

# What Symptoms and How Long? An Interpretable AI Approach for Depression Detection in Social Media

Junwei Kuang  
Beijing Institute of Technology  
[kuangjw@bit.edu.cn](mailto:kuangjw@bit.edu.cn)

Jiaheng Xie  
University of Delaware  
[jxie@udel.edu](mailto:jxie@udel.edu)

Zhijun Yan  
Beijing Institute of Technology  
[yanzhijun@bit.edu.cn](mailto:yanzhijun@bit.edu.cn)

## Abstract

*Depression is the most prevalent and serious mental illness, which induces grave financial and societal ramifications. Depression detection is key for early intervention to mitigate those consequences. Such a high-stake decision inherently necessitates interpretability. Although a few depression detection studies attempt to explain the decision, these explanations misalign with the clinical depression diagnosis criterion that is based on depressive symptoms. To fill this gap, we develop a novel Multi-Scale Temporal Prototype Network (MSTPNet). MSTPNet innovatively detects and interprets depressive symptoms as well as how long they last. Extensive empirical analyses show that MSTPNet outperforms state-of-the-art depression detection methods. This result also reveals new symptoms that are unnoted in the survey approach. We further conduct a user study to demonstrate its superiority over the benchmarks in interpretability. This study contributes to IS literature with a novel interpretable deep learning model for depression detection in social media.*

**Keywords:** social media mining, depression detection, prototype learning, multi-scale, interpretability

## 1. Introduction

Depression is one of the most prevalent mental disorders (WHO, 2017), and brought significant societal and financial consequences. Approximately 280 million people suffer from depression worldwide, accounting for 3.8% of the world's population (Murray, 2022). More than one million people worldwide commit suicide due to depression annually, on par with the number of deaths from cancer (WHO, 2017). The economic toll linked to depression increased from \$236.6 billion to \$326.2 billion during 2010-2018 in the United States (Greenberg et al., 2021) and is projected to be the world's leading economic burden by 2030 (WHO, 2017). While many effective depression

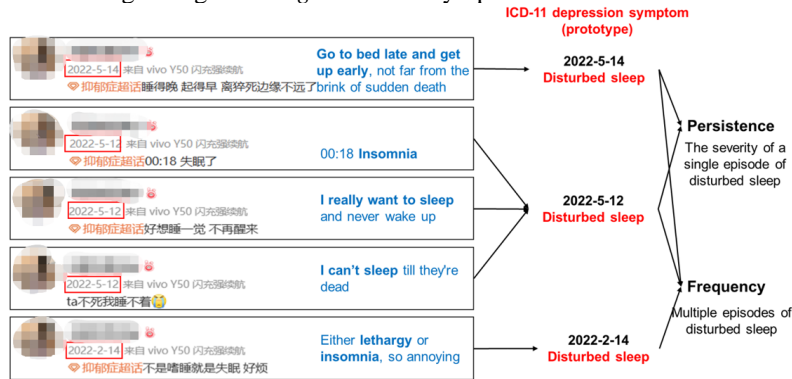
treatments exist, more than 70% of patients do not seek treatments at the early stages of depression due to stigmatization around depression (Schomerus et al., 2009; Shen et al., 2017; Wang et al., 2007). To address this gap and avoid preventable ramifications, depression detection is key to early intervention and providing social support and treatment (Picardi et al., 2016).

While surveys remain the primary source of depression detection (Kroenke et al., 2001), social media unleashes the unprecedented potential to expand its reach. Moreover, depressed patients are more willing to communicate on social media compared to offline (Naslund et al., 2016). Many scholars develop depression detection models on social media for early intervention (Chau et al., 2020; Liu et al., 2022). Although achieving satisfying performance, most of these studies rely on black-box methods, which results in limited applicability and potential risk in high-stake scenarios such as healthcare decision-making (Chiong et al., 2021; Zogan et al., 2022). To overcome the non-interpretable dilemma, a few depression detection studies attempt to explain why users are classified as depressed based on the importance score or attention weights of interpretable inputs such as words in a post (Cheng & Chen, 2022). However, existing interpretable models depart from clinical depression diagnosis criteria and receive compromised trust from end users. To tackle their limitations, there has recently been a rising interest in utilizing symptoms for interpreting depression detection. Pioneering studies have shown the potential benefits of improving accuracy, generalizability, and interpretability (Nguyen et al., 2022; Zhang, Chen, Wu, et al., 2022). Therefore, our research objective is to develop an *interpretable depression detection model in social media based on symptom-based depression diagnostic criteria*.

The symptom-based interpretable methods for depression detection can be categorized into dictionary-based, similarity-based, and classification-based (Shen et al., 2017). The core of these methods is to identify depressive symptoms from user-generated posts on social media. However, these methods still face three

limitations. First, prior methods only identify pre-defined symptoms. However, depressive symptoms may evolve over time. Second, previous methods rely on domain-specific knowledge, which require significant labor costs and suffer from poor generalizability. Third, extant methods focus on *what symptoms* users present while neglecting *how long* these

symptoms last (Kroenke et al., 2001). Fortunately, user-generated posts on social media can reveal such “how long” aspects of depressive symptoms. As shown in Figure 1, the user reported the disturbed sleep symptom numerous times, ranging from Feb 14 to May 14. Certain periods (e.g., May 12 to May 14) show denser symptom mentions than others.



