

# Sentiment-guided Transformer with Severity-aware Contrastive Learning for Depression Detection on Social Media

Tianlin Zhang<sup>1</sup> Kailai Yang<sup>1</sup> Sophia Ananiadou<sup>1,2</sup>

<sup>1</sup>National Centre for Text Mining, Department of Computer Science,  
The University of Manchester, United Kingdom

<sup>2</sup>Alan Turing Institute, United Kingdom  
{tianlin.zhang,kailai.yang}@postgrad.manchester.ac.uk  
sophia.ananiadou@manchester.ac.uk

## Abstract

Early identification of depression is beneficial to public health surveillance and disease treatment. There are many models that mainly treat the detection as a binary classification task, such as detecting whether a user is depressed. However, identifying users' depression severity levels from posts on social media is more clinically useful for future prevention and treatment. Existing severity detection methods mainly model the semantic information of posts while ignoring the relevant sentiment information, which can reflect the user's state of mind and could be helpful for severity detection. In addition, they treat all severity levels equally, making the model difficult to distinguish between closely-labeled categories. We propose a sentiment-guided Transformer model, which efficiently fuses social media posts' semantic information with sentiment information. Furthermore, we also utilize a supervised severity-aware contrastive learning framework to enable the model to better distinguish between different severity levels. The experimental results show that our model achieves superior performance on two public datasets, while further analysis proves the effectiveness of all proposed modules.

## 1 Introduction

According to the World Health Organization (WHO)<sup>1</sup>, depression is a common mental disorder and a leading cause of disability worldwide. It is estimated that around 350 million people suffer from depression and over 70% of people do not receive timely treatment (Olfson et al., 2016). Therefore, if depression is detected at an early stage, it could be useful for disease intervention and treatment. In addition, with social media becoming increasingly popular, people often share their feelings and mental states through social media platforms such as

Twitter and Reddit, making social media posts critical for mental illness detection, including depression detection (De Choudhury et al., 2013; Skaik and Inkpen, 2020; Chiong et al., 2021).

Recent advances in natural language processing (NLP) have played an increasingly essential role in supporting the analysis of user-generated contents from social media, including mental illness surveillance (Skaik and Inkpen, 2020; Chancellor and De Choudhury, 2020; Ríssola et al., 2021). There are many methods leveraging NLP technologies for automated depression detection. However, existing depression detection methods mainly treat detection as a binary classification task due to the limitation of existing annotated datasets. According to the Beck's Depression Inventory (BDI) (Beck et al., 1961), not all cases of depression are the same; depression could be classified as: *minimal*, *mild*, *moderate* and *severe* (see Figure 1). For our methods to be translational, detecting the severity level of depression is expected to be more important for further prevention and treatment (Hollon et al., 2002). For instance, evidence-based psychological therapies could be appropriate for mild or moderate depression (Naseem et al., 2022), while users with severe level might be provided more resource-intensive interventions (Desrochers and Houck, 2014).

Previous psychological studies (Rude et al., 2004; Molendijk et al., 2010) have examined the correlation between sentiment and depression, illustrating that sentiment can reflect the user's state of mind and help to identify depressed individuals (Babu and Kanaga, 2022; Zhang et al., 2023). For example, a depressed user may post a comment "Today was a really bad day. I had no energy and thought about suicide all day" with a negative sentiment, while non-depressed individuals tend to express a positive or neutral attitude, such as "idk why I'm doing this. I guess I just am. Maybe it will help someone else". However, existing sever-

<sup>1</sup><https://www.who.int/health-topics/depression>

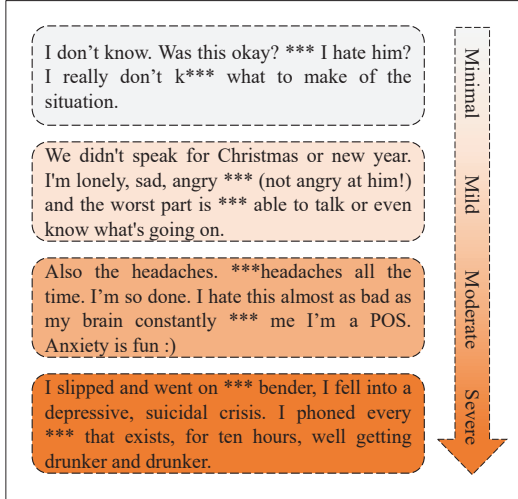


Figure 1: Example of posts sorted in order of increasing severity level (minimal → mild → moderate → severe). To prevent misuse, all posts have been obfuscated and paraphrased.

ity detection methods mainly focus on modeling the semantic information of posts while ignoring the relevant sentiment information. Additionally, traditional classification models with cross-entropy loss treat all severity levels equally, making it challenging to distinguish between closely-labeled categories. For example, ‘*moderate*’ and ‘*severe*’ are semantically closer in severity level than ‘*mild*’ and ‘*severe*’, making them more difficult to distinguish. To tackle the above challenges, we propose a sentiment-guided Transformer model with severity-aware contrastive learning for depression severity detection. We first leverage post encoders MentalRoBERTa and SentiLARE (Ke et al., 2020) to obtain the semantic as well as sentiment hidden features and then design a sentiment-guided Transformer model to better fuse these two features. For the second challenge, we implement a supervised severity-aware contrastive learning framework, which helps the model adaptively differentiate the weights between different negative samples based on the label information. We validate our method on two depression severity detection datasets. The experimental results show that our model consistently outperforms other baselines, such as DepressionGCN (Naseem et al., 2022) and MentalRoBERTa (Ji et al., 2022). Our main contributions can be summarized as follows:

(1) We propose a sentiment-guided Transformer fusion model to fuse sentiment information into the depression severity detection task. This network

can better combine semantic information and sentiment information compared with existing PLMs-based models.

(2) We implement a supervised severity-aware contrastive learning framework to enable our model to better capture class-specific features and distinguish between different severity levels.

(3) We conduct experiments and further analysis on mainstream public datasets. The experimental results illustrate that our model achieves competitive performance compared with others and each of the proposed modules is effective in depression detection.

## 2 Related Work

With the fast-growing numbers of social media users, using user-generated posts to detect depression manually is no longer practical. Automated NLP and text mining technologies provide new opportunities for mental health analysis (Ríssola et al., 2021; Zhang et al., 2022).

