

Revisiting Paraphrase Question Generator using Pairwise Discriminator

Badri N. Patro¹, Dev Chauhan², Vinod K. Kurmi¹, Vinay P. Namboodiri²

¹ *Department of Electrical Engineering, Indian Institute of Technology Kanpur, India*

² *Department of Computer Science and Engineering, Indian Institute of Technology Kanpur, India*

Github Link: <https://github.com/dev-chauhan/PQG-pytorch>

Abstract

In this paper, we propose a method for obtaining sentence-level embeddings. While the problem of securing word-level embeddings is very well studied, we propose a novel method for obtaining sentence-level embeddings. This is obtained by a simple method in the context of solving the paraphrase generation task. If we use a sequential encoder-decoder model for generating paraphrase, we would like the generated paraphrase to be semantically close to the original sentence. One way to ensure this is by adding constraints for true paraphrase embeddings to be close and unrelated paraphrase candidate sentence embeddings to be far. This is ensured by using a sequential pair-wise discriminator that shares weights with the encoder that is trained with a suitable loss function. Our loss function penalizes paraphrase sentence embedding distances from being too large. This loss is used in combination with a sequential encoder-decoder network. We also validated our method by evaluating the obtained embeddings for a sentiment analysis task. The proposed method results in semantic embeddings and outperforms the state-of-the-art on the paraphrase generation and sentiment analysis task on standard datasets. These results are also shown to be statistically significant.

Keywords: Sentiment Analysis, GAN, VQA, Paraphrase, Gquestion generation, LSTM, Pairwise, Adversarial learning, Discriminator

Email address: (badri, devgiri, vinodkk, vinaypn)@iitk.ac.in (Badri N. Patro¹, Dev Chauhan², Vinod K. Kurmi¹, Vinay P. Namboodiri²)

1. Introduction

The problem of obtaining a semantic embedding for a sentence that ensures that the related sentences are closer and unrelated sentences are farther lies at the core of understanding languages. This is a too challenging task to obtaining and improving embedding for input text sequence. This would be relevant for a wide variety of machine reading comprehension and related tasks, such as sentiment analysis. Towards this problem, we propose a supervised method that uses a sequential encoder-decoder framework for paraphrase generation. The task of generating paraphrases is closely related to the task of obtaining semantic sentence embeddings. In our approach, we aim to ensure that the generated paraphrase embedding should be close to the corresponding true sentence and far from unrelated sentences. The embeddings so obtained help us to obtain state-of-the-art results for paraphrase generation task.

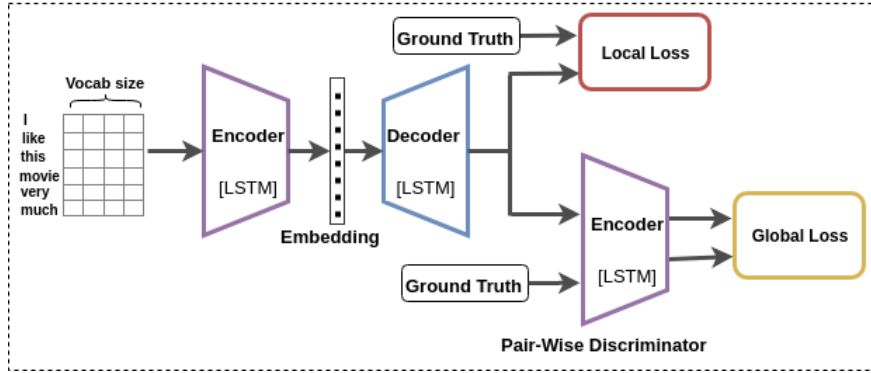


Figure 1: Pairwise Discriminator based Encoder-Decoder for Paraphrase Generation: This is the basic outline of our model which consists of an LSTM encoder, decoder and discriminator. Here the encoders share the weights. The discriminator generates discriminative embeddings for the Ground Truth-Generated paraphrase pair with the help of ‘global’ loss. Our model is jointly trained with the help of a ‘local’ and ‘global’ loss which we describe in section 3.

In this work, we proposed a pair-wise loss function for the task of paraphrase question generator, which will bring similar structure sentences close to each other as compared to the dissimilar type of sentences. In this work, we use local cross-entropy loss to generate each word in the sentence and global pair-wise discriminator loss to capture complete sentence structure in the given paraphrase sentences. Our

model consists of a sequential encoder-decoder that is further trained using a pair-wise discriminator. The encoder-decoder architecture has been widely used for machine translation and machine comprehension tasks. In general, the model ensures a ‘local’ loss that is incurred for each recurrent unit cell. It only ensures that a particular word token is present at an appropriate place. This, however, does not imply that the whole sentence is correctly generated. To ensure that the whole sentence is correctly encoded, we make further use of a pair-wise discriminator that encodes the whole sentence and obtains an embedding for it. We further ensure that this is close to the desired ground-truth embeddings while being far from other (sentences in the corpus) embeddings. This model thus provides a ‘global’ loss that ensures the sentence embedding as a whole is close to other semantically related sentence embeddings. This is illustrated in Figure 1. We further evaluate the validity of the sentence embeddings by using them for the task of sentiment analysis. We observe that the proposed sentence embeddings result in state-of-the-art performance for both these tasks. In this work, we use standard datasets like Quora Question Pair (QQP) dataset for paraphrase question generation task and Stanford Sentiment Treebank (SST) dataset for sentiment analysis task. In table-1, we show the various state of the art methods and its various attributes for paraphrase generation task. We observe that the contribution of various losses in the different state of the art method. In our method, we use pair-wise discriminator loss and show the improvement over the adversarial loss, KL-divergence loss, reinforcement loss.

Our contributions are as follows:

- We propose a model for obtaining sentence embeddings for solving the paraphrase generation task using a pair-wise discriminator loss added to an encoder-decoder network.
- We show that these embeddings can also be used for the sentiment analysis task.
- We validate the model using standard datasets (QQP and SST) with a detailed comparison with state-of-the-art methods and also ensure that the results are statistically significant.

Methods	Base Model	Adversarial	Loss	Task	Dataset
Seq-to-Seq [1]	deterministic	✗	CE	Paraphrase	COCO
Attention [2]	deterministic	✗	CE	Paraphrase	COCO
Residual LSTM [3]	deterministic	✗	CE	Paraphrase	COCO
VAE [4]	probabilistic	✗	CE, KL	Paraphrase	QQP, COCO
RbM-SL [5]	deterministic	✗	CE, RL	Paraphrase	QQP, Twitter
VAE-M [6]	probabilistic	✗	CE, KL	Paraphrase	QQP, COCO
EDL (Ours)	deterministic	✗	CE	Paraphrase, SA	QQP, SST
EDLPG (Ours)	deterministic	✓	CE, AD, PA	Paraphrase, SA	QQP
EDLPGS (Ours)	deterministic	✓	CE, AD, PA	Paraphrase, SA	QQP
EDLP (Ours)	deterministic	✗	CE, PA	Paraphrase, SA	QQP, SST
EDLPS (Ours)	deterministic	✗	CE, PA	Paraphrase, SA	QQP, SST

Table 1: Overview of various Paraphrase Question Generation (PQG) methods and their various properties. AD: Adversarial, CE: Cross Entropy, PA: Pairwise SA: Sentiment Analysis, RL: Reinforcement Learning, KL: KL divergence Loss Learning

2. Related Work

Given the flexibility and diversity of natural language, it has been a challenging task to represent text efficiently. There have been several hypotheses proposed for representing the same. [7, 8, 9] proposed a distribution hypothesis to represent words, i.e., words which occur in the same context have similar meanings. One popular hypothesis is the bag-of-words (BOW) or Vector Space Model [10], in which a text (such as a sentence or a document) is represented as the bag (multiset) of its words. [11] proposed an extended distributional hypothesis and [12, 13] proposed a latent relation hypothesis, in which a pair of words that co-occur in similar patterns tend to have similar semantic relation. Word2Vec [14, 15, 16] is also a popular method for representing every unique word in the corpus in a vector space. Here, the embedding of every word is predicted based on its context (surrounding words). NLP researchers have also proposed phrase-level, and sentence-level representations [17, 18, 19, 20, 15]. [21, 22, 23, 24, 25] have analyzed several approaches to represent sentences and phrases by a weighted average of all the words in the sentence, combining the word vectors in an order given by a parse tree of a sentence and by using matrix-vector operations. The primary issue with BOW models and weighted averaging of word vectors is the loss of semantic meaning of the words, the parse tree approaches can only work for sentences because of its dependence on sentence parsing mechanism. [26, 27] proposed a method to obtain a

vector representation for paragraphs and use it for some text-understanding problems like sentiment analysis and information retrieval.

