

# Building and Using Personal Knowledge Graph to Improve Suicidal Ideation Detection on Social Media

Lei Cao, Huijun Zhang, and Ling Feng *Senior Member, IEEE*

**Abstract**—A large number of individuals are suffering from suicidal ideation in the world. There are a number of causes behind why an individual might suffer from suicidal ideation. As the most popular platform for self-expression, emotion release, and personal interaction, individuals may exhibit a number of symptoms of suicidal ideation on social media. Nevertheless, challenges from both data and knowledge aspects remain as obstacles, constraining the social media-based detection performance. Data implicitness and sparsity make it difficult to discover the inner true intentions of individuals based on their posts. Inspired by psychological studies, we build and unify a high-level suicide-oriented knowledge graph with deep neural networks for suicidal ideation detection on social media. We further design a two-layered attention mechanism to explicitly reason and establish key risk factors to individual's suicidal ideation. The performance study on microblog and Reddit shows that: 1) with the constructed personal knowledge graph, the social media-based suicidal ideation detection can achieve over 93% accuracy; and 2) among the six categories of personal factors, *post*, *personality*, and *experience* are the top-3 key indicators. Under these categories, *posted text*, *stress level*, *stress duration*, *posted image*, and *ruminant thinking* contribute to one's suicidal ideation detection.

**Index Terms**—Suicidal ideation detection, social media, personal knowledge graph, social interaction.

## I. INTRODUCTION

SUICIDAL ideation is a problem which is on the rise across all the nations in the world. According to different reports and statistical figures, the lifetime prevalence of suicidal ideation of the entire world population is around 9 per cent [1], and this figure is significantly higher for individuals in the age group of 18 to 25 years old. There are a number of causes behind why an individual might suffer from suicidal ideation, and involvement of personality traits in susceptibility to suicide has been investigated for a long time. However, due to the limitations of mass group reaching and the diversity of conceptual and methodological approaches, the extent of their independent contributions remain is still have been difficult to establish. The risk factors for the immediate precursor to suicidal ideation are still not well-known especially in developing countries [1]. In any circumstance, for those with suicidal ideation, the earlier they can be detected, the better the chance of suicide prevention [2].

The authors are with the Department of Computer Science and Technology, Centre for Computational Mental Healthcare, Research Institute of Data Science, Tsinghua University, Beijing, China.  
E-mail: cao-l17@mails.tsinghua.edu.cn, zhang-hj17@mails.tsinghua.edu.cn, fengling@mail.tsinghua.edu.cn

In the literature, psychologists have developed a number of suicide risk measurements (such as Suicide Probability Scale [3], Adult Suicide Ideation Questionnaire [4], Suicidal Affect-Behavior-Cognition Scale [5], etc.) to assess individual's suicidal ideation. As this kind of methods requires people to either fill in a subjective questionnaire or participate in a professional interview, it is only applicable to a small group of people. For those who are suffering but tend to hide inmost thoughts and refuse to seek helps from others, the approach cannot work [6, 7, 8, 9].

With social media (like Twitter, online forums, and microblogs) becoming an integral part of daily lives nowadays, more and more people go to social media for information acquisition, self-expression, emotion release, and personal interaction. The large-scale, low-cost, and open advantages of social media offer us the unprecedented opportunity to capture individual's symptoms and traits related to suicidal ideation [10, 11, 12, 13, 14, 15, 16]. Nevertheless, analysis and detection of individual's suicidal ideation through social media are not trivial, facing a number of challenges from the following two aspects.

### A. Data-Aspect Challenges

Data implicitness and data sparsity are two critical data-aspect problems, challenging social media-based solutions for a long time.

**Data Implicitness.** Due to the unique equality, freedom, fragmentation, and individuality characteristics of social media, users linguistic and visual expressions on social media are implicit, reserved, and even anti-real. To illustrate, let's compare users' normal posts with their commenting posts in a hidden *microblog tree hole*. Here, a microblog tree hole refers to a microblog space, whose author has committed suicide, and other microblog users usually with suicidal ideation tend to share and post their inner real feelings and thoughts as comments under the last post of the passed one. One such a tree hole on Sina Weibo (the largest microblog platform in China) has collected over 1,700,000 commenting posts from more than 160,000 suicide attempts under the last post of microblog user called *Zoufan*, who committed suicide due to depression on March 17, 2012.

Figure 1 lists two users' example posts on microblog. From their normal posts, we can sense a joyful feeling. However, if reading their commenting posts in the hidden tree hole, we are touched by their severe hopelessness and suicidal thoughts.

User	Normal Posts	Hidden Tree Hole Posts
1	<p>2018-06-15 12:31:53 “What to buy after getting the salary?”</p> <p>2018-06-18 21:45:22 “I sang a song today. It's ugly. Everyone laughed at me. Who cares?”</p> <p>2018-06-19 08:12:13 “Ohh! Today is a new day. Good things happen one by one.”</p>	<p>2018-06-17 17:11:05 “In fact, I am dead, I am a worthless person. My work is a failure.”</p> <p>2018-06-18 19:31:27 “No place to go, timid, no matter where I go, I am still like this! Incurable! The place I want to go to heaven.”</p> <p>2018-06-19 05:45:00 “No meaning to live. I am like a corpse, who can understand me?”</p>
2	<p>2019-05-25 22:15:10 “My birthday! Thanks for the happy moment!”</p> <p>2019-05-26 07:00:00 “Breakfast is delicious!”</p>	<p>2019-05-25 23:55:15 “Today is my birthday, but I seem to have passed away.”</p> <p>2019-05-26 01:29:11 “I want to kill myself, living with my parents without a job. So hard to find a satisfactory one.”</p>

Fig. 1: Two users' normal posts versus their hidden tree hole posts on Sina Weibo.

To further examine the differences between users' normal posts and their commenting posts in the hidden tree hole, we crawled and analyzed 3,652 individuals' posts from May 1, 2018 to April 30, 2019 on Sina Weibo. As Table I shows, their posting contents were quite different. The hidden tree hole posts focused more on oneself through more self-concern words (like “I”, “me”, “my”, “mine”, “am”, etc.) rather than others-concern words, and they used more suicide-related words than the normal posts, demonstrating users' willingness and directives to express their inner suicidal thoughts in the hidden tree hole. On the contrary, the users were reluctant to show their feelings in the open normal posts, and they used more others-concern (such as “they”, “their”, “she”, “he”, etc.) words than the tree hole posts. Such phenomena challenge the suicide risk detection through the normal posts.

**Data Sparsity.** Apart from data implicitness, subjective to habits, characters, emotions, etc., some people may not want to express themselves actively on social media. Particularly, when immersed in hopelessness, loneliness, and suicidal thoughts, people tend to be self-enclosed and have very low motivations to write or leave something on social media.

To illustrate, we identified 3,652 users which made both normal posts and hidden tree hole posts on Sina Weibo, and their tree hole posts contained suicide-related words. We hypothesized that this group of users have suicidal ideation. For comparison purpose, we also randomly crawled another 3,677 active users, who only posted normal blogs on microblog, and

TABLE I  
STATISTICS OF USERS' NORMAL POSTS VS. HIDDEN TREE HOLE POSTS ON SINA WEIBO BASED ON 3,652 USERS FROM MAY 1, 2018 TO APRIL 30, 2019.

Perspective	Normal Posts	Hidden Tree Hole Posts
Prop. of posts containing <i>self-concern</i> words	43%	50%
Prop. of posts containing <i>others-concern</i> words	8%	6%
Prop. of posts containing <i>suicide-related</i> words	31%	95%
Prop. of <i>self-concern</i> words per post	4%	9%
Prop. of <i>others-concern</i> words per post	8%	1%
Prop. of <i>suicide-related</i> words per post	0.02%	0.3%
Total post number	252,901	190,087

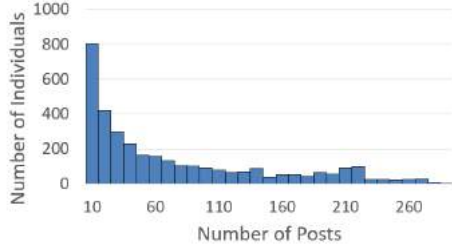
their blogs contained no suicide-related words at all. We called this group of users *ordinary users*. From Figure 2, we can find that users with suicidal thoughts made much less normal posts than ordinary users. Overall, the majority of this group of users had less than 1 normal post per week, verifying the serious data sparsity problem.