Early researchers (Islam et al., 2018; Trifan et al., 2020) extracted linguistic features (e.g., Linguistic Inquiry and Word Count (LIWC), Bag-of-Words) and statistical features (e.g., n-gram, Term Frequency-Inverse Document Frequency (TF-IDF)) from posts, then employed them in traditional machine learning methods such as SVM (Tadesse et al., 2019), Logistic Regression (Martínez-Castaño et al., 2018) and Random Forest (Cacheda et al., 2019) to detect depression. Recently, deep learning-based models have been used for depression detection, allowing models to automatically capture valuable features without feature engineering. Some works have shown that CNN (Orabi et al., 2018) and RNN (Wu et al., 2020), including Long Short Term Memory (LSTM) (Ghosh and Anwar, 2021) and Gated Recurrent Unit (GRU) (Sekulic and Strube, 2019) are effective. Furthermore, PLMs such as BERT (Owen et al., 2020), RoBERTa (Dinu and Moldovan, 2021), (Ji et al., 2023) are also widely utilized for many mental illness detection tasks due to strong context modelling ability. In particular, Ji et al. (2022) trained the domain-specific language model MentalRoBERTa for mental healthcare, which introduces the domain-related knowledge and outperforms basic PLMs on several mental health datasets. The depression severity identification task is becoming increasingly attractive due to the practical clinical implications, so some studies have emerged

(Losada et al., 2019; Kayalvizhi et al., 2022). Although the severity level is measured on an ordinal scale, current methods treat it as a traditional multi-classification task. Naseem et al. (2022) first reformulate depression identification as an ordinal classification problem via ordinal loss and employ Graph Convolutional Neural Network (GCN) to capture contextual information, which achieves state-of-the-art performance on depression severity datasets.

Sentiment analysis is one of the key technologies in NLP, which aims to analyse people’s sentiments or opinions expressed in texts (usually classified as positive, negative and neutral). Sentiment analysis is widely applied to reviews, social media, newswires, and medical informatics (Yue et al., 2019; Colón-Ruiz and Segura-Bedmar, 2020; Yang et al., 2023b). Recently, PLMs-based models have achieved competitive performances in sentiment analysis. For example, Yin et al. (2020) propose SentiBERT, a variant of BERT model that effectively captures compositional sentiment semantics. Ke et al. (2020) present SentiLARE with a context sentiment mechanism and knowledge integration, which facilitates sentiment understanding.

Contrastive learning is a recent technique for enhancing the performance of models by learning different representations of contrasting samples. The main purpose is to generate negative pairs using data augmentation and minimise the contrastive loss of the positive pairs. Similar methods are transferred to NLP tasks, like ConSERT (Yan et al., 2021) for better unsupervised sentence representation learning. Supervised contrastive learning can leverage label information for representation learning in a supervised setting (Gao et al., 2021), which improves the ability of the model to differentiate between samples with different labels. Recent work also uses supervised contrastive learning for depression detection tasks (Yang et al., 2022; Wang et al., 2022), illustrating its effectiveness.

### 3 Methodology

#### 3.1 Task definition

We focus on detecting the depression severity level of users by analysing posts on social media. Formally, given a post  $P = \{s_1, s_2, \dots, s_i\}$ , where  $s_i$  is the  $i$ -th sentence of the post, our goal is to classify each post into a corresponding depression severity level label  $y \in Y$ , where  $Y$  is the list of increasing depression severity levels. For example, according

to the definition of different datasets described in Section 4,  $Y$  is  $\{minimal, mild, moderate, severe\}$  or  $\{not\ depressed, moderately\ depressed, severely\ depressed\}$ . This task can be formalized as a multi-class text classification problem.

#### 3.2 Depression severity detection model

Our model architecture (Figure 2) consists of a post encoder, a sentiment-guided Transformer and a supervised severity-aware contrastive learning component. We will introduce the details of each component as follows.

##### 3.2.1 Post encoder

With the huge success of PLMs in several NLP tasks, PLMs are widely used as sentence encoder for various tasks, including for mental health detection (Zhang et al., 2022). In particular, Ji et al. (2022) trained a domain-specific language model MentalRoBERTa for mental healthcare, which introduces domain-related knowledge and outperforms basic PLMs on several mental health datasets. We exploited MentalRoBERTa as the semantic encoder. Since we require sentiment information to help us to better understand a user’s state of mind, we also adopt the pre-trained language model SentiLARE (Ke et al., 2020) that achieves good performance on various sentiment analysis tasks, as the sentiment encoder to obtain the sentiment hidden representation of each sentence. Specifically, given a post  $P = \{s_1, s_2, \dots, s_i\}$ , we have obtained  $i$  split sentences and each sentence is  $s = \{[CLS], x_1, x_2, \dots, x_j\}$ , where  $[CLS]$  is a special token for marking the start of the sentence and  $x_j$  is the  $j$ -th word. Then, we feed each sentence sequence into MentalRoBERTa and SentiLARE to get the two embeddings of the  $[CLS]$  as the corresponding sentence embeddings. We denote the process as:

$$\begin{aligned} H^m &= \{H_1^m, H_2^m, \dots, H_i^m\} \\ &= Mental\_Encoder(\{s_1, s_2, \dots, s_i\}) \end{aligned} \quad (1)$$

$$\begin{aligned} H^s &= \{H_1^s, H_2^s, \dots, H_i^s\} \\ &= Senti\_Encoder(\{s_1, s_2, \dots, s_i\}) \end{aligned} \quad (2)$$

where *Mental\_Encoder* and *Senti\_Encoder* denote the MentalRoBERTa encoder and SentiLARE encoder respectively;  $H_i$  denotes the  $i$ -th sentence embedding;  $H^m$  and  $H^s \in \mathbb{R}^{n \times d}$  are encoded input embeddings,  $n$  denotes the number of sentences, and  $d$  is the hidden dimension of the encoder.