Many language models have been proposed for obtaining better text embeddings in machine translation [1, 28, 29, 30], question generation [31], dialogue generation [32, 33, 34], document summarization [35], text generation [36, 37, 38, 39, 40, 41] and question answering [42, 43]. For paraphrase generation task, Prakash *et al.* [3] have generated paraphrases using stacked residual LSTM based network. Hasan *et al.* [44] proposed a encoder-decoder framework for this task. Gupta *et al.* [4] explored a VAE approach to generate paraphrase sentences using recurrent neural networks. Li *et al.* [5] used reinforcement learning for paraphrase generation task. Very recently, Yang *et al.* [6] has proposed another variational method for generating paraphrase questions.

In our previous work [45], we propose a pairwise discriminator based method to generation paraphrase questions. In this work, we extend our previous work by analyzing other variants of our model, like adversarial learning (EDLPG), as described in section-4.1.2. In section-4.1.3, we compare our method with the latest state of the art methods. In this work, we visualize the performance of different variants of our model over various epochs, as in section-4.1.4. We provide more qualitative results in both paraphrase question generation task and sentiment analysis task in section-4.1.5 and section-4.2.5. In section-4.1.1, we provide more detail about the QQP dataset, and we provide a few examples of this dataset in table-5.

3. Method

In this paper, we propose a text representation method for sentences based on an encoder-decoder framework using a pairwise discriminator for paraphrase generation and then fine-tune these embeddings for sentiment analysis tasks. Our model is an extension of *seq2seq* [1] model for learning better text embeddings.

3.1. Overview

Task: In the paraphrase generation problem, given an input sequence of words $X = [x_1, \dots, x_L]$, we need to generate another output sequence of words $Y = [q_1, \dots, q_T]$ that

has the same meaning as X . Here L and T are not fixed constants. Our training data consists of M pairs of paraphrases $\{(X_i, Y_i)\}_{i=1}^M$ where X_i and Y_i are the paraphrase of each other.

Our method consists of three modules, as illustrated in Figure 2: first is a Text Encoder which consists of LSTM layers, second is LSTM-based Text Decoder, and the last one is an LSTM-based Discriminator module. These are shown respectively in part 1, 2, 3 of Figure 2. Our network with all three parts is trained end-to-end. The weight parameters of encoder and discriminator modules are shared. Instead of taking a separate discriminator, we shared it with the encoder so that it learns the embedding based on the ‘global’ as well as ‘local’ loss. After training, at test time, we used encoder to generate feature maps and pass it to the decoder for generating paraphrases. These text embeddings can be further used for other NLP tasks, such as sentiment analysis.

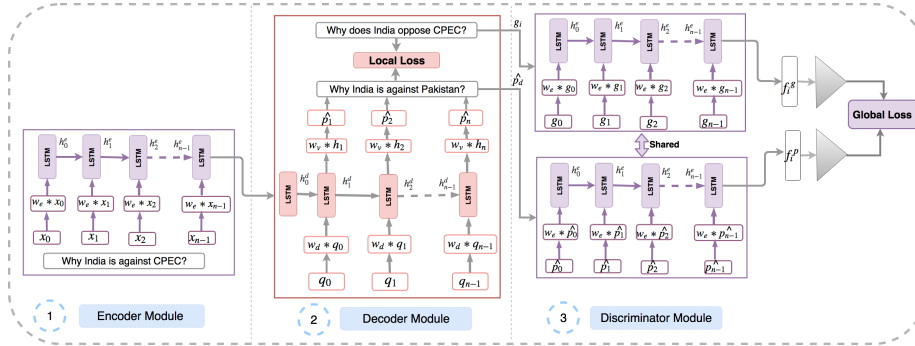


Figure 2: This is an overview of our model. It consists of 3 parts: 1) LSTM-based Encoder module which encodes a given sentence, 2) LSTM-based Decoder Module which generates natural language paraphrases from the encoded embeddings and 3) LSTM-based pairwise Discriminator module which shares its weights with the Encoder module and this whole network is trained with local and global loss.

3.2. Encoder-LSTM

We use an LSTM-based encoder to obtain a representation for the input question X_i , which is represented as a matrix in which every row corresponds to the vector representation of each word. We use a one-hot vector representation for every word

and obtain a word embedding c_i for each word using a Temporal CNN [46, 47] module that we parameterize through a function $G(X_i, W_e)$ where W_e are the weights of the temporal CNN. Now this word embedding is fed to an LSTM-based encoder which provides encoding features of the sentence. We use LSTM [48] due to its capability of capturing long term memory [47]. As the words are propagated through the network, the network collects more and more semantic information about the sentence. When the network reaches the last word (L_{th} word), the hidden state h_L of the network provides a semantic representation of the whole sentence conditioned on all the previously generated words (q_0, q_1, \dots, q_t) . Question sentence encoding feature f_i is obtained after passing through an LSTM which is parameterized using the function $F(C_i, W_l)$ where W_l are the weights of the LSTM. This is illustrated in part 1 of Figure 2.

3.3. Decoder-LSTM

The role of decoder is to predict the probability for a whole sentence, given the embedding of input sentence (f_i). RNN provides a nice way to condition on previous state value using a fixed length hidden vector. The conditional probability of a sentence token at a particular time step is modeled using an LSTM as used in machine translation [1]. At time step t , the conditional probability is denoted by $P(q_t|f_i, q_0, \dots, q_{t-1}) = P(q_t|f_i, h_t)$, where h_t is the hidden state of the LSTM cell at time step t . h_t is conditioned on all the previously generated words $(q_0, q_1, \dots, q_{t-1})$ and q_t is the next generated word.

Generated question sentence feature $\hat{p}_d = \{\hat{p}_1, \dots, \hat{p}_T\}$ is obtained by decoder LSTM which is parameterized using the function $D(f_i, W_{dl})$ where W_{dl} are the weights of the decoder LSTM. The output of the word with maximum probability in decoder LSTM cell at step k is input to the LSTM cell at step $k + 1$ as shown in Figure 2. At $t = -1$, we are feeding the embedding of input sentence obtained by the encoder module. $\hat{Y}_i = \{\hat{q}_0, \hat{q}_1, \dots, \hat{q}_{T+1}\}$ are the predicted question tokens for the input X_i . Here, we are using \hat{q}_0 and \hat{q}_{T+1} as the special START and STOP token respectively. The predicted question token (\hat{q}_i) is obtained by applying Softmax on the probability distribution \hat{p}_i . The question tokens at different time steps are given by the following

equations where LSTM refers to the standard LSTM cell equations:

$$\begin{aligned}
d_{-1} &= \text{Encoder}(f_i) \\
h_0 &= \text{LSTM}(d_{-1}) \\
d_t &= W_d * q_t, \forall t \in \{0, 1, 2, \dots, T-1\} \\
h_{t+1} &= \text{LSTM}(d_t, h_t), \forall t \in \{0, 1, 2, \dots, T-1\} \\
\hat{p}_{t+1} &= W_v * h_{t+1} \\
\hat{q}_{t+1} &= \text{Softmax}(\hat{p}_{t+1}) \\
\text{Loss}_{t+1} &= \text{loss}(\hat{q}_{t+1}, q_{t+1})
\end{aligned} \tag{1}$$

Where \hat{q}_{t+1} is the predicted question token and q_{t+1} is the ground truth one. In order to capture local label information, we use the Cross Entropy loss which is given by the following equation:

$$L_{local} = \frac{-1}{T} \sum_{t=1}^T q_t \log P(\hat{q}_t | q_0, \dots, q_{t-1}) \tag{2}$$

Here T is the total number of sentence tokens, $P(\hat{q}_t | q_0, \dots, q_{t-1})$ is the predicted probability of the sentence token, q_t is the ground truth token.