**Figure 1. Typical Posts Disclosing Depressive Symptoms in Social Media**

The abovementioned limitations motivate us to develop a novel interpretable depression detection method. Following the computational design science paradigm and prior IS research on health analytics (Yu et al., 2023), we propose and rigorously evaluate a novel interpretable model, Multi-Scale Temporal Prototype Network (MSTPNet). MSTPNet is built upon an emergent stream of case-based interpretable models, prototype learning (Ming et al., 2019), which interprets the prediction for new inputs by comparing them with a few learned prototypes. In this study, typical posts disclosing depressive symptoms can be recognized as prototypes. To consider how long the symptoms last, MSTPNet modifies standard prototype learning methods by devising two novel layers: a temporal segmentation layer that eliminates the negative effects of irrelevant and redundant posts on symptom identification to facilitate period-level analysis (i.e., “What symptoms did the user suffer in a period?”), and a multi-scale temporal prototype layer that captures the temporal distribution of symptoms. In practice, our method can be implemented in social media to detect depressed patients and interpret their temporal symptoms. When implementing intervention, platform

managers need to combine human intelligence to judge rather than rely entirely on artificial intelligence.

## 2. Literature Review

Social media-based depression detection is broadly classified into traditional machine learning, black-box deep learning, and interpretable deep learning. The traditional machine learning-based depression detection model mostly relies on effective input features (Li et al., 2019). Table 1 summarizes the traditional machine learning-based method for depression detection.

However, these studies have shown unsatisfactory predictive power mainly because hand-crafted features and traditional machine learning models are not complex enough to capture high-level interactions between features. Black-Box deep learning methods have demonstrated significantly higher predictive power in depression detection (Malhotra & Jindal, 2022). These improvements have benefited from the development of embedding techniques and the utilization of various neural network architectures. Table 2 summarizes recent black-box deep learning-based depression detection methods in social media.

**Table 1. Traditional Machine Learning Methods in Depression Detection**

Reference	Dataset	Sample (depression/non)	Input features	Methods
Choudhury et al. (2013)	Twitter	476 (171/305)	Emotion, Depression language, Language style	SVM
Tsugawa et al. (2015)	Twitter	209 (81/128)	Emotions, Linguistic style, Topic, social Network	LDA, SVM
Chen et al. (2018)	Twitter	1200 (600/600)	Emotion swings, LIWC	SVM, RF
Chau et al. (2020)	Blog	804 (274/530)	N-gram, Lexicon based, LIWC	SVM, Rule-based, GA
Chiong et al. (2021)	Twitter	2804 (1402/1402)	N-gram	SVM, DT, NB, KNN,

**Table 2. Black-Box Deep Learning Methods in Depression Detection**

Reference	Dataset	Sample (depression/non)	Input features	Methods
Orabi et al. (2018)	Twitter	899 (327/572)	Text	CNN/RNN
Chiu et al. (2021)	Instagram	520 (260/260)	Text, Image, Posting time	LSTM with temporal weighting
Ghosh and Anwar (2021)	Twitter	6562 (1402/5160)	Text	LSTM
Zogan et al. (2022)	Twitter	4800 (2500/2300)	Text, Image	HAN
Kour and Gupta (2022)	Twitter	1681 (941/740)	Text	CNN + Bi-LSTM

**Table 3. Interpretable Deep Learning Methods in Depression Detection**

Reference	Type	Method	Usage	Explanations
Adarsh et al. (2023)	Approximation	LIME	Post-hoc	Important raw inputs
Cheng and Chen (2022)	Attention	Attention	Intrinsic	Important raw inputs
Zogan et al. (2022)	Attention	HAN	Intrinsic	Important raw inputs
Shen et al. (2017)	Symptom	Dictionary-based	Post-hoc	Predicted symptoms
Zhang, Chen, Wu, et al. (2022)	Symptom	Classification-based	Post-hoc	Predicted symptoms
Zhang, Chen, Mengyue Wu, et al. (2022)	Symptom	Similarity-based	Post-hoc	Predicted symptoms
Our study	Symptom	Similarity-based	Intrinsic	More symptoms, and how long

Despite their satisfying performance, their lack of interpretability limits their applicability in high-stake decision-making scenarios (Rudin, 2019). Interpretable deep learning methods refer to deep learning methods that provide a certain explanation (Li et al., 2022). Table 3 summarizes and contrasts recent interpretable deep learning-based methods and our study in social media-based depression detection.

Symptom-based interpretable deep learning methods align with clinical depression criteria, but still face two limitations. First, they generally require high labor costs and only identify pre-defined symptoms, neglecting new symptoms unnoted in offline depression screening questionnaires in the online setting. Second, symptom-based interpretable methods focus only on the type of depressive symptoms users suffer, neglecting how long these symptoms last, which is equally critical for a clinical depression diagnosis. These limitations motivate us to develop a novel interpretable depression detection method that is capable of discovering

depressive symptoms in a data-driven manner while capturing how long these symptoms last.

We resort to an emergent interpretable model paradigm that is closely related to our task: prototype learning. Prototype learning methods learn prototypes that have clear semantic meanings, and intrinsic explanations are generated based on the comparison between input and each prototype (Nauta et al., 2021). Chen et al. (2019) originally propose ProtoPNet, which explains the contribution of prototypical parts of the predicted image by comparing the learned prototypes. Multiple prototype learning variants have also been proposed for various tasks. Typical posts disclosing depressive symptoms can be recognized as prototypes in our study. By calculating how similar a user’s posts are to these prototypes, this user’s depressive symptoms can be inferred, which serves as a natural interpretation mechanism. Table 4 contrasts major prototype learning methods with our method.