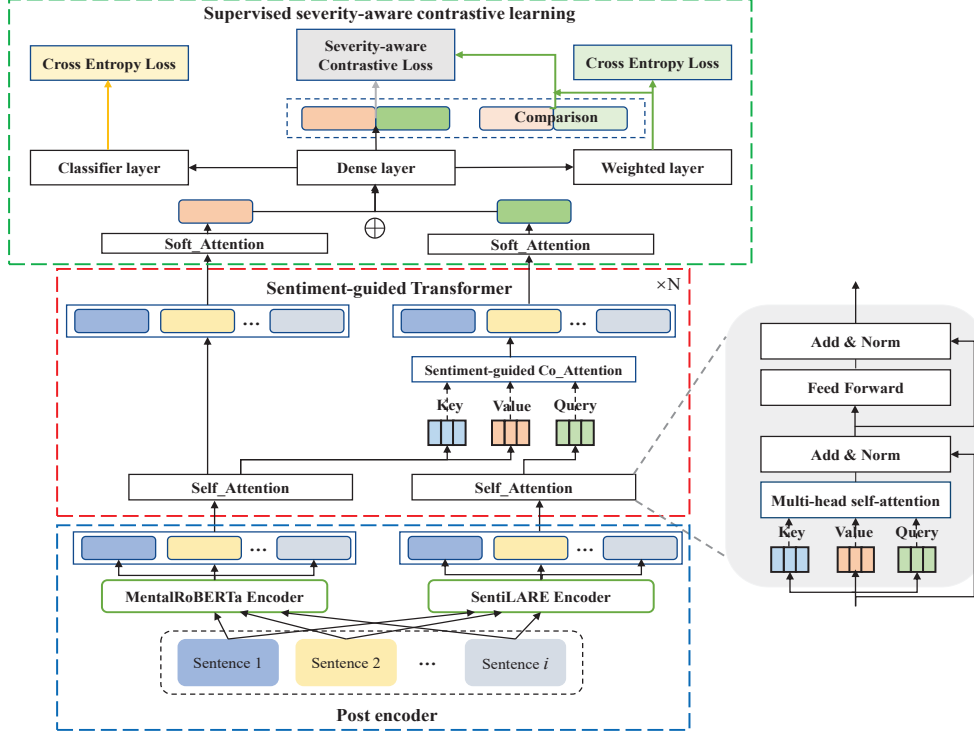


Figure 2: The overall architecture of our proposed model. The model contains three main components: post encoder, sentiment-guided Transformer and supervised severity-aware contrastive learning. The symbol  $\oplus$  denotes vector concatenation. The internal architecture of Self\_Attention module is shown in light grey block. More details about our model are introduced in the main text.

### 3.2.2 Sentiment-guided Transformer

The architecture of the Sentiment-guided Transformer is illustrated in the red-dotted box of Figure 2, which is composed of a stack of  $N$  blocks with their own training parameters. Each block contains a Self\_Attention module and a Sentiment-guided Co\_Attention module. For one social media post, there are two post embeddings generated as inputs  $\{H^m; H^s\}$  for sentiment-guided Transformer.

Firstly, the two embeddings are fed into Self\_Attention module separately. This module is the most important component of Google’s Transformer (Vaswani et al., 2017) that learns the sequence’s representation consisting of a multi-head self-attention mechanism (MHA), a feed-forward network and two layers of normalization (Ba et al., 2016) with residual connection (He et al., 2016). The self-attention mechanism can be described as the relationships between a query and a set of key–value pairs to compute an output, where query, key, value and output are vectors. For semantic embedding, given the matrix  $H^m$ , we can get a query matrix  $Q \in \mathbb{R}^{n \times d}$ , a key matrix  $K \in \mathbb{R}^{n \times d}$  and a value matrix  $V \in \mathbb{R}^{n \times d}$ . The self-attention

function is expressed as:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (3)$$

where  $n$  is the length of the post,  $d$  is the embedding hidden dimension,  $softmax$  denotes softmax operation,  $\sqrt{d}$  is a scaling factor.  $Q = K = V = H^m$ .

To enable the model to jointly learn contextual information from different representation sub-spaces, the three matrices  $Q, K, V$  are multiplied with the  $W_i^Q, W_i^K, W_i^V$  respectively, and the MHA with  $h$  heads is used. Then, the results of all attention heads are concatenated together, and a weight matrix  $W^O$  is used to obtain the final output of the encoder:

$$MHA(Q, K, V) = concat(head_1, \dots, head_h)W^O \quad (4)$$

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V) \quad (5)$$

where parameter weight matrices  $W_i^Q \in \mathbb{R}^{d \times d_h}$ ,  $W_i^K \in \mathbb{R}^{d \times d_h}$ ,  $W_i^V \in \mathbb{R}^{d \times d_h}$  and  $W^O \in \mathbb{R}^{d \times d}$ ,  $d_h = d/h$ .

After Self\_Attention module, we will get the semantic-level output  $SA^m$  and the sentiment-level output  $SA^s$ . In order to obtain contextual sentiment-guided sentence representation, we leverage the Sentiment-guided Co\_Attention module to fuse sentiment information with the corresponding semantic information. In contrast to the Self\_Attention, the inputs of Sentiment-guided Co\_Attention are from two different embeddings.

We switch the query  $Q^s$  from the Self\_Attention output  $SA^s$  (the sentiment-level post hidden representation), the key  $K^m$  and the value  $V^m$  from  $SA^m$  (the semantic-level post hidden representation). The calculation process is the same as the multi-head self-attention mechanism:

$$Co\_A(Q^s, K^m, V^m) = softmax\left(\frac{Q^s(K^m)^T}{\sqrt{d}}\right)V^m \quad (6)$$

### 3.2.3 Supervised severity-aware contrastive learning

After all the Sentiment-guided Transformer blocks are computed, we obtain the final block's outputs, including the semantic hidden feature representation  $o^m = \{o_1^m, o_2^m, \dots, o_i^m\}$  and the sentiment-guided hidden feature representation  $o^s = \{o_1^s, o_2^s, \dots, o_i^s\}$ . Considering the different importance of each sentence to the post, we leverage Soft\_Attention mechanism (Yang et al., 2016) to get semantic-level post embedding  $p^m$  and sentiment-level post embedding  $p^s$ :

$$u_i = \tanh(W_p o_i + b_p) \quad (7)$$

$$\alpha_i = \frac{\exp(u_i^T u_p)}{\sum_i \exp(u_i^T u_p)} \quad (8)$$

$$p = \sum_i \alpha_i o_i \quad (9)$$

Where  $W_p$  is a weight matrix,  $b_p$  is a bias vector, and  $u_p$  is a post context vector which is randomly initialized. Then, we concatenate these two embeddings  $p^f = \text{concat}(p^m, p^s)$  as the fusion post embedding that retains contextual semantic and sentiment information.