Moreover, due to violent emotional fluctuations, individuals with suicidal ideation tend to delete their previous posts related to suicidal ideation, striving to hide the true inner intentions. We counted the number of normal posts from the 7,329 users (3,652 with suicidal ideation and 3,677 ordinary users) in one year from 2018.05.01 to 2019.04.30 twice (one on 2019.04.30, and the other on 2019.08.31), and discovered

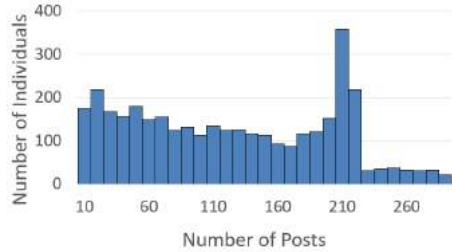
TABLE II

STATISTICS OF DELETED NORMAL POST NUMBERS BASED ON 7,329 USERS (3,652 WITH SUICIDAL IDEATION AND 3,677 NORMAL USERS) FROM MAY 1, 2018 TO APRIL 30, 2019 ON SINA WEIBO.

	# Users with Suicidal Ideation	# Ordinary Users
Number of users who deleted their previous normal posts after 3 months	1,673	326
Proportion of users who deleted their previous normal posts after 3 months	45.81%	8.87%
Proportion of deleted normal posts after 3 months	35.70%	3.95%



(a) Users with suicidal ideation



(b) Ordinary users

Fig. 2: Distributions of normal post numbers from 7,329 users (3,652 with suicidal ideation and 3,677 ordinary users) from May 1, 2018 to April 30, 2019 on Sina Weibo.

that (1) around 45.81% of the users with suicidal ideation deleted their previous normal posts after 3 months, while only 8.87% of ordinary users did it. (2) around 35.70% of normal posts were deleted by the users with suicidal ideation. For ordinary users, the figure was just about 3.95%. In other words, individuals with suicidal feelings tend to delete much more normal posts than ordinary ones (Table II).

The above data-aspect challenges hinder the detection performance greatly.

**Our Work.** To address the data-aspect challenges, we drew inspirations from previous psychological studies that individuals tend to find and follow one's own kind in the networks, looking for similar people, comfort each other, find a way to get rid of the difficulties, or even commit suicide together [17]. In such an emotional sharing and resonance process, negative suicidal emotion and influence could be developed and fast propagated among potentially suicidal individuals [18]. As in Figure 3, we can notice similar negative feelings from user 1 and his two followed friends. In our crawled 7,329 users, each user follows 5 neighbour users (as friends) on average.

Motivated by the findings, we leveraged individual's friends information via the follower-following relationship on microblogs to enrich users data and cope with data implicitness

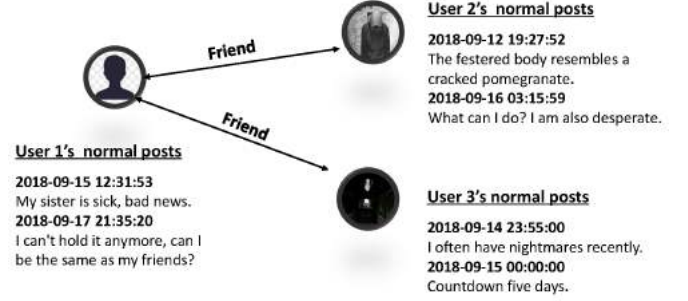


Fig. 3: Examples of normal posts by user 1 and his two followed users on Sina Weibo.

and sparsity problems in suicidal ideation detection.

### B. Knowledge-Aspect Challenges

Although there has been a large body of research addressing the task of suicidal ideation based on social media, the detection performance remains to be constrained due to the disparity between the informal language used by social media users and the concepts defined by domain experts in medical knowledge bases [19, 20, 21]. Moreover, methods based on domain knowledge have been successful in many fields [22, 23, 24]. To fill the gap, recently, [19] incorporated domain specific knowledge into a learning framework to predict the severity of suicide risk for a Redditor user. Specifically, it employed two medical lexicons (TwADR-L and AskaPatient) to map social media contents to medical concepts [25], and identified contents with negative emotions based on anonymized and annotated suicide notes available from Informatics for Integrating Biology and the Bedside (i2b2) challenge. It further developed a suicide risk severity lexicon using medical knowledge bases and suicide ontology to detect cues relevant to suicidal thoughts and actions, and used language modeling, medical entity recognition and normalization and negation detection to create a gold standard dataset of 500 Redditors developed by four practicing psychiatrists [19].

The work [19] proved that utilizing domain specific knowledge is a promising approach for suicidal ideation analysis. In reality, there could be a large number of causes behind why an individual might suffer from suicidal ideation. It is of extreme importance that the suicidal ideation detection method should be able to recognize the warning signs or symptoms around the individual who is suffering, and involve them to improve the detection performance. However, few contribution measurements of different indicators in relation to individual's *personal information* (display image on social media, gender,

age, etc.), *personality* (pursuing perfection, ruminant thinking, or interpersonal sensitivity), *experience* (stress or previous suicide attempt), negative emotion and emotion transition pattern, and social interaction impact have been established for suicidal ideation detection. In the literature, so far there has no social media-based *personal* suicide-oriented knowledge framework being developed and incorporated into suicidal ideation analysis.

**Our Work.** Beyond looking at individual's series of posting contents and social interactions on social media [26, 27], we built a structured personal suicide-oriented knowledge graph, and instantiated the personal knowledge graph by applying deep learning techniques [28] to 7,329 microblog users (3,652 with suicidal ideation and 3,677 without suicidal ideation) and 500 Redditors (108 without suicidal ideation, 99 with suicide indicator, 171 with suicidal ideation, 77 with suicide behavior and 45 with actual attempt).

We further built a two-layered attention mechanism (property attention and neighbour attention) to propagate model messages through the designed personal knowledge graph, and adaptively tune the weights (contributions) of different risk factors (properties and neighbour users) to individuals' suicidal ideation.

Finally, we incorporated and unified the achieved personal suicide-oriented knowledge graph with graph neural networks to facilitate suicidal ideation detection. Our experimental results on microblog and Reddit showed that with the personal suicide-oriented knowledge graph and the two attention mechanisms, suicidal ideation detection can achieve over 93% of accuracy and F1-measure. Among the six categories of personal factors, *posts*, *personality*, and *experience* are the top-3 key indicators. Under these categories, *posted texts*, *stress level*, *stress duration*, *posted image*, and *ruminant thinking* contribute the most to one's suicidal ideation detection.

In summary, the contributions of the study are three-fold.

- We build and unify a high-level suicide-oriented knowledge graph with graph neural networks for suicidal ideation detection on social media.
- We design a two-layered attention mechanism that explicitly reasons and establishes key risk factors to individual's suicidal ideation.
- We deliver a social media-based suicide-oriented knowledge graph. Together with the labeled dataset (<https://github.com/bryant03/Sina-Weibo-Dataset>) with 3,652 (3,677) users with (without) suicidal ideation from Sina Weibo, they could facilitate further suicide-related studies in the computer science and psychology fields.

## II. LITERATURE REVIEW ON DETECTING SUICIDAL IDEATION

### A. Traditional Questionnaire-based Suicide Risk Assessment

Traditional methods rely on questionnaires and face-to-face diagnosis to assess whether one is in the risk of suicide [29]. Some widely used psychological measurements include Suicide Probability Scale (SPS) [3], Depression Anxiety Stress Scales-21 (DASS-21) [30, 31], Adult Suicide Ideation Questionnaire [4], Suicidal Affect-Behavior-

Cognition Scale [5], functional Magnetic Resonance Imaging (fMRI) signatures [32], and so on. While this kind of approaches are professional and effective, they require respondents to either fill in a questionnaire or participate in a professional interview, constraining its touching to suicidal people who have low motivations to seek help from professionals [6, 7, 8, 9]. A recent study found out that taking a suicide assessment may bring negative effect to individuals with depressive symptoms [33].

### B. Suicide Risk Detection from Social Media

Recently, detection of suicide risk from social media is making great progress due to the advantages of reaching massive population, low-cost, and real-time [34]. [35] reported that suicidal users tend to spend a lot of time online, have great likelihood of developing online personal relationships, and great use of online forums.

**Suicide Risk Detection from Suicide Notes.** [36] built a suicide note classifier used machine learning techniques, which performs better than human psychologists in distinguishing fake online suicide notes from real ones. [37] hunted suicide notes based on lexicon-based keyword matching on MySpace.com (a popular site for adolescents and young adults, particularly sexual minority adolescents with over 1 billion registered users worldwide) to check whether users have an intent to commit suicide.

