

# RMHIDD: A Reddit Mental Health Intervention Dialogue Dataset

Anonymous EMNLP submission

## Abstract

In the modern human society, mental health is one of the most critical concern. Over past many years a large proportion of population has been affected with serious mental disorders. People with mental illness require effective mental health intervention and treatment as early as possible to decrease the chances of any further mental defilement. In this paper, we present RMHIDD, a new dialogue corpus for automated mental health intervention. The dataset is consists of over 200K Reddit posts collected from 18 different sub-Reddit groups with each post consisting of sequential conversation between the users. On this dataset, we also trained various models for dialogue generation task, namely-'Seq2Seq', 'BART' and 'DialoGPT'. In our analysis we found that the BART model outperformed other models with a higher Perplexity score of 19.7. We also found that the DialoGPT model surpasses other models on various machine translation evaluation metrics. The results generated from various language models were promising and showed the possibility of building automated mental health intervention.

## 1 Introduction

Mental health is one of the most serious global concerns. In the last few years, there has been a huge increase in the number of people affected by some kind of mental disorder. A report from World Health Organization (WHO)<sup>1</sup> states that 1 out of every 4 people in the world is affected by some kind of mental disorder in their different stages of life. According to WHO's report on depression, it indicates that around 322 million people all around the globe have been affected by depression and this accounts for 18% growth in the total cases from 2005 to 2015. Other major mental illnesses

<sup>1</sup>Mental health action plan 2013 - 2020, available at [https://www.who.int/mental\\_health/](https://www.who.int/mental_health/)

such as anxiety and bipolar disorder which has affected around 264 million people and 60 million people worldwide respectively. Despite increased awareness about mental health conditions and it's management, a report from WHO shows that in every 4 people 3 of them who are suffering from serious mental illness lack timely treatment which pushes them into a major serious mental disorder state. This is a fearful condition and in many cases due to the lack of timely medical help people with mental disorders tend to commit suicide. Every year around 800 thousand people die because of suicide(Organization et al., 2017). Mental health systems are underdeveloped and are not sufficient enough to reach every person who's in need. The basic measures for the prevention of mental illness are psychological intervention and oral consultation. Due to insufficient medical facilities(Jacob and Patel, 2014), the majority of people remain deprived of much-needed treatment and support. Approximately 45% of total world population is living in countries where for every 100K people lesser than 1 psychiatrist is available<sup>2</sup>. Moreover, a report from WHO states that approximately 76% to 85% of people having a mental illness and living in countries with medium and low income do not get the necessary treatment. For high-income countries, it ranges between 35% to 50%. A combination of various factors such as social stigma (Barney et al., 2006) against mentally ill people in the society, unwillingness, or hesitation in asking for help/support, resource scarcity is few reasons behind mismanagement of mental health conditions. Additionally, For the past few years due to the rise in popularity of social media platforms, millions of people are using these platforms to either provide or receive mental health support. Through these online mediums, people express their feelings more

<sup>2</sup>Available at [http://www.who.int/mental\\_health/evidence/atlas](http://www.who.int/mental_health/evidence/atlas)

freely and seek help without any hesitation. One of the most popular online social media platform for sharing mental disorder experiences is Reddit.

Reddit consists of various subject-specific communities called 'subreddits' where people post their thoughts (topic) and other users can reply through commenting on the post. In this paper, we collected<sup>3</sup> over 200K human-generated posts from various mental disorders help subreddits in a nested way in order to preserve the sequence of conversations (dialogues) generated between the users on a Reddit post. The major **motivation** behind creating and publishing the RMHIDD dialogue corpus is to enable researchers and scientist all around the world to utilize latest advancements in natural language processing and understanding and develop innovative automated mental health intervention tools (such as intelligent chatbots, e-therapy, e-screening, detecting and predicting mental disorders through dialogues) for addressing and solving the issue of mental disorder in our society. We also performed the task of automated dialogue generation that involves generating helpful/supportive dialogues based on the user's mental illness. We used various language models namely - 'Seq2Seq(Sutskever et al., 2014)', 'BART(Lewis et al., 2019)' and 'DialoGPT(Zhang et al., 2019)' and compared their performance on our proposed dataset. In our experiment, we trained the Seq2Seq model whereas the weights of BART and DialoGPT models were fine-tuned on our dataset. We also performed a comparative study of all the models by evaluating them using various automatic evaluation metrics. The responses generated by the dialogue generation models were very promising and demonstrated the potential and application of natural language processing in the field of automated healthcare systems.

The rest of the paper is organized in the following way. In the section 2 we have summarised the related work done in the field of analyzing user-generated data for mental disorders. In Section 3 we discuss the data collection method and involved steps and in section 4, we discuss various language models used in experiment. Section 5 contains the sequential experiment setup. Section 6 presents analysis and discuss of experiment results and finally in section 7 we conclude the paper.