Recent works (Gao et al., 2021; Yang et al., 2022; Wang et al., 2022) utilize supervised contrastive learning to capture the class-specific information, enabling sentences with the same label to be cohesive and those with different labels to be mutually exclusive. In addition, following (Suresh and

Ong, 2021), we design a supervised severity-aware contrastive learning framework to help the model differentiate the more difficult negatives.

To ensure that there are at least two samples of a class in one batch for contrastive learning, according to (Yang et al., 2022), we copy each post embedding to produce augmented samples with the batch size will changing from  $N_b$  to  $2N_b$ . The traditional supervised contrastive loss is given by:

$$L_{CL} = \sum_{i=1}^{2N_b} -\frac{1}{|J|} \sum_{j \in J} \log \frac{\exp(h_i^f \cdot h_j^f / \tau)}{\sum_{k \in I \setminus i} \exp(h_i^f \cdot h_k^f / \tau)} \quad (10)$$

$$h_i^f = L2 - \text{Normaliz}(\text{Dense\_layer}(p_i^f)) \quad (11)$$

where  $J = \{j : j \in I, y_j = y_i \wedge j \neq i\}$ ,  $i$  is the sample index,  $h_i^f$  is the L2-normalized vector of post hidden representation obtained from an encoder (dense layer in Figure 2), and  $\tau$  indicates the temperature hyper-parameter.

We can see that traditional supervised contrastive learning treats both negative samples (the comparison in Figure 2) the same. However, in the depression severity detection task, the labels contain an inherent ordinal nature. For example, the sample with the label *severe* is semantically closer to the sample with the label *moderate* than the one with the label *mild* in the representation space. Therefore, to distinguish better the closely-labeled examples, we introduce supervised severity-aware contrastive loss for adaptively weighting the samples based on the relationships between the severity labels.

In order to obtain the adaptive weights for different comparison samples, following (Suresh and Ong, 2021), we use a dual-model strategy as shown in Figure 2. The post hidden representations are fed into a weighted layer that is same as the classifier layer, and the output is optimised using cross-entropy loss  $L_w$ . Then, by using softmax function to obtain prediction probabilities, we can get the confidence scores of the sample that are classified to different classes:

$$hh_i^f = \text{weighted\_layer}(h_i^f) \quad (12)$$

$$w_i = \{w_{i,c}\}_{c=1}^C = \frac{\exp(hh_{i,c}^f)}{\sum_{c=1}^C \exp(hh_{i,c}^f)} \quad (13)$$

where  $C$  is the number of the classes,  $w_{i,c}$  denotes the confidence that  $i$ -th sample is classified to class  $c$ ,  $hh_{i,c}^f$  is the  $c$ -th value of  $hh_i^f$ .

The  $w_i$  is taken as weighting vector to weight the pair-wise similarity values, and the severity-aware contrastive loss is as follows:

$$L_s = \sum_{i=1}^{2N_b} -\frac{1}{|J|} \sum_{j \in J} \log \frac{w_{i,y_i} \exp(h_i^f \cdot h_j^f / \tau)}{\sum_{k \in I \setminus i} w_{i,y_k} \exp(h_i^f \cdot h_k^f / \tau)} \quad (14)$$

Where  $w_{i,y_i}$  represents the weight between the  $i$ -th sample and its corresponding severity label  $y_i$ . We can find that the loss assign higher weights to confusable negatives based on the model’s confidence scores.

We combine the training of depression severity detection and the supervised severity-aware contrastive learning task in a multi-task learning framework. The framework is jointly optimised using the following training losses: (1) the classification output’s cross-entropy loss  $L_c$ ; (2) the severity-aware contrastive loss  $L_s$ ; (3) the weighted layer output’s cross-entropy loss  $L_w$ . We formalized the overall training loss as follows:

$$L_{overall} = \alpha(L_c + L_w) + (1 - \alpha)L_s \quad (15)$$

where  $\alpha$  is a tunable weight parameter.

## 4 Experiments and analysis

### 4.1 Datasets and Implementation Details

To evaluate and compare our method with related works, we conducted experiments on two publicly available depression severity detection datasets: DsD (Naseem et al., 2022) and DepSign (Kayalvizhi and Thenmozhi, 2022). The statistical details and some examples of these two datasets are shown in Table 1 and Table 2.

**Depression severity Dataset (DsD):** The dataset is collected from the social media platform Reddit and annotated according to the disorder annotation scheme (Mowery et al., 2015), where each post is labeled with one of four severity levels:  $\{minimal, mild, moderate, severe\}$  following the depression rating scale by (Zimmerman et al., 2013).

**Depression signs detection Dataset (DepSign):** The dataset is from shared task and also collected from Reddit. The posts belong to the following subreddits groups, "r/depression, r/mental health". The data is annotated by two domain experts according

to the signs of depression, and each post is labeled with the severity level of signs of depression:  $\{not\ depressed, moderate, severe\}$ . Since the shared task does not make the test set publicly available, we merge the train set and valid set as the whole dataset and then construct train/valid/test sets.

We perform pre-processing steps on the datasets to clean the texts. We use the *ekphrasis*<sup>2</sup> library that is a processing tool for texts from social media to correct spelling, remove hashtags and normalize some special words (URLs, emails, digits, user names, elongated words, etc.). We further use *emoj*<sup>3</sup> to replace emoticons and emojis with the associated words, then leverage *NLTK*<sup>4</sup> to split the post.