**Suicide Risk Detection from Community Forums.** [38] applied textual sentiment analysis and summarization techniques to users' posts and posts' comments in a Chinese web forum in order to identify suicide expressions. [39] examined online forums in Japan, and discovered that the number of communities which a user belongs to, the intransitivity, and the fraction of suicidal neighbors in the social network contributed the most to suicidal ideation. [40] built a logistic regression framework to analyze Reddit users' shift tendency from mental health sub-communities to a suicide support sub-community. heightened self-attentional focus, poor linguistic coherence and coordination with the community, reduced social engagement and manifestation of hopelessness, anxiety, impulsiveness and loneliness in shared contents are distinct markers characterizing these shifts. Based on the suicide lexicons detailing suicide indicator, suicidal ideation, suicide behavior, and suicide attempt, [11] built four corresponding semantic clusters to group semantically similar posts on Reddit and questions in a questionnaire together, and used the clusters to assess the aggregate suicide risk severity of a Reddit post. [19] used Reddit as the unobtrusive data source to conduct knowledge-aware assessment of severity of suicide risk for early intervention. Its performance study showed that convolutional neural network (CNN) provided the best performance due to the discriminative features and use of domain-specific knowledge resources, in comparison to SVM-L that has been used in the state-of-the-art tools over similar datasets. [41] introduced a new evaluation measure to move beyond suicide risk classification to a paradigm in which prioritization is the focus. The proposed joint ranking approach outperformed the logistic regression under the new evaluation measure. [42]

used document embeddings generated by transfer learning to classify users at suicide risk or not. [43] employed an LSTM-CNN combined model to evaluate and compare to other classification models.

**Suicide Risk Detection from Blogs.** [44] used search keywords and phrases relevant to suicide risk factors to filter potential suicide-related tweets, and observed a strong correlation between Twitter-derived suicide data and real suicide data, showing that Twitter can be viewed as a viable tool for real-time monitoring of suicide risk factors on a large scale. The correlation study between suicide-related tweets and suicidal behaviors was also conducted based on a cross-sectional survey [45], where participants answered a self-administered online questionnaire, containing questions about Twitter use, suicidal behaviour, depression and anxiety, and demographic characteristics. The survey result showed that Twitter logs could help identify suicidal young Internet users.

Based on eight basic emotion categories (joy, love, expectation, anxiety, sorrow, anger, hate, and surprise), [46] examined three accumulated emotional traits (i.e., emotion accumulation, emotion covariance, and emotion transition) as the special statistics of emotions expressions in blog streams for suicide risk detection. A linear regression algorithm based on the three accumulated emotional traits was employed to examine the relationship between emotional traits and suicide risk. The experimental result showed that by combining all of three emotion traits together, the proposed model could generate more discriminative suicidal prediction performance.

Natural language processing and machine learning techniques, such as Latent Dirichlet Allocation (LDA), Logistic Regression, Random Forest, Support Vector Machine, Naive Bayes, Decision Tree, etc., were applied to identify users' suicidal ideation based on their linguistic contents and online behaviors on Sina Weibo [47, 48, 49, 50, 12] and Twitter [51, 52, 53, 54, 55, 56, 57]. Deep learning based architectures like Convolutional Neural Network (CNN), Recurrent Neural Network (RNN), Long Short-Term Memory Neural Network (LSTM), etc., were also exploited to detect users' suicide risk on social media [13, 14, 58, 15]. [16] detected users' change points in emotional well-being on Twitter through a martingale framework, which is widely used for change detection in data streams.

Very recently, [59] designed suicide-oriented word embeddings to capture the implicit suicidal expression and proposed a suicide risk detection model based on LSTM and ResNet. [27] applied the machine learning algorithms to detect suicidal profiles in Twitter based on features extracted from users accounts and tweeting behaviors. [26, 60] proposed a multipronged approach and implemented different neural network models such as sequential models and graph convolutional networks, which were trained on textual content shared in Twitter, tweeting histories of the users and social network formed between different users posting about suicidality.

Different from the previous work, in this study, we looked beyond individual's posting contents and social interactions on social media, and incorporated the constructed personal knowledge graph into suicidal ideation detection. The performance study showed that the proposed method can achieve over 93%

of accuracy and F1-measure.

### III. FRAMEWORK

#### A. Building the Personal Suicide-Oriented Knowledge Graph

We define the terminology for representing the nodes and edges of the personal suicide-oriented knowledge graph, and then convert the extracted data from individuals' social media accounts to instantiate such a personal knowledge graph. Inspired by the psychological investigation into the predictors of suicide risk [61, 62, 63, 64, 65, 66, 67], we extracted and analyzed individuals' relevant social media behaviors from the following six perspectives (i.e., *personal information*, *personality*, *experience*, *post behavior*, *emotion expression*, and *social interaction*), as shown in Figure 4.

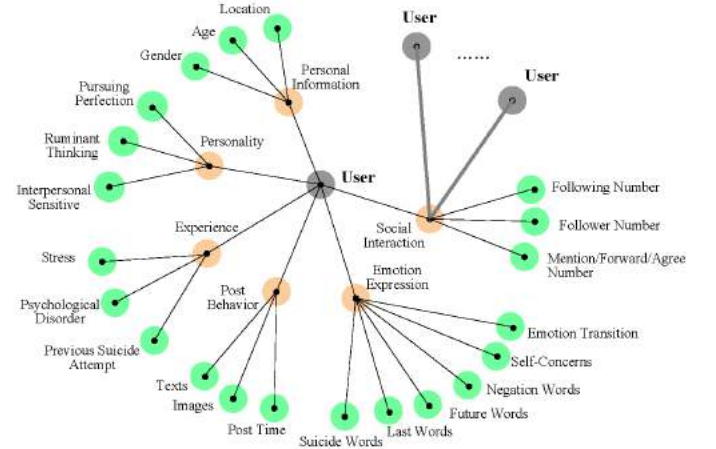


Fig. 4: The ontology of the social media-based suicide-oriented knowledge graph.

1) *Personal Information (Gender, Age, Location)*: Based on the previous study result that suffering women are three times more likely than suffering men to attempt suicide [61], we captured user  $u$ 's personal details including gender, age, and location from his microblog account, and involved them into the personal knowledge graph.

We used a 3-dimensional vector to describe user's gender information.  $Gender(u) = (1,0,0)$ ,  $(0,1,0)$ , or  $(0,0,1)$ , representing *female*, *male*, or *unknown*, respectively.

We used a scalar from zero to one to represent user's age. The scalar is obtained by dividing the user's age by the maximal user age in the collected dataset.  $Age(u) = \frac{u.age}{max\_user\_age}$ , where  $max\_user\_age=65$ .

Like user's gender expression, we used a vector of length 8 to represent user  $u$ 's location.  $Loc(u) = (1,0,0,\dots,0)$ ,  $(0,1,0,\dots,0)$ ,  $(0,0,1,\dots,0)$ ,  $\dots$ , or  $(0,0,0,\dots,1)$ . Each element of the vector represents a certain geographic direction. It takes value 1 if the user is in the corresponding location, and 0 otherwise. Totally 8 location directions (i.e., *east*, *south*, *north*, *south-west*, *north-west*, *middle*, *north-east*, plus an *unknown* location) were considered.

2) *Experience (Stress, Psychological Disorder, Previous Suicide Attempt)*: Suicidal ideation often occurs when an individual is no longer able to cope with some kind of difficult or overwhelming situation [62]. Some of the commonly

agreed causes of suicidal ideation include situations where an individual experiences continuous stressful periods of high stress levels and different stress categories [63]. To detect user's stress (stress periods, stress levels, and stress categories) within the year, we applied the algorithm [68] to user's posting behaviors, and captured a series of stressful periods  $\mathcal{S}(u)$ , each of which is of the form  $s=(s_p, s_l, s_c)$ , where  $s \in \mathcal{S}(u)$ ,  $s_p$  is a temporal stress period which should be over five days,  $s_l$  is the stress level 1 (*weak stress*) or 2 (*strong stress*), and  $s_c$  is the stress category, which can be *study*, *work*, *family*, *interpersonal relation*, *romantic relation*, or *self-cognition*. We aggregated user  $u$ 's stressful periods into the total number of stress periods, the average stress level, and the number of different stress categories that user  $u$  experienced throughout the year.

$$Stress(u)=(sNum(u), sLevel(u), sCatNum(u)),$$

where  $sNum(u) = |\mathcal{S}(u)|$ ,  $sLevel(u) = \frac{\sum_{s \in \mathcal{S}(u)} s_l}{|\mathcal{S}(u)|}$ , and  $sCatNum(u) = |\{s_c \mid s \in \mathcal{S}(u)\}|$ .

