

# Identifying Nuances of Multi-Task Learning for Bengali and English Emotional Texts

Anonymous EMNLP submission

## Abstract

In this paper, we present a multi-task learning (MTL) model to classify sentiment and emotion in Bengali and English languages. For this multi-task learning work, different Bengali and English datasets were collected from publicly available sources and developed two MTL models by utilizing pre-trained mBERT and MuRIL models. Our proposed MTL model outperforms their corresponding standalone classifiers with an average F1-score of 0.5728 (+0.041) and 0.7590 (+0.046) for Bengali sentiment, emotion and English sentiment, emotion classification tasks respectively.

## 1 Introduction

Bengali is the 6<sup>th</sup> most popular language in the world spoken by over 200 million peoples<sup>1</sup>. Also, it is the second most spoken language Indian sub-continent and the national and most widely spoken language in Bangladesh. Furthermore, with the popularity of social media and the internet, the number of Bengali language-spoken users significantly increased in the past few years. As of 2023, the total internet users in Bangladesh were 66.94 million and among them, 44.7 million were social media users which is 26% of the total population (Kemp, 2023). In addition, the number of internet users increased by 22.33 million between 2021 and 2023 (Kemp, 2021).

Over the decades, with the advancement of machine learning and deep learning techniques, NLP methods can efficiently find sentiments and emotions in social media and other texts not only in English languages but also in low-resource languages such as Bengali. However, most of the research focused on only learning one task: either sentiment classification or emotion analysis. But to find both sentiment and emotion in a sentence or text, we have to execute two separate models which may increase overhead.

<sup>1</sup><https://salc.uchicago.edu/language-study/bengali>

Multi-task learning (MTL) as the name suggests, is a machine-learning technique that is capable of learning and handling multiple tasks at the same time. Researches show that, in the majority of cases, MTL models perform significantly better than their corresponding standalone classifiers for similar kinds of tasks.

In this present article, we focused on developing a multi-task learning framework for sentiment and emotion classification for Bengali and English texts and analyzed the performances of MTL models with their corresponding standalone classifiers. The main contributions of this paper can be summarized as follows:

- We presented two MTL schemes by using transformer-based pre-trained multilingual-BERT (mBERT) (Devlin et al., 2019) and MuRIL (Khanuja et al., 2021) models and compared the performances with standalone classifiers for Bengali and English languages.
- Our proposed MTL models provide a superior result than their corresponding standalone classifiers for both Bengali and English languages.

## 2 Related Work

The concept of MTL was first proposed by (Caruana, 1997). Since then, MTL approaches have been efficiently used in different domains of computer science including NLP.

(Liu et al., 2016) proposed three LSTM-based MTL frameworks for text classification where two models were unidirectional LSTM based and the third model was bidirectional LSTM based with two classification heads in each MTL framework. The authors evaluated their models for some popular NLP datasets such as ‘SST-1’, ‘IMBD’, etc. and their proposed MTL framework provide better performances over single-task learning.

078 An adversarial MTL framework was proposed  
079 by (Liu et al., 2017) primarily using LSTMs and  
080 the authors achieved a better performance in their  
081 proposed adversarial MTL model on 16 different  
082 datasets.

083 (Majumder et al., 2019), (Tan et al., 2023) and  
084 (Savini and Caragea, 2020) proposed MTL frame-  
085 works for sentiment and sarcasm classification.  
086 (Majumder et al., 2019) used GRU-based archi-  
087 tecture and attention mechanism to classify sen-  
088 timent and sarcasm whereas (Tan et al., 2023)  
089 and (Savini and Caragea, 2020) used BiLSTM  
090 in their study. In addition, (Savini and Caragea,  
091 2020) used a non-contextual pre-trained embed-  
092 ding FastText (Bojanowski et al., 2016), which  
093 elevates their performance. Another sentiment  
094 and sarcasm analysis MTL work was proposed by  
095 (El Mahdaouy et al., 2021) using the pre-trained  
096 BERT (Devlin et al., 2019) model where the au-  
097 thors focused only on Arabic languages.

