

# Intelligent Classification of Educational Questions Using XLNet-CNN and Deep Contextual Embeddings.

## *Abstract*

The automated classification of examination questions according to Bloom's Taxonomy (BT) assists question setters in developing high-quality assessments by accurately categorising questions into cognitive levels. While most previous studies in this area have employed traditional machine learning methods, relatively few have explored deep learning-based approaches. Contextual embeddings, in particular, have shown effectiveness across various natural language processing tasks. This study aims to evaluate a hybrid optimal pre-trained contextual word embedding technique, XLNet, combined with a Convolutional Neural Network (CNN) model tailored for BT-based question classification. To this end, the study examines the performance of the proposed XLNet+ CNN model with state-of-the-art models. Experimental results indicate that the XLNet + CNN model offers performance comparable to existing models, with the added advantage of higher precision at higher cognitive levels.

**Keywords:** XLNet, CNN, NLP, Bloom, Education technology, Deep learning, word embedding

## 1. INTRODUCTION

The classification of examination questions by cognitive level is an important part of educational assessment, which ensures that students' learning outcomes are effectively measured. Bloom's Taxonomy is a widely accepted framework in educational theory that provides a structured approach to categorising learning objectives into three primary domains, viz., *cognitive*, *affective*, and *psychomotor*. Among these, the cognitive domain deals with different levels of thinking skills, from basic knowledge recall to higher-order analysis, synthesis, and evaluation. These hierarchical levels provide guidelines for educators to design balanced assessments that align with students' cognitive abilities [1].

Instructors often rely on Bloom's Taxonomy to formulate exam questions. However, manually constructing questions that accurately reflect these levels can be time-consuming and subjective and might also lead to inconsistencies in assessing students' skills. Therefore, there has been increased interest in automating the classification of exam questions based on Bloom's Taxonomy [2].

Despite the potential benefits of automation, most studies have focused on traditional machine learning techniques and only a few have explored deep learning approaches [3]. Deep learning models based on pre-trained word embeddings have been proven successful in various natural language processing tasks, such as automating question classification [4].

37 This research work aims to address the challenge of automating the classification of exam  
38 questions by developing a deep learning model that leverages the pre-trained contextualised  
39 word embeddings XLNet, in conjunction with Convolutional Neural Networks (CNNs). CNNs  
40 are known for their ability to extract key features from text, and are well-suited for this task as  
41 they focus on identifying patterns within the question text that correspond to specific cognitive  
42 levels of Bloom's Taxonomy.

43 The paper is structured as follows: Section 2 presents the study's motivation and contribution.  
44 Related work is presented in Section 3. The methodology is discussed in Section 4. Section 5  
45 discusses the results, and Section provides the conclusion and future scope

46

## 47 **2 MOTIVATION AND CONTRIBUTION**

48

49 In today's fast-growing world, there is a requirement for high-quality educational assessments.  
50 Designing such assessments that align with Bloom's taxonomy presents significant  
51 challenges. Because it requires educators to carefully make  
52 the question paper that assesses various levels of student understanding according to the action  
53 verbs. Therefore, manually mapping the question papers to Bloom's taxonomy is difficult.

54 Given the requirements of higher education institutions and advancements in NLP and DL, there  
55 is a strong need for automated solutions to help teachers objectively classify examination  
56 questions [5]. This process of automation can save teachers' effort, improve consistency, and help  
57 to align learning outcomes and assessment strategies. The current work proposes a novel  
58 automated system for classifying examination questions based on Bloom's Taxonomy using  
59 advanced DL techniques. The key contributions are as follows:

60 1. To classify examination questions based on Bloom's Taxonomy, with the help of a novel  
61 hybrid deep learning approach that integrates contextual embeddings from XLNet with CNN.

62 2. The paper also compares the proposed model with the existing state-of-the-art models.

63 3. The proposed framework has been evaluated on the dataset and shows comparable  
64 performance.