Furthermore, people who are suffering from psychological disorders like depression or bipolar disorder are also more prone to suffering from suicidal ideation. People who previously tried to attempt suicide also tend to have a higher suicidal risk than those who did not. Like [69], we set  $Disorder(u)=1$ , if and only if user  $u$  posted something on microblog like “*I have depression/bipolar disorder ...*”, and 0 otherwise.  $Attempt(u)=1$ , if and only if user  $u$  posted something on microblog like “*I've committed suicide #number# times ...*”, and 0 otherwise.

3) *Personality (Pursuing Perfection, Ruminant Thinking, Interpersonal Sensitive)*: Involvement of personality factors in susceptibility to suicidality has been the subject of research since the 1950s. Psychological findings showed that personality dimensions are significantly associated with suicide-related behaviors, and specific personality factors impact suicide uniquely for each gender [64]. Such personality types like *pursuing perfection*, *ruminant thinking*, and *interpersonal sensitivity* could be useful markers of suicide risk [65].

To assess the extent of one's *pursuing perfection* characteristic, we identified 36 potential perfectionist users from the “Perfectionism” subreddit (<https://www.reddit.com/r/perfectionism/>), whose posts contained sentences like “*like many perfectionists, I feel ...*”, “*I definitely have unhealthy perfectionism*”, and so on. We calculated top 100 high-frequency words from their posts. We remove stop words and words representing things like examination, and acquire such words like “*no enough*”, “*perfect*”, “*loser*”, “*failure*”, “*imperfect*”, “*compete*”, and so on. We took them as seed words, and enrich the perfection-related lexicon with synonymous or similar words. We also made a ruminant thinking-related lexicon consisting of very destructive negative words such as “*regret*”, “*repent*”, “*rue*”, “*penitent*”, “*self-blame*”, etc. Table III lists typical perfection-related and ruminant thinking-related words that we used in the study.

We measured user's pursuing perfection and ruminant thinking based on the mean proportion of perfection-related and ruminant-related words per post, respectively. Let  $P(u)$  denote the set of posts made by user  $u$ . For each post  $p \in P(u)$ , as-

TABLE III  
TYPICAL PERFECTION-RELATED AND RUMINANT  
THINKING-RELATED WORDS.

Type	Words
Perfection-Related	perfect, perfectionist, perfection, stupid, wrong, upset, mood, struggle, goal, realistic, not enough, lose, loser, failure, imperfect, compete, force, negative, disappoint, despair, letdown, useless, always, never, still, upset, problem
Ruminant Thinking-Related	regret, repent, rue, penitent, confess, hate, self-blame, grievance, complaint, injustice, unfairness, growl, discontent, lose, loser, owe, sinner, die

sume function  $pwordNum(p, u)$  and  $rwordNum(p, u)$  return the number of perfection-related and ruminant-related words in  $p$ , respectively.

$$Perfect(u) = \frac{1}{|P(u)|} \sum_{p \in P(u)} \frac{pwordNum(u, p)}{|p|}$$

$$Ruminat(u) = \frac{1}{|P(u)|} \sum_{p \in P(u)} \frac{rwordNum(u, p)}{|p|}$$

We measured user  $u$ 's interpersonal sensitivity by counting how many times the user experienced stressful periods in the stress category of *inter-personnel relation*.

$$Sensitive(u) =$$

$$|\{s \mid s \in \mathcal{S}(u) \wedge s_c = \text{"interpersonal relation"}\}|$$

where  $\mathcal{S}(u)$  is a set of stressful periods detected from user  $u$ 's posts, and  $s=(s_p, s_l, s_c)$  is one of it in  $\mathcal{S}(u)$ .

4) *Post Behavior (Text, Image, Post Time)*: An individual may post texts and images anytime at will. These linguistic and/or visual contents, as well as the posting time, may reveal personal thoughts or feelings. Since Bert [70] and ResNet [71] are common used pre-trained models which perform well on textual feature extraction [72, 73] and visual feature extraction [74, 75]. For each post, we encoded its linguistic and visual contents into a 768-dimensional and a 300-dimensional vector, respectively, through a pre-trained Bert model and a 34-layer ResNet. When the image was missing, we filled in with a default null image, which is usually displayed at the position for image wanted in applications. When the post contains multiple images, we took the average feature vector of the images as the visual content representation of the post. We mapped the user's absolute post time into the hour from 0 to 23. Assume user  $u$  made totally  $n$  posts  $post_1(text_1, image_1, ptime_1), \dots, post_n(text_n, image_n, ptime_n)$ , where  $text_i \in \mathbb{R}^{768 \times 1}$ ,  $image_i \in \mathbb{R}^{300 \times 1}$ , and  $ptime_i \in \{0, 1, 2, \dots, 23\}$  for  $1 \leq i \leq n$ .

To further learn user  $u$ 's posts representation, we concatenated  $(text_i, image_i, ptime_i)$  as a 1069-dimensional vector  $postvec_i$ , where  $postvec_i = text_i || image_i || ptime_i \in \mathbb{R}^{1 \times 1069}$ , and then applied LSTM to extract textual and visual information from  $postvec_1, \dots, postvec_n$ . Let the hidden dimension size of LSTM be 300.

$$out_i, h_i = LSTM(postvec_i, h_{i-1}).$$

The last output of LSTM was  $out_n \in \mathbb{R}^{1 \times 300}$ . To reduce the computational complexity of subsequent calculations, a fully



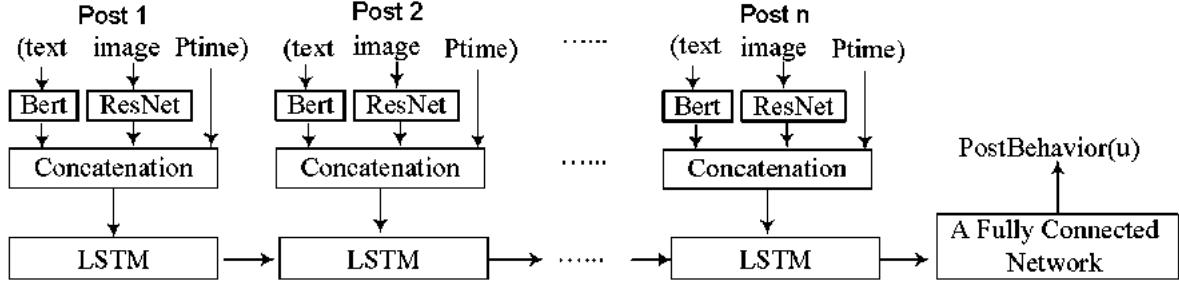


Fig. 5: The framework for learning user  $u$ 's post behavior.

connected layer was used to generate a 30-dimensional vector as the value of user's *post behavior* property:

$$PostBehavior(u) = ReLU(out_n W_0 + b_0) \in \mathbb{R}^{1 \times 30},$$

where  $out_n \in \mathbb{R}^{1 \times 300}$ ,  $W_0 \in \mathbb{R}^{300 \times 30}$ , and  $b_0 \in \mathbb{R}^{1 \times 30}$ .

Details of learning user  $u$ 's post behavior representation are illustrated in Figure 5.

5) *Emotion Expression (Suicide Words, Last Words, Future Words, Negation Words, Self-Concerns, Emotion Transition)*: Suicides tend to express rather than repress their despair suicidal feelings, less future words, more negation and self-referenced first-person pronouns such as *I*, *me*, *my*, *myself*, and *mine*, etc. in their posts [66]. The existence of last words (which express regret, guilty, wishes towards their families and friends, or funeral), particularly in those recent posts, is also recognized as an important observable risk factor of suicidal ideation [76]. From the Chinese Suicide Dictionary [77] and Chinese Linguistic Inquiry and Word Count [78], we extracted 586 about suicide, 125 about last words, 86 about future words, 665 about negation words, and 36 words/phrases about self-concern. We calculated the mean word frequencies (i.e., mean word proportions over the total words per post) as the property values of  $SuicideProp(u)$ ,  $LastWordProp(u)$ ,  $FutureProp(u)$ ,  $NegProp(u)$ , and  $SelfProp(u)$  for user  $u$ . Here, only two most recent weeks' posts were considered in the computation of  $LastWordProp(u)$ .