098 (Singh et al., 2022) proposed an MTL architec-  
099 ture for emoji, sentiment and emotion analysis by  
100 using the xlm-RoBERTa-base (Liu et al., 2019)  
101 and their proposed multi-task learning classifiers  
102 provide better performances than standalone clas-  
103 sifiers. The authors also analyse sentiment and  
104 emotion intensities in their studies along with the  
105 classification tasks.

106 An MTL framework was for sentiment, emo-  
107 tion, target analysis (targeting a specific commu-  
108 nity such as black people, women, LGBT etc.),  
109 hate speech and offensive language classification  
110 by (Del Arco et al., 2021) using the BERT (Devlin  
111 et al., 2019) models.

112 In this present article, we proposed a multi-task  
113 learning framework for sentiment and emotion in  
114 Bengali and English text. As per our literature, this  
115 type of work is new and not widely explored in the  
116 context of the Bengali language.

### 117 3 Dataset

118 For this MTL work, we prepared two separate  
119 datasets: one was annotated with three sentiment  
120 labels (positive, negative and neutral) and another  
121 was annotated with six emotion labels (anger, fear,  
122 happy, sad, disgust and surprise).

123 For the emotion dataset, we collected 6314 sam-  
124 ples from the ‘BanglaEmotion’ dataset (Rahman,  
125 Md Aatur, 2020). The ‘BanglaEmotion’ dataset  
126 was prepared from the users’ comments from two  
127 different Facebook groups in Bangladesh and an-

128 notated each comment with one of six emotion la-  
129 bels: angry, fear, happy, sad, disgust and surprise.

130 Keeping these six emotion labels in mind, we  
131 extended this dataset to another 6314 English  
132 texts collected from ‘GoEmotion’ (Demszky et al.,  
133 2020) and ‘emotion\_dataset’ (Saravia et al., 2018).

134 For the Sentiment dataset, the Bengali texts  
135 were collected from the ‘SentNoB’ dataset (Islam  
136 et al., 2021). The ‘SentNoB’ dataset was pre-  
137 pared from the social media users’ comments on  
138 news and videos with a sample size of around 15K,  
139 annotated each comment with one of three senti-  
140 ment labels: positive, negative, and neutral. How-  
141 ever, we considered only 9519 samples from this  
142 dataset.

143 After that, the subset of the ‘SentNoB’ dataset  
144 was extended to the English texts, collected from  
145 the ‘tweet\_sentiment\_multilingual’ dataset (Barbi-  
146 eri et al., 2022) with around 3.03K records.

147 Next, 10% of Bengali and English texts were  
148 split out from the sentiment and emotion dataset  
149 and preserved for testing purposes. The data dis-  
150 tributions for the sentiment and emotion datasets  
151 are provided in Figures 1 and 2 respectively.

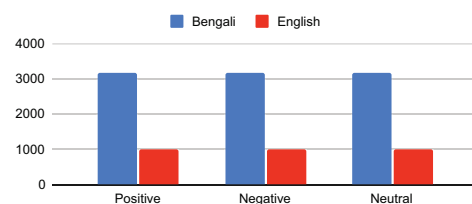


Figure 1: Distribution of Bengali and English lan-  
guages in Sentiment dataset

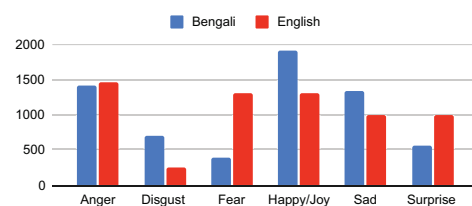


Figure 2: Distribution of Bengali and English lan-  
guages in Emotion dataset

### 152 4 Methodology

153 This section discusses the methodologies of our  
154 proposed work. We aim to develop a multi-task  
155 learning framework that can classify sentiments  
156 and emotions at a time time in Bengali and En-  
157 glish texts. To do that we took the help of the

