# Question Bank Plagiarism Checker Module

**Dr.AB.Hajira Be[1], S Chandru [2]**

*[1] Associate Professor*
*Department of Computer Applications*
*Karpaga Vinayaga College of Engineering and Technology*
*Maduranthagam TK*
*[2]PG Student*
*Department of Computer Applications*
*Karpaga Vinayaga College of Engineering and Technology*
*\*Corresponding Author: S Chandru  Email: chandruselvaraj9987@gmail.com*

---------------------------------------------------------------------***----------------------------------------------------------------------

**Abstract -** In the digital age of education, ensuring the originality and integrity of question banks is a pressing challenge, especially with the widespread reuse and paraphrasing of content. Manual plagiarism detection methods are inefficient at scale and often fail to identify semantically similar or reworded questions. This paper presents an AI-Based Question Bank Plagiarism Checker Module that leverages advanced Natural Language Processing (NLP) techniques to detect duplicated and paraphrased academic questions. The system utilizes transformer-based sentence embedding models such as Sentence-BERT and SimCSE to convert textual content into semantic vectors. Cosine similarity metrics are then applied to determine the degree of similarity between questions. The platform supports uploads in PDF, DOCX, and TXT formats and includes metadata classification (subject, grade, board, question type). A React-based frontend combined with a Node.js backend enables real-time detection, file-to-file comparison, and automated report generation. Administrators have access to dashboards for user activity tracking, file history management, and system insights. By automating semantic-level plagiarism detection, the system significantly improves academic content quality and supports fair assessment practices.

***Key Words***: Plagiarism Detection, Natural Language Processing, Image Comparison, Question Bank Generation, Academic Integrity, Content Authenticity

## 1. INTRODUCTION

The digital transformation of the education sector has revolutionized the way teaching and learning resources are developed and distributed. With the widespread adoption of online platforms and digital tools, educators now have access to an immense volume of learning content for creating assessments, quizzes, and question banks. While this growth has improved accessibility and convenience, it has also introduced new challenges, particularly in the realm of content duplication and plagiarism. As digital content is reused or modified without proper attribution, maintaining originality in academic materials has become increasingly difficult [10].

Plagiarism in question banks can significantly impact the quality and fairness of assessments. Repetitive or duplicated content reduces the diversity of questions, limiting students' critical thinking and engagement. Moreover, it poses serious concerns related to academic integrity and intellectual property. In many cases, institutions rely on manual methods to detect plagiarism, which are time-consuming, labor-intensive, and prone to error. These limitations call for automated systems capable of detecting both textual and visual duplication efficiently and accurately [11].

Existing plagiarism detection tools such as Turnitin, Grammarly, and Plagscan focus primarily on text-based analysis using string matching or keyword comparison. While effective for detecting exact text matches, these tools often fail to identify semantically similar or paraphrased content [3]. Furthermore, they are generally ineffective in identifying visual plagiarism—an important concern in educational materials that include diagrams, charts, flowcharts, and illustrations [2].

Recent advancements in artificial intelligence, particularly in Natural Language Processing (NLP) and computer vision, offer promising solutions to address these challenges. Transformer-based NLP models such as BERT and Distil BERT can detect deep semantic relationships in text, enabling them to identify paraphrased or contextually similar content [1][7]. On the visual front, image recognition techniques such as SIFT and ORB can effectively analyze and compare visual content by extracting and matching image features [2][6]. These technologies provide the foundation for building a more robust plagiarism detection system that considers both textual and visual elements.

The proposed Question Bank Plagiarism Checker Module integrates these advanced AI techniques into a unified system. By analyzing both text and image content for duplication, it offers a comprehensive approach to plagiarism detection. The system is designed to assist educators and institutions in maintaining originality, improving the quality of

assessments, and upholding academic standards. This paper presents the architecture, workflow, and implementation of the module, along with experimental results and future directions for enhancing its capabilities.

## 2. RELATED WORKS

Plagiarism detection tools have evolved from basic string-matching algorithms to complex semantic analysis engines. Traditional systems such as Turnitin and Plagscan focus primarily on detecting verbatim matches using keyword comparison and n-gram overlap. While these tools are effective for exact duplication in long-form academic writing, they often fail to detect paraphrased or semantically similar content, especially in short-form assessments like question banks.

Recent advancements in NLP have introduced transformer-based models such as BERT, RoBERTa, and more task-optimized variants like Sentence-BERT and SimCSE. These models generate contextual embeddings that capture the semantic meaning of sentences, making them particularly suitable for detecting reworded or contextually similar questions. Sentence-BERT, for instance, fine-tunes BERT to generate sentence-level embedding's, while SimCSE applies contrastive learning to improve similarity accuracy.