<sup>3</sup>We will make the data publicly available

## 2 Related Work

A lot of works have been done that is based on collecting user-generated online data and utilizing it for the analysis and creating insights into mental disorder in people. The author in (De Choudhury and De, 2014) presented a study on characteristics shown by patients with mental illness on social media platforms (Reddit) such as self-disclosure, anonymity, and how it affects the social support received by the patients. The author in papers (Gkotsis et al., 2016), (Park and Conway, 2018) investigated the linguistic characteristic specifically present in the Reddit posts concerned with mental health and illness. In this paper(Gkotsis et al., 2017) the author performed the analysis of Reddit posts and proposed a deep learning-based detection and classification Reddit posts in 11 fine-grained classes of mental disorder. Another paper (Thorstad and Wolff, 2019) also presented an automated mental disorder detection model trained using clinical subreddits which focuses on lexical features in user-generated data to detect the mental illness present in the user. This work presented an automated system for targeted mental health intervention based on user-generated data. In the paper (Shen and Rudzicz, 2017) the author builds a dataset consisting of anxiety-related user-generated posts. The author also applied topic analysis, vector embeddings, emotional norms, and N-gram language modeling for generating features to classify posts in anxiety levels. In paper (Abd Yusof et al., 2018), the author developed lexical features for depression classification tasks and created a dataset using LiveJournal<http://www.livejournal.com> to evaluate feature effectiveness. In the paper (Wongkoblap et al., 2018), the author investigated the relationship between depression and life satisfaction using Facebook user's data and also presented a multilevel predictive model for finding depression in users. Number of researches focusing on identifying depression, anxiety, suicide and bipolarity in social media networks has been done such as (Murrieta et al., 2018), (Lee et al., 2018), (Chen et al., 2018), (Leis et al., 2019), (Wolohan et al., 2018), (Wongkoblap et al., 2019), (Gruda and Hasan, 2019), (Sahota and Sankar, 2020), (Baba et al., 2019). Lot of work focusing on automated healthcare facilities has been done(Liliana Laranjo, 2018). In the papers (Lucas et al., 2017), (Philip et al., 2017), (Tanaka et al., 2017), authors used sequence based step-by-step

guided conversation models. Recently neural network based medical dialogue generation models were also proposed. In paper (Wei et al., 2018), author utilized reinforcement learning and developed task oriented dialogue system for automatic medical diagnosis. Paper (Xu et al., 2019) proposed a knowledge-routed relational dialogue system.

Despite plenty of existing research work and resources available on mental disorders, the count of people affected with mental illness rises sharply each year. Most of these works attempt to comprehend the user’s action and behavior over social media platforms and develop methods/models to detect the degree of mental disorder among the people. In this paper, we focused on creating a dataset that comprises instances of dialogues between mental disorder help seekers and support providers, collected from social media platforms where people are free to express themselves. Through our work, we want to take a step forward towards developing automated mental health intervention systems that would be readily available to the people suffering from mental disorders.

### 3 Dataset

The prime source for the collection of our data was Reddit<sup>4</sup>. Reddit is basically a social media website comprised of multiple distinct online communities called subreddits. These are topic-specific forums dedicated to a single topic (e.g., depression, relationship advice, anxiety, etc.) where a user creates a topic, expresses themselves and other users can comment or vote for other comments on that post. We scrapped posts from various mental health management, advice, or support providing subreddits.

- **Mental disorder subreddits** :  
r/depression, r/anxiety,  
r/stress, r/BipolarReddit
- **Advice Support subreddits** :  
r/therapy, r/depression\_help,  
r/Anxietyhelp, r/SuicideWatch,  
r/relationship\_advice,  
r/offmychest,  
r/askatherapist,  
r/relationships
- **Motivating Uplifting subreddits** :  
r/TheMixedNuts, r/MadeMeSmile,  
r/FreeCompliments,

<sup>4</sup>Available at <https://old.reddit.com/>

r/UpliftingNews,  
r/DecidingToBeBetter  
r/GetMotivated,

In the table 1, we have described the number of subreddits used for collecting posts along with the number of dialogues, utterances and the average token length per dialogue. We collected all the content posted in the duration of one year (2, July 2019 to 2, July 2020) on the above-mentioned subreddits. Using PRAW API<sup>5</sup> to extract all the content in a nested manner to conserve the sequence of dialogues (topic and comments) in the post. While scrapping dialogues from the posts we removed unwanted bot auto-generated texts by skipping those lines. We also normalized the scrapped dialogues by removing the curse words. We also removed the posts with no comments, in total, we extracted around 200K online user-generated Reddit posts containing instances of dialogues between the users. Table 2 presents the extracted dialogue example from our collected dataset.

### 4 Methods

In this section, we gave an overview of various state-of-the-art and well-established dialogue generation models. For our experiment, we used 3 deep learning encoder-decoder based models, i.e., Seq2Seq (Wu et al., 2016), DialoGPT (Zhang et al., 2019) and BART (Lewis et al., 2019). We trained the Seq2Seq model on our dataset, whereas for DialoGPT and BART we fine-tuned these models on our dataset.