**Table 4. Existing Prototype Learning Methods vs. Our Method**

Reference	Method	Novelty	Input	TD*
Chen et al. (2019)	ProtoPNet	Prototype for image classification	An image	No
Hase et al. (2019)	HPNet	Hierarchical prototype	An image	No
Ming et al. (2019)	ProSeNet	Prototype for text classification	A piece of text	No
Zhang et al. (2020)	TapNet	Attentional prototype	A time series of ECG	No
Nauta et al. (2021)	ProtoTree	Prototype and decision tree	An image	No
Trinh et al. (2021)	DPNet	Dynamic prototype	A clip of a video	No
Deng et al. (2022)	K-HPN	Pairwise prototype	A piece of text	No
<b>Our study</b>	<b>MSTPNet</b>	<b>Multi-scale temporal prototypes</b>	<b>A sequence of text documents</b>	<b>Yes</b>

\* TD stands for “Temporal Distribution”, which includes frequency and persistence of prototypes at period level

The majority of prototype learning methods focus on static subjects, such as an image and a piece of text. When applied to our study, these methods only consider whether depressive symptoms appear, neglecting how long each symptom last (Chen et al.,

2019). While a few prototype learning methods process dynamic subjects such as video, these methods focus on directly identifying complex prototypes with temporal properties, rather than analyzing the temporal distribution of prototypes after identifying

them. Our method aims to incorporate the temporal distribution of symptoms into the prototype learning method to effectively capture how long depressive symptoms last to improve the predictive power and interpretability.

### 3. The MSTPNet Approach

Figure 2 shows the architecture of MSTPNet, which features four building blocks. The feature learning layer aims to represent each post as an embedding vector with a fixed length and rich semantic meaning. Different from analyzing each post independently, our proposed temporal segmentation layer assigns posts into different periods based on the

semantic similarity and time interval between posts, which facilitate period-level analysis. Instead of learning complex dynamic prototypes (e.g., “long-term disturbed sleep”) directly, our proposed multi-scale temporal prototype layer breaks the task down into two parts. We first infer depressive symptoms (e.g., “disturbed sleep”) in each period by comparing posts with learned prototypes, and then explicitly measure the frequency (e.g., the proportion of periods where disturbed sleep appears) and persistence (e.g., the number of continuous periods where disturbed sleep all appears) of each symptom. Based on the above interpretable temporal measurement of each symptom, the classification layer classifies a user into depression or non-depression categories.

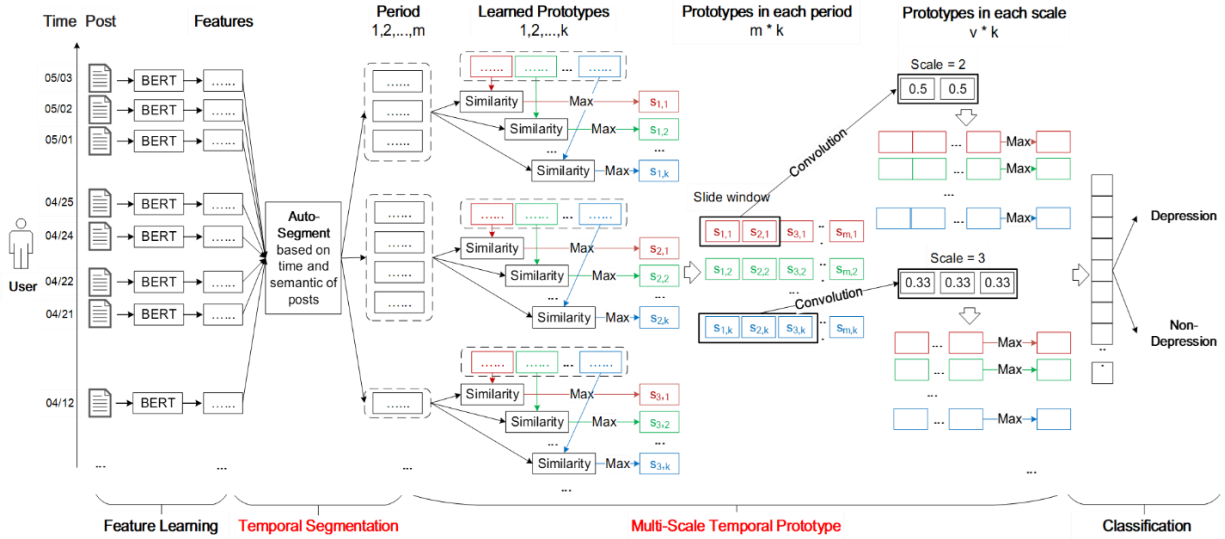


Figure 2. MSTPNet Architecture

To learn an effective representation for each post, we deploy a feature learning layer using the cutting-edge pre-trained language model BERT (Devlin et al., 2019). Specifically, for a post  $X_i$ :

$$H_i = BERT(X_i) \quad (1)$$

Our temporal segmentation layer builds upon a bottom-up hierarchical clustering algorithm (Shetty & Singh, 2021) to segment the social media posts  $u = (H_1, t_1; H_2, t_2; \dots; H_n, t_n)$  into  $m$  periods  $u = (C_1, C_2, \dots, C_m)$ ,  $C_i = (H_{i,1}, t_{i,1}; H_{i,2}, t_{i,2}; \dots; H_{i,l}, t_{i,l})$ . The key to segmentation methods is the distance measurement between different posts. We propose a new measurement that combines both semantic similarity and the time interval between posts in Formula (2), (3), and (4).

$$sim_{sem}^{i,j} = \frac{H_i \cdot H_j}{\|H_i\| \cdot \|H_j\|} \quad (2)$$

$$sim_{time}^{i,j} = \exp\left(-\frac{|t_i - t_j|}{w_d}\right) \quad (3)$$

$$sim^{i,j} = w_a * sim_{time}^{i,j} + (1 - w_a) * sim_{sem}^{i,j} \quad (4)$$

In each iteration, the method calculates the similarity between each pair of segments, and then merges the most similar pair into a new segment, until the time distance between the two segments exceeds the pre-defined length  $h$  of periods. The remaining clusters  $(C_1, C_2, \dots, C_m)$  are the segmentation results, where  $C_i$  is the  $i$ -th segment of the focal user, and  $X_{i,j}$  is the  $j$ -th post in the segment  $C_i$ .

$$s_{m,k} = \max_{j=1,2,\dots,l} \exp\left(-\|H_{m,j} - p_k\|^2\right) \quad (5)$$