## 65 **3. RELATED WORK**

66 With advancements in Natural Language Processing and Deep Learning, text classification tasks  
67 have gained significant prominence. Any text classification task relies on word embeddings,  
68 which transform textual data into numerical representations [6]. The literature suggests that  
69 earlier studies employed traditional Machine Learning (ML) approaches using feature selection  
70 techniques. However, with the emergence of deep learning models, neural network-based  
71 methods such as Convolutional Neural Networks (CNNs) have come to prominence in sentence  
72 classification tasks [7]. Furthermore, Yoon K (2014) highlighted that CNNs can automatically  
73 learn hierarchical features from text with minimal feature engineering [8]. CNNs are also

74 computationally efficient compared to sequential models such as LSTMs, which makes them  
75 suitable for large-scale applications [9].

76 In recent years, transformer-based models such as RoBERTa and XLNet have gained  
77 prominence. This is because they can generate contextualised word embeddings. These models  
78 are pre-trained on large corpora and can capture semantic and syntactic relationships more  
79 effectively than traditional embedding techniques like Word2Vec or GloVe [10].

80 Studies have shown that transformer-based embeddings significantly improve performance in  
81 various NLP tasks[11].In the context of education, many recent studies have explored the use of  
82 pre-trained models to classify examination questions according to Bloom's Taxonomy. These  
83 approaches emphasise the importance of selecting appropriate embedding techniques to capture  
84 contextual meaning and improve model effectiveness [12].

85 Despite these advancements,there are many research gaps. Many existing studies rely on either  
86 traditional embeddings or standalone deep learning models. The literature shows limited  
87 exploration of hybrid architectures that combine transformer-based embeddings with CNNs [13].

88 To address these limitations, recent research trends focus on integrating contextual embeddings  
89 from transformer models with efficient feature extraction mechanisms such as CNNs. Such  
90 hybrid approaches aim to leverage the strengths of both models, i.e., contextual understanding  
91 from transformers and spatial feature extraction from CNNs,to achieve higher accuracy and  
92 better generalisation.In this context, the present work builds upon prior studies by combining  
93 XLNet embeddings with CNN-based architectures to develop an efficient and accurate system  
94 for classifying examination questions according to Bloom's Taxonomy and comparing its  
95 performance with the state-of-the-art RoBERTa+CNN model as proposed by [12]

96

## 97 **4. METHODOLOGY**

98 This section presents the methodology for implementing the proposed framework.

99

### 100 **4.1. Dataset**

101 In this study, five datasets were utilised, comprising both previously established datasets and one  
102 collected specifically for this research. The dataset collected contained 1,200 questions, all of  
103 which were labelled by educators. After merging all five datasets, a comprehensive dataset of  
104 2,522 questions was formed.Ultimately, the proposed models were trained and evaluated  
105 exclusively on the combined dataset as given in [12] . Figure 1 presents the percentage of the  
106 number of questions against each Bloom's category.Figure 2 presents word clouds of various  
107 classes containing action verbs.