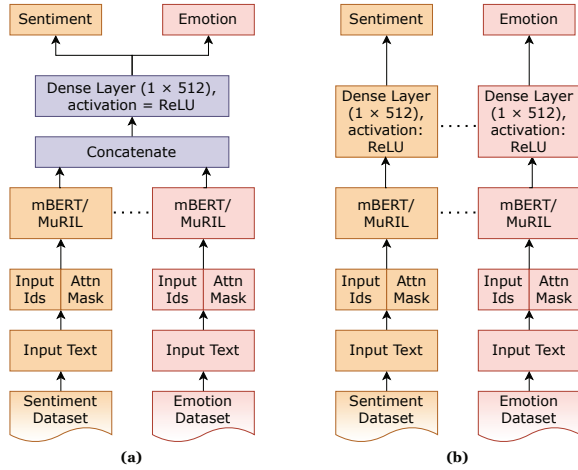


Figure 3: Proposed MTL models. (a) MTL architecture with a shared dense layer, (b) MTL architecture with task-specific dense layers.

pre-trained multilingual BERT (mBERT) (Devlin et al., 2019) and MuRIL (Multilingual Representations for Indian Languages) (Khanuja et al., 2021) models and developed the model architecture. The proposed MTL models are provided in Figure 3.

**Tokenization:** Before going to the actual training process, the training input text were converted into a sequence of tokens. For this tokenization process, we used the pre-trained mBERT and MuRIL tokenizers for their corresponding models with a maximum sequence length of 256. Each tokenizer returns a sequence of input ids and attention masks which were fed into the pre-trained models as depicted in Figure 3.

**Model Selection:** For this MTL work, we had chosen mBERT (Devlin et al., 2019) and MuRIL (Khanuja et al., 2021) models. The reason behind this, the mBERT is a transformer-based model and was trained on a total of 104 languages including Bengali and English. So, it can learn both Bengali and English contexts in a sentence or text. On the other hand, the MuRIL follows the BERT-base architecture and was trained on 17 Indian languages and this model provides better results than mBERT in different benchmark datasets. In addition, since MuRIL is specifically trained on 17 Indian languages, it may learn Indian languages (such as Bengali) in more better way than mBERT.

Next, after passing the input ids and attention masks to the mBERT/MuRIL models the pooled output of the last layer was fed into a Dense layer. In this stage, two experiments were performed: one is a shared dense layer, where the pooled out-

put from two mBERT/MuRIL models was concatenated and then fed into a single dense layer of 512 neurons. We say this as  $MTL_{shared}$ . In the second approach, instead of concatenating, the pooled output of two mBERT/MuRIL models was fed into two separate dense layers of 512 neurons. We say this as  $MTL_{task-specific}$ . All the dense layers used ReLU as their activation function.

**Classification:** For the  $MTL_{shared}$  approach the output of the shared dense layer was fed into the sentiment and emotion output layers (Figure 3(a)) and for the  $MTL_{task-specific}$  approach the outputs of the separate dense layers fed into each sentiment and emotion output layers (Figure 3(b)) with 3 and 6 hidden units for sentiment and emotion classification heads respectively. All the output layers used the softmax as their activation function.

**Training:** Before beginning the training process we randomly split the training dataset into 9:1 ratio where 90% data were used for training and 10% data were used as validation split.

We trained our proposed model up to 4 epochs with a learning rate of  $2e-5$  and for the multi-task loss function, the SparseCategoricalCrossentropy loss function was used and monitored the loss for the validation split of the training dataset.

$$L_{total} = L_{sentiment} + L_{emotion}$$

Where  $L_{sentiment}$  and  $L_{emotion}$  represent the loss for sentiment and emotion tasks. The hyperparameters that were used to train the model are provided in Table 1.

Parameter	Value
Dropout	0.1
Loss function	SparseCategoricalCrossentropy
Optimizer	Adam (Kingma and Ba, 2014)
Learning rate	$2e-5$
Epoch	4
Batch size	8

Table 1: Hyperparameters used to train the model

## 5 Experiment and Result

### 5.1 Experimental Setup

All the experiments were executed using the libraries of ‘TensorFlow’ and ‘Keras’ and the pre-trained tokenizers and models were imported from the ‘HuggingFace’.