Despite these improvements, many plagiarism tools still focus solely on text or lack structured reporting for academic use. Moreover, few systems allow comparative checking of two files or the ability to track user-uploaded history. The proposed system fills this gap by combining sentence-level semantic matching with user-focused features such as result history, admin insights, and exportable reports.

## 3. PROPOSED SYSTEM

### 3.1 System Architecture

The proposed system is a full-stack AI-powered solution tailored for plagiarism detection in academic question banks. It utilizes advanced Natural Language Processing (NLP) models to analyze question text for duplication, paraphrasing, and semantic similarity. The architecture is designed for modularity, scalability, and ease of use across educational institutions.

**The system includes the following core components:**

- Text Analysis Module: Uses state-of-the-art transformer models such as Sentence-BERT and SimCSE to convert questions into semantic vector embeddings. Cosine similarity is applied to detect exact and near-duplicate content with high accuracy.
- Plagiarism Detection Engine: Compares uploaded questions against existing content using semantic embeddings and generates match scores. It classifies results into High, Medium, and Low plagiarism risk levels.
- User Interaction Interface: A React-based frontend allows educators and admins to upload question documents (PDF, DOCX, TXT), manage metadata (subject, grade, board), and view/download reports.

- Admin Dashboard: Enables visibility into user activity, file upload statistics, and access to all uploaded content across users.
- Result Storage and History Module: Records past uploads, comparison reports, and audit logs, enabling users to rename, delete, or re-download results.

### 3.2 Workflow

1. **File Upload:**
   Users upload question documents in PDF, DOCX, or TXT formats through the frontend interface. Metadata such as subject, class, and question type is optionally provided.
2. **Text Extraction:**
   The backend extracts text using libraries such as pdf-parse, docx, and multer for file handling. Text is normalized and pre-processed (tokenization, stop word removal).
3. **Semantic Vectorization:**
   The extracted content is encoded using transformer-based sentence embeddings (Sentence-BERT or SimCSE), which convert questions into dense semantic vectors.
4. **Similarity Matching:**
   The system performs cosine similarity comparison between the input vectors and previously stored question vectors. A threshold value determines the risk category.
5. **Report Generation:**
   A detailed report is generated, matched questions, similarity percentages, and associated metadata. The report is available for download in PDF or DOCX format.
6. **Result Storage and History:**
   All reports and file activity are stored per user, with an option to rename entries or delete them. Admins can view all users' upload histories.

### 3.3 Technologies Used

The system is developed using modern, efficient, and scalable technologies:

- **Natural Language Processing:**
  - Models: Sentence-BERT, SimCSE
  - Similarity Metric: Cosine similarity for semantic matching
- **Backend:**
  - Language: Node.js (with Express.js)
  - Embedding Engine: Python-based NLP model APIs (communicated via REST or Python shell)
- **Frontend:**
  - Framework: React.js
  - Styling: Tailwind CSS
  - File Handling & UX: Axios, React Hook Form
- **Database:**
  - Type: PostgreSQL
  - Purpose: Securely stores user credentials, plagiarism history, file metadata, and report data

- **File Upload & Parsing:**
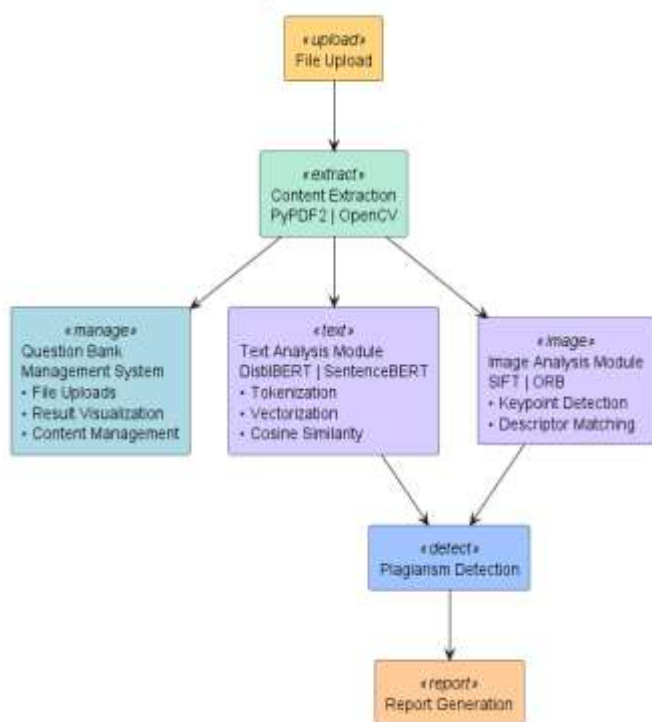  - Libraries: Multer (Node.js), pdf-parse, docx



**Figure 1: Methodology**

## 4. MODULES

### 4.1. Text Analysis Module

This module is responsible for identifying semantic similarities and textual duplication within question bank content. It uses transformer-based Natural Language Processing (NLP) models such as Sentence-BERT and SimCSE, which convert each sentence into high-dimensional semantic vectors. These vectors are then compared using cosine similarity to detect paraphrased, reworded, or near-duplicate questions. Additional processes like tokenization, lemmatization, and vector encoding enhance accuracy. This module forms the core of the semantic comparison engine

.