Apart from the specific words, we drew the inspiration from the study [46] that there is a significant difference in emotion transition between suicide users and non-suicide users. That is, love emotion is more regularly followed by anxiety or sorrow emotion among the suicide people, while love is more regularly followed by joy emotion among the non-suicide people. Hereby, as in [46], we considered eight emotions (*love*, *joy*, *expectation*, *anxiety*, *sorrow*, *anger*, *hate*, and *surprise*), and counted two types of emotion transitions from love to joy ( $love\text{-}joy(u)$ ) or from love to anxiety/sorrow ( $love\text{-}anxiety/sorrow(u)$ ) within the ten most recent continuous posts as the property value of emotion transition:  $EmotionTrans(u) = (love\text{-}joy(u), love\text{-}anxiety/sorrow(u))$ .

For user  $u$ , if there existed a pair of posts, where the first one contained one or more love lexicons, and the second later post contained one or more joy lexicons based on the Chinese emotion lexicons DUTIR (<http://ir.dlut.edu.cn/>), the

value of  $love\text{-}joy(u)$  was increased by 1. The value of  $love\text{-}anxiety/sorrow(u)$  was computed in the same way.

6) *Interaction (Following Number, Follower Number, Mention/Forward/Agree Number, Neighbour Users)*: Social isolation is a significant and reliable predictor of suicidal ideation. When trapped in trouble, if one could get the support, understanding, and acceptance from family, friends and peers, his/her stress and follow-up negative emotions could be offset, which could probably avoid extreme suicidal behavior. In reality, desperate people tend to have a weak social network, and thus get weak social supports [67]. Here, one's online social engagement is reflected from his following/follower number, mention/forward/agree number, and neighbour users. Specially, in order not to be affected by "Zombie fan" (fake fans in social media), we filtered out fans who posted less than three original posts in the last six months when counting the follower number.

$Interact(u) = (FollowingNum(u), FollowerNum(u), IntActNumb(u), Users(u))$ , where the first three elements take non-negative integer values, and  $Users(u)$  is a set of neighbour users which user  $u$  follows.

Formally, we describe the personal suicide-oriented knowledge graph as  $G = \{V, E\}$ , where  $V$  and  $E$  are node set and edge set respectively. A user node can have multiple property nodes, and may link to another user node via the following relationship on the social media. The weight of a social interaction-user edge implies the suicidal emotion influence between the two users on the social media, and the weight of a user-property edge represents the contribution of the property to the suicidal ideation. These weights will be computed during the following learning process.

Finally, the statistic of knowledge graph is shown in Table IV.

## B. Suicidal Ideation Detection based on the Personal Knowledge Graph

We designed a property attention mechanism and neighbour attention mechanism to reason the key factors related to suicidal ideation, and discover the suicidal emotion influence among the neighbour users on the social media.

1) *Property Attention*: Excepting user  $u$ 's neighbour users, we concatenated all the instantiated properties from  $u$ 's personal knowledge graph, and obtained a 61-dimensional property vector  $P_u = (p_1, p_2, \dots, p_n) \in \mathbb{R}^{61 \times 1}$ , where  $n=61$ .

TABLE IV  
STATISTICS OF USERS' PERSONAL SUICIDE-ORIENTED KNOWLEDGE GRAPHS ON WEIBO.

Category	Properties	Users with Suicidal Ideation	Ordinary Users without Suicidal Ideation
Personal Information	Gender ( <i>male,female,unknown</i> )	(21.3%, 78.5%, 0.2%)	(50.6%, 42.3%, 7.1%)
	Age on average	25.8	28.3
	Location ( <i>east, south, north, south-west, north-west, middle, north-east, unknown</i> )	(18.8%, 10.5%, 7.4%, 6.6%, 3.0%, 7.0%, 3.3%, 43.5%)	(25.0%, 15.0%, 15.7%, 5.8%, 2.9%, 7.9%, 5.3%, 22.3%)
Personality	Pursuing Perfection	0.25%	0.17%
	Ruminant Thinking	0.086%	0.052%
	Interpersonal Sensitive	1.3	1.0
Experience	Number of Stressful Periods	2.1	1.8
	Psychological Disorder	0.1%	0.03%
	Previous Attempt	3.5%	0.1%
Post Behavior	Number of Texts	252,901	491,130
	Number of Images	93,461	260,667
Emotion Expression	Proportion of Suicide Words Per Post	0.034%	0.016%
	Proportion of Last Words Per Post	0.0013%	0.000023%
	Proportion of Future Words Per Post	0.031%	0.045%
	Proportion of Negation Words Per Post	0.012%	0.009%
	Proportion of Self-Concern Words Per Post	0.079%	0.029%
	Emotion Transition ( <i>love-joy, love-anxiety/sorrow</i> )	(0.1, 0.5)	(0.3, 0.1)
Social Interaction	Following Number	207.0	378.1
	Follower Number	566.9	1515.3
	Mention/Forward/Agree Number	3.4	10.9
	Number of Neighbour Users	4.3	5.1

To compute the significance of different properties, we enforced  $P_u$  with a property attention vector  $\alpha \in \mathbb{R}^{1 \times 61}$  and acquired a new property vector  $P'_u \in \mathbb{R}^{61 \times 1}$ :

$$P'_u = P_u \times \alpha^T, \\ \alpha = \text{softmax}((P_u)^T W_1 + b_1),$$

where  $W_1 \in \mathbb{R}^{61 \times 61}$  and  $b_1 \in \mathbb{R}^{1 \times 61}$  are two trainable parameters.

2) *Neighbour Attention*: We imposed neighbour influences upon user  $u$ 's property representation  $P'_u$ . Considering different neighbour influences upon user  $u$ 's suicidal ideation, for instance, some people may easily be influenced by elders, and else may trust their peers more, we adopted a neighbour attention mechanism similar to the graph attention mechanism [28].

To reduce computational complexity, we firstly figured out user  $u$ 's initial hidden state  $h_u \in \mathbb{R}^{1 \times 60}$  from his property representation  $P'_u$  through a fully connected layer:

$$h_u = \tanh((P'_u)^T W_2 + b_2),$$

where  $W_2 \in \mathbb{R}^{61 \times 60}$  and  $b_2 \in \mathbb{R}^{1 \times 60}$  are trainable parameters.

Let  $\mathbb{N}_u$  be the set of direct neighbour users linked with user  $u$  on the social media. For each neighbour pair  $(u, \tilde{u})$  (where  $\tilde{u} \in \mathbb{N}_u$ ), we obtained the attention coefficient  $c_{u,\tilde{u}} \in \mathbb{R}^{1 \times 1}$  through another fully connected layer with concatenate operation  $\parallel$ :

$$c_{u,\tilde{u}} = \tanh((h_u \parallel h_{\tilde{u}})^T W_3 + b_3),$$

where  $W_3 \in \mathbb{R}^{120 \times 1}$  and  $b_3 \in \mathbb{R}^{1 \times 1}$  are two trainable parameters.

In this way, we obtained a vector of user  $u$ 's neighbour attention coefficients:

$$C_u = [c_{u,1}, c_{u,2}, \dots, c_{u,|\mathbb{N}_u|}] \in \mathbb{R}^{1 \times |\mathbb{N}_u|}$$

where  $|\mathbb{N}_u|$  is the number of  $u$ 's direct neighbours.

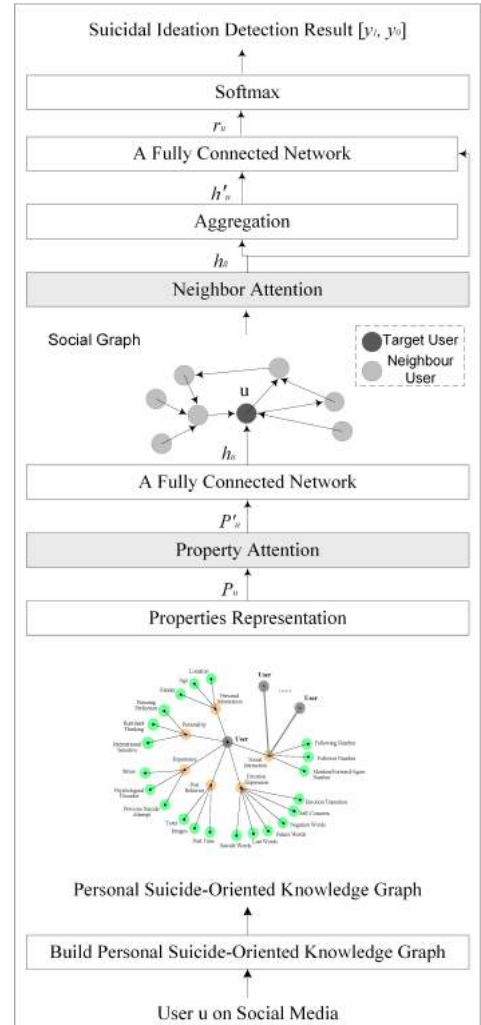


Fig. 6: Architecture of the personal knowledge graph based suicidal ideation detection.