To evaluate the performances of proposed MTL models, for sentiment classification, we passed

Task	Model	STL				MTL <sub>shared</sub>				MTL <sub>task-specific</sub>			
		Bengali		English		Bengali		English		Bengali		English	
		Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Sentiment	mBERT	0.6296	0.6331	0.6142	0.6125	0.6227	0.6148	0.6142	0.6124	0.6458	0.6482	0.6296	0.6271
	MuRIL	0.6609	0.6639	0.6265	0.6278	<b>0.6806</b>	<b>0.6803</b>	0.6420	0.6411	0.6655	0.6683	<b>0.6852</b>	<b>0.6858</b>
Emotion	mBERT	0.4873	0.3319	0.8244	0.7985	0.5056	0.3485	<b>0.8560</b>	<b>0.8323</b>	0.5015	0.3514	0.8354	0.8182
	MuRIL	0.5142	0.3998	0.8339	0.8115	<b>0.557</b>	<b>0.4654</b>	0.8528	0.8215	0.5348	0.4139	0.8528	0.8342

Table 2: Performance comparison of STL vs proposed scheme of MTLs for both Bengali and English languages on the test dataset. (All the F1-scores provided here are the macro F1-scores.)

separately Bengali and English sentiment test data and compared the accuracy and F1-scores with standalone classifiers (STL), and the same was done for the emotion classification task.

## 5.2 Result

The results for STL, MTL<sub>shared</sub> and MTL<sub>task-specific</sub> models are provided in Table 2 for both Bengali and English language. From Table 2 we can see that our proposed MTL frameworks outperform to their corresponding STL models for both Bengali and English test datasets. For Bengali sentiment classification, we see a performance improvement (F1-score) of 2.41%, as well as for English sentiment classification, the MTL<sub>task-specific</sub> model shows an improvement (F1-score) of 8.46% than their corresponding STL model. Additionally, our sentiment classification result provides a better F1-score than the best result provided by (Islam et al., 2021)<sup>2</sup>.

For the emotion classification task, the MTL models show a performance improvement (F1-score) of 14.15% and 4.06% with respect to the STL models for Bengali and English emotion test datasets respectively. In addition, our MTL<sub>shared</sub> Bengali emotion classification result provides an improvement in accuracy by 4.88% and an improvement in F1-score by 28.58% than the best accuracy and F1-score provided by (Rahman and Seddiqui, 2019)<sup>3</sup>.

Furthermore, It can also be observed that the MuRIL models provide better performance than the mBERT models for the Bengali language. This is because the MuRIL model trained on specifically 17 Indian languages, so it can learn Bengali contexts in more better way than the mBERT model which was trained on 104 languages.

Additionally, if we closely observe the results,

<sup>2</sup>We considered only a subset of the full dataset with sample size around 9.5K, the author’s original dataset was with a sample size of around 15K.

<sup>3</sup>The authors recorded 0.5298 and 0.3324 as their best accuracy and F1-score respectively.

we can see that, the MTL<sub>shared</sub> approach learns emotions in a better way than the MTL<sub>task-specific</sub> approach. On the other hand for sentiment classification results, the MTL<sub>shared</sub> gives a better performance for Bengali languages. However, in the context of English languages, the MTL<sub>shared</sub> model failed to provide a good result whereas the MuRIL-based MTL<sub>task-specific</sub> model provides superior results than MTL<sub>shared</sub> and STL approaches.

## 6 Conclusion

In this paper, we proposed an MTL framework for sentiment and emotion classification in Bengali and English languages by transformer-based pre-trained models mBERT and MuRIL and our proposed MTL models outperform their corresponding STL models in both Bengali and English languages. Also, for the Bengali language, the MuRIL-based MTL models perform better than the mBERT-based MTL models.

In future, we’ll expand our existing dataset to make a more robust model. Additionally, we’ll consider other Indian languages such as Hindi and CodeMixed texts (mixing of Bengali and English languages) in future works.