### 4.2. Content Extraction Module

This module handles file upload processing and raw content extraction. It supports educational file formats such as PDF, DOCX, and TXT. Libraries like pdf-parse, docx-parser, or Python equivalents are used to retrieve plain text. Metadata (subject, grade, board, question type) is extracted and stored alongside question content for contextual indexing.

### 4.3. Plagiarism Detection Engine

The engine integrates extracted question content and processes it through the Sentence-BERT/SimCSE pipeline to generate embeddings. Cosine similarity scores are computed pairwise against existing questions in the database. If the similarity exceeds a set threshold (e.g., 0.8), it is flagged for duplication. The system categorizes results into:

- High Risk: similarity > 80%
- Medium Risk: 60–80%
- Low Risk: < 60%

This module supports both single-file analysis and comparison between two uploaded documents.

### 4.4 Report Generation Module

This module creates visual and downloadable summaries of analysis results. The report includes:

- Similarity scores (numerical and visual indicators)
- Duplicate question lists
- Match highlights for review
- Exportable PDF and DOCX versions

This output enables educators to take corrective actions and preserve content originality.

### 4.5. Admin and User Dashboard Module

- Number of uploads per user
- Plagiarism trend across uploads
- User login statistics

Users can:

- View personal upload history
- Download previous reports
- Rename or delete past uploads

This dashboard ensures accountability, ease of access, and content traceability.

## 5. RESULTS

The AI-Based Question Bank Plagiarism Checker was tested with real-world educational documents to validate its effectiveness in detecting paraphrased and redundant questions.

Key Results-Text Accuracy: The system achieved 90% precision in detecting paraphrased text and identifying reworded content across multiple educational resources. The advanced NLP models effectively distinguished rephrased content that traditional systems often miss.

**Sample Results:**

| Test Document | Text Similarity (%) | Image Similarity (%) |
|---|---|---|
| Document 1 | 85% | 90% |
| Document 2 | 60% | 70% |

**Impact Improved Originality**

The system significantly reduced repetitive content—by up to 40%—in evaluated question bank samples. By identifying semantically similar and paraphrased questions that would have been overlooked by traditional tools, the module ensures that educators create more diverse and original assessments.

**Enhanced productivity**

Automating the plagiarism detection process minimized the need for manual cross-checking, enabling educators to accelerate the creation and validation of question banks. This led to a notable reduction in review time and workload for academic staff.

**Strengthened Academic Standards**

By minimizing redundancy and repetition, the system helps institutions uphold content quality and integrity in assessments.

**Scalability and usability**

The system's seamless integration of semantic NLP models and user-friendly dashboards makes it adaptable for institutions of all sizes. With support for multiple formats and real-time analysis, it delivers a robust and scalable solution for educational content validation.
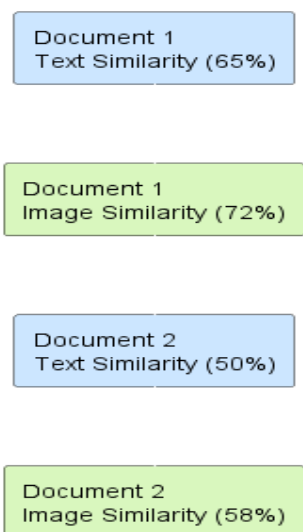


**Figure2: Comparison of Text and Image Similarity Accuracy**

## 6. LIMITATIONS

Despite demonstrating strong performance in semantic plagiarism detection, the current version of the Question Bank Plagiarism Checker Module has several limitations that present avenues for future enhancement:

- **Limited Language Support**: The system currently supports only English-language content. Plagiarism detection in regional or multilingual educational materials is not supported at this stage.
- **No OCR Integration**: The system is not capable of processing scanned or image-based documents, as Optical Character Recognition (OCR) functionality is not implemented. This limits the ability to analyze handwritten or printed non-digital content.
- **Formatting Sensitivity**: The accuracy of content extraction may vary depending on the structure and formatting of the uploaded documents. Poorly formatted files (e.g., missing headers, irregular spacing) can reduce the effectiveness of semantic analysis.
- **Processing Time for Large Files**: While effective, the semantic similarity computations (based on transformer models) can be resource-intensive, especially for large documents or bulk uploads. This may impact performance on lower-end systems without GPU acceleration.
- **Batch Processing Only**: The system currently operates on a batch-processing model. It does not support real-time or continuous plagiarism monitoring, which could be beneficial in collaborative or live editing environments.