108

Proportion of Each Category

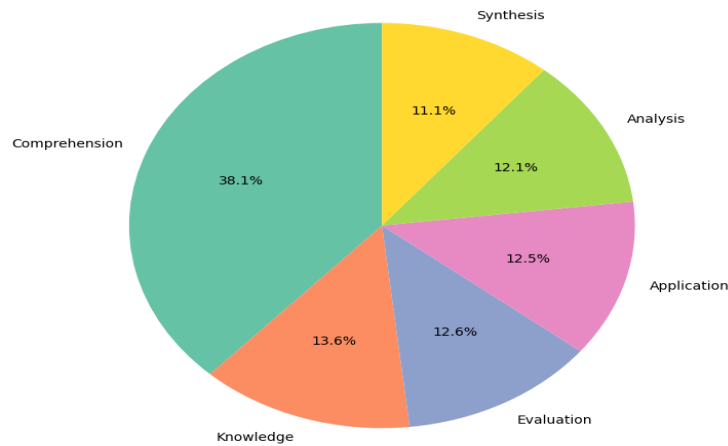
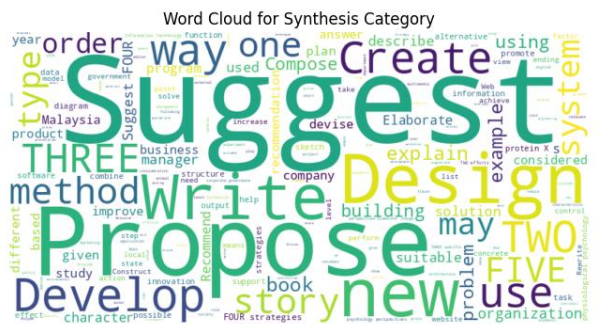
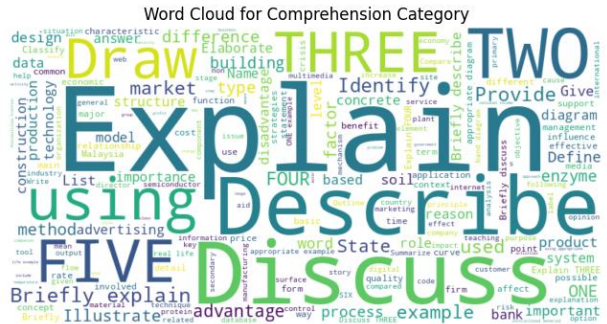


Figure 1. Percentage of Number of questions in each class in the dataset

109  
110  
111  
112



113

114  
115  
116  
117  
118

## Figure 2 Word cloud of various classes

### 4.2 Proposed Framework (XLNet+CNN)

Algorithm 1 below provides the key details of the proposed framework.

#### **Input:**

*Dataset of examination questions*

#### **Output:**

*Predicted Bloom's Taxonomy category*

#### **Step 1: Data Loading and Preparation**

1. *Load dataset from Excel file.*
2. *Assign numerical labels to Bloom's Taxonomy categories.*
3. *Preprocess text data (cleaning, normalization).*
4. *Tokenize text using XLNet tokenizer:*
  - *Convert text into input IDs*
  - *Generate attention masks*
  - *Set fixed sequence length (e.g., 256 tokens)*

#### **Step 2: Tokenization and Embedding Generation**

5. *Store tokenized outputs (input IDs and attention masks) in arrays.*
6. *Pass tokenized inputs to XLNet to generate contextual embeddings.*

#### **Step 3: Model Architecture Design**

7. *Input XLNet embeddings into Convolutional Neural Networks (CNN):*
  - *Apply Conv1D layers to extract features*
  - *Apply Max-Pooling layers to reduce dimensionality*
8. *Flatten the output of CNN layers.*
9. *Pass flattened output to Dense layer with ReLU activation.*
10. *Add output layer with Softmax activation to classify into six categories.*

#### **Step 4: Training and Validation**

11. *Split dataset into training and validation sets.*
12. *Compile model using:*
  - *Optimizer: Adam*
  - *Loss function: Categorical Cross-Entropy*