## 7 Limitations

Our proposed work also has some limitations. Firstly, the dataset size in this experiment is relatively small (around 12K samples per dataset). Secondly, we have considered only two multilingual models: mBERT and MuRIL. There are also more available multilingual models such as xlm-RoBERTa (Liu et al., 2019) or IndicBERT (Kakwani et al., 2020) etc., and we’ll explore them in the future. Also, not considering the ‘large’ models such as ‘mBERT-large’ or ‘MuRIL-large’ is one of the limitations of this work. Thirdly, we have considered only Bengali and English languages for this MTL work. And fourth, we only performed our experiments with a batch size of 8 and did not perform the experiments with a higher



307	batch size (16, 32, 64, etc.) due to resource limita-	Simon Kemp. 2021. <a href="#">Digital in Bangladesh: All the</a>	362
308	tions.	<a href="#">statistics you need in 2021</a> DataReportal Global	363
		<a href="#">Digital Insights</a> .	364
309	<b>References</b>	Simon Kemp. 2023. <a href="#">Digital 2023: Bangladesh</a>	365
310	Francesco Barbieri, Luis Espinosa Anke, and Jose	<a href="#">DataReportal Global Digital Insights</a> .	366
311	Camacho-Collados. 2022. <a href="#">XLM-T: Multilingual</a>	Simran Khanuja, Diksha Bansal, Sarvesh Mehtani,	367
312	<a href="#">language models in Twitter for sentiment analysis</a>	Savya Khosla, Atreyee Dey, Balaji Gopalan,	368
313	<a href="#">and beyond</a> . In <i>Proceedings of the Thirteenth Lan-</i>	Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja	369
314	<i>guage Resources and Evaluation Conference</i> , pages	Nagipogu, Shachi Dave, Shruti Gupta, Subhash	370
315	258–266, Marseille, France. European Language Re-	Chandra Bose Gali, Vish Subramanian, and Partha	371
316	sources Association.	Talukdar. 2021. <a href="#">Muril: Multilingual representations</a>	372
317	Piotr Bojanowski, Edouard Grave, Armand Joulin,	<a href="#">for indian languages</a> .	373
318	and Tomas Mikolov. 2016. Enriching word vec-	Diederik P. Kingma and Jimmy Ba. 2014. <a href="#">Adam: A</a>	374
319	tors with subword information. <i>arXiv preprint</i>	<a href="#">method for stochastic optimization</a> . <i>arXiv (Cornell</i>	375
320	<i>arXiv:1607.04606</i> .	<i>University)</i> .	376
321	Rich Caruana. 1997. <a href="#">Multitask Learning</a> . <i>Machine</i>	Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2016.	377
322	<i>learning</i> , 28(1):41–75.	<a href="#">Recurrent neural network for text classification with</a>	378
323	Flor Miriam Plaza Del Arco, Sercan Halat, Sebastian	<a href="#">multi-task learning</a> . <i>arXiv (Cornell University)</i> ,	379
324	Padó, and Roman Klinger. 2021. <a href="#">Multi-Task Learn-</a>	pages 2873–2879.	380
325	<a href="#">ing with Sentiment, Emotion, and Target Detection</a>	Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2017.	381
326	<a href="#">to Recognize Hate Speech and Offensive Language</a> .	<a href="#">Adversarial multi-task learning for text classifica-</a>	382
327	<i>arXiv (Cornell University)</i> .	<a href="#">tion</a> . In <i>Proceedings of the 55th Annual Meeting</i>	383
328	Dorottya Demszky, Dana Movshovitz-Attias, Jeong-	<i>of the Association for Computational Linguistics</i>	384
329	woo Ko, Alan Cowen, Gaurav Nemade, and Sujith	<i>(Volume 1: Long Papers)</i> , pages 1–10, Vancouver,	385
330	Ravi. 2020. <a href="#">Goemotions: A dataset of fine-grained</a>	Canada. Association for Computational Linguistics.	386
331	<a href="#">emotions</a> .	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	387
332	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	388
333	Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of</a>	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	389