We use grid search to explore the parameters. The parameter settings are as follows: the number of Sentiment-guided Transformer blocks  $N \in \{1, 2, 4, 8\}$ , the number of hidden dimension  $H \in \{128, 256, 512, 1024\}$ , the number of heads  $h \in \{2, 4, 8, 16\}$ , dropout  $\delta \in \{0.1, 0.2, 0.3, 0.4\}$ , learning rate  $lr \in \{5 \times 10^{-4}, 1 \times 10^{-4}, 5 \times 10^{-5}, 1 \times 10^{-5}\}$ , the optimization  $O \in \{Adam, AdamW\}$ , the batch size  $b \in \{8, 16, 32\}$ , tunable weight parameter  $\alpha \in \{0.2, 0.5, 0.8\}$ . The optimal hyperparameters used for both datasets are  $N = 4, H = 1024, h = 8, \delta = 0.2, lr = 5 \times 10^{-5}, O = AdamW, b = 16, \alpha = 0.5$ . We train the model on an Nvidia Tesla V100 GPU with 16GB of memory for 100 epochs and apply early stopping with patience of 20 epochs. We use Huggingface’s Transformers library<sup>5</sup> as pre-trained language models. For consistency, we use the same experimental settings and 10-fold cross-validation for all models. Our reported results are averaged across all folds on both datasets.

### 4.2 Evaluation metrics

To better evaluate our model from the point of view of depression severity assessment, we choose the evaluation metrics Graded Precision (GP), Graded Recall (GR) and FScore used by (Naseem et al., 2022; Gaur et al., 2019). They altered the formulation of False Positives (FP) and False Negatives (FN). FP is changed as the proportion of the number of times the predicted severity level ( $l^p$ ) is higher than the actual severity level ( $l^a$ ). FN is the propor-

<sup>2</sup><https://pypi.org/project/ekphrasis/>

<sup>3</sup><https://pypi.org/project/emoj/>

<sup>4</sup><https://www.nltk.org/>

<sup>5</sup><https://huggingface.co/models>

Stats/Datasets	DsD	DepSign
Number of posts	3,553	13,387
Average number of sentences per post	4.79	9.25
Average number of tokens per post	87.91	154.55
Label	{ <i>minimal, mild, moderate, severe</i> }	{ <i>not depressed, moderate, severe</i> }

Table 1: Datasets Statistics

Dataset	Examples	Label
DsD	It was created by a friend of the h**d. She’s in for a long, hard road after she gets done with this f**l. If anyone can help. Please do.	minimal
	I have NOONE to talk too. stress I’ve taken d**bt classes. Yes I’ve tried grounding tech**es and get frustrated after a b**ch don’t work. Cannot Afford a pyscologist Now I want to drink till I’m drunk a**n and my head sounds as though there is yelling when I’m the only one home. At least I have a drs ap**t coming up soon so I can b**g it up.	mild
	Try and learn me**on. I thought it wasn’t my scene but I got an app called H**e and it surprisingly helps. I’ve been close to killing myself 3 times.	moderate
DepSign	Youre all tough. : Youre all s**g people, youve gone t**h so much and made it this far.	not depressed
	Somone else Feeling like 2020 will be t**e last year on earth b**e even wen your hammerd your Feeling like a moron thats depressed?	moderate
	Words can’t d**e how bad I feel right now : I just w**t to fall asleep forever.	severe

Table 2: Some examples of the two datasets. The posts have been obfuscated and paraphrased for user privacy.

tion of the number of times the predicted severity level  $l^p$  is lower than the actual severity level  $l^a$ .

$$FP = \frac{\sum_{i=1}^{N^T} I(l_i^p > l_i^a)}{N^T}, FN = \frac{\sum_{i=1}^{N^T} I(l_i^a > l_i^p)}{N^T} \quad (16)$$

where  $\Delta(l_i^p, l_i^a)$  denotes the difference between predicted  $l^p$  and actual  $l^a$  depression severity level for post  $P_i$ ,  $N^T$  is the size of test set.

The precision and recall are reformulated as GP and GR since FP and FN contain the comparison between  $l^p$  and  $l^a$ :

$$GP = \frac{TP}{TP + FP}, GR = \frac{TP}{TP + FN} \quad (17)$$

$$FScore = \frac{2GP \cdot GR}{GP + GR} \quad (18)$$

where TP is the True Positives, FScore is the F1-score calculated by GP and GR.

### 4.3 Overall results

We compare our model with recent deep learning-based methods to show the overall performance. Table 3 shows the experimental results of our model and the baselines. The comparative approaches are as follows:

- **DepressionNet** (Hamad et al., 2021) : a text summarization model for depression detection.

Model	DsD			DepSign		
	GP	GR	FS	GP	GR	FS
DepressionNet	80.0	70.0	78.0	-	-	-
DepressionGCN	<b>95.0</b>	75.0	85.0	84.2	88.3	86.2
BERT	92.6	80.2	85.9	83.7	89.3	86.4
MentalBERT	93.0	80.8	86.5	84.3	89.9	87.0
MentalRoBERTa	94.1	80.4	86.7	84.4	<b>90.3</b>	87.2
Our model	93.1	<b>82.7</b>	<b>87.5</b>	<b>88.1</b>	90.2	<b>89.1</b>

Table 3: Experimental performance comparison for different models. We highlight top-1 values in bold.

- **DepressionGCN** (Naseem et al., 2022) : a recently proposed depression severity detection method combining GCN and Bi-LSTM structures.
- **BERT** (Devlin et al., 2019) : a widely used basic pre-trained language model.
- **MentalBERT** and **MentalRoBERTa** (Ji et al., 2022) : two domain specific PLMs pre-trained on mental health corpus. Moreover, MentalRoBERTa is the current state-of-the-art model on this task.

As shown in Table 3, we notice that PLMs-based models, especially mental health-related PLMs, offer a general advantage with over 85% performance and outperform DepressionGCN on both datasets, which shows the advantage of context modelling for PLMs. Compared with these baselines, our model achieves a new state-of-the-art on both datasets (87.5% / 89.1% of FScore), with 0.8% and 1.9% improvements over the MentalRoBERTa.

In summary, the results prove that the proposed model achieves better performance on the depression severity detection task.

#### 4.4 Ablation studies

In this section, we perform ablation studies of our model on both DsD and DepSign datasets to illustrate the effectiveness of each module of our proposed model. Specifically, we conduct module ablation experiment, sentiment fusion ablation experiment and contrastive ablation experiment.