- *Metric: Accuracy*

*13. Train model for a fixed number of epochs (e.g., 5 epochs).*

*14. Store training and validation accuracy and loss.*

**Step 5: Evaluation and Classification**

*15. Evaluate model performance on validation dataset.*

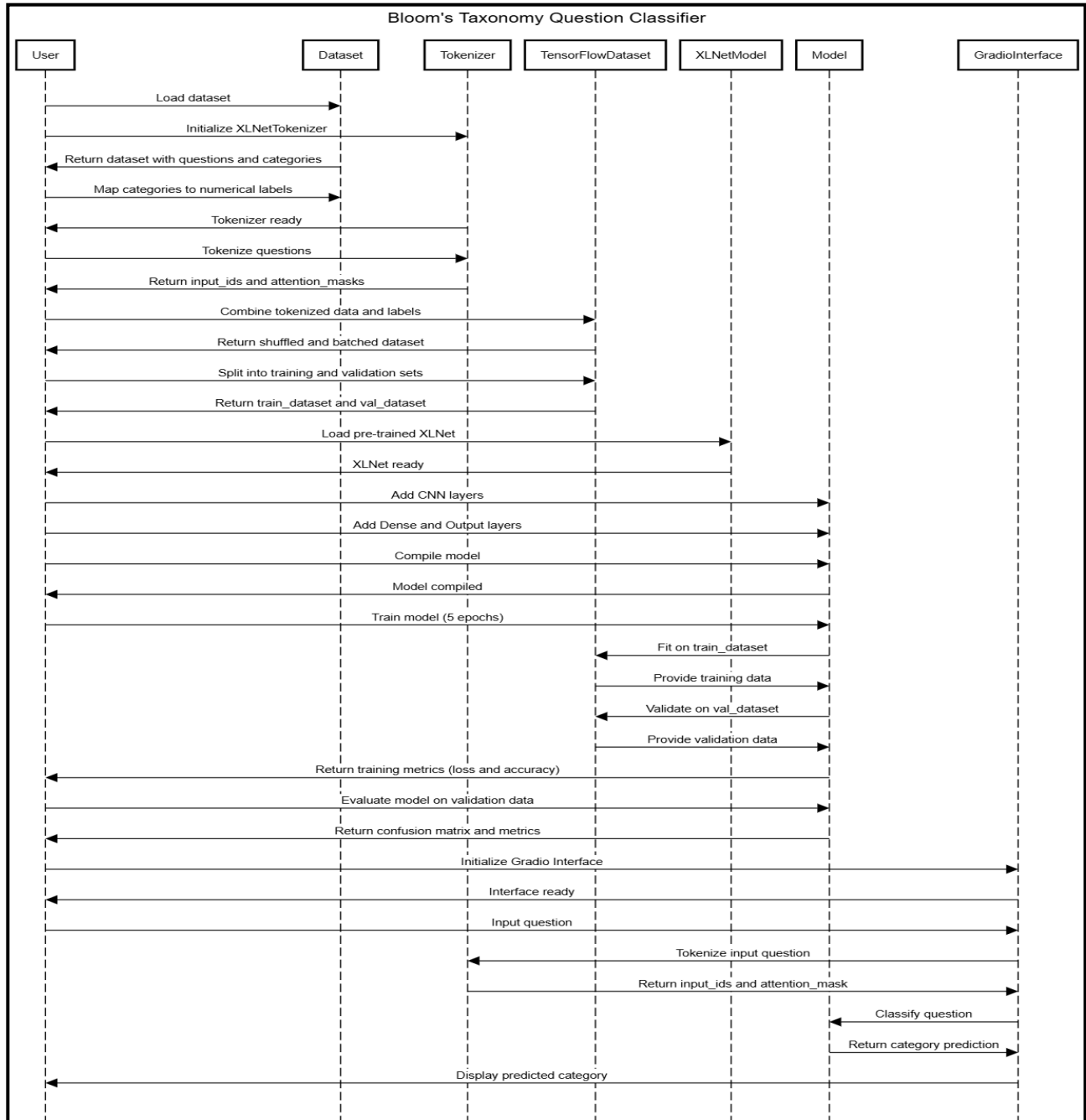
*16. Visualize accuracy and loss trends.*

*17. Define classification function:*

- *Tokenize new question using XLNet tokenizer*
- *Pass through trained model*
- *Obtain Softmax probabilities*
- *Select highest probability class*

