, "Handwritten Text Recognition using Deep Learning," 2020 International Conference on Recent Trends in Electronics, Information, Communication & Technology (RTEICT), Bangalore, India, 2020, pp. 388-392, doi: 10.1109/RTEICT49044.2020.9315679.