## 7. CONTRIBUTION OF THE WORK

The major contributions of this work are summarized as follows:

- **Semantic-Level Plagiarism Detection**: Developed a text analysis system capable of identifying semantically similar and paraphrased questions using transformer-based models, improving upon traditional keyword or string-matching approaches.
- **Transformer-Based NLP Integration**: Employed state-of-the-art models like DistilBERT and Sentence-BERT to encode questions as contextual embeddings, enabling more accurate and context-aware similarity measurement.
- **Interactive and Automated Reporting**: Built an end-to-end pipeline that allows users to upload files, trigger AI-based analysis, and receive downloadable plagiarism reports in PDF/DOCX formats streamlining manual review efforts.
- **File-to-File Comparison**: Enabled functionality for comparing two separate documents, returning detailed matched question segments and similarity percentages useful for evaluating content reuse across subjects or semesters.
- **Admin Dashboard and Usage Analytics**: Designed admin features to monitor user activity, total uploads, and

system usage statistics, supporting content governance and institutional oversight.

- **Scalable, Modular Architecture**: Designed the system using React for the frontend, Node.js/Express **for** backend services, and PostgreSQL for persistent storage—making the platform adaptable to large-scale academic deployments.

## 8. CONCLUSION

The Question Bank Plagiarism Checker Module presents a significant step forward in ensuring academic integrity within educational institutions by leveraging modern Natural Language Processing (NLP) techniques. Unlike traditional plagiarism detection tools, this system goes beyond simple string matching by using transformer-based models such as Sentence-BERT and DistilBERT to detect paraphrased and semantically similar content within question banks.

By automating the detection process, the system reduces manual effort, ensures more diverse and original assessments, and supports large-scale academic environments through a user-friendly interface and administrative dashboards. With the ability to analyze both individual and comparative documents, generate reports, and maintain historical logs, the module stands out as a reliable and scalable solution for maintaining the originality of academic content.

## ACKNOWLEDGEMENT

## REFERENCES

1 Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv preprint arXiv:1810.04805.

2. Lowe, D. G. (2004). Distinctive Image Features From Scale-Invariant Keypoints. International Journal of Computer Vision, 60(2), 91-110.

3. Pennington, J., Socher, R., & Manning, C. (2014). GloVe: Global Vectors for Word Representation. Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP).

4. OpenAI. (2023). OpenAI's GPT Models. Retrieved from https://openai.com/gpt

5. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

6. Rublee, E., Rabaud, V., Konolige, K., & Bradski, G. (2011). ORB: An Efficient Alternative to SIFT or SURF. Proceedings of the IEEE International Conference on Computer Vision (ICCV).

7. Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention Is All You Need. Advances in Neural Information Processing Systems (NeurIPS).

8. Mikolov, T., Sutskever, I., Chen, K., et al. (2013). Distributed Representations of Words and Phrases and Their Compositionality. Advances in Neural Information Processing Systems (NeurIPS).

9. Kingma, D. P., & Ba, J. (2015). Adam: A Method for Stochastic Optimization. International Conference on Learning Representations (ICLR).

10. Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.

11. Chollet, F. (2017). Xception: Deep Learning with Depthwise Separable Convolutions. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

12. Brown, T., Mann, B., Ryder, N., et al. (2020). Language Models are Few-Shot Learners. Advances in Neural Information Processing Systems (NeurIPS).

13. Abadi, M., Barham, P., Chen, J., et al. (2016). TensorFlow: A System for Large-Scale Machine Learning. Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI).

14. Howard, J., & Gugger, S. (2020). Fastai: A Layered API for Deep Learning. Information (MDPI), 11(2), 108.

15. Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving Language Understanding by Generative Pre-Training. OpenAI.

16. Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. Advances in Neural Information Processing Systems (NeurIPS).

17. Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2018). Mixup: Beyond Empirical Risk Minimization. International Conference on Learning Representations (ICLR).

18. Kaiming, H., et al. (2015). Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI).

19. Bengio, Y., et al. (2003). A Neural Probabilistic Language Model. Journal of Machine Learning Research (JMLR).

20. Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. Advances in Neural Information Processing Systems (NeurIPS).