For a given dialogue having an alternating sequence of utterances between the users, we decided to take two utterances, (i)  $D_1$  (person A issue) the content of the main topic created by a user and (ii)  $D_2$  (person B response) the comment with the highest number of upvotes. So, for each dialogue, we created a pair of utterances, i.e.  $\{D_1, D_2\}$  which is used for training all of our dialogue generation models. Given an input  $D_1$ , the dialogue generation model outputs  $D_2$ .

#### 4.1 Seq2Seq

We utilized an encoder-decoder framework for dialogue generation tasks. Following the original architecture proposed by the author (Wu et al., 2016) for machine translation tasks, we build an LSTMs based deep seq2seq model with attention. The

<sup>5</sup><https://github.com/praw-dev/praw>

Subreddit Group	#Subreddits	#Dialogue	#Utterances	Avg Tokens/Dialogue
Mental disorder subreddits	4	100,492	246,077	187
Advice Support subreddits	8	83,065	172,688	191
Motivating Uplifting subreddits	6	41,533	95,102	183
Total	18	225,090	513,867	516

Table 1: Data Statistics: We listed 3 subreddit group along with the their associated total number of dialogues, utterances and average tokens per dialogue present in the dataset.

model takes  $D_1$  as input and outputs  $D_2$  as the generated dialogue. Each of the encoder and decoder was consist of 2 LSTM and 1 BiLSTM layer. The input was first passed through 2 LSTM layers, followed by a single BiLSTM layer which generated the latent representation. Similarly, in the decoder, we applied the same 2 LSTM layers with the final BiLSTM layer as the decoding layer.

For a given training set  $S$ , we intend to make the log probability of the output sequences  $T$  maximum where the given input sequences  $S$  given (Sutskever et al., 2014).

$$\frac{1}{|S|} \sum_{(T,S) \in \mathcal{S}} \log_p(T|S) \quad (1)$$

Once the training is done, according to the LSTM the most probable output sequence is produced:

$$T^c = \underset{T}{\operatorname{argmax}} p(T|S) \quad (2)$$

To obtain final predictions, in the decoder we used the softmax layer and performed decoding using beam search<sup>6</sup>. Finally, the obtained outputs were passed into the loss function, and parameters were updated through backpropagation. Adam (Kingma and Ba, 2014) optimizer was used in the model.

## 4.2 DialoGPT

In the paperrad2m018ipriving, the author proposed a transformer based language model- GPT. For a given token sequence  $x_1, \dots, x_n$ , in a language model the probability over sequence was defined as:  $p(x_1, \dots, x_n) = \prod_{i=2}^n p(x_i|x_1, \dots, x_{i-1})$ , where historical sequences are used for predicting the next token. In case of GPT, transformer decoder was used to define  $p(x_i|x_1, \dots, x_{i-1})$ . The decoder consists of stacked self-attention feed-forward layers (each accompanied by normalization layer) for

<sup>6</sup><https://google.github.io/seq2seq/nmt/decoding-with-beam-search>

encoding  $x_1, \dots, x_{i-1}$  and which was then used to predict  $x_i$ . In the case of GPT-2 (Radford et al., 2019) which was an improvement over GPT, the normalization layer was moved to each of the sub-blocks input. An extra normalization layer was added after the last self-attention block.

For our experiment, we used DialoGPT (Zhang et al., 2019) which was a GPT-2 based model trained on a very large corpus consisting of English Reddit dialogues. The corpus was consist of 147,116,725 instances of dialogues, collected over a period of 12 years. The model takes the dialogue utterances history  $S$  and ground truth response  $T = x_1, \dots, x_n$ , the DialoGPT model aims at maximizing the probability:  $p(T|S) = p(x_1|S) \prod_{i=2}^n p(x_i|S, x_1, \dots, x_{i-1})$ , where the transformer model defines the conditional probabilities. Through a maximum mutual information (MMI) function (Li et al., 2015), the model also gets penalized for generating uninteresting responses. In our experiment, we used *DialoGPT<sub>small</sub>* with 117 million weight parameters.

## 4.3 BART

BART (Lewis et al., 2019) is a denoising autoencoder that tries to rebuild a corrupted document by performing masked token prediction with the help of bidirectional encoding methods and generates text regressively for natural language generation tasks using a masked attention mechanism. The mask attention mechanism enables the BART model to train on sequence from left to right, generating texts based on the left part of the sequence.