3.4. Discriminative-LSTM

The aim of the Discriminative-LSTM is to make the predicted sentence embedding f_i^p and ground truth sentence embedding f_i^g indistinguishable as shown in Figure 2. Here we pass \hat{p}_d to the shared encoder-LSTM to obtain f_i^p and also the ground truth sentence to the shared encoder-LSTM to obtain f_i^g . The discriminator module estimates a loss function between the generated and ground truth paraphrases. Typically, the discriminator is a binary classifier loss, but here we use a global loss, similar to [41] which acts on the last hidden state of the recurrent neural network (LSTM). The main objective of this loss is to bring the generated paraphrase embeddings closer to its ground truth paraphrase embeddings and farther from the other ground truth paraphrase embeddings (other sentences in the batch). Here our discriminator network ensures that the generated embedding can reproduce better paraphrases. We are using the idea of sharing discriminator parameters with encoder network, to enforce learning of

embeddings that not only minimize the local loss (cross entropy), but also the global loss.

Suppose the predicted embeddings of a batch is $e_p = [f_1^p, f_2^p, \dots, f_N^p]^T$, where f_i^p is the sentence embedding of i^{th} sentence of the batch. Similarly ground truth batch embeddings are $e_g = [f_1^g, f_2^g, \dots, f_N^g]^T$, where N is the batch size, $f_i^p \in \mathcal{R}^d$ $f_i^g \in \mathcal{R}^d$. The objective of global loss is to maximize the similarity between predicted sentence f_i^p with the ground truth sentence f_i^g of i^{th} sentence and minimize the similarity between i^{th} predicted sentence, f_i^p , with j^{th} ground truth sentence, f_j^g , in the batch. The loss is defined as

$$L_{global} = \sum_{i=1}^N \sum_{j=1}^N \max(0, ((f_i^p \cdot f_j^g) - (f_i^p \cdot f_i^g) + 1)) \quad (3)$$

Gradient of this loss function is given by

$$\left(\frac{dL}{de_p} \right)_i = \sum_{j=1, j \neq i}^N (f_j^g - f_i^g) \quad (4)$$

$$\left(\frac{dL}{de_g} \right)_i = \sum_{j=1, j \neq i}^N (f_j^p - f_i^p) \quad (5)$$

3.5. Cost function

Our objective is to minimize the total loss, that is the sum of local loss and global loss over all training examples in QQP dataset. The total loss is:

$$L_{total} = \frac{1}{M} \sum_{i=1}^M (L_{local} + L_{global}) \quad (6)$$

Where M is the total number of examples, L_{local} is the cross entropy loss, L_{global} is the global loss.

Dataset	# Train	# Validation	# Test
QQP-I	50k	5.2k	30k
QQP-II	100k	5.2k	30k

Table 2: Statistics of QQP dataset

Model	BLEU_1	BLEU_2	BLEU_3	BLEU_4	ROUGE_L	METEOR	CIDEr
EDP	0.0012	0.0000	0.0000	0.0000	0.00519	0.0056	0.0001
EDG	0.0012	0.0000	0.0000	0.0000	0.0019	0.0071	0.0001
EDPG	0.0015	0.0000	0.0000	0.0000	0.0023	0.0071	0.0002
EDL	0.4162	0.2578	0.1724	0.1219	0.4191	0.3244	0.6189
EDLPG	0.4152	0.2569	0.1725	0.1223	0.4168	0.3235	0.6081
EDLPGS	0.4177	0.2554	0.1695	0.1191	0.4206	0.3244	0.6468
EDLP	0.4370	0.2785	0.1846	0.1354	0.4399	0.3305	0.8723
EDLPS	0.4754	0.3160	0.2249	0.1672	0.4781	0.3488	1.0949

Figure 3: Analysis of variants of our proposed method on QQP Dataset as mentioned in section 4.1.2. Summary results from different models for **100k** dataset. Here L and P refer to the Local and Pairwise discriminator loss and S represents the parameter sharing between the discriminator and encoder module. G represents adversarial loss. ED represents Encoder-Decoder network. As we can see that our proposed method EDLPS clearly outperforms the other ablations on all metrics and detailed analysis is present in section 4.1.2.

Model	BLEU_1	BLEU_2	BLEU_3	BLEU_4	ROUGE_L	METEOR	CIDEr
EDG	0.0012	0.0000	0.0000	0.0000	0.0019	0.007	0.0001
EDPG	0.0012	0.0000	0.0000	0.0000	0.0019	0.0071	0.0001
EDL	0.3877	0.2336	0.1532	0.1067	0.3913	0.3133	0.4550
EDLPG	0.3823	0.2281	0.1487	0.1028	0.3847	0.3113	0.4322
EDLPGS	0.3956	0.2373	0.1552	0.1077	0.3997	0.3156	0.4945
EDLP	0.4159	0.2511	0.1683	0.1188	0.4079	0.3200	0.5431
EDLPS	0.4553	0.2981	0.2105	0.1560	0.4583	0.3421	0.9690

Figure 4: Analysis of variants of our proposed method on QQP Dataset as mentioned in section 4.1.2. Summary results from different models for **50k** dataset. Here L and P refer to the Local and Pairwise discriminator loss and S represents the parameter sharing between the discriminator and encoder module. G represents adversarial loss. ED represents Encoder-Decoder network. As we can see that our proposed method EDLPS clearly outperforms the other ablations on all metrics and detailed analysis is present in section 4.1.2.