Beyond one period or all periods, we focus on a specified number of consecutive periods to measure the frequency and persistence of depressive symptoms. The number of consecutive periods is called “scale” in this study. Different scales enable analysis at different granularity and can capture comprehensive clues to detect depression. Therefore, we employ multiple scales (i.e., multi-scale) with different sizes for period-level analysis, which is conceptually similar to the filters with different sizes

in CNN to analyze image data. The difference is that our “filters” are not learnable parameters but are explicitly set to get the average value over continuous periods, which is easy for humans to understand. Specifically, let  $W$  be the set of scales used in our model, and  $w_j$  denotes the size of the  $j$ -th scale. We calculate the existence (i.e., average similarity) of depressive symptoms in each pair with a window length of the scale, and then take the highest value as the existence (i.e.,  $g_{j,k}$ ) of the depressive symptom  $k$  on the scale  $w_j$ , as shown in Formula (6).

$$g_{j,k} = \max_{m=1,2,\dots,M-w_j+1} \frac{1}{w_j} \sum_{m=1}^{m+w_j-1} s_{m,k} \quad (6)$$

We let  $G = (g_{1,1}, g_{1,2}, \dots, g_{1,K}, \dots, g_{j,K}, \dots, g_{J,K})$ , where  $J$  is the number of scales. The classification layer computes the probability of depression given all  $g_{j,k}$  ( $G$ ) of a focal user in formulas (7) and (8):

$$\hat{y}_i = \frac{\exp(z_i)}{\sum_{s=0}^1 \exp(z_s)} \quad (7)$$

$$Z = QG \quad (8)$$

Following Ming et al. (2019), the loss function of MSTPNet to be minimized is defined based on the

binary cross-entropy (CE) loss with four additional regularization terms. Specifically,

$$Loss = CE + \lambda_c R_c + \lambda_e R_e + \lambda_d R_d + \lambda_{l_1} \|Q\|_1 \quad (9)$$

$$CE = \sum_{(x,y) \in D} y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \quad (10)$$

where  $\lambda_c$ ,  $\lambda_e$ ,  $\lambda_d$  and  $\lambda_{l_1}$  are hyperparameters that determine the weight of the regularizations.

## 4. Empirical Analysis

We use the WU3D, an annotated dataset regarding depression detection in a Chinese social media platform (Wang et al., 2022). It contains chronological sequences of posts from 10,325 depressed users and a random control group of 22,245 users. We set imbalance ratios as 1:8, which approximates the ratio of adults with depression risk in China (Fu et al., 2023). We split this dataset into 60% for training, 20% for validation, and 20% for test. We set  $K=70$ ,  $h=15$ ,  $w_a=0.4$ , and the scales contain 1, 2, 3, 5, 8, 12, 16 and 20. The evaluation results are reported in Tables 5. MSTPNet outperforms benchmark models in F1 and accuracy and outperforms interpretable deep learning in all metrics.

**Table 5. SOTA Methods vs. Our Method**

Category	Models	F1	Precision	Recall	Accuracy
Traditional	Yang et al. (2020)	0.412***	0.887***	0.268***	0.918***
Machine	Chen et al. (2018)	0.508***	0.891***	0.355***	0.926***
Learning	Chiong et al. (2021)	0.380***	0.363***	0.399***	0.861***
Methods	Chau et al. (2020)	0.706***	0.604***	0.850	0.924***
Black-Box	Orabi et al. (2018)	0.828***	0.878***	0.785	0.965**
Deep	Chiu et al. (2021)	0.822**	0.936*	0.732**	0.966**
Learning	Ghosh and Anwar (2021)	0.820**	0.921**	0.741***	0.965***
Methods	Naseem et al. (2022)	0.826**	0.963	0.723**	0.967*
Interpretable	Cheng and Chen (2022)	0.806***	0.910**	0.723**	0.963**
Deep	Zogan et al. (2022)	0.795***	0.840***	0.754*	0.958**
Learning	Ming et al. (2019)	0.816**	0.929*	0.729***	0.965**
Methods	Chen et al. (2019)	0.735***	0.876***	0.633***	0.951**
	Trinh et al. (2021)	0.675***	0.774***	0.598***	0.938***
	<b>MSTPNet</b>	<b>0.851</b>	<b>0.957</b>	<b>0.766</b>	<b>0.971</b>

Note: \*p < 0.05; \*\*p < 0.01; \*\*\*p < 0.001

**Table 6. Ablation Studies**

Model	F1	Precision	Recall	Accuracy
MSTPNet (Ours)	0.851	0.957	0.766	0.971
MSTPNet removing clustering layer	0.801***	0.923***	0.702***	0.962***
MSTPNet removing MS using Max	0.760***	0.868***	0.676***	0.954***
MSTPNet removing MS using Mean	0.690***	0.846***	0.583***	0.944***

Note: \*p < 0.05; \*\*p < 0.01; \*\*\*p < 0.001

We further perform ablation studies to show their effectiveness as shown in Table 6. We remove the temporal segmentation layer to validate the effectiveness of period-level analysis. We also replace the multi-scale temporal prototype (MS) layer with a common prototype learning layer. We test two options: the maximum or mean existence strength of prototype. Our MSTPNet provides a level of

interpretability that is absent in other interpretable deep models. Figure 3 provides a visual comparison of different types of interpretation. The approximation-based explanation mainly reveals the importance score and direction of interpretable input (e.g., words) linked to the final output, as shown in Figure 3(a). At a finer level, attention-based explanation enables the end user to attend to important words within a post, and

important posts for depression detection, as shown in Figure 3(b). The above two explanations are common in many scenarios but fall short of interpreting

depression detection because they depart from the clinical depression diagnosis criterion that is based on depressive symptoms.