A softmax function was then applied to compute a series of scores  $\mathcal{B}_u = [\beta_{u,1}, \beta_{u,2}, \dots, \beta_{u,|\mathbb{N}_u|}] \in \mathbb{R}^{1 \times |\mathbb{N}_u|}$  to represent the influences of neighbour users  $\mathbb{N}$  upon user  $u$ . The higher the score is, the bigger the suicidal influence is.

$$\mathcal{B}_u = \text{softmax}(C_u).$$

By means of  $\mathcal{B}_u$ , we aggregated information from neighbour users and updated  $u$ 's hidden state  $h_u$  accordingly into  $h'_u \in \mathbb{R}^{1 \times 60}$ :

$$h'_u = \sigma\left(\sum_{\tilde{u} \in \mathbb{N}_u} \beta_{u,\tilde{u}} h_{\tilde{u}} + h_u\right).$$

Once obtained, we utilized a fully connected layer to get the user's final representation  $r_u \in \mathbb{R}^{1 \times 60}$ :

$$r_u = \tanh(h'_u W_4 + b_4),$$

where  $W_4 \in \mathbb{R}^{60 \times 60}$  and  $b_4 \in \mathbb{R}^{1 \times 60}$ .

3) *Suicidal Ideation Detection*: Finally, a fully connected layer and softmax function were applied to detect suicidal ideation of user  $u$ :

$$[y_1, y_0] = \text{softmax}(r_u W_5 + b_5),$$

where  $y_1, y_0$  represent the possibility of user  $u$  with or without suicidal ideation,  $W_5 \in \mathbb{R}^{60 \times 2}$  and  $b_5 \in \mathbb{R}^{1 \times 2}$  are trainable parameters.

#### IV. EXPERIMENTS

##### A. Datasets

In this study, we evaluate our personal suicide-oriented knowledge graph based suicidal ideation detection methods on two types of datasets: microblog posts from Sina Weibo and forum posts from Reddit.

1) *Sina Weibo Dataset*: On March 17, 2012, a user whose screen name was "Zoufan" left his last word on Sina Weibo and then committed suicide. Since then over 160,000 users gathered here and posted more than 1,700,000 comments, and the numbers keep growing today, forming a large microblog tree hole. The majority of these commenting posts spoke out the posters' minds, disclosing their tragic experiences, hopeless thoughts, and even plans of suicide.

We crawled all the users' commenting posts from May 1, 2018 to April 30, 2019, and computed each user's expression degrees related to suicide according to the Chinese social media-based suicide dictionary [66]. The dictionary lists 2168 suicide-related words and phrases, falling into 13 categories (suicidal thought, self-harm, psychache, mental illness, hopeless feeling, somatic complain, self-regulation, negative personality, stress, trauma/hurt, talking about others, shame/guilt, anger/hostility). Each word/phrase is assigned a weight from 1 to 3, indicating its binding degree with suicide risk.

Assume user  $u$  posted a set of comments  $P(u)$  in the tree hole. We define  $u$ 's expression degree related to suicide by counting his total suicide-related weighted words/phrases in  $P(u)$ .

$$\begin{aligned} \text{UserDegree}(u) &= \sum_{post \in P(u)} \text{postDegree}(\text{post}) \\ \text{postDegree}(\text{post}) &= \sum_{w \in \text{post}} \text{weight}(w) \end{aligned}$$

where  $\text{weight}(w)$  is the weight of word/phrase  $w$  in  $p$  according to the suicide dictionary [66]. A word/phrase not in the suicide dictionary has weight 0.

We ranked and selected top 4,000 users who had high expression degrees related to suicide. After that, we recruited four PhD researchers to manually scan the 4,000 users and keep 3,652 users, who clearly mentioned suicide thought, suicide plan, or self-injury at least five times in different days. For example, if a user expressed clear suicidal thoughts like "At this moment, I especially want to die. I feel very tired. I really want to be free.", "Heh, even Weibo can't give me a reason to keep on going." or "I don't want to do anything for the last five days of my life." more than 5 times in different days, then we label him/her at suicide risk. Finally, we labeled these 3,652 users with suicidal ideation at suicide risk.

As comparison, we randomly selected 3,677 ordinary users on Sina Weibo, subject to the following three conditions. Firstly, they were active users, each making over 100 normal posts on the open microblog during the same temporal period. Secondly, neither of their posts contained any suicide-related word or phrase. Thirdly, they did not make any commenting post in any hidden tree hole. In the study, we regarded the ordinary users without suicidal ideation. If two users had a following-follower relationship on Sina Weibo, we built an edge linking the two neighbour users. Table V illustrates the data collected from Sina Weibo.

The training set, validation set, and test set contained 6,129, 600, and 600 users, who were exclusively and randomly chosen from the total 7,329 microblog users, respectively. Users with suicidal ideation and ordinary users without suicidal ideation occupied nearly half in each set.

2) *Reddit Dataset*: The Reddit dataset was crawled from the popular online forum Reddit, and contained 500 users in total. There were five categories of users, i.e., Supportive, Suicide Indicator, Suicidal Ideation, Suicidal Behavior, and Actual Attempt with user number of 108, 99, 171, 77, and 45, respectively. From Supportive to Actual Attempt, the risk of suicide gradually increases. We split the Reddit dataset into training set, validation set, and test set with the number of 303, 99, and 98, respectively.

##### B. Experimental Setup

We compared the performance of our personal suicide-oriented knowledge graph based method with the following four methods.

- **CNN** [19]: A Convolutional Neural Network takes embeddings of user posts as input.
- **LSTM+Attention** [15]: An self-attention mechanism based Long Short-Term Memory model which captures the contextual information between suicide-related words and others.
- **SDM** [59]: An hierarchical attention network based on Long Short-Term Memory module and ResNet module.
- **Text+History+Graph** [26]: A multipronged approach includes bidirectional Long Short-Term Memory module and Graph Neural Network. The textual content shared in Twitter, the historical tweeting activity of the users and

TABLE V  
STATISTIC OF DATA COLLECTED FROM SINA WEIBO.

	# Users	# Normal Posts	# Normal Posts with Images	# Neighbour Users to be followed
Users with suicidal ideation	3,652	252,901	93,461	4.3
Ordinary users without non-suicide	3,677	491,130	260,667	5.1
Total users	7,329	744,031	354,128	4.7

social network formed between different users posting about suicidality are concerned to identify and explore users' suicidal ideation.

As Reddit dataset only contained posts from each user, during the experiment, we deleted the ResNet module and the user profile module of the SDM method, the Social Graph module of the Text+History+Graph method, as well as the Neighbour Attention module, Personal Information property and Social Interaction property of our KG-based method. We took Bert contextual word embeddings [70] as our word embeddings.

Five metrics (accuracy, precision, recall, F1-measure, macro-F1-measure) were adopted in the performance evaluation.

$$Precision = \frac{TP}{TP+FP}, \quad Recall = \frac{TP}{TP+FN}$$

$$Accuracy = \frac{TP+TN}{TP+FP+TN+FN}$$

$$F1-measure = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

where true positive (TP), false positive (FP), true negative (TN), and false negative (FN) are defined in the following confusion matrix (Table VI).

As toward Reddit dataset there were more than two classes (actually 5), we employed a commonly used multi-class macro-F1-measure, which calculates F1-measure for each class, and then finds their un-weighted mean.

TABLE VI  
CONFUSION MATRIX USED TO EVALUATE THE PERFORMANCE.

Detected \ Actual	Positive	Negative
Positive	True Positive (TP)	False Negative (FN)
Negative	False Positive (FP)	True Negative (TN)

### C. Effectiveness of Involving Personal Suicide-Oriented Knowledge Graphs in Suicidal Ideation Detection

Table VII shows that, with the considerations of visual information from posts and users' profile information, the SDM method achieved better performance compared with that of CNN method and LSTM-Attention method. After using the social graph of users, Text+History+Graph method outperformed the above methods slightly in accuracy and F1-measure. Furthermore, involving personal suicide-oriented knowledge graphs into suicidal ideation detection could improve the detection accuracy, F1-measure, precision, and recall by 2.18%, 1.88%, 1.9%, 1.87% over the Text+History+Graph method, demonstrating the effectiveness of the personal knowledge graph based method.

TABLE VII  
PERFORMANCE COMPARISON ON THE SINA WEIBO DATASET.