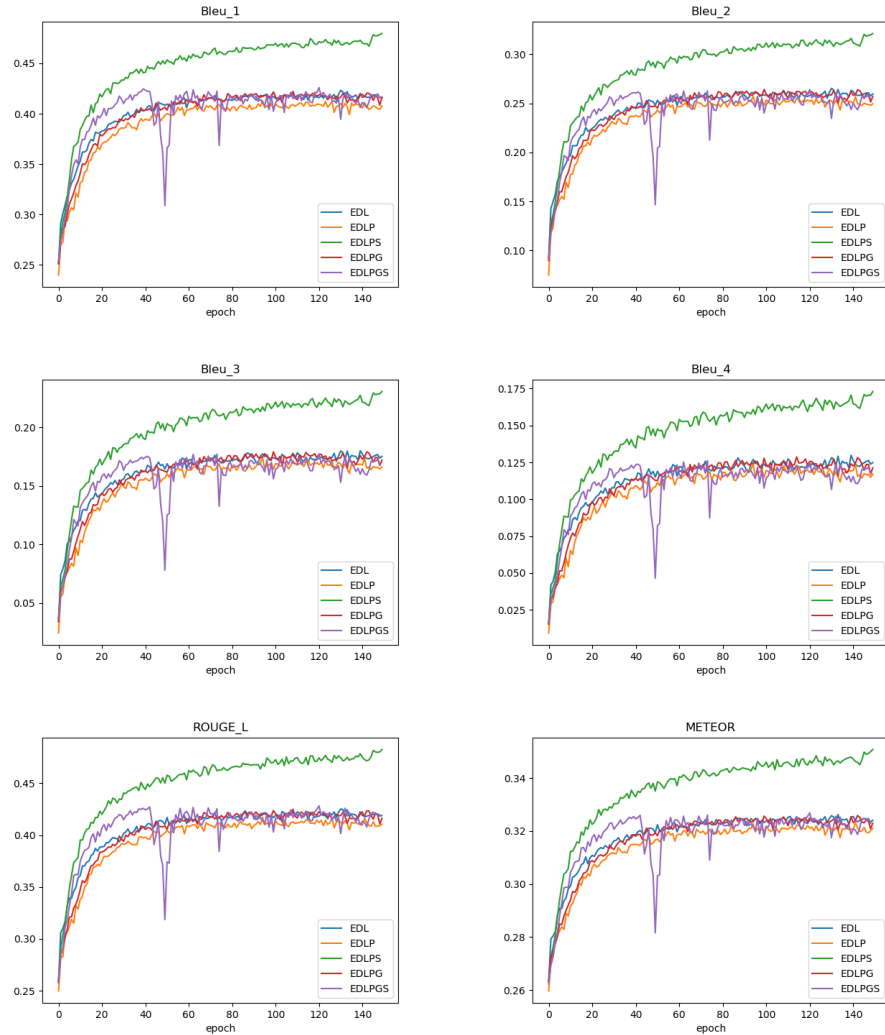


Table 3: This figure shows performance score of proposed models across various epochs for QQP-II (100k) dataset. The EDLPS model gives best performance among all the models across all the performance scores.

4. Experiments

We perform experiments to better understand the behavior of our proposed embeddings. To achieve this, we benchmark Encoder-Decoder with shared discriminator using Local-Global Pairwise loss (EDLPS) embeddings on two text understanding problems, Paraphrase Generation, and Sentiment Analysis. We use the Quora Question Pairs

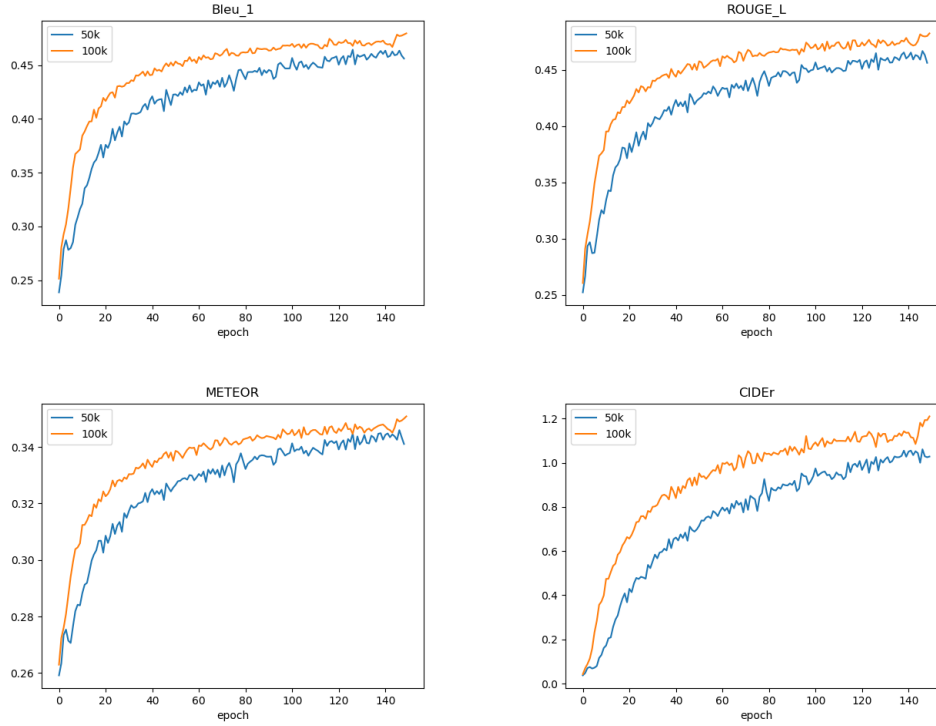


Figure 5: This figure shows performance of the model on QQP-I (50k) and QQP-II (100k) datasets. We observe that on 100k dataset model performs better than 50k dataset over all the scores.

(QQP) dataset ¹ for paraphrase generation and Stanford Sentiment Treebank (SST) dataset [26] for sentiment analysis. In this section, we describe the different datasets, experimental setup, and results of our experiments.

4.1. Paraphrase Generation Task

Paraphrase generation is an important problem in many NLP applications such as question answering, information retrieval, information extraction, and summarization. It involves the generation of similar meaning sentences. Paraphrase generation depends on how much meaningful information can be captured through the encoding scheme.

¹website: <https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs>

Dataset	Model	BLEU1 (\uparrow)	METEOR (\uparrow)	TER (\downarrow)
50K	VAE-S [4]	11.9	17.4	69.4
	VAE-SVG-eq [4]	17.4	21.4	61.9
	Seq2Seq + Att	26.0	20.3	-
	Residual LSTM	27.3	22.3	-
	RbM-SL [5]	35.8	28.1	-
	VAE-M[6]	37.1	24.0	61.4
	VAE-B[6]	38.2	22.5	56.6
	EDLP(Ours)	41.5	32.0	51.0
	EDLPS(Ours)	45.5	34.2	50.8
100K	VAE-S [4]	17.5	21.6	67.1
	VAE-SVG-eq [4]	22.9	24.7	55.0
	Seq2Seq + Att	36.5	26.2	-
	Residual LSTM	37.3	28.1	-
	RbM-SL [5]	43.5	32.8	-
	VAE-B[6]	45.0	31.4	55.5
	EDLP(Ours)	43.7	33.0	48.3
	EDLPS(Ours)	47.7	34.8	47.5

Table 4: Analysis of Baselines and State-of-the-Art methods for paraphrase generation on Quora dataset. As we can see clearly that our model outperforms the state-of-the-art methods by a significant margin in terms of BLEU and TER scores. Detailed analysis is present in section 4.1.3. A lower TER score is better whereas for the other metrics, a higher score is better. The Best results across baseline models are in bold letter. We put “-” for no results mentioned.

4.1.1. Dataset

We use the newly released Quora Question Pairs (QQP) dataset for this task. The QQP dataset is a newly release dataset for the paraphrase benchmark so far. The main principle objective to built for this dataset is to consider a different question for each logically distinct question. There is a total of 400k question pairs present the dataset, out of which, 149k are potential paraphrase. We split this dataset into two different ways, such as QQP-I and QQP-II result in two experiments. We have trained our model using

Id	Qid1	Qid2	Question1	Question2	Is
578	1154	1155	How do I manage time for studies?	How do you manage time between work and study?	0
591	1180	1181	How do I be a boyfriend?	How does a girl get a boyfriend?	0
715	1426	1427	How magnets are made?	What are magnets made of?	0
603	1204	1205	How do I find the phenotypic ratio?	What is a phenotype ratio?	0
592	1182	1183	How is time travel possible?	Do you think time travel is possible?	1
705	1406	1407	How can I consult a good free online astrologer?	Are there any good free online astrologers?	1
726	1448	1449	What is the meaning and purpose to life?	What is the exact meaning of life?	1
755	1505	1506	What is the ultimate way to serve humanity?	How can one serve humanity?	1

Table 5: This table shows sample example are present in this dataset.Each pair have question id (Qid1), its paraphrase question id (Qid2) and the ground truth label (Is) whether it is paraphrase or not. Is tends for “Is Duplicate”.