For this transformer-based dialogue system, we create a BART language model wrapper which includes the API of the BART-large model from hugging face-transformers. This pretrained model has 400M trainable parameters with 6 encoding and decoding layers in each block, 16 attention heads both at the encoding and decoding layer. We represent each encoder layer as an *Encoder(.)* which outputs the hidden state of the respective layer. We

	Dialogue	
400		450
401	creator_id : I am , stuck .	451
402	creator_id : I deal with social anxiety , and lately things have been worse than ever	452
403	and I don ' t know what I should do .	453
404	creator_id : To give you an idea with what I 've been feeling ; I 've always disliked myself, but	454
405	lately it's become a real hatred .	455
406	creator_id : I see any little thing about myself and I feel disgusted , and angry , I can't	456
407	even take photos of myself or feel comfortable when others do because I know	457
408	when I see it I ' ll feel repulsed ..	458
409	creator_id : And whenever my friends try to make plans I feel unmotivated , and afraid , and I	459
410	usually make up some stupid lie to get out of things .	460
411	creator_id : And I can ' t make plans because I don ' t want to come off as clingy or whatever,	461
412	and it's really frustrating .	462
413	creator_id : Things aren't getting better, but they're not getting worse, it's like this, numbing pain	463
414	been going on for so long it's frustrating and I'm sick of it .	464
415	creator_id : If you know what I'm talking about, if you've felt this before,	465
416	tell me what to do next.	466
417		467
418	Commenter_id(1) : I have felt literally the exact same type of way you are describing .	468
419	Commenter_id(1) : One thing I would suggest is to stop preparing for the anxiety to come .	469
420	Commenter_id(1) : Sometimes we have a tendency to constantly prepare for " war " which keeps	470
421	us in this exhausting loop of hypervigilance .	471
422	Commenter_id(1) : Here is an article that has really helped me to stop letting my triggers have control	472
423	of when my anxiety pops up. Commenter_id(1) : Please read <a href="https://www.thatanxietyguy.com/">https://www.thatanxietyguy.com /</a>	473
424		474
425	Commenter_id(2) : Keep a journal ... write down everything you feel, describe it in as much	475
426	detail as you can .	476
427	Commenter_id(2) : If you have trauma in your past , write about it , try to make written connections	477
428	between what you feel now and other events in your past where you felt the same .	478
429	Commenter_id(2) : For me , my anxiety was due to past abuse , so as an adult I became an	479
430	approval seeker , validation seeker and people pleaser in an attempt to gain certainty,	480
431	safety and self esteem from my environment .	481
432	Commenter_id(2) : It's a hard road , but you need to cross the bridge of "I Don't Give a Crap"...	482
433	easier said than done but you need to realize that you don't need validation from others, you give it	483
434	to yourself, give yourself permission for everything you do think or say, you don't need to control	484
435	whether people and pleased with you, you don't owe anyone an explanation for being who you are,	485
436	feeling what you feel, wanting what you want... the hatred you feel for yourself could possibly be	486
437	linked to how you perceived a parent felt about you or treated you when you were younger.	487
438	Commenter_id(2) : Please read my other posts ... they may be helpful.	488
439		489
440		490
441		491
442		492
443		493
444		494
445		495
446		496
447		497
448		498
449		499

Table 2: An example dialogue in the RMHIDD dataset. creator\_id represents the username of the post(issue) creator, Commenter\_id represents the commented username. Here on one post have two comments

feed the encoder block of the BART model with a set of input id's from the dialogue history  $Q$ . Let the input for the first encoder layer be  $h_e^0$ . The  $h_e^0$  is converted into an embedding matrix which passes through the  $1^{st}$  layer's encoder function yielding a hidden state for the first layer. This step is repeated for each  $l^{th}$  layer, where  $l \in \{1, \dots, 12\}$ . We get a hidden state  $h_e$  for every  $l^{th}$  layer by applying the  $Encoder(\cdot)$  function as shown in equation(3). The final  $12^{th}$  layer of encoder block output its hidden state  $h_e^{12}$  which is utilized by the hidden state of the decoder layer for sequential decoding.

$$h_e^l = Encoder(h_e^{l-1}) \quad (3)$$

$$h_d^l = Decoder(h_d^{l-1} \cdot h_e^{12}) \quad (4)$$

Next we feed the set of target-response  $T = \{x_1, x_2, \dots, x_n\}$  to the decoder block. Similar to encoder block we represent each decoding layer as the  $Decoder(\cdot)$  function, which generates hidden state  $h_d$  for each decoder layer. Again, let the decoder's input be  $h_d^0$ . We feed the model's decoder block with decoder's input ids along with the hidden state of  $12^{th}$  encoder layer. With the help of the decoder block function it generates the hidden state  $h_d$  for each layers as shown in equation(4).