##### 4.4.1 Module ablation

To evaluate the contribution of different modules in our model, we perform a module ablation study. ‘w/o contrastive’ denotes the model without severity-aware contrastive learning module. ‘w/o SGT’ denotes the model without sentiment-guided Transformer module. ‘w/o all’ removes all modules while keeping MentalRoBERTa encoder. The results are shown in Table 4. We note that our complete model achieves the best performance while the performance degrades when components are removed. In detail, e.g., with DepSign, the performance of FScore decreases by 0.6% when contrastive learning is removed and the FScore drops by 1.5% when SGT is removed. The model without SGT consistently performs worse than the model without contrastive, which indicates the importance of sentiment information in depression detection. If the two modules are both removed, the model degrades to the basic MentalRoBERTa encoder and the FScore drops by 1.9%. The results above demonstrate that each proposed module works effectively and the performance is improved when these components are combined together.

Model	DsD			DepSign		
	GP	GR	FS	GP	GR	FS
Our model	93.1	<b>82.7</b>	<b>87.5</b>	<b>88.1</b>	90.2	<b>89.1</b>
w/o contrastive	93.2	81.8	87.1	86.8	<b>90.3</b>	88.5
w/o SGT	<b>94.3</b>	80.6	86.9	85.2	90.2	87.6
w/o all	94.1	80.4	86.7	84.4	<b>90.3</b>	87.2

Table 4: The results of module ablation study. We highlight top-1 values in bold.

##### 4.4.2 Sentiment fusion ablation

In order to demonstrate the effectiveness of our proposed sentiment fusion strategy (sentiment-guided Transformer), we compare it with other methods. ‘SA(m)’ means only using hidden representations (m) from MentalRoBERTa encoder combining Self\_Attention module. ‘SA(m)+s’ denotes the

concatenation that the representations from SA(m) are concatenated with a sentiment label, which is a common fusion strategy to add external features. We leverage SenticNet 6 (Cambria et al., 2020) to obtain the sentiment label of each sentence and the concatenation is implemented at the classification layer. ‘SA(m)+SA(s)’ denotes the concatenation of semantic representations from SA(m) and sentiment representations from SA(s). ‘SGT(m+s)’ denotes our sentiment-guided Transformer. In addition, we also explore whether emotion information is valuable for depression severity detection, thus we fuse semantic information and emotion information using the emotion-guided Transformer module (EGT). We use the pre-trained model j-hartmann-roberta<sup>6</sup> as the post emotion encoder, which is trained on various emotion classification datasets and achieves good performances.

The results are shown in Table 5. We notice that adding sentiment information improves the performance of the model regardless of the fusion strategy used. In particular, our proposed sentiment-guided Transformer performs better than others since the fusion structure can better integrate semantic and sentiment hidden information, while only using label information is not sufficient. We observe that the emotion-guided Transformer also improves the performance on two datasets, by only 0.1% FScore on DsD, and the improvements are both smaller than the sentiment-guided Transformer. A possible reason is that emotion classification is a more difficult task than sentiment classification, resulting in lower accuracy of sentiment classification (66% F1-score on emotion task than 80%-90% F1-score on sentiment task), which will bring the noise to the depression detection.

Model	DsD			DepSign		
	GP	GR	FS	GP	GR	FS
SA(m)	94.1	80.4	86.7	84.4	90.3	87.2
SA(m)+s	92.5	<b>82.1</b>	86.9	84.3	<b>90.4</b>	87.2
SA(m)+SA(s)	92.7	81.8	86.9	86.3	90.0	88.1
SGT(m+s)	<b>93.2</b>	81.8	<b>87.1</b>	<b>86.8</b>	90.3	<b>88.5</b>
EGT(m+e)	92.7	81.7	86.8	86.5	90.1	88.3

Table 5: The results of sentiment fusion ablation study. We highlight top-1 values in bold.

##### 4.4.3 Contrastive learning ablation

We also conducted a contrastive learning ablation experiment. ‘Base’ means our proposed model

<sup>6</sup><https://huggingface.co/j-hartmann/emotion-english-distilroberta-base>



without a contrastive learning module. ‘+CL’ denotes the base model with supervised contrastive learning. ‘+SCL’ is the base model with our supervised severity-aware contrastive learning. Table 6 shows the performance comparison of the case study. The results demonstrate that adding supervised contrastive learning improves detection performance. In detail, on both datasets, the FS-core of using severity-aware contrastive learning is 0.3% higher than that of traditional supervised learning, which proves the effectiveness of our designed severity-aware contrastive learning.

Model	DsD			DepSign		
	GP	GR	FS	GP	GR	FS
Base	<b>93.2</b>	81.8	87.1	86.8	90.3	88.5
+CL	92.4	82.5	87.2	86.8	<b>90.9</b>	88.8
+SCL	93.1	<b>82.7</b>	<b>87.5</b>	<b>88.1</b>	90.2	<b>89.1</b>

Table 6: The results of contrastive learning ablation study. We highlight top-1 values in bold.

#### 4.5 Case study

To validate the effectiveness of our model, we provide three cases from the test results of MentalRoBERTa and our model in Appendix. A for further investigation. In the first case, MentalRoBERTa fails to detect the depression severity, and classifies it as *mild*. Meanwhile, our model successfully predicts the label due to the use of sentiment information (positive). In the second case, both models are able to make accurate predictions, indicating semantic information’s importance, and sentiment information does not bring additional errors. In the third case, although our model fails to detect the severity, the predicted label (*moderate*) is closer to the true label (*severe*) than the prediction of MentalRoBERTa (*mild*), which also reveals that the negative sentiment in this post is helpful for predicting depression severity.

#### 4.6 Ethical considerations

Although the datasets are publicly available, we also follow strict ethical protocols (Benton et al., 2017) of sharing datasets. To prevent misuse and protect user privacy, all examples shown in this paper have been obfuscated and paraphrased according to the moderate disguise scheme (Bruckman, 2002).

### 5 Conclusion

In this paper, we propose a depression severity detection model based on sentiment-guided

Transformer and severity-aware contrastive learning. The sentiment-guided Transformer can efficiently fuse semantic and sentiment information from user’s posts. Then we use a severity-aware contrastive learning framework, which can fully leverage label information for capturing severity-specific features and helping the model distinguish closely-labeled categories. The experimental results show that our model performs better than recent models and achieves superior performance on two publicly available datasets. Further analysis determines the effectiveness of all proposed modules.