Method	Acc.	F1.	Prec.	Rec.
CNN	86.77%	84.47%	84.19%	84.77%
LSTM+Attention	88.89%	88.10%	88.56%	87.65%
SDM	91.35%	90.97%	90.11%	91.85%
Text+History+Graph	91.56%	91.81%	91.85%	91.77%
KG-based method	<b>93.74%</b>	<b>93.69%</b>	<b>93.75%</b>	<b>93.64%</b>

TABLE VIII  
PERFORMANCE COMPARISON ON THE REDDIT DATASET.

Method	Acc.	macro.F1.	Prec.	Rec.
CNN	52.31%	52.61%	53.05%	52.19%
LSTM+Attention	55.12%	54.49%	53.16%	55.89%
SDM	60.58%	60.44%	60.31%	60.58%
Text+History+Graph	59.61%	59.02%	58.42%	59.64%
KG-based method	<b>65.92%</b>	<b>65.93%</b>	<b>65.35%</b>	<b>66.21%</b>

Since five-class classification is a harder task than two-class classification, and the modalities of the Reddit dataset are less than that of the Sina Weibo dataset, the overall performance of all the methods dropped, compared with that on the Sina Weibo dataset, as shown in Table VIII. Without the contribution of social graph, the performance of the Text+History+Graph method dropped to 59.61% and 59.02% in accuracy and F1-measure, respectively, which were worse than the SDM method. With the help of personal suicide-oriented knowledge graph, the KG-based method still outperformed the rest methods.

### D. Effectiveness of Property and Neighbour Attention Mechanisms

Property attention and neighbour attention are two crucial modules of our personal knowledge graph based method. They aim to reason about the key properties and find the most influential neighbour users to user's suicidal ideation, reflected from their associated property weights and neighbour weights, respectively. To investigate the effectiveness of the two attention mechanisms, we replaced each learnt attention weight with an average or random value in the range of [0,1], respectively. Table IX shows that the two attention mechanisms were effective, achieving the best performance compared to the alternative average and random weight solutions, and the average solution was better than the random solution.

It is worth mentioning here that the property+neighbour attention strategy drew inspirations from the Graph Attention networks (GAT) [28], which incorporates the neighbor attention mechanism. This also contributes to its best performance

TABLE IX  
PERFORMANCE OF PROPERTY AND NEIGHBOUR  
ATTENTION MECHANISMS.

Method	Acc.	F1.	Prec.	Rec.
avg-weight-property	92.55%	93.23%	93.55%	92.91%
avg-weight-neighbour	92.17%	93.02%	93.17%	92.88%
random-weight-property	90.12%	89.56%	89.49%	89.63%
random-weight-neighbour	91.86%	92.16%	92.55%	91.77%
property+neighbour attention	<b>93.74%</b>	<b>93.69%</b>	<b>93.75%</b>	<b>93.64%</b>

TABLE X  
PERFORMANCE OF USING DIFFERENT GRAPH NEURAL  
NETWORK (GNN) MODELS.

Method	Acc.	F1.	Prec.	Rec.
GCN	92.26%	92.69%	92.86%	92.54%
GraphSAE	92.49%	91.16%	90.89%	91.43%
GAT-inspired (property+neighbour attention)	<b>93.74%</b>	<b>93.69%</b>	<b>93.75%</b>	<b>93.64%</b>

than those of the other two well-known GNN models (GCN [79] and GraphSAGE [80]), which are commonly used for neighbour information aggregation. As illustrated in Table X, after replacing the neighbour attention module with the GCN method or GraphSAGE method, there are obvious declines about more than 1.25%, in accuracy and 1% in F1-measure on the Weibo dataset, since the neighbour attention mechanism can derive the different influence from one's different neighbour users.

#### E. Contributions of Social Neighbours and Personal Suicide-Oriented Knowledge Graph to Suicidal Ideation Detection

The effectiveness of the personal suicide-oriented knowledge graph prompted us to further investigate the impact of different personal factors upon suicidal ideation detection. We used *information gain* to assess the change of information amount involving a certain personal factor (property) based on *entropy* and *conditional entropy*. A bigger information gain implies a bigger significance of the factor.

Formally, let  $F$  be a category, and  $f$  be a vector value in the domain of category  $F$ , i.e.,  $f \in \text{Dom}(F)$ . Let  $y \in \{y_1, y_0\}$  denote the detection result (with/without suicidal ideation).  $\text{Prob}(\cdot)$  represents probability and  $H(\cdot)$  is the entropy.

$$\text{InfoGain}(y|F) = H(y) - H(y|F)$$

$$\text{where } H(y) = - \sum_{y \in \{y_1, y_0\}} \text{Prob}(y) \log \text{Prob}(y),$$

$$H(y|F) = \sum_{f \in \text{Dom}(F)} \text{Prob}(F=f) H(y|F=f),$$

$$H(y|f) = \sum_{y \in \{y_1, y_0\}} \text{Prob}(y|F=f) \log \text{Prob}(y|F=f).$$

We firstly considered the six categories of personal properties (*Personal Information*, *Personality*, *Experience*, *Post Behaviors*, *Emotion Expression*, and *Social Interaction*) in one's personal suicide-oriented knowledge graph, and then delved into the concrete properties of the key categories.

To ease the computation, for a property which takes continuous values (like *Age*, *Pursuing Perfection*, *Ruminant Thinking*, *Interpersonal Sensitive*, *Stress Duration*, *Stress Level*, *Stress Category*, *Suicide Words*, *Last Words*, *Future Words*, *Negation*

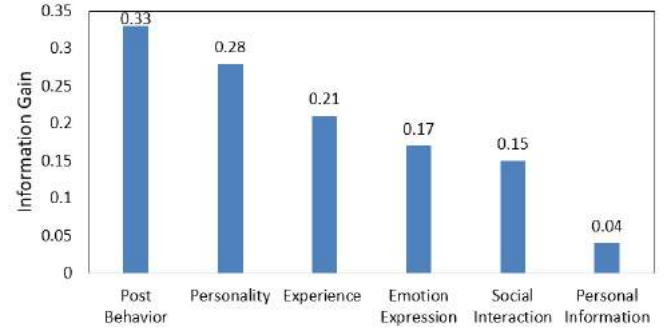


Fig. 7: Information gains of the top-3 categories.

*Words*, *Self-Concerns Words*, *Following Number*, *Follower Number*, and *Mention/Forward/Agree Number*), we used the mean to divide all its values into two classes, and transformed the domain of the property into  $\{0, 1\}$ . The domains of discrete property values (like *Location*, *Gender* and *Emotion Transition*) remained unchanged.

We mapped the *Texts* property values under the *Post Behavior* category to three classes based on the emotional polarities given by SnowNLP (<https://github.com/isnowfy/snownlp>). Assume user  $u$  wrote a sequence of texts  $\text{Texts}(u)$  in his posts.

$$\text{Class}(\text{Texts}(u)) = \begin{cases} 0 & \text{if } EP(\text{Texts}(u)) \leq -0.3; \\ 1 & \text{if } -0.3 < EP(\text{Texts}(u)) < 0.3; \\ 2 & \text{if } EP(\text{Texts}(u)) \geq 0.3 \end{cases}$$

$$\text{where } EP(\text{Texts}(u)) = \frac{\sum_{t \in \text{Texts}(u)} \text{SnowNLP}(t)}{|\text{Texts}(u)|}.$$

In a similar way, we mapped the *Images* property values under the *Post Behavior* category to four classes based on brightness and proportion of warm color, calculated by RGB. Assume user  $u$  posted a sequence of images  $\text{Images}(u)$ .

$$\text{Class}(\text{Images}(u)) = \begin{cases} 0 & \text{if } B(\text{Images}(u)) < 0.5 \wedge W(\text{Images}(u)) < 0.5; \\ 1 & \text{if } B(\text{Images}(u)) < 0.5 \wedge W(\text{Images}(u)) \geq 0.5; \\ 2 & \text{if } B(\text{Images}(u)) \geq 0.5 \wedge W(\text{Images}(u)) < 0.5; \\ 3 & \text{if } B(\text{Images}(u)) \geq 0.5 \wedge W(\text{Images}(u)) \geq 0.5 \end{cases}$$

$$\text{where } B(\text{Images}(u)) = \frac{\sum_{i \in \text{Images}(u)} \text{RGB-Bright}(i)}{|\text{Images}(u)|}, \text{ and}$$

$$W(\text{Images}(u)) = \frac{\sum_{i \in \text{Images}(u)} \text{RGB-Warm}(i)}{|\text{Images}(u)|}.$$

Figure 7 shows the average information gain of each category in the detection of suicidal ideation. From the result presented, we can see that all the listed categories brought positive impact on suicidal detection, and the top-3 categories were *Post Behavior*, *Personality*, and *Experience*, whose information gains were 0.33, 0.28, and 0.21, respectively.