In the BART's language model wrapper we have included a linear layer, which generates output to-

500	<b>Creator id:</b>	550
501	what is wrong with me, and should i seek help ? in the 22 years i've been alive, i've had	551
502	depression for most of it for various reasons. i'm currently in my last year of college and	552
503	currently have the worst bout i've ever experienced. i'm very close to finishing, yet i find it	553
504	hard to focus on schoolwork because i often have my mind clouded by ideas that none of this	554
505	is worth it, and that i shouldn't bother trying . as a result, it's next to impossible to focus on	555
506	homework which leads to procrastination, and class time is hard to engage in. i find myself	556
507	unmotivated and feel trapped in a spiral that will lead to inevitable failure. this is a problem	557
508	that has persisted throughout my college career but has hit me harder now more than ever.	558
509	i know these thoughts are not true, but it still affects me nonetheless .i often end up stressing	559
510	because i keep shirking my work and thus continue to put it off. i don't believe i'm on a path	560
511	to self- destruction, but i don't want this problem to affect my life once i graduate . is this	561
512	some problem i need addressed or is it just me being lazy? what should i do ?	562
513	<b>Ground truth Response:</b>	563
514	it is totally fine to have some concern about your future since you are about to finish college. you	564
515	should see a therapist as most colleges offer therapy for free in college counseling centers. you've	565
516	essentially already paid for it as part of your tuition fees. you should relax and take some	566
517	break from the college. finding new hobbies, making new friends can help you.	567
518	in your case, it is best to talk to someone who is an expert in this field. all the best for the future.	568
519	<b>Seq2Seq:</b>	569
520	please be careful and happy. you should take a break and therapy is good for you. go to a doctor.	570
521	all the best.	571
522	<b>DialoGPT:</b>	572
523	sorry to hear about it. sleep little bit also do exercise daily. if you feel sick you should go	573
524	to a doctor. seek for support from expert. it is all right. please read my article. thanks	574
525	<b>BART:</b>	575
526	i can understand that college is difficult. just be confident i suggest you to take the therapist help and get	576
527	professional help from a family. make new friends and talk to your friends. college counselors are	577
528	good for you. don't give up all the best.	578

Table 3: Generated responses from various models on a test dialogue

kens probabilities(logits) by applying a normalized exponential function(softmax). This output helps in determining the words within a sequence. Our fine-tuned model aims at maximizing the likelihood as stated in equation(5) by training  $\theta$  parameters on minimizing cross-entropy of BART model as stated in equation(6)

$$P(T|Q) = P(x_1|Q) \prod_{i=2}^n P(x_i|Q, x_1, \dots, x_{i-1}) \quad (5)$$

$$\mathcal{L}_{xe}(\theta) = -\log P_{\theta}(T|Q) = -\sum_{t=1}^N \log P_{\theta}(y_t|y_{t:t-1}, Q) \quad (6)$$

## 5 Experiment

In this section, we elaborated on the data preprocessing steps involved in structuring the dataset. We also discussed the hyperparameter setting and optimization strategies used for training the dialogue generation model.

### 5.1 Data Preprocessing

For our experiment, due to computational limitations, we trained and evaluated all of our models on

a randomly collected subset dataset with the size of 50,000 dialogues. As described above<sup>4</sup>, each instance of dialogue in the dataset we structured them into a pair of utterances('issue' and response'). As shown in the table 3, *Creator\_id* (input utterance) and *Ground truth Response* (ground-truth response utterance) were used for training the models. For a given pair of utterance  $\{D_1, D_2\}$  in the dataset, we removed all the emojis, unnecessary symbols, and characters. We also replaced the most common abbreviations of words with their original form. We corrected the words with the misspelling. All the unwanted extra spaces in the utterances were removed. For the experiment we divided the dataset into 3 parts: train/validate/train, The distribution of dataset across data was 70%/20%/10% respectively. Hyperparameters were fine-tuned using the validation data.

### 5.2 Experiment Settings

**Seq2Seq :** In our Seq2Seq model, we used an embedding layer (trainable matrix) with dimension size of 128. The 2 LSTM layer present in the encoder and decoder was consist of 128 cells. Each of the forwarding and backward LSTM cell in the

	Seq2Seq	DialoGPT	BART
Perplexity	225.3	27.2	<b>19.7</b>
NIST-4	0.60	<b>1.82</b>	1.56
BLEU-2	2.16%	<b>9.19%</b>	7.38%
BLEU-4	1.53%	<b>2.83%</b>	1.97%
METEOR	3.71%	<b>8.60%</b>	7.53%

Table 4: Performance score of various models on automatic evaluation metrics.

single BiLSTM layer also consisted of 128 cells. The training was done using the batch size of 16 with the max input and output sequence length set to 400 and 100 respectively. We trained our model for 64 epochs with an initial learning rate of 0.001. At the dense layer, we applied a dropout with a probability of 0.2, and the beam size was set to 3.

**BART** : For our experiment, we have used  $BART_{base}$ . We used the Huggingface transformer  $BART_{base}$  model<sup>7</sup> provided by the Facebook. As specified in the paper(Lewis et al., 2019), for BART model we followed the given fine-tuning parameters and train the model for 5 epochs with batch size 64. We use the Adam(Kingma and Ba, 2014) optimizer with the linear warm-up scheduler and an initial learning rate of  $4e-5$ . We fed the encoder with noised input tokens of length 400 with its respective attention mask tokens done by the Byte-pair-encoding tokenizer. Similarly we tokenize decoder input with token length of 100 and feed decoder with its respective attention mask. The model trains both the encoder and decoder architect jointly so we get a score logit from the model. We train the model by calculating the cross entropy loss with label smoothing(factor = 0.1) from the logits. Based on validation score we save the model weights and use it on the test dataset evaluation.