*18. Map predicted label to Bloom's Taxonomy level.*

UNDER PEER REVIEW IN IJAR



119

120 **Figure 3: Sequence diagram of the proposed framework implementation**

121 Hyperparameter tuning for the proposed framework have been done as represented in Table 2

122 **Table 2: Hyperparameters for XLNetCNN:**

123

Hyperparameter	Value
Text Preprocessing	Tokenizer: XLNetTokenizer
Max Sequence Length	256
Input IDs	Shape: (batch size, 256)

<b>Attention Masks</b>	Shape: (batch_size, 256)
<b>Model Type</b>	XLNet (xlnet-base-cased)
<b>Embedding Layer</b>	TFXLNetModel
<b>Conv1D Layer Filters</b>	128
<b>Conv1D Kernel Size</b>	3
<b>Conv1D Activation</b>	ReLU
<b>MaxPooling1D Pool Size</b>	2
<b>Second Conv1D Filters</b>	64
<b>Second Conv1D Kernel Size</b>	3
<b>Second Conv1D Activation</b>	ReLU
<b>Fully Connected Layer</b>	Dense (256 units, ReLU)
<b>Output Layer</b>	Dense (6 units, Softmax)
<b>Optimizer</b>	Adam (learning_rate=1e-5)
<b>Loss Function</b>	Categorical Crossentropy
<b>Metrics</b>	Categorical Accuracy
<b>Batch Size</b>	16

124

## 125 5. RESULT AND DISCUSSION

126 The results of this study highlight the significant impact of word embedding techniques on the  
 127 classification of examination questions into Bloom's Taxonomy (BT) categories. **XLNet+CNN**  
 128 model was trained and validated on a dataset of 2,522 questions, with performance evaluated  
 129 using Accuracy and F1 Score metrics. The classification report is shown in Figure 4, and the  
 130 confusion matrix in Figure 5.

131

Classification Report:				
	precision	recall	f1-score	support
Knowledge	0.92	0.80	0.86	60
Comprehension	0.84	0.95	0.89	212
Application	0.79	0.63	0.70	60
Analysis	0.94	0.78	0.85	58
Synthesis	0.76	0.88	0.82	48
Evaluation	0.94	0.84	0.89	58
accuracy			0.85	496
macro avg	0.87	0.81	0.83	496
weighted avg	0.86	0.85	0.85	496