50k paraphrase question pair for QQP-I and 100k for QQP-II and validate our model performance on a validation set of 5k question pair. Finally, we evaluated the model performance on a test set of 30k question pair as pointed out in [5]. The question pairs having the binary value 1 are the ones which are the paraphrase of each other, and the others are duplicate questions. The QQP dataset has IDs for each question in the pair, the full text for each question, and a binary value that indicates whether the questions in the pair are truly a duplicate of each-other. Wherever the binary value is 1, the question in the pair are not identical; they are rather paraphrases of each other. So, we choose all such question pairs with binary value 1. There are a total of 149K such questions. We mention our data split statistic in table-2. Some examples of Quora Question Paraphrase pairs are provided in Table 5.

S.No	Original Question	Ground Truth Paraphrase	Generated Paraphrase
1	Is university really worth it?	Is college even worth it?	Is college really worth it?
2	Why India is against CPEC?	Why does India oppose CPEC?	Why India is against Pakistan?
3	How can I find investors for my tech startup?	How can I find investors for my startup on Quora?	How can I find investors for my startup business?
4	What is your view/opinion about surgical strike by the Indian Army?	What world nations think about the surgical strike on POK launch pads and what is the reaction of Pakistan?	What is your opinion about the surgical strike on Kashmir like?
5	What will be Hillary Clinton's strategy for India if she becomes US President?	What would be Hillary Clinton's foreign policy towards India if elected as the President of United States?	What will be Hillary Clinton's policy towards India if she becomes president?

Table 6: Examples of Paraphrase generation on Quora Dataset. We observe that our model is able to understand abbreviations as well and then ask questions on the basis of that as is the case in the second example.

4.1.2. Ablation Analysis

We experimented with different variations for our proposed method. We start with a baseline model, which we take as a simple encoder and decoder network with only the local loss (EDL) as proposed by Sutskever *et al.* [1]. Further, we have experimented with encoder-decoder and a discriminator network with only global loss (EDP) to distinguish the ground truth paraphrase with the predicted one. Another variation of our model is used both the global and local loss (EDLP). The discriminator is the same as our proposed method, only the weight sharing is absent in this case. Another variation of our model is EDLPG, which has a discriminator and is trained alternately between local and global loss, similar to GAN training. Finally, we make the discriminator share weights with the encoder and train this network with both the losses (EDLPS). The analyses are given in table 3. Among the ablations, the proposed EDLPS method works way better than the other variants in terms of BLEU and METEOR metrics by achieving an improvement of 7% and 3% in the scores respectively over the baseline method for

QQP-I (50k) dataset and an improvement of 6% and 2% in the scores respectively for QQP-II (100k) dataset.

4.1.3. Baseline and state-of-the-art Method Analysis

There has been relatively less work on this dataset, and the only work which we came across was that of [4]. We experimented with a simple Encoder-Decoder (EDL) framework [29, 1] for this task and chose that as our baseline. In this method, we use a LSTM encoder to encode the questions and then a LSTM decoder to generate the paraphrase question. The results of these variations are present in table 3. We compare our result with five baseline models² such as variational auto-encoder (VAE-SVG-eq)[4] model, Seq2Seq + attention [2] model, Residual LSTM [3] model, a deep reinforcement learning approach (RbM-SL)[5] and another VAE based generative architecture approach (VAE-B) [6] as provided in table 4 We further compare our best method EDLPS model with VAE-B[6], which is the current state-of-the-art on the QQP dataset. As we can see from the table that we achieve a significant improvement of 7.3% in BLEU1 score compare with VAE-B [6] , 6.1% in METEOR score compare with RbM-SL [5] and 5.8% in TER score (A lower TER score is better) compare with VAE-B [6] for 50K dataset and similarly 2.7% in BLEU1 score compare with VAE-B [6], 2.4% in METEOR score compare with RbM-SL [5] and 7.5% in TER score compare with VAE-SVG-eq [4] for 100K dataset.

4.1.4. Performance visualisation for PQG models

In figure-3, we show performance score for various models such as Encoder-Decoder model with local cross-entropy loss (EDL), Encoder-Decoder model with local and pairwise loss (EDLP), Encoder-Decoder model with local and pairwise loss along with shared encoder and discriminator network (EDLPS), Encoder-Decoder model with local, pairwise and adversarial loss (EDLPG), Encoder-Decoder model with local, pairwise and adversarial loss along with shared encoder and discriminator network (EDLPGS). We have shown the performance score of each model across every epoch.

²we report same baseline results as mentioned in [6]

We observe that the score EDLPS method performs best across all scores (BLEU, ROUGE, METEOR, and CIDEr). We also observe that the EDLPGS method performs some peculiar behavior on different epochs. In figure-5, we visualize the performance of various score across the datasets that is QQP-I (50k) and QQP-II (100k) dataset. We observe that in QQP-II, the performance best over all the scores across all epochs.

4.1.5. Qualitative results for Paraphrase Generation

This table 6 contains few paraphrase questions generated by our model along with the original question and the ground truth paraphrase questions. In table 7, we provide some more examples of the paraphrase generation task. Our model is also able to generate sentences that capture higher-level semantics like in the last example of table 7.

4.1.6. Experimental Protocols for Paraphrase Generation

We follow the experimental protocols and evaluation methods, as mentioned in [4] for the Quora Question Pairs (QQP) dataset. We also followed the dataset split mentioned in [5] to calculate the accuracies on a different test set and provide the results on our project webpage. We trained our model end-to-end using local loss (cross entropy loss) and global loss. We have used RMSPROP optimizer to update the model parameter and found these hyperparameter values to work best to train the Paraphrase Generation Network: learning rate = 0.0008, batch size = 150, $\alpha = 0.99$, $\epsilon = 1e - 8$. We have used learning rate decay to decrease the learning rate on every epoch by a factor given by:

$$\text{Decay_factor} = \exp\left(\frac{\log(0.1)}{a * b}\right)$$

where $a = 1500$ and $b = 1250$ are set empirically.

4.1.7. Statistical Significance Analysis

We have analyzed statistical significance [49] for our proposed embeddings against different ablations and the state-of-the-art methods for the paraphrase generation task. The Critical Difference (CD) for Nemenyi [50] test depends upon the given α (confidence level, which is 0.05 in our case) for average ranks and N (number of tested datasets). If the difference in the rank of the two methods lies within CD, then they are

S.No	Original Question	Ground Truth Paraphrase	Generated Paraphrase
1	How do I add content on Quora?	How do I add content under a title at Quora?	How do I add images on Quora ?
2	Is it possible to get a long distance ex back?	Long distance relationship: How to win my ex-gf back?	Is it possible to get a long distance relationship back ?
3	How many countries are there in the world? Thanks!	How many countries are there in total?	How many countries are there in the world ? What are they ?
4	What is the reason behind abrupt removal of Cyrus Mistry?	Why did the Tata Sons sacked Cyrus Mistry?	What is the reason behind firing of Cyrus Mistry ?
5	What are some extremely early signs of pregnancy?	What are the common first signs of pregnancy? How can I tell if I'm pregnant? What are the symptoms?	What are some early signs of pregnancy ?
6	How can I improve my critical reading skills?	What are some ways to improve critical reading and reading comprehension skills?	How can I improve my presence of mind ?

Table 7: Examples of Paraphrase generation on Quora Dataset.

not significantly different, otherwise, they are statistically different. Figure 6 visualizes the post hoc analysis using the CD diagram. From the figure, it is clear that our embeddings work best, and the results are significantly different from the state-of-the-art methods.