**DialoGPT** : We used  $DialoGPT_{small}$ (Zhang et al., 2019) and fine-tuned the model on our dataset. The fine-tuning was done for 5 epochs and the batch size was set as 64. The token length for the encoder and decoder was set to 400 and 100 respectively. Similar to the BART model, we used Adam(Kingma and Ba, 2014) optimizer with along with liner learning schedule. The initial learning rate was decided to be  $4e-5$ . During training cross-entropy loss was calculated with the label smoothing factor of 0.1.

<sup>7</sup>Available at <https://huggingface.co/models>

## 6 Results and Discussion

We used five performance evaluation matrix to compare the performance of all dialogue generation models. We calculated the perplexity score, METEOR (Lavie and Agarwal, 2007), BLEU-n (Papineni et al., 2002) score( $n = 2, 4$ ), NIST-n (Dodington, 2002) score for  $n = 4$ . For machine translation NIST, METEOR and BLEU are very frequently used evaluation metrics. They compute the similarity by matching the n-grams ground-truth and the model’s generated response. In BLEU n-gram precision is calculated by adding equal weight whereas NIST also calculates the informativeness of a particular n-gram and penalizes the non-informative n-grams. Through Perplexity, we calculated and compared the smoothness and quality of produced responses.

In the table 4, we have summarized the performance result of all the dialogue generation model. The following are the observation we can take from the table. Firstly, the overall performance of pretrained language models was superior to the un-trained Seq2Seq model. The reason behind this was the advantage of transfer learning through which pretrained models effectively leverages the knowledge extracted from the large data. Secondly, out of all three models, BART achieved the lowest perplexity score of 19.7, whereas  $DialoGPT_{small}$  and Seq2Seq achieved a score of 27.2 and 225.3 respectively. The biggest advantage of  $BART_{base}$  was that it was trained on much bigger and diverse data in a way to reconstruct the texts from the corrupted documents, which therefore enhanced and increased BART capabilities as compared to other models. Seq2Seq model scored the highest perplexity which was on an average 89% more than the large pretrained models( $DialoGPT_{small}$  and  $BART_{base}$ ). The third observation that could be made was on the machine translation benchmark scores such as METEOR, BLEU, and NIST, the best performance was given by the  $DialoGPT_{small}$  model. Since the

*DialoGPT<sub>small</sub>* model was pretrained on a large Reddit dialogue dataset which gave the model more contextual understanding for handling our dataset and as a result, more related and relevant dialogue n-grams were produced by the model. With the advantage of large pretraining, the *BART<sub>base</sub>* model surpassed the Seq2Seq model on all the machine translation benchmark scores. In the table 3, we have provided the generated responses from all the dialogue generation models on an example dialogue from the test dataset. On average the generated dialogues length form, various models was approximately 50.

## 7 Conclusion and Future Work

In this paper, we have presented a mental health intervention dialogue dataset. We collected a large number of mental disorder related user-generated data from online platform. Using our dataset, we conducted a systematic analysis of various state-of-the-art dialogue generation language models as an attempt to develop automated mental health intervention system. In our study, we discovered that the large pretrained model (*DialoGPT<sub>small</sub>* and *BART<sub>base</sub>*) performed better than the un-trained model (Seq2Seq) on the task of dialogue response generation. The results obtained from various models were very promising and shows the potential of developing automated mental health intervention system in future. We believe that this dataset would enable computer scientist to design and develop more sophisticated, intelligent and feasibly available advance mental health intervention systems such as chatbots, that would help millions of people. In future we aim at extending our current work by collecting a large scale user-generated multilingual mental health dialogue dataset. Through this we would be able to develop a multilingual intervention systems that would not be restricted to single language.

## References

Noor Fazilla Abd Yusof, Chenghua Lin, and Frank Guerin. 2018. Assessing the effectiveness of affective lexicons for depression classification. In *International Conference on Applications of Natural Language to Information Systems*, pages 65–69. Springer.

Takahiro Baba, Kensuke Baba, and Daisuke Ikeda. 2019. Detecting mental health illness using short comments. In *International Conference on Ad-*

*vanced Information Networking and Applications*, pages 265–271. Springer.

Lisa J Barney, Kathleen M Griffiths, Anthony F Jorm, and Helen Christensen. 2006. Stigma about depression and its impact on help-seeking intentions. *Australian & New Zealand Journal of Psychiatry*, 40(1):51–54.

Xuetong Chen, Martin D Sykora, Thomas W Jackson, and Suzanne Elayan. 2018. What about mood swings: Identifying depression on twitter with temporal measures of emotions. In *Companion Proceedings of the The Web Conference 2018*, pages 1653–1660.