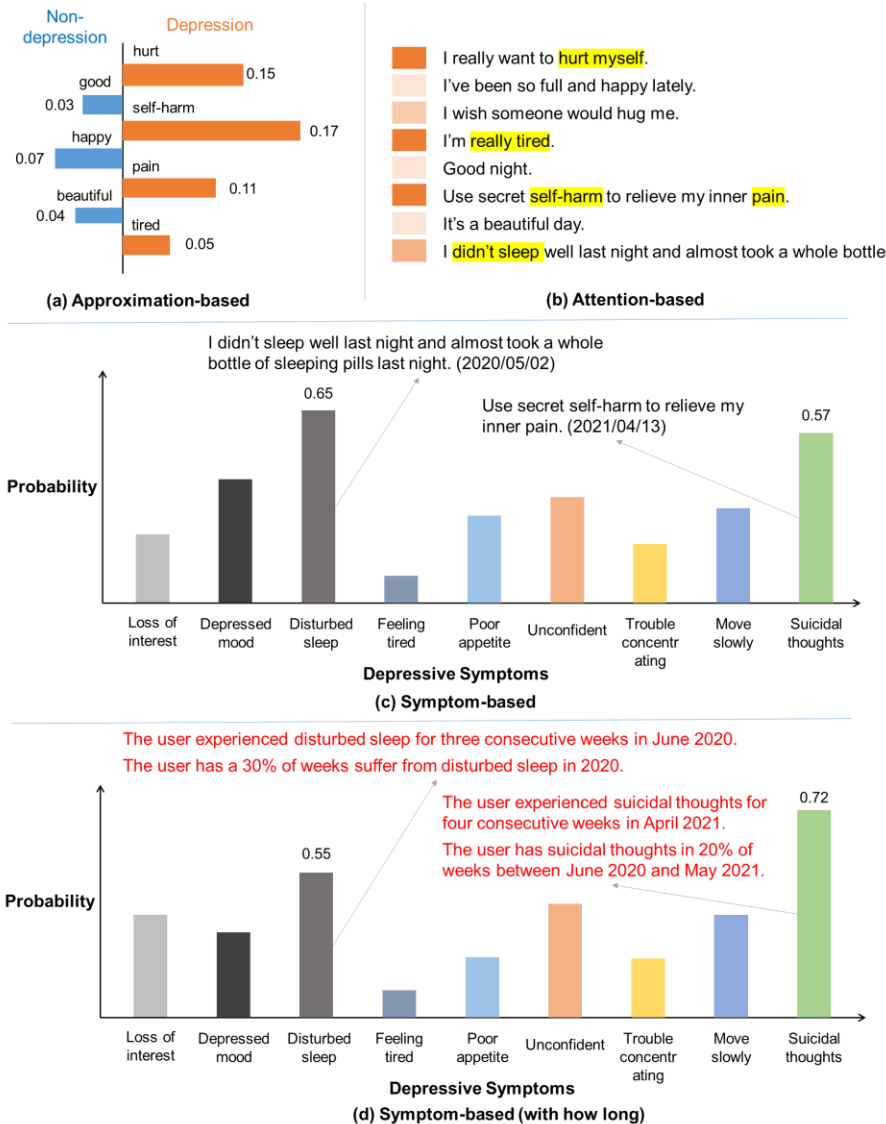


Figure 3. Visual comparison of different types of model interpretability

Figure 3(c) shows a symptom-based explanation that contains the strength of each depressive symptom, which is generally measured by the maximum similarities between user-generated posts and each symptom. The posts with the highest similarity are presented as evidence, which provides the end user to judge the credibility of the interpretation.

Our interpretation is also symptom-based, but the innovation is that we combine the duration of symptoms to more accurately measure the existence strength of each depressive symptom, as shown in Figure 3(d). The interpretation is capable of capturing what symptoms users suffer and how long these

symptoms last, which aligns with the clinical depression diagnosis criterion. Moreover, our MSTPNet is based on prototype learning and can show new depressive symptoms rather than just pre-defined.

Table 7 shows the five most salient symptom prototypes that our model learned and their responding posts as interpretations. For each prototype, we match it with the symptoms defined by ICD-11. We also found new symptoms unnoted in the previous literature, such as share depression-related content. For clarity, we highlighted the keywords related to symptoms in each post.



**Table 7. Symptom Prototypes**

Symptom Prototype	ICD-11/New
Don't conceive when you are emotionally depressed. Once you are anxious, depressed, or have a heavy mental burden.	Depressed mood
I can't seem to find anyone to talk to except on Weibo. It's so sad that I can't keep going. No one cares about whether I'm okay or not. No one cares about my condition.	fatigue or low energy
Unhappy, unhappy, too depressed, too depressed. I said I was depressed, and he said, do you want to commit suicide? I don't know if cutting my hands twice in a row counts as suicide	Loss of interests Suicidal thoughts
Standing at the window and wanting to jump off, I'm really sad, I can't do anything, I'm so tired, I'm counting the time, I'm leaving tonight. I hope there's no pain in the other world	Suicidal thoughts
Different from typical depression, patients with smiling depression do not stay at home every day and do not socialize with others, but have good social functions	<b>Share depression-related content</b>

Note: The yellow highlight is only used in the paper for demonstration purposes

To validate that capturing the temporal distribution of symptoms improves users' trust and perceived helpfulness in our model, we conduct a user study. We recruited 92 volunteers and informed the participants that they would be assigned an interpretable ML model to predict depression using social media data. We randomly selected one user sample that the model classified as depressed and show the participants how the model interpreted this classification. Since our main methodological contribution is designing the temporal distribution of the prototype, we design two randomized groups: one to present MSTPNet's interpretation, and the other to present the interpretation without considering the

temporal distribution which is in line with the state-of-the-art prototype learning methods. The only difference between the two groups is that MSTPNet considers the temporal distribution.