### Limitations

There are also several potential limitations. First, the severity of the identification of depression is subjective, which inevitably leads to annotation bias, and we cannot verify the actual diagnosis. Thus, our model is not intended to be used as a psychiatric diagnosis tool but an estimate of depression level for users, which can be utilized to direct intervention and treatment for non-clinical use. Second, the datasets used are collected from a single social media platform (Reddit) and are imbalanced. To train a more robust and effective model, future works should collect more precise data from other social media and increase collaboration with clinicians to ensure the quality of the data. Despite these limitations, we believe that our work will facilitate depression severity detection. In future work, we would like to explore other fusion strategies to better integrate various information and assess the performance of ChatGPT (Yang et al., 2023a) on the depression severity detection task.

### Acknowledgements

This work is supported in part by the Alan Turing Institute and funds from MRC MR/R022461/1 and the project JPNP20006 from New Energy and Industrial Technology Development Organization (NEDO).

### References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Nirmal Varghese Babu and E Kanaga. 2022. Sentiment analysis in social media data for depression detection

- using artificial intelligence: A review. *SN Computer Science*, 3(1):1–20.
- Aaron T Beck, Calvin H Ward, Mock Mendelson, Jeremiah Mock, and John Erbaugh. 1961. An inventory for measuring depression. *Archives of general psychiatry*, 4(6):561–571.
- Adrian Benton, Glen Coppersmith, and Mark Dredze. 2017. Ethical research protocols for social media health research. In *Proceedings of the first ACL workshop on ethics in natural language processing*, pages 94–102.
- Amy Bruckman. 2002. Studying the amateur artist: A perspective on disguising data collected in human subjects research on the internet. *Ethics and Information Technology*, 4(3):217–231.
- Fidel Cacheda, Diego Fernandez, Francisco J Novoa, Victor Carneiro, et al. 2019. Early detection of depression: social network analysis and random forest techniques. *Journal of medical Internet research*, 21(6):e12554.
- Erik Cambria, Yang Li, Frank Z Xing, Soujanya Poria, and Kenneth Kwok. 2020. Senticnet 6: Ensemble application of symbolic and subsymbolic ai for sentiment analysis. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pages 105–114.
- Stevie Chancellor and Munmun De Choudhury. 2020. Methods in predictive techniques for mental health status on social media: a critical review. *NPJ digital medicine*, 3(1):1–11.
- Raymond Chiong, Gregorius Satia Budhi, Sandeep Dhakal, and Fabian Chiong. 2021. A textual-based featuring approach for depression detection using machine learning classifiers and social media texts. *Computers in Biology and Medicine*, 135:104499.
- Cristóbal Colón-Ruiz and Isabel Segura-Bedmar. 2020. Comparing deep learning architectures for sentiment analysis on drug reviews. *Journal of Biomedical Informatics*, 110:103539.
- Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz. 2013. Predicting depression via social media. In *Seventh international AAAI conference on weblogs and social media*.
- John Edward Desrochers and Gail M Houck. 2014. *Depression in children and adolescents: Guidelines for school practice*. National Association of School Nurses Silver Spring, MD.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Anca Dinu and Andreea-Codrina Moldovan. 2021. Automatic detection and classification of mental illnesses from general social media texts. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 358–366.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910.
- Manas Gaur, Amanuel Alambo, Joy Prakash Sain, Ugur Kursuncu, Krishnaprasad Thirunarayan, Ramakanth Kavuluru, Amit Sheth, Randy Welton, and Jyotishman Pathak. 2019. Knowledge-aware assessment of severity of suicide risk for early intervention. In *The world wide web conference*, pages 514–525.
- Shreya Ghosh and Tarique Anwar. 2021. Depression intensity estimation via social media: a deep learning approach. *IEEE Transactions on Computational Social Systems*, 8(6):1465–1474.
- Zogan Hamad, Razzak Imran, Mohammad Shoaib Jameel, and Xu Guandong. 2021. Depressionnet: A novel summarization boosted deep framework for depression detection on social media. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 133–142. ACM (Association for Computing Machinery).
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Steven D Hollon, Michael E Thase, and John C Markowitz. 2002. Treatment and prevention of depression. *Psychological Science in the public interest*, 3(2):39–77.
- Md Rafiqul Islam, Muhammad Ashad Kabir, Ashir Ahmed, Abu Raihan M Kamal, Hua Wang, and Anwaar Ulhaq. 2018. Depression detection from social network data using machine learning techniques. *Health information science and systems*, 6(1):1–12.
- Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria. 2022. Mentalbert: Publicly available pretrained language models for mental healthcare. In *Proceedings of the Language Resources and Evaluation Conference*, pages 7184–7190.
- Shaoxiong Ji, Tianlin Zhang, Kailai Yang, Sophia Ananiadou, Erik Cambria, and Jörg Tiedemann. 2023. Domain-specific continued pretraining of language models for capturing long context in mental health. *arXiv preprint arXiv:2304.10447*.