Munmun De Choudhury and Sushovan De. 2014. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. In *Eighth international AAAI conference on weblogs and social media*.

George Doddington. 2002. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of the second international conference on Human Language Technology Research*, pages 138–145.

George Gkotsis, Anika Oellrich, Tim Hubbard, Richard Dobson, Maria Liakata, Sumithra Velupillai, and Rina Dutta. 2016. The language of mental health problems in social media. In *Proceedings of the Third Workshop on Computational Linguistics and Clinical Psychology*, pages 63–73.

George Gkotsis, Anika Oellrich, Sumithra Velupillai, Maria Liakata, Tim JP Hubbard, Richard JB Dobson, and Rina Dutta. 2017. Characterisation of mental health conditions in social media using informed deep learning. *Scientific reports*, 7:45141.

Dritjon Gruda and Souleiman Hasan. 2019. Feeling anxious? perceiving anxiety in tweets using machine learning. *Computers in Human Behavior*, 98:245–255.

K Stanly Jacob and Vikram Patel. 2014. Classification of mental disorders: a global mental health perspective. *The Lancet*, 383(9926):1433–1435.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Alon Lavie and Abhaya Agarwal. 2007. Meteor: An automatic metric for mt evaluation with high levels of correlation with human judgments. In *Proceedings of the second workshop on statistical machine translation*, pages 228–231.

Kyung Sang Lee, Hyewon Lee, Woojae Myung, Gil-Young Song, Kihwang Lee, Ho Kim, Bernard J Carroll, and Doh Kwan Kim. 2018. Advanced daily prediction model for national suicide numbers with social media data. *Psychiatry investigation*, 15(4):344.