The first part of the user study collects seven control variables: age, education, gender, computer literacy, deep learning literacy, trust in AI, and medical literacy (Osborne et al., 2013). This user study passed randomization checks. The summary statistics of the final participants and randomization *p*-values are reported in Tables 8, 9, and 10. As shown in Table 10, such large *P*-value shows that there is no significant difference in the user background between the treatment group and control group.

**Table 8. Summary Statistics (Categorical)**

Variable	Category	Count	Variable	Category	Count
Age	18 and lower	1	Education	College freshman	9
	18 – 24	21		College junior	1
	25 – 34	32		College senior	3
	35 - 44	1		Master	21
Gender	Female	29	Doctorate	21	
	Male	26			

**Table 9. Summary Statistics (Continuous)**

Variable	Min	1st Qu.	Median	Mean	3rd Qu.	Max
Computer Literacy	1.000	2.000	3.000	2.727	3.000	4.000
Deep Learning Literacy	0.000	2.000	2.000	2.236	3.000	4.000
Trust in AI	1.000	2.000	3.000	2.527	3.000	4.000
Health Literacy	1.250	2.750	3.000	3.059	3.250	4.000

**Table 10. Randomization Checks**

	Age	Education	Gender	Computer Literacy	Deep Learning Literacy	Trust in AI	Medical Literacy
P-value	0.433	0.534	0.921	0.488	0.788	0.944	0.894

Then, we show the corresponding interpretation to each group separately, as shown in Tables 11 and 12. Subsequently, we ask the participants to rate their trust and perceived helpfulness of the given model, which are the common interpretability measurements for ML models (Xie et al., 2023). The measurement scale of trust is adopted from Chai et al. (2011), and Cronbach's Alpha is 0.901 for this scale, suggesting excellent reliability. The measurement scale of

helpfulness is adopted from Adams et al. (1992). These scales also add an attention check question ("Please just select neither agree nor disagree"). After removing the participants who failed the attention check or the manipulation check, 55 participants remain in the final analyses. Table 13 shows the mean of trust and perceived helpfulness for each of the two groups. The participants' trust in our model (mean = 2.556) is significantly higher than the baseline model

(mean = 1.659,  $p < 0.001$ ). The participants' perceived helpfulness in our model (mean = 2.913) is

significantly higher than the baseline model (mean = 2.567,  $p = 0.023$ ).

**Table 11. The Interpretation of the Baseline Model**

What Symptoms	Evidence Posts
Suicidal thoughts	The first year of high school was insulted to the point of doubting life and wanting to <b>commit suicide</b>
Disturbed sleep	I <b>didn't sleep well last night</b> and almost took a whole bottle of sleeping pills last night

**Table 12. The Interpretation of Our MSTPNet**

What symptoms	How Long		Persistence	Evidence Posts
	Frequency	Evidence Posts		
Suicidal thoughts	0.5 (times each week)	Use secret <b>self-harm</b> to relieve your inner pain. (04/13) The first year of high school was insulted to the point of doubting life and wanting to <b>commit suicide</b> . (04/21) I'm <b>leaving tonight</b> , I hope there's no pain in the other world. (06/08)	2 (weeks)	Use secret <b>self-harm</b> to relieve your inner pain. (04/13)  The first year of high school was insulted to the point of doubting life and wanting to <b>commit suicide</b> . (04/21)
Disturbed sleep	0.3 (times each week)	I <b>didn't sleep well last night</b> and almost took a whole bottle of sleeping pills last night (05/02)	1 (weeks)	I <b>didn't sleep well last night</b> and almost took a whole bottle of sleeping pills last night. (05/02)

Note: The red color is only used in the paper for demonstration purposes. The user study uses black color.

**Table 13. Interpretability Comparison Between MSPTNet and Baseline**

Interpretability measurement	MSTPNet (mean)	Baseline (mean)	MSTPNet (std)	Baseline (std)	P-value
Trust	2.556	1.659	0.749	0.677	< 0.001
Perceived Helpfulness	2.913	2.567	0.677	0.641	0.023

After the participants rate the trust and perceived helpfulness for the given model, we then show them the interpretation of the other model as a comparison. We ask them to choose a model interpretation that they trust more. 45 of 55 (81.8%) participants chose MSTPNet over the baseline. The above user study results prove that by interpreting the symptoms and their temporal distribution, MSTPNet improves users' trust and perceived helpfulness in our model, which offers empirical evidence for our contribution.

## 5. Discussion and Conclusion

We propose a novel interpretable deep learning method to detect and interpret depression based on what symptoms the user has and how long these related symptoms last. We conduct extensive evaluations to demonstrate the superior predictive power of our method over state-of-the-art benchmarks and showcase its interpretation of detection. Furthermore, through a user study, we show that our method outperforms these benchmarks in terms of interpretability. Our study contributes to the extant literature in three aspects. First, we develop computational methods to solve business and societal

problems and aims to make methodological contributions. In this regard, our proposed MSTPNet is a novel prototype learning method that processes a sequence of user-generated posts and interprets detection based on symptom and their temporal distribution. Second, our study contributes to healthcare IS research with a novel interpretable deep learning method that detects depression using social media data and interprets its detection. Third, our study contributes to medical research with new symptoms unnoted in the previous literatures.

Our study establishes a few generalized design principles: (1) A temporal segmentation module could facilitate period-level analysis and mitigate the effect of redundant and irrelevant information; (2) It's cost-effective and flexible to explicitly separate a complex task into two related simple tasks; (3) Showing the temporal distribution of prototypes could improve interpretability and boost the trust and perceived helpfulness. These design principles prescribe how to predict and interpret the hidden state of a user from a sequence of user-related data, such as social media posts, electric health records.

This study can be generalized to social media-based user-level detection or prediction, and



interpretable prediction where the active period, frequency, and persistence of the prototype are essential. For social media platforms, our method can be deployed in the platform. Platforms can recommend online resources such as articles and videos and offer resources for treatments.