4.2. Sentiment Analysis with Stanford Sentiment Treebank (SST) Dataset

4.2.1. Dataset

This dataset consists of sentiment labels for different movie reviews and was first proposed by [51]. [26] extended this by parsing the reviews to subphrases and then

Model	Error Rate (Fine-Grained)
Naive Bayes [26]	59.0
SVMs [26]	59.3
Bigram Naive Bayes [26]	58.1
Word Vector Averaging [26]	67.3
Recursive Neural Network [26]	56.8
Matrix Vector-RNN [26]	55.6
Recursive Neural Tensor Network [26]	54.3
Paragraph Vector [27]	51.3
EDD-LG(shared) (Ours)	35.6

Table 8: Performance of our method compared to other approaches on the Stanford Sentiment Treebank Dataset. The error rates of other methods are reported in [27]

fine-graining the sentiment labels for all the phrases of movie reviews using Amazon Mechanical Turk. The labels are classified into 5 sentiment classes, namely {Very Negative, Negative, Neutral, Positive, Very Positive}. This dataset contains a total of 126k phrases for the training set, 30k phrases for the validation set, and 66k phrases for the test set.

4.2.2. Tasks and Baselines

In [26], the authors propose two ways of benchmarking. We consider the 5-way fine-grained classification task where the labels are {Very Negative, Negative, Neutral, Positive, Very Positive}. The other axis of variation is in terms of whether we should label the entire sentence or all phrases in the sentence. In this work, we only consider labeling all the phrases. [26] apply several methods to this dataset, and we show their performance in table 8.

4.2.3. Sentiment Visualization of the Sentence

Li *et al.*[52] have proposed a mechanism to visualize language features. We conducted a toy experiment for our EDD-LG(shared) model. We provide visualization

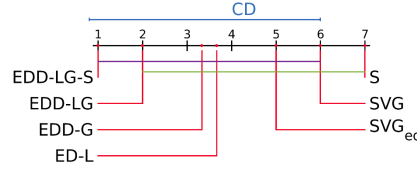


Figure 6: The mean rank of all the models on the basis of BLEU score are plotted on the x-axis. Here EDD-LG-S refers to our EDD-LG shared model and others are the different variations of our model described in section 4.1.2 and the models on the right are the different variations proposed in [4]. Also the colored lines between the two models represents that these models are not significantly different from each other. $CD=5.199, p=0.0069$

of different parts of the sentence on which our model focuses while predicting the sentiment in figure 7. In figure- 7 represents the saliency heat map for EDD-LG(shared) model sentiment analysis. We obtained 60-dimensional feature maps for each word present in the target sentence. The heat map captures the measure of the influence of the sentimental decision. In the heat map, each word of a sentence (from top to bottom, first word at the top) represents its contribution to making the sentimental decision. For example, in the first image in 7, the word ‘comic’ contributed more (2nd word, row 10-20). Similarly, in the second image, first, second, and third (‘A’, ‘wildly’, ‘funny’) words have more influence on making this sentence have a positive sentiment.

4.2.4. Experimental Protocols

For the task of Sentiment analysis, we are using a similar method of performing the experiments as used by [26]. We treat every subphrase in the dataset as a separate sentence and learn their corresponding representations. We then feed these to a logistic regression to predict the movie ratings. During inference time, we used a method similar to [27] in which we freeze the representation of every word and use this to construct a representation for the test sentences which are then fed to a logistic regression for predicting the ratings. In order to train a sentiment classification model, we have used RMSPROP, to optimize the classification model parameter and we found these hyperparameter values to be working best for our case: learning rate = 0.00009, batch size = 200, $\alpha = 0.9$, $\epsilon = 1e - 8$.

Phrase ID	Phrase	Sentiment
162970	The heaviest, most joyless movie	Very Negative
159901	Even by dumb action-movie standards, Ballistic : Ecks vs. Sever is a dumb action movie.	
158280	Nonsensical, dull “cyber-horror” flick is a grim, hollow exercise in flat scares and bad acting	
159050	This one is pretty miserable, resorting to string-pulling rather than legitimate character development and intelligent plotting.	
157130	The most hopelessly monotonous film of the year, noteworthy only for the gimmick of being filmed as a single unbroken 87-minute take.	
156368	No good jokes, no good scenes, barely a moment	Negative
157880	Although it bangs a very cliched drum at times	
159269	They take a long time to get to its gasp-inducing ending.	
157144	Noteworthy only for the gimmick of being filmed as a single unbroken 87-minute	
156869	Done a great disservice by a lack of critical distance and a sad trust in liberal arts college bumper sticker platitudes	
221765	A hero can stumble sometimes.	Neutral
222069	Spiritual rebirth to bruising defeat	
218959	An examination of a society in transition	
221444	A country still dealing with its fascist past	
156757	Have to know about music to appreciate the film’s easygoing blend of comedy and romance	
157663	A wildly funny prison caper.	Positive
157850	This is a movie that’s got oodles of style and substance.	
157879	Although it bangs a very cliched drum at times, this crowd-pleaser’s fresh dialogue, energetic music, and good-natured spunk are often infectious.	
156756	You don’t have to know about music to appreciate the film’s easygoing blend of comedy and romance.	
157382	Though of particular interest to students and enthusiast of international dance and world music, the film is designed to make viewers of all ages, cultural backgrounds and rhythmic ability want to get up and dance.	
162398	A comic gem with some serious sparkles.	Very Positive
156238	Delivers a performance of striking skill and depth	
157290	What Jackson has accomplished here is amazing on a technical level.	
160925	A historical epic with the courage of its convictions about both scope and detail.	
161048	This warm and gentle romantic comedy has enough interesting characters to fill several movies, and its ample charms should win over the most hard-hearted cynics.	

Table 9: Examples of Sentiment classification on test set of kaggle competition dataset.

4.2.5. Results

We report the error rates of different methods in table 8. We can clearly see that the performance of bag-of-words or bag-of-n-grams models (the first four models in the table) is not up to the mark and instead the advanced methods (such as Recursive Neural Network [26]) perform better on sentiment analysis task. Our method outperforms

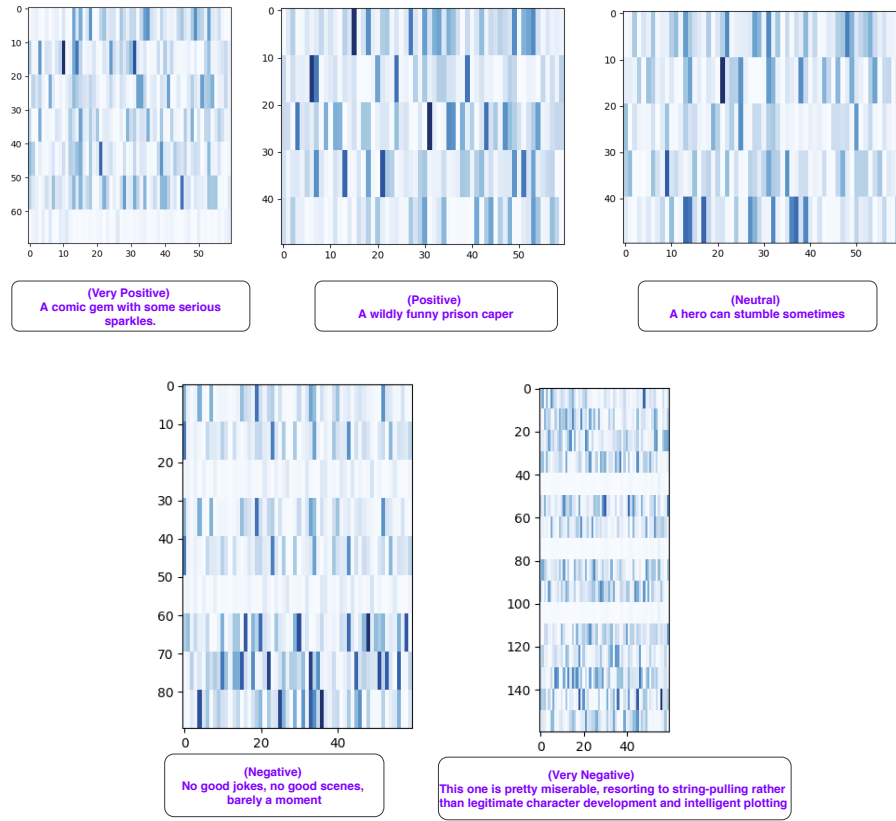


Figure 7: These are the visualisations for the sentiment analysis for some examples and we can clearly see that our model focuses on those words which we humans focus while deciding the sentiment for any sentence. In the second image, ‘wildly’ and ‘funny’ are emphasised more than the other words.

all these methods by an absolute margin of 15.7% which is a significant increase considering the rate of progress on this task. We have also uploaded our models to the online competition on Rotten Tomatoes dataset ³ and obtained an accuracy of 62.606% on their test-set of 66K phrases.