750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799

- 800 Angela Leis, Francesco Ronzano, Miguel A Mayer, Alec Radford, Jeffrey Wu, Rewon Child, David Luan, 850  
801 Laura I Furlong, and Ferran Sanz. 2019. Detecting Dario Amodei, and Ilya Sutskever. 2019. Language 851  
802 signs of depression in tweets in spanish: behavioral models are unsupervised multitask learners. *OpenAI 852  
803 and linguistic analysis. Journal of medical Internet Blog, 1(8):9. 853  
804 research, 21(6):e14199. 854*
- 805 Mike Lewis, Yinhan Liu, Naman Goyal, Mar- 855  
806 jan Ghazvininejad, Abdelrahman Mohamed, Omer 856  
807 Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. 857  
808 Bart: Denoising sequence-to-sequence pre-training 858  
809 for natural language generation, translation, and 859  
810 comprehension. *arXiv preprint arXiv:1910.13461. 860*
- 811 Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, 861  
812 and Bill Dolan. 2015. A diversity-promoting objec- 862  
813 tive function for neural conversation models. *arXiv 863  
814 preprint arXiv:1510.03055. 864*
- 815 Huong Ly Tong Ahmet Baki Kocaballi Jessica Chen 865  
816 Rabia Bashir Didi Surian Blanca Gallego Farah Ma- 866  
817 grabi Annie YS Lau Liliana Laranjo, Adam G Dunn. 867  
818 2018. Conversational agents in healthcare: a sys- 868  
819 tematic review. *Journal of the American Medical 869  
820 Informatics Association. 870*
- 821 Gale M Lucas, Albert Rizzo, Jonathan Gratch, Ste- 871  
822 fan Scherer, Giota Stratou, Jill Boberg, and Louis- 872  
823 Philippe Morency. 2017. Reporting mental health 873  
824 symptoms: breaking down barriers to care with vir- 874  
825 tual human interviewers. *Frontiers in Robotics and 875  
826 AI, 4:51. 876*
- 827 Julissa Murrieta, Christopher C Frye, Linda Sun, 877  
828 Linh G Ly, Courtney S Cochancela, and Elizabeth V 878  
829 Eikley. 2018. # depression: Findings from a litera- 879  
830 ture review of 10 years of social media and depres- 880  
831 sion research. In *International Conference on Infor- 881  
832 mation, pages 47–56. Springer. 882*
- 833 World Health Organization et al. 2017. Depression and 883  
834 other common mental disorders: global health es- 884  
835 timates. Technical report, World Health Organiza- 885  
836 tion. 886
- 837 Kishore Papineni, Salim Roukos, Todd Ward, and Wei- 887  
838 Jing Zhu. 2002. Bleu: a method for automatic eval- 888  
839 uation of machine translation. In *Proceedings of the 889  
840 40th annual meeting of the Association for Computa- 890  
841 tional Linguistics, pages 311–318. 891*
- 842 Albert Park and Mike Conway. 2018. Harness- 892  
843 ing reddit to understand the written-communication 893  
844 challenges experienced by individuals with mental 894  
845 health disorders: Analysis of texts from mental 895  
846 health communities. *Journal of medical Internet re- 896  
847 search, 20(4):e121. 897*
- 848 Pierre Philip, Jean-Arthur Micoulaud-Franchi, Patricia 898  
849 Sagaspe, Etienne De Sevin, Jérôme Olive, Stéphanie 899  
850 Bioulac, and Alain Sauteraud. 2017. Virtual human 900  
851 as a new diagnostic tool, a proof of concept study 901  
852 in the field of major depressive disorders. *Scientific 902  
853 reports, 7(1):1–7. 903*
- 904 Punet KC Sahota and Pamela L Sankar. 2020. Bipolar 904  
854 disorder, genetic risk, and reproductive decision- 905  
855 making: A qualitative study of social media discus- 906  
856 sion boards. *Qualitative health research, 30(2):293– 907  
857 302. 908*
- 909 Judy Hanwen Shen and Frank Rudzicz. 2017. De- 909  
859 tecting anxiety through reddit. In *Proceedings of 910  
860 the Fourth Workshop on Computational Linguistics 911  
861 and Clinical Psychology—From Linguistic Signal to 912  
862 Clinical Reality, pages 58–65. 913*
- 914 Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. 914  
863 Sequence to sequence learning with neural networks. 915  
864 In *Advances in neural information processing sys- 916  
865 tems, pages 3104–3112. 917*
- 918 Hiroki Tanaka, Hideki Negoro, Hidemi Iwasaka, and 918  
866 Satoshi Nakamura. 2017. Embodied conversational 919  
867 agents for multimodal automated social skills train- 920  
868 ing in people with autism spectrum disorders. *PLoS 921  
869 one, 12(8):e0182151. 922*
- 923 Robert Thorstad and Phillip Wolff. 2019. Predicting 923  
872 future mental illness from social media: A big-data 924  
873 approach. *Behavior research methods, 51(4):1586– 925  
874 1600. 926*
- 927 Zhongyu Wei, Qianlong Liu, Baolin Peng, Huaixiao 927  
875 Tou, Ting Chen, Xuan-Jing Huang, Kam-Fai Wong, 928  
876 and Xiang Dai. 2018. Task-oriented dialogue sys- 929  
877 tem for automatic diagnosis. In *Proceedings of the 930  
878 56th Annual Meeting of the Association for Computa- 931  
879 tional Linguistics (Volume 2: Short Papers), pages 932  
880 201–207. 933*
- 934 JT Wolohan, Misato Hiraga, Atreyee Mukherjee, Zee- 934  
881 shan Ali Sayyed, and Matthew Millard. 2018. De- 935  
882 tecting linguistic traces of depression in topic- 936  
883 restricted text: Attending to self-stigmatized depres- 937  
884 sion with nlp. In *Proceedings of the First Interna- 938  
885 tional Workshop on Language Cognition and Com- 939  
886 putational Models, pages 11–21. 940*
- 941 Akkapon Wongkoblaph, Miguel A Vadillo, and Vasa 941  
887 Curcin. 2018. A multilevel predictive model for 942  
888 detecting social network users with depression. In 943  
889 *2018 IEEE International Conference on Healthcare 944  
890 Informatics (ICHI), pages 130–135. IEEE. 945*
- 946 Akkapon Wongkoblaph, Miguel A Vadillo, and Vasa 946  
892 Curcin. 2019. Modeling depression symptoms from 947  
893 social network data through multiple instance learn- 948  
894 ing. *AMIA Summits on Translational Science Pro- 949  
895 ceedings, 2019:44. 950*
- 951 Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V 951  
896 Le, Mohammad Norouzi, Wolfgang Macherey, 952  
897 Maxim Krikun, Yuan Cao, Qin Gao, Klaus 953  
898 Macherey, et al. 2016. Google’s neural machine 954  
899 955

900 translation system: Bridging the gap between hu- 950  
901 man and machine translation. *arXiv preprint* 951  
902 *arXiv:1609.08144*. 952

903 Lin Xu, Qixian Zhou, Ke Gong, Xiaodan Liang, Jian- 953  
904 heng Tang, and Liang Lin. 2019. End-to-end 954  
905 knowledge-routed relational dialogue system for au- 955  
906 tomatic diagnosis. In *Proceedings of the AAAI Con-* 956  
907 *ference on Artificial Intelligence*, volume 33, pages 957  
908 7346–7353. 958

909 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, 959  
910 Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing 960  
911 Liu, and Bill Dolan. 2019. Dialogpt: Large-scale 961  
912 generative pre-training for conversational response 962  
913 generation. *arXiv preprint arXiv:1911.00536*. 963  
914 964  
915 965  
916 966  
917 967  
918 968  
919 969  
920 970  
921 971  
922 972  
923 973  
924 974  
925 975  
926 976  
927 977  
928 978  
929 979  
930 980  
931 981  
932 982  
933 983  
934 984  
935 985  
936 986  
937 987  
938 988  
939 989  
940 990  
941 991  
942 992  
943 993  
944 994  
945 995  
946 996  
947 997  
948 998  
949 999