## Acknowledgements

This study was supported by National Natural Science Foundation of China (No: 721110107003, 71872013).

## References

- Adams, D. A., Nelson, R. R., & Todd, P. A. (1992). Perceived usefulness, ease of use, and usage of information technology: A replication. *MIS Quarterly*, 16(2), 227-247.
- Adarsh, V., Arun Kumar, P., Lavanya, V., & Gangadharan, G. R. (2023). Fair and Explainable Depression Detection in Social Media. *Information Processing & Management*, 60(1), 103168. <https://doi.org/https://doi.org/10.1016/j.ipm.2022.103168>
- Chai, S., Das, S., & Rao, H. R. (2011). Factors Affecting Bloggers' Knowledge Sharing: An Investigation Across Gender. *Journal of Management Information Systems*, 28(3), 309-342. <https://doi.org/10.2753/mis0742-1222280309>
- Chau, M., Li, T. M. H., Wong, P. W. C., Xu, J. J., Yip, P. S. F., & Hsinchun, C. (2020). FINDING PEOPLE WITH EMOTIONAL DISTRESS IN ONLINE SOCIAL MEDIA: A DESIGN COMBINING MACHINE LEARNING AND RULE-BASED CLASSIFICATION. *MIS Quarterly*, 44(2), 933-956. <https://doi.org/10.25300/MISQ/2020/14110>
- Chen, C., Li, O., Tao, D., Barnett, A., Rudin, C., & Su, J. K. (2019). This looks like that: deep learning for interpretable image recognition. Proceedings of the 33rd International Conference on Neural Information Processing Systems,
- Chen, X., Sykora, M., Jackson, T. W., & Elayan, S. (2018). What about Mood Swings: Identifying Depression on Twitter with Temporal Measures of Emotions. Proceedings of the The Web Conference 2018,
- Cheng, J. C., & Chen, A. L. (2022). Multimodal time-aware attention networks for depression detection. *Journal of Intelligent Information Systems*, 59(2), 319-339.
- Chiong, R., Budhi, G. S., Dhakal, S., & Chiong, F. (2021). A textual-based featuring approach for depression detection using machine learning classifiers and social media texts. *Computers in Biology and Medicine*, 135(8), 104499. <https://doi.org/https://doi.org/10.1016/j.combiomed.2021.104499>
- Chiu, C. Y., Lane, H. Y., Koh, J. L., & Chen, A. L. (2021). Multimodal depression detection on instagram considering time interval of posts. *Journal of Intelligent Information Systems*, 56, 25-47.
- Choudhury, M. D., Counts, S., & Horvitz, E. (2013). Social media as a measurement tool of depression in populations. Proceedings of the 5th Annual ACM Web Science Conference, Paris, France.
- Deng, S., Zhang, N., Chen, H., Tan, C., Huang, F., Xu, C., & Chen, H. (2022). Low-resource extraction with knowledge-aware pairwise prototype learning. *Knowledge-Based Systems*, 235, 107584. <https://doi.org/https://doi.org/10.1016/j.knosys.2021.107584>
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Proceedings of NAACL-HLT,
- Fu, X., Zhang, K., Zhang, X., & Chen, Z. (2023). *Report on National Mental Health Development in China (2021-2022)*. S. S. A. Press.
- Ghosh, S., & Anwar, T. (2021). Depression intensity estimation via social media: a deep learning approach. *IEEE Transactions on Computational Social Systems*, 8(6), 1465-1474.
- Greenberg, P. E., Fournier, A.-A., Sisitsky, T., Simes, M., Berman, R., Koenigsberg, S. H., & Kessler, R. C. (2021). The Economic Burden of Adults with Major Depressive Disorder in the United States (2010 and 2018). *Pharmaco Economics*, 39(6), 653-665. <https://doi.org/10.1007/s40273-021-01019-4>
- Hase, P., Chen, C., Li, O., & Rudin, C. (2019). Interpretable Image Recognition with Hierarchical Prototypes. Proceedings of the AAAI Conference on Human Computation and Crowdsourcing,
- Kour, H., & Gupta, M. K. (2022). An hybrid deep learning approach for depression prediction from user tweets using feature-rich CNN and bi-directional LSTM. *Multimedia Tools and Applications*, 81(17), 23649-23685.
- Kroenke, K., Spitzer, R. L., & Williams, J. B. W. (2001). The PHQ-9 - Validity of a brief depression severity measure. *Journal of General Internal Medicine*, 16(9), 606-613. <https://doi.org/10.1046/j.1525-1497.2001.016009606.x>
- Li, X., Xiong, H., Li, X., Wu, X., Zhang, X., Liu, J., Bian, J., & Dou, D. (2022). Interpretable deep learning: Interpretation, interpretability, trustworthiness, and beyond. *Knowledge and Information Systems*, 64(12), 3197-3234.
- Li, X. W., Zhang, X., Zhu, J., Mao, W. D., Sun, S. T., Wang, Z. H., Xia, C., & Hu, B. (2019). Depression recognition using machine learning methods with different feature generation strategies. *Artificial Intelligence in Medicine*, 99, 101696, Article 101696. <https://doi.org/10.1016/j.artmed.2019.07.004>
- Liu, D. X., Feng, X. L., Ahmed, F., Shahid, M., & Guo, J. (2022). Detecting and Measuring Depression on Social Media Using a Machine Learning Approach: Systematic Review. *Jmir Mental Health*, 9(3), e27244, Article e27244. <https://doi.org/10.2196/27244>
- Malhotra, A., & Jindal, R. (2022). Deep learning techniques for suicide and depression detection from online social media: A scoping review. *Applied Soft Computing*, 130, 109713.