We provide 5 examples for each sentiment in table 9. We can see clearly that our proposed embeddings are able to get the complete meaning of smaller as well as larger sentences. For example, our model classifies ‘Although it bangs a very cliched drum

³website: www.kaggle.com/c/sentiment-analysis-on-movie-reviews

at times’ as Negative and ‘Although it bangs a very cliched drum at times, this crowd-pleaser’s fresh dialogue, energetic music, and good-natured punk are often infectious.’ as positive showing that it is able to understand the finer details of language. Some more examples of our model for the Sentiment analysis task on the SST dataset in table 10. The link for the code is provided here ⁴.

4.3. Quantitative Evaluation

We use automatic evaluation metrics which are prevalent in machine translation domain: BLEU [53], METEOR [54], ROUGE-n [55] and Translation Error Rate (TER) [56]. These metrics perform well for Paraphrase generation task and also have a higher correlation with human judgments [57, 58]. BLEU uses n-gram precision between the ground truth and the predicted paraphrase. considers exact match between reference whereas ROUGE considers recall for the same. On the other hand, METEOR uses stemming and synonyms (using WordNet) and is based on the harmonic mean of unigram-precision and unigram-recall. TER is based on the number of edits (insertions, deletions, substitutions, shifts) required to convert the generated output into the ground truth paraphrases and quite obviously a lower TER score is better whereas other metrics prefer a higher score for showing improved performance. We provided our results using all these metrics and compared it with existing baselines.

5. Conclusion

In this paper we have proposed a sentence embedding using a sequential encoder-decoder with a pairwise discriminator. We have experimented with this text embedding method for paraphrase generation and sentiment analysis. We also provided experimental analysis which justifies that a pairwise discriminator outperforms the previous state-of-art methods for NLP tasks. We also performed ablation analysis for our method, and our method outperforms all of them in terms of BLEU, METEOR and TER scores. We plan to generalize this to other text understanding tasks and also extend the same idea in vision domain.

⁴Code: <https://github.com/dev-chauhan/PQG-pytorch>

6. Acknowledgment

We acknowledge the help provided by our DeITA Lab members and our family who have supported us in our research activity.

Phrase ID	Phrase	Sentiment
156628 157078 159749 163483 163882 164436	The movie is just a plain old monster a really bad community theater production of West Side Story Suffers from rambling , repetitive dialogue and the visual drabness endemic to digital video . lapses quite casually into the absurd It all drags on so interminably it 's like watching a miserable relationship unfold in real time . Your film becomes boring , and your dialogue is n't smart	Very Negative
156567 156689 157730 157695 158814 159281 159632 159770	It would be hard to think of a recent movie that has worked this hard to achieve this little fun A depressing confirmation There 's not enough here to justify the almost two hours. a snapshot of a dangerous political situation on the verge of coming to a head It is ridiculous , of course A mostly tired retread of several other mob tales. We are left with a superficial snapshot that , however engaging , is insufficiently enlightening and inviting . It 's as flat as an open can of pop left sitting in the sun .	Negative
156890 160247 160754 160773 201255 201371 221444 222102	liberal arts college bumper sticker platitudes the movie 's power as a work of drama Schweig , who carries the film on his broad , handsome shoulders to hope for any chance of enjoying this film also examining its significance for those who take part those who like long books and movies a country still dealing with its fascist past used to come along for an integral part of the ride	Neutral
157441 157879 157663 157749 157806 157850	the film is packed with information and impressions . Although it bangs a very cliched drum at times , this crowd-pleaser 's fresh dialogue , energetic music , and good-natured spunk are often infectious. A wildly funny prison caper. This is one for the ages. George Clooney proves he 's quite a talented director and Sam Rockwell shows us he 's a world-class actor with Confessions of a Dangerous Mind . this is a movie that 's got oodles of style and substance .	Positive
157742 160562 160925 161048 161459 162398 162779 163228	Kinnear gives a tremendous performance . The film is painfully authentic , and the performances of the young players are utterly convincing . A historical epic with the courage of its convictions about both scope and detail. This warm and gentle romantic comedy has enough interesting characters to fill several movies , and its ample charms should win over the most hard-hearted cynics . is engrossing and moving in its own right A comic gem with some serious sparkles . a sophisticated , funny and good-natured treat , slight but a pleasure Khouri then gets terrific performances from them all .	Very Positive

Table 10: Examples of Sentiment classification on test set of kaggle dataset.

7. Reference

References

- [1] I. Sutskever, O. Vinyals, Q. V. Le, Sequence to sequence learning with neural networks, in: *Advances in neural information processing systems*, 2014, pp. 3104–3112.
- [2] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate, *arXiv preprint arXiv:1409.0473*.
- [3] A. Prakash, S. A. Hasan, K. Lee, V. Datla, A. Qadir, J. Liu, O. Farri, Neural paraphrase generation with stacked residual lstm networks, *arXiv preprint arXiv:1610.03098*.
- [4] A. Gupta, A. Agarwal, P. Singh, P. Rai, A deep generative framework for paraphrase generation, *arXiv preprint arXiv:1709.05074*.
- [5] Z. Li, X. Jiang, L. Shang, H. Li, Paraphrase generation with deep reinforcement learning, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 3865–3878.
- [6] Q. Yang, D. Shen, Y. Cheng, W. Wang, G. Wang, L. Carin, et al., An end-to-end generative architecture for paraphrase generation, in: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 3123–3133.
- [7] Z. S. Harris, Distributional structure, *Word* 10 (2-3) (1954) 146–162.
- [8] J. R. Firth, A synopsis of linguistic theory, 1930-1955, *Studies in linguistic analysis*.
- [9] M. Sahlgren, The distributional hypothesis, *Italian Journal of Disability Studies* 20 (2008) 33–53.