334	<a href="#">deep bidirectional transformers for language under-</a>	<a href="#">Roberta: A robustly optimized BERT pretraining ap-</a>	390
335	<a href="#">standing</a> . In <i>Proceedings of the 2019 Conference</i>	<a href="#">proach</a> . <i>CoRR</i> , abs/1907.11692.	391
336	<i>of the North American Chapter of the Association</i>	Navonil Majumder, Soujanya Poria, Haiyun Peng, Niy-	392
337	<i>for Computational Linguistics: Human Language</i>	ati Chhaya, Zhaoxia Wang, and Alexander Gel-	393
338	<i>Technologies, Volume 1 (Long and Short Papers)</i> ,	bukh. 2019. <a href="#">Sentiment and sarcasm classification</a>	394
339	pages 4171–4186, Minneapolis, Minnesota. Associ-	<a href="#">with multitask learning</a> . <i>IEEE Intelligent Systems</i> ,	395
340	ation for Computational Linguistics.	34(3):38–43.	396
341	Abdelkader El Mahdaouy, Abdellah El Mekki, Ka-	Md. Ataur Rahman and Md. Hanif Seddiqui. 2019.	397
342	bil Essefar, Nabil El Mamoun, Ismail Berrada, and	<a href="#">Comparison of classical machine learning ap-</a>	398
343	Ahmed Khoumsi. 2021. <a href="#">Deep multi-task model for</a>	<a href="#">proaches on bangla textual emotion analysis</a> .	399
344	<a href="#">sarcasm detection and sentiment analysis in Arabic</a>	Rahman, Md Ataur. 2020. <a href="#">Banglaemotion: A bench-</a>	400
345	<a href="#">language</a> . In <i>Proceedings of the Sixth Arabic Natu-</i>	<a href="#">mark dataset for bangla textual emotion analysis</a> .	401
346	<i>ral Language Processing Workshop</i> , pages 334–339,	Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang,	402
347	Kyiv, Ukraine (Virtual). Association for Computa-	Junlin Wu, and Yi-Shin Chen. 2018. <a href="#">CARER: Con-</a>	403
348	tional Linguistics.	<a href="#">textualized affect representations for emotion recog-</a>	404
349	Khondoker Ittehadul Islam, Sudipta Kar, Md Saiful Is-	<a href="#">nition</a> . In <i>Proceedings of the 2018 Conference on</i>	405
350	lam, and Mohammad Ruhul Amin. 2021. <a href="#">SentNoB:</a>	<i>Empirical Methods in Natural Language Processing</i> ,	406
351	<a href="#">A dataset for analysing sentiment on noisy Bangla</a>	pages 3687–3697, Brussels, Belgium. Association	407
352	<a href="#">texts</a> . In <i>Findings of the Association for Computa-</i>	for Computational Linguistics.	408
353	<i>tional Linguistics: EMNLP 2021</i> , pages 3265–3271,	Edoardo Savini and Cornelia Caragea. 2020. <a href="#">A multi-</a>	409
354	Punta Cana, Dominican Republic. Association for	<a href="#">task learning approach to sarcasm detection (student</a>	410
355	Computational Linguistics.	<a href="#">abstract)</a> . <i>Proceedings of the AAAI Conference on</i>	411
356	Divyanshu Kakwani, Anoop Kunchukuttan, Satish	<i>Artificial Intelligence</i> , 34(10):13907–13908.	412
357	Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M.	Gopendra Vikram Singh, Dushyant Singh Chauhan,	413
358	Khapra, and Pratyush Kumar. 2020. IndicNLPsuite:	Mauajama Firdaus, Asif Ekbal, and Pushpak Bhat-	414
359	Monolingual Corpora, Evaluation Benchmarks and	tacharyya. 2022. <a href="#">Are emoji, sentiment, and emotion</a>	415
360	Pre-trained Multilingual Language Models for In-		
361	dian Languages. In <i>Findings of EMNLP</i> .		

416 Friends? a multi-task learning for emoji, sentiment,  
417 and emotion analysis. In *Proceedings of the 36th Pa-*  
418 *cific Asia Conference on Language, Information and*  
419 *Computation*, pages 166–174, Manila, Philippines.  
420 Association for Computational Linguistics.

421 Yik Yang Tan, Chee-Onn Chow, Jeevan Kanesan,  
422 Joon Huang Chuah, and YongLiang Lim. 2023. *Sen-*  
423 *timent Analysis and Sarcasm Detection using Deep*  
424 *Multi-Task Learning*. *Wireless Personal Communi-*  
425 *cations*, 129(3):2213–2237.