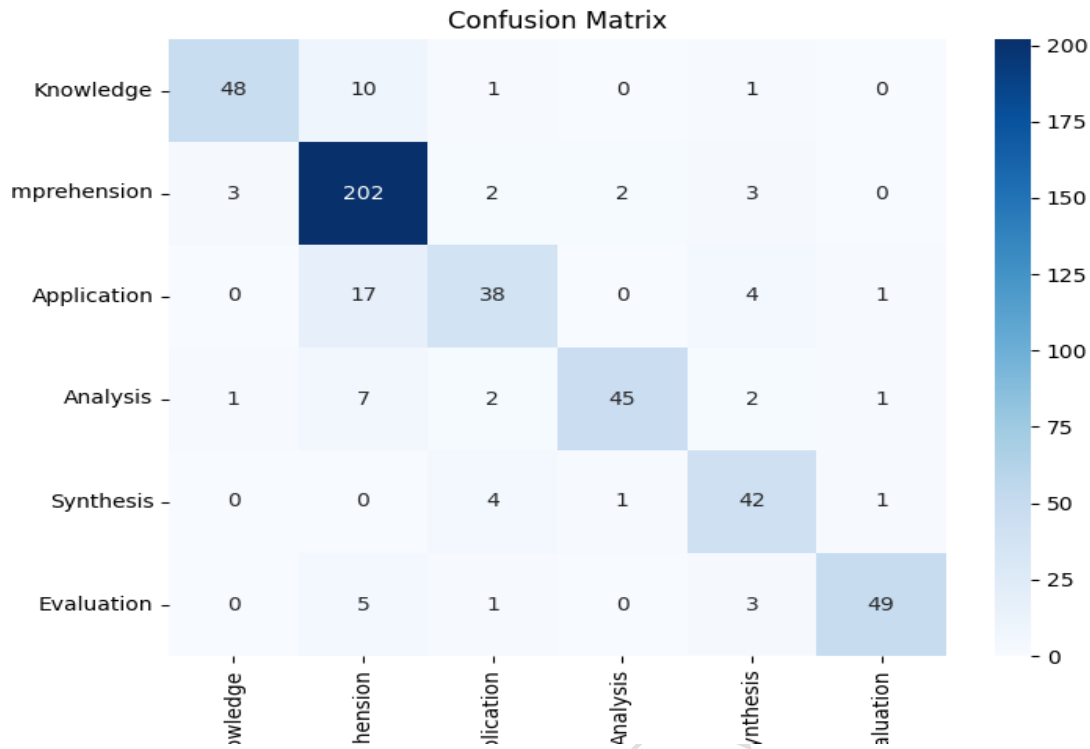
132

133

134

135

**Figure 4 Classification Report of the proposed Framework**



**Figure 5: Confusion matrix of the proposed Framework**

### 5.1. Comparative Performance Analysis

The performance of the proposed framework (XLNET +CNN) has been compared with the state-of-the-art best RoBERTa +CNN as proposed by. Table 3 presents the results below, and Table 4 presents the class-wise performance of both models in terms of precision, recall, and F1 Score.

**Table 3 Comparative Performance of the proposed framework**

Technique	Accuracy	F1-score
RobertaCnn (Gani,M.)	86.2%	86.1%
xlne+CNN (Proposed Model)	85.63%	84.03%

**Table 4: Cla-s wise performance**

Class	XLNet + CNN (P / R / F1)	RoBERTa + CNN (P / R / F1)
Knowledge	0.92 / 0.80 / 0.86	0.91 / 0.86 / 0.88
Comprehension	0.84 / 0.95 / 0.89	0.89 / 0.93 / 0.91
Application	0.79 / 0.63 / 0.70	0.85 / 0.72 / 0.78
Analysis	0.94 / 0.78 / 0.85	0.90 / 0.87 / 0.88

Class	XLNet + CNN (P / R / F1)	RoBERTa + CNN (P / R / F1)
Synthesis	0.76 / 0.88 / 0.82	0.72 / 0.93 / 0.81
Evaluation	0.94 / 0.84 / 0.89	0.86 / 0.75 / 0.80

146 The comparative analysis of the XLNet + CNN and RoBERTa + CNN models indicates that both  
147 approaches exhibit comparable performance in classifying examination questions according to  
148 Bloom's Taxonomy. While the RoBERTa + CNN model shows a slightly higher overall  
149 accuracy and F1-score, the difference in performance is marginal, suggesting that both models  
150 are equally effective for this task. XLNet + CNN model demonstrates certain advantages at the  
151 class level. It achieves higher precision in higher-order cognitive categories such as Analysis and  
152 Evaluation. Additionally, XLNet shows competitive performance in handling the complex  
153 semantic structures of higher-order thinking-skills questions due to its permutation-based  
154 contextual learning mechanism. Overall, the findings suggest that RoBERTa + CNN performs  
155 marginally better in terms of Accuracy and F1 Score. However, the XLNet + CNN model offers  
156 comparable performance while providing the added advantage of higher precision at higher  
157 cognitive levels.

158

## 159 **6. CONCLUSION AND FUTURE WORK:**

160 This study evaluated the performance of the proposed XLNet + CNN model for automatically  
161 classifying examination questions by Bloom's Taxonomy and compared it with the state-of-the-  
162 art RoBERTa + CNN model. The results indicated that both models achieved comparable  
163 performance. Integrating automated question-generation systems with educational paradigms  
164 advances intelligent assessment systems in education.