- [10] G. Salton, A. Wong, C.-S. Yang, A vector space model for automatic indexing, *Communications of the ACM* 18 (11) (1975) 613–620.
- [11] D. Lin, P. Pantel, Dirt@ sbt@ discovery of inference rules from text, in: *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2001, pp. 323–328.
- [12] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, R. Harshman, Indexing by latent semantic analysis, *Journal of the American society for information science* 41 (6) (1990) 391.
- [13] P. D. Turney, M. L. Littman, Measuring praise and criticism: Inference of semantic orientation from association, *ACM Transactions on Information Systems (TOIS)* 21 (4) (2003) 315–346.
- [14] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, *arXiv preprint arXiv:1301.3781*.
- [15] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- [16] Y. Goldberg, O. Levy, word2vec explained: Deriving mikolov et al.’s negative-sampling word-embedding method, *arXiv preprint arXiv:1402.3722*.
- [17] J. Mitchell, M. Lapata, Composition in distributional models of semantics, *Cognitive science* 34 (8) (2010) 1388–1429.
- [18] F. M. Zanzotto, I. Korkontzelos, F. Fallucchi, S. Manandhar, Estimating linear models for compositional distributional semantics, in: *Proceedings of the 23rd International Conference on Computational Linguistics*, Association for Computational Linguistics, 2010, pp. 1263–1271.
- [19] A. Yessenalina, C. Cardie, Compositional matrix-space models for sentiment analysis, in: *Proceedings of the Conference on Empirical Methods in Natural*

- Language Processing, Association for Computational Linguistics, 2011, pp. 172–182.
- [20] E. Grefenstette, G. Dinu, Y. Zhang, M. Sadrzadeh, M. Baroni, Multi-step regression learning for compositional distributional semantics, in: Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013)–Long Papers, 2013, pp. 131–142.
 - [21] R. Socher, C. C. Lin, C. Manning, A. Y. Ng, Parsing natural scenes and natural language with recursive neural networks, in: Proceedings of the 28th international conference on machine learning (ICML-11), 2011, pp. 129–136.
 - [22] Y. Kim, Convolutional neural networks for sentence classification, arXiv preprint arXiv:1408.5882.
 - [23] R. Lin, S. Liu, M. Yang, M. Li, M. Zhou, S. Li, Hierarchical recurrent neural network for document modeling, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015, pp. 899–907.
 - [24] W. Yin, H. Schütze, B. Xiang, B. Zhou, Abcnn: Attention-based convolutional neural network for modeling sentence pairs, arXiv preprint arXiv:1512.05193.
 - [25] N. Kalchbrenner, E. Grefenstette, P. Blunsom, A convolutional neural network for modelling sentences, arXiv preprint arXiv:1404.2188.
 - [26] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, C. Potts, Recursive deep models for semantic compositionality over a sentiment treebank, in: Proceedings of the 2013 conference on empirical methods in natural language processing, 2013, pp. 1631–1642.
 - [27] Q. Le, T. Mikolov, Distributed representations of sentences and documents, in: International Conference on Machine Learning, 2014, pp. 1188–1196.
 - [28] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio, Learning phrase representations using rnn encoder–decoder for statistical machine translation (2014) 1724–1734.

- [29] O. Vinyals, Q. Le, A neural conversational model, arXiv preprint arXiv:1506.05869.
- [30] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al., Google’s neural machine translation system: Bridging the gap between human and machine translation. corr abs/1609.08144 (2016).
- [31] X. Du, J. Shao, C. Cardie, Learning to ask: Neural question generation for reading comprehension 1 (2017) 1342–1352.
- [32] L. Shang, Z. Lu, H. Li, Neural responding machine for short-text conversation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Vol. 1, 2015, pp. 1577–1586.
- [33] J. Li, W. Monroe, A. Ritter, M. Galley, J. Gao, D. Jurafsky, Deep reinforcement learning for dialogue generation, arXiv preprint arXiv:1606.01541.
- [34] J. Li, W. Monroe, T. Shi, A. Ritter, D. Jurafsky, Adversarial learning for neural dialogue generation, arXiv preprint arXiv:1701.06547.
- [35] A. M. Rush, S. Chopra, J. Weston, A neural attention model for abstractive sentence summarization, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015, pp. 379–389.
- [36] Y. Zhang, Z. Gan, K. Fan, Z. Chen, R. Henao, D. Shen, L. Carin, Adversarial feature matching for text generation, arXiv preprint arXiv:1706.03850.
- [37] Z. Hu, Z. Yang, X. Liang, R. Salakhutdinov, E. P. Xing, Toward controlled generation of text, in: International Conference on Machine Learning, 2017, pp. 1587–1596.
- [38] L. Yu, W. Zhang, J. Wang, Y. Yu, Seqgan: Sequence generative adversarial nets with policy gradient., in: AAAI, 2017, pp. 2852–2858.

- [39] J. Guo, S. Lu, H. Cai, W. Zhang, Y. Yu, J. Wang, Long text generation via adversarial training with leaked information, arXiv preprint arXiv:1709.08624.
- [40] X. Liang, Z. Hu, H. Zhang, C. Gan, E. P. Xing, Recurrent topic-transition gan for visual paragraph generation, CoRR, abs/1703.07022 2.
- [41] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, H. Lee, Generative adversarial text to image synthesis, arXiv preprint arXiv:1605.05396.
- [42] J. Yin, X. Jiang, Z. Lu, L. Shang, H. Li, X. Li, Neural generative question answering, in: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, AAAI Press, 2016, pp. 2972–2978.
- [43] Y. Miao, L. Yu, P. Blunsom, Neural variational inference for text processing, in: International Conference on Machine Learning, 2016, pp. 1727–1736.
- [44] S. A. Hasan, B. Liu, J. Liu, A. Qadir, K. Lee, V. Datla, A. Prakash, O. Farri, Neural clinical paraphrase generation with attention, in: Proceedings of the Clinical Natural Language Processing Workshop (ClinicalNLP), 2016, pp. 42–53.
- [45] B. N. Patro, V. K. Kurmi, S. Kumar, V. Namboodiri, Learning semantic sentence embeddings using sequential pair-wise discriminator, in: Proceedings of the 27th International Conference on Computational Linguistics, 2018, pp. 2715–2729.
- [46] X. Zhang, J. Zhao, Y. LeCun, Character-level convolutional networks for text classification, in: Advances in neural information processing systems, 2015, pp. 649–657.
- [47] H. Palangi, L. Deng, Y. Shen, J. Gao, X. He, J. Chen, X. Song, R. Ward, Deep sentence embedding using long short-term memory networks: Analysis and application to information retrieval, IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP) 24 (4) (2016) 694–707.
- [48] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural computation 9 (8) (1997) 1735–1780.

- [49] J. Demšar, Statistical comparisons of classifiers over multiple data sets, *Journal of Machine learning research* 7 (Jan) (2006) 1–30.
- [50] D. Fišer, T. Erjavec, N. Ljubešić, Janes v0. 4: Korpus slovenskih spletnih uporabniških vsebin, *Slovenščina* 2 (4) (2016) 2.
- [51] B. Pang, L. Lee, Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales, in: *Proceedings of the 43rd annual meeting on association for computational linguistics*, Association for Computational Linguistics, 2005, pp. 115–124.
- [52] J. Li, X. Chen, E. Hovy, D. Jurafsky, Visualizing and understanding neural models in nlp, in: *Proceedings of NAACL-HLT*, 2016, pp. 681–691.
- [53] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, Bleu: a method for automatic evaluation of machine translation, in: *Proceedings of the 40th annual meeting on association for computational linguistics*, Association for Computational Linguistics, 2002, pp. 311–318.
- [54] S. Banerjee, A. Lavie, Meteor: An automatic metric for mt evaluation with improved correlation with human judgments, in: *Proc. of ACL workshop on Intrinsic and Extrinsic Evaluation measures for Machine Translation and/or Summarization*, Vol. 29, 2005, pp. 65–72.
- [55] C.-Y. Lin, Rouge: A package for automatic evaluation of summaries, in: *Text summarization branches out: Proceedings of the ACL-04 workshop*, 2004.
- [56] M. Snover, B. Dorr, R. Schwartz, L. Micciulla, J. Makhoul, A study of translation edit rate with targeted human annotation, 2006.
- [57] N. Madnani, J. Tetreault, M. Chodorow, Re-examining machine translation metrics for paraphrase identification, in: *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, 2012, pp. 182–190.

- [58] S. Wubben, A. v. d. Bosch, E. Krahmer, Paraphrasing headlines by machine translation: Sentential paraphrase acquisition and generation using google news, LOT Occasional Series 16 (2010) 169–183.