- <https://doi.org/https://doi.org/10.1016/j.asoc.2022.109713>
- Ming, Y., Xu, P., Qu, H., & Ren, L. (2019). Interpretable and Steerable Sequence Learning via Prototypes. Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.
- Murray, C. J. L. (2022). The Global Burden of Disease Study at 30 years. *Nature Medicine*, 28(10), 2019-2026. <https://doi.org/10.1038/s41591-022-01990-1>
- Naseem, U., Dunn, A. G., Kim, J., & Khushi, M. (2022). Early identification of depression severity levels on reddit using ordinal classification. Proceedings of the ACM Web Conference 2022.
- Naslund, J. A., Aschbrenner, K. A., Marsch, L. A., & Bartels, S. J. (2016). The future of mental health care: peer-to-peer support and social media. *Epidemiology and psychiatric sciences*, 25(2), 113-122.
- Nauta, M., van Bree, R., & Seifert, C. (2021). Neural prototype trees for interpretable fine-grained image recognition. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Nguyen, T., Yates, A., Ziriky, A., Desmet, B., & Cohan, A. (2022). Improving the Generalizability of Depression Detection by Leveraging Clinical Questionnaires. 60th Annual Meeting of the Association for Computational Linguistic.
- Orabi, A. H., Buddhitha, P., Orabi, M. H., & Inkpen, D. (2018). Deep Learning for Depression Detection of Twitter Users. Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic.
- Osborne, R. H., Batterham, R. W., Elsworth, G. R., Hawkins, M., & Buchbinder, R. (2013). The grounded psychometric development and initial validation of the Health Literacy Questionnaire (HLQ). *BMC Public Health*, 13(1), 658. <https://doi.org/10.1186/1471-2458-13-658>
- Picardi, A., Lega, I., Tarsitani, L., Caredda, M., Matteucci, G., Zerella, M. P., Miglio, R., Gigantesco, A., Cerbo, M., Gaddini, A., Spandonaro, F., Biondi, M., & The Set-Dep, G. (2016). A randomised controlled trial of the effectiveness of a program for early detection and treatment of depression in primary care. *Journal of Affective Disorders*, 198, 96-101. <https://doi.org/10.1016/j.jad.2016.03.025>
- Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5), 206-215.
- Schomerus, G., Matschinger, H., & Angermeyer, M. C. (2009). The stigma of psychiatric treatment and help-seeking intentions for depression. *European Archives of Psychiatry and Clinical Neuroscience*, 259(5), 298-306. <https://doi.org/10.1007/s00406-009-0870-y>
- Shen, G., Jia, J., Nie, L., Feng, F., Zhang, C., Hu, T., Chua, T.-S., & Zhu, W. (2017). Depression Detection via Harvesting Social Media: A Multimodal Dictionary Learning Solution. IJCAI.
- Shetty, P., & Singh, S. (2021). Hierarchical clustering: a survey. *International Journal of Applied Research*, 7(4), 178-181.
- Trinh, L., Tsang, M., Rambhatla, S., & Liu, Y. (2021). Interpretable and trustworthy deepfake detection via dynamic prototypes. Proceedings of the IEEE/CVF winter conference on applications of computer vision.
- Tsugawa, S., Kikuchi, Y., Kishino, F., Nakajima, K., Itoh, Y., & Ohsaki, H. (2015). Recognizing depression from twitter activity. Proceedings of the 33rd annual ACM conference on human factors in computing systems.
- Wang, P. S., Angermeyer, M., Borges, G., Bruffaerts, R., Chiu, W. T., de Girolamo, G., Fayyad, J., Gureje, O., Haro, J. M., Huang, Y. Q., Kessler, R. C., Kovess, V., Levinson, D., Nakane, Y., Browne, M. A. O., Ormel, J. H., Posada-Villa, J., Aguilar-Gaxiola, S., Alonso, J., . . . Conso, W. H. O. W. M. H. S. (2007). Delay and failure in treatment seeking after first onset of mental disorders in the World Health Organization's World Mental Health Survey Initiative. *World Psychiatry*, 6(3), 177-185. <Go to ISI>://WOS:000252131300016
- Wang, Y. D., Wang, Z. Y., Li, C. H., Zhang, Y. L., & Wang, H. Z. (2022). Online social network individual depression detection using a multitask heterogeneous modality fusion approach. *Information Sciences*, 609, 727-749. <https://doi.org/10.1016/j.ins.2022.07.109>
- WHO. (2017). *Depression and Other Common Mental Disorders: Global Health Estimates*. WHO.
- Xie, J., Chai, Y., & Liu, X. (2023). Unbox the Blackbox: Predict and Interpret YouTube Viewership Using Deep Learning. *Journal of Management Information Systems*, 40(2), 541-579.
- Yang, X., McEwen, R., Ong, L. R., & Zihayat, M. (2020). A big data analytics framework for detecting user-level depression from social networks. *International Journal of Information Management*, 54, 102141.
- Yu, S., Chai, Y., Samtani, S., & Liu, H. (2023). Motion Sensor-Based Fall Prevention for Senior Care: A Hidden Markov Model with Generative Adversarial Network (HMM-GAN) Approach. *Information Systems Research*, Online.
- Zhang, X., Gao, Y., Lin, J., & Lu, C.-T. (2020). Tapnet: Multivariate time series classification with attentional prototypical network. Proceedings of the AAAI Conference on Artificial Intelligence.
- Zhang, Z., Chen, s., Mengyue Wu, & Zhu, K. Q. (2022). Psychiatric Scale Guided Risky Post Screening for Early Detection of Depression. Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence.
- Zhang, Z., Chen, S., Wu, M., & Zhu, K. (2022). Symptom Identification for Interpretable Detection of Multiple Mental Disorders on Social Media. Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing.
- Zogan, H., Razzak, I., Wang, X. Z., Jameel, S., & Xu, G. D. (2022). Explainable depression detection with multi-aspect features using a hybrid deep learning model on social media. *World Wide Web-Internet and Web Information Systems*, 25(1), 281-304. <https://doi.org/10.1007/s11280-021-00992-2>