- S Kayalvizhi, Thenmozhi Durairaj, Bharathi Raja Chakravarthi, et al. 2022. Findings of the shared task on detecting signs of depression from social media. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 331–338.
- S Kayalvizhi and D Thenmozhi. 2022. Data set creation and empirical analysis for detecting signs of depression from social media postings. *arXiv preprint arXiv:2202.03047*.
- Pei Ke, Haozhe Ji, Siyang Liu, Xiaoyan Zhu, and Minlie Huang. 2020. Sentilare: Sentiment-aware language representation learning with linguistic knowledge. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6975–6988.
- David E Losada, Fabio Crestani, and Javier Parapar. 2019. Overview of erisk 2019 early risk prediction on the internet. In *International Conference of the Cross-Language Evaluation Forum for European Languages*, pages 340–357. Springer.
- Rodrigo Martínez-Castaño, Juan C Pichel, David E Losada, and Fabio Crestani. 2018. A micromodule approach for building real-time systems with python-based models: Application to early risk detection of depression on social media. In *European Conference on Information Retrieval*, pages 801–805. Springer.
- Marc L Molendijk, Lotte Bamelis, Arnold AP van Emmerik, Arnoud Arntz, Rimke Haringsma, and Philip Spinhoven. 2010. Word use of outpatients with a personality disorder and concurrent or previous major depressive disorder. *Behaviour Research and Therapy*, 48(1):44–51.
- Danielle L Mowery, Craig Bryan, and Mike Conway. 2015. Towards developing an annotation scheme for depressive disorder symptoms: A preliminary study using twitter data. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 89–98.
- Usman Naseem, Adam G Dunn, Jinman Kim, and Matloob Khushi. 2022. Early identification of depression severity levels on reddit using ordinal classification. In *Proceedings of the ACM Web Conference 2022*, pages 2563–2572.
- Mark Olfson, Carlos Blanco, and Steven C Marcus. 2016. Treatment of adult depression in the united states. *JAMA internal medicine*, 176(10):1482–1491.
- Ahmed Hussein Orabi, Prasadith Buddhitha, Mahmoud Hussein Orabi, and Diana Inkpen. 2018. Deep learning for depression detection of twitter users. In *Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic*, pages 88–97.
- David Owen, Jose Camacho-Collados, and Luis Espinosa Anke. 2020. Towards preemptive detection of depression and anxiety in twitter. In *Proceedings of the Fifth Social Media Mining for Health Applications Workshop & Shared Task*, pages 82–89.
- Esteban A Ríssola, David E Losada, and Fabio Crestani. 2021. A survey of computational methods for online mental state assessment on social media. *ACM Transactions on Computing for Healthcare*, 2(2):1–31.
- Stephanie Rude, Eva-Maria Gortner, and James Pennebaker. 2004. Language use of depressed and depression-vulnerable college students. *Cognition & Emotion*, 18(8):1121–1133.
- Ivan Sekulic and Michael Strube. 2019. Adapting deep learning methods for mental health prediction on social media. *W-NUT 2019*, 14(162.2):322.
- Ruba Skaik and Diana Inkpen. 2020. Using social media for mental health surveillance: a review. *ACM Computing Surveys (CSUR)*, 53(6):1–31.
- Varsha Suresh and Desmond Ong. 2021. Not all negatives are equal: Label-aware contrastive loss for fine-grained text classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4381–4394.
- Michael M Tadesse, Hongfei Lin, Bo Xu, and Liang Yang. 2019. Detection of depression-related posts in reddit social media forum. *IEEE Access*, 7:44883–44893.
- Alina Trifan, Dave Semeraro, Justin Drake, Radek Bukowski, and José Luís Oliveira. 2020. Social media mining for postpartum depression prediction. In *Digital Personalized Health and Medicine*, pages 1391–1392. IOS Press.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wei-Yao Wang, Yu-Chien Tang, Wei-Wei Du, and Wen-Chih Peng. 2022. Nycu\_twd@ It-edi-acl2022: Ensemble models with vader and contrastive learning for detecting signs of depression from social media. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 136–139.
- Min Yen Wu, Chih-Ya Shen, En Tzu Wang, and Arbee LP Chen. 2020. A deep architecture for depression detection using posting, behavior, and living environment data. *Journal of Intelligent Information Systems*, 54(2):225–244.
- Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. Consert: A contrastive framework for self-supervised sentence representation transfer. In *Proceedings of the 59th Annual*

*Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5065–5075.

Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, and Sophia Ananiadou. 2023a. On the evaluations of chatgpt and emotion-enhanced prompting for mental health analysis. *arXiv preprint arXiv:2304.03347*.

Kailai Yang, Tianlin Zhang, and Sophia Ananiadou. 2022. A mental state knowledge-aware and contrastive network for early stress and depression detection on social media. *Information Processing & Management*, 59(4):102961.

Kailai Yang, Tianlin Zhang, and Sophia Ananiadou. 2023b. Disentangled variational autoencoder for emotion recognition in conversations. *arXiv preprint arXiv:2305.14071*.

Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pages 1480–1489.

Da Yin, Tao Meng, and Kai-Wei Chang. 2020. Sentibert: A transferable transformer-based architecture for compositional sentiment semantics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3695–3706.

Lin Yue, Weitong Chen, Xue Li, Wanli Zuo, and Minghao Yin. 2019. A survey of sentiment analysis in social media. *Knowledge and Information Systems*, 60(2):617–663.

Tianlin Zhang, Annika M Schoene, Shaoxiong Ji, and Sophia Ananiadou. 2022. Natural language processing applied to mental illness detection: a narrative review. *NPJ digital medicine*, 5(1):1–13.

Tianlin Zhang, Kailai Yang, Shaoxiong Ji, and Sophia Ananiadou. 2023. Emotion fusion for mental illness detection from social media: A survey. *Information Fusion*, 92:231–246.

Mark Zimmerman, Jennifer H Martinez, Diane Young, Iwona Chelminski, and Kristy Dalrymple. 2013. Severity classification on the hamilton depression rating scale. *Journal of affective disorders*, 150(2):384–388.

## A Case study

The three cases from the test results of Mental-ROBERTa and our mode are shown in Table 7.

Post	Golden	MentalRoBERTa	our model
Now I know everyone is going to maybe think it's from the suboxone but it or at least used to be ** of the things that helped my anxiety. If I admitted m**f I am afraid they'd admit me and m**e me go cold turkey. I am unemployed ** looking for employment. I live with my grandparents.	minimum	mild	minimum
I'm worried I have a blood clot or something that gets aggravated when I wear them. I ** want to be okay and to have a good t**e on this trip, but I've been so out of it, ** I'm at my wit's end. Right now I'm lying down and I feel blood rushing all through my h**, and bulging of blood vessels a**d my nose. I'm extremely sleep deprived and woozy but I'm scared to go to sleep and am in pain. I'm so scared.	moderate	moderate	moderate
Did I mention my parents are religious? I don't know if ** is normal for religious people to treat. Whenever I tell them I'm terrified of being homeless ** tell me I'm a "acting like a baby" and "get over it" my parents parents did not treat them this way. They're basically mad because they(I guess ** or 40 years ago were different when they were my age? Because ** said they both lived on their own at ** and that they find it creepy I'm ** and they ** it creepy being around me).	severe	mild	moderate

Table 7: We provide three cases from the datasets we use. The posts, the golden labels and the predicted labels of MentalRoBERTa and our model are shown.