165 This study has certain limitations, including a relatively small dataset and class imbalance, which  
166 may affect the model's generalizability and bias the performance toward dominant  
167 categories. This approach hasn't considered the contextual metadata, as it was limited to textual  
168 features. The future work can focus on enhancing the robustness of the proposed approach. This  
169 can be done using larger and more diverse datasets. Also, incorporating explainable AI  
170 techniques will improve transparency and trust in the system.

171

## 172 **REFERENCES**

- 173 [1] Adams, N. E. (2015). Bloom's taxonomy of cognitive learning objectives. *Journal of the Medical*  
174 *Library Association: JMLA*, 103(3), 152.
- 175 [2] Banujan, K., Kumara, S., Prasanth, S., & Ravikumar, N. (2023). Revolutionising Educational  
176 Assessment: Automated Question Classification Using Bloom's Taxonomy and Deep Learning  
177 Techniques--A Case Study on Undergraduate Examination Questions. *International Journal of Education*  
178 *and Development using Information and Communication Technology*, 19(3), 259-278.
- 179 [3] Li, Q., Peng, H., Li, J., Xia, C., Yang, R., Sun, L., ... & He, L. (2022). A survey on text classification:  
180 From traditional to deep learning. *ACM Transactions on Intelligent Systems and Technology*  
181 *(TIST)*, 13(2), 1-41.
- 182 [4] Asudani, D. S., Nagwani, N. K., & Singh, P. (2023). Impact of word embedding models on text  
183 analytics in deep learning environment: a review. *Artificial intelligence review*, 56(9), 10345-10425.
- 184 [5] Alqahtani, T., Badreldin, H. A., Alrashed, M., Alshaya, A. I., Alghamdi, S. S., Bin Saleh, K., ...  
185 & Albekairy, A. M. (2023). The emergent role of artificial intelligence, natural learning processing, and  
186 large language models in higher education and research. *Research in social and administrative*  
187 *pharmacy*, 19(8), 1236-1242.

- 188 [6] Minaee, S., Kalchbrenner, N., Cambria, E., Nikzad, N., Chenaghlu, M., & Gao, J. (2021). Deep  
189 learning--based text classification: a comprehensive review. *ACM computing surveys (CSUR)*, 54(3), 1-  
190 40.
- 191 [7] Soni, S., Chouhan, S. S., & Rathore, S. S. (2023). TextConvoNet: A convolutional neural network  
192 based architecture for text classification. *Applied Intelligence*, 53(11), 14249-14268.
- 193 [8] Lee, Y., Yoon, S., & Jung, K. (2018, October). Comparative studies of detecting abusive language on  
194 twitter. In *Proceedings of the 2nd workshop on abusive language online (ALW2)* (pp. 101-106).
- 195 [9] Shen, L., Sun, Y., Yu, Z., Ding, L., Tian, X., & Tao, D. (2024). On efficient training of large-scale deep  
196 learning models. *ACM Computing Surveys*, 57(3), 1-36.
- 197 [10] Areshey, A., & Mathkour, H. (2024). Exploring transformer models for sentiment classification: A  
198 comparison of BERT, RoBERTa, ALBERT, DistilBERT, and XLNet. *Expert Systems*, 41(11), e13701.
- 199 [11] Naseem, U., Razzak, I., Musial, K., & Imran, M. (2020). Transformer based deep intelligent  
200 contextual embedding for twitter sentiment analysis. *Future Generation Computer Systems*, 113, 58-69.
- 201 [12] Gani, M. O., Ayyasamy, R. K., Sangodiah, A., & Fui, Y. T. (2023). Bloom's Taxonomy-based exam  
202 question classification: The outcome of CNN and optimal pre-trained word embedding  
203 technique. *Education and Information Technologies*, 28(12), 15893-15914.
- 204 [13] Wu, H., Liu, Y., & Wang, J. (2020). Review of Text Classification Methods on Deep  
205 Learning. *Computers, Materials & Continua*, 63(3).
- 206

UNDER PEER REVIEW