We further looked into all the categories, and examined the contributions from their respective properties. As shown in Figure 8, the top-5 key properties were *posted text*, *stress level*, *stress duration*, *image*, and *ruminant thinking*.

To be more illustrative in the detection of suicidal ideation, we conducted a further feature selection experiment. From the

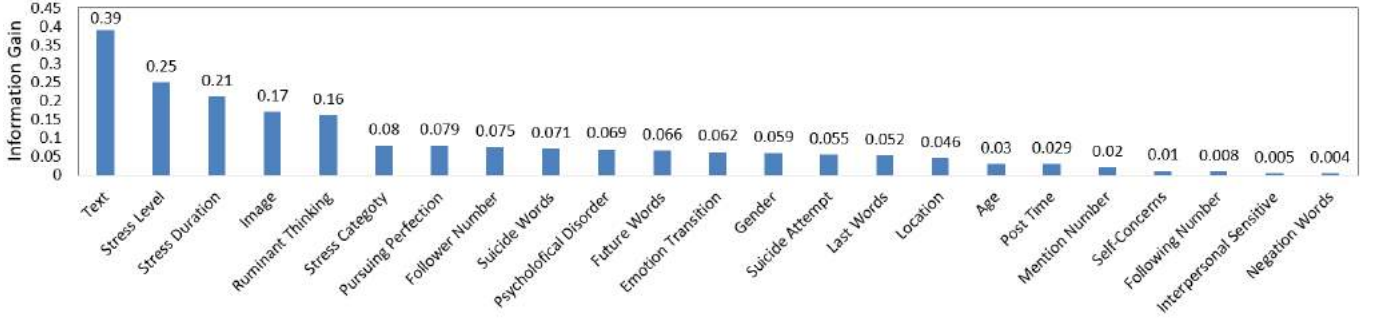


Fig. 8: Information gains of all the properties from personal suicide-oriented knowledge graph.

TABLE XI  
DERIVED FIVE TEST DATASETS CONTAINING DIFFERENT PROPORTIONS OF USERS WITH ANTI-REAL POSTS.

	$TD_1$	$TD_2$	$TD_3$	$TD_4$	$TD_5$
# Users with suicidal ideation	43	100	200	300	257
# with anti-real posts	43 (100%)	43 (43%)	43 (21.5%)	43 (14.3%)	0 (0%)
# without anti-real posts	0	57	157	257	257
# Ordinary users (random selection)	43	100	200	300	257

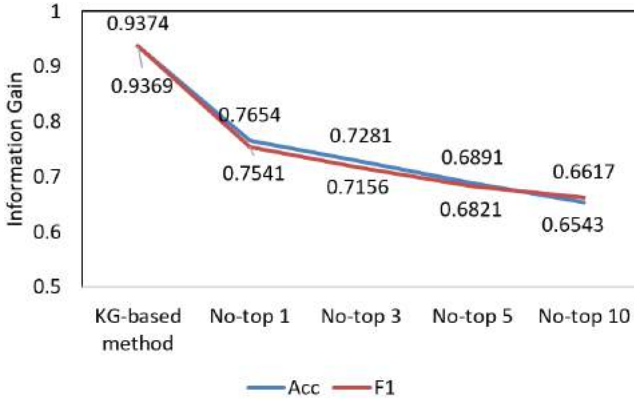


Fig. 9: Results of the feature selection experiment through removing the top- $x$  key properties in information gain (where  $x=1,3,5,10$ ).

result presented in Figure 9, we note that after removing the top- $x$  key properties in information gain (where  $x=1,3,5,10$ ), both accuracy and F1-measure kept declining, verifying the effectiveness of the key properties in suicidal risk analysis.

#### F. Impact of Data Noise on Suicidal Ideation Detection Performance

Due to the less restrictive and free-style nature of the social media, the presence of data noise is quite common, and may have a misleading effect on the experimental results. To examine the impact of data noise on the suicidal ideation detection performance, we consider two types of data noise: (1) **anti-real posts** (that is, users didn't fill in real information of themselves on the social media.) (2) **a few posts** (that is, users made only a few posts on the social media), and conducted experiments on a few datasets, each possessing different levels of data noise.

TABLE XII  
DETECTION PERFORMANCE ON THE FIVE DERIVED TEST DATASETS CONTAINING DIFFERENT PROPORTIONS OF USERS WITH ANTI-REAL POSTS.

Method	$TD_1$	$TD_2$	$TD_3$	$TD_4$	$TD_5$
CNN	65.75%	73.73%	82.97%	86.77%	90.12%
LSTM+Attention	67.27%	76.54%	84.23%	88.89%	92.12%
SDM	74.98%	82.31%	87.59%	91.35%	94.03%
Text+History+Graph	75.61%	83.59%	88.64%	91.56%	94.06%
KG-based method	<b>80.08%</b>	<b>88.11%</b>	<b>90.76%</b>	<b>93.74%</b>	<b>95.45%</b>
–On suicidal users	66.51%	82.33%	88.06%	93.97%	97.58%
–On ordinary users	93.65%	93.89%	93.46%	93.51%	93.51%

1) *Users with Anti-Real Posts*: Thanks to the existence of the tree hole, we could glimpse some anti-real sentiment expression from a user's normal open post as follows. For a user who made posts in both the tree hole and on the open microblog platform, if (1) he expressed despair, depression, anxiety, or suicidal ideation in the tree hole, but showed uplift feelings on the open posts within the 24 hours; and (2) the above situation happened more than twice within the latest one month, we annotated the user with anti-real posts.

Our test Sina Weibo dataset contains 300 users with suicidal ideation and 300 ordinary users. From the 300 users with suicidal ideation, we identified totally 43 users with anti-real posts. By varying the proportions of users with suicidal ideation in the test dataset, we derived 5 test datasets with different levels of data noise. In each derived test dataset (Table XI), as ordinary users made no posts in the tree hole, we randomly picked up the same number of ordinary users from the 300 ordinary test users as the control group.

Table XII lists the performance of all the methods on the five different test datasets. Our knowledge graph based method outperforms all the other methods in the four metrics. The less proportions of the users with anti-real posts, the higher

TABLE XIII  
DERIVED TEST DATASET CONTAINING USERS WITH LESS THAN 5 POSTS THROUGHOUT THE YEAR.

	$TD_6$
# Users with suicidal ideation	88
# Ordinary users	88

TABLE XIV  
DETECTION PERFORMANCE ON THE DERIVED TEST DATASET CONTAINING USERS WITH ONLY A FEW POSTS.

Method	Acc.	F1.	Prec.	Rec.
CNN	72.35%	72.38%	72.12%	72.65%
LSTM+Attention	74.56%	74.91%	74.98%	74.84%
SDM	85.49%	85.36%	85.46%	85.26%
Text+History+Graph	86.88%	86.52%	86.49%	86.55%
KG-based method	<b>89.73%</b>	<b>89.49%</b>	<b>89.92%</b>	<b>89.49%</b>
–On suicidal users	92.04%	90.00%	88.04%	92.04%
–On ordinary users	87.50%	89.53%	91.67%	87.50%

the detection performance. When the proportion of the users with suicidal ideation and anti-real posts dropped to 21.5%, 14.3%, and 0%, the detection accuracy of the knowledge graph based method on the suicidal group reached 88.06%, 93.97%, and 97.58%, respectively. The detection performance for the ordinary group remained unchanged, around 93.65% in accuracy and 93.89% in F1-measure.

2) *Users with Only a Few Posts*: From the total 600 test users (300 ordinary users and 300 users with suicidal ideation), we filtered out 88 users with suicidal ideation and another 88 ordinary users without suicidal ideation, and all of them had less than 5 posts throughout the year (Table XIII).

It is not surprising that the performance of all the trained models (in Table XIV) dropped, given such a small number of test posts. Among the five detection methods, our knowledge graph based method achieved the highest overall accuracy of 89.73% on all the users, 92.04% on the suicidal group, and 87.50% on the ordinary group.

3) *Experiments on a Small-Sized New Dataset*: Apart from examining the impact of data noise through the six datasets derived from our Weibo test dataset, we also built a small-sized new dataset, containing 85 suicide users (who had passed away as verified by news from 2012 to 2014) and another 85 ordinary users randomly sampled from our test dataset. It is expected that the data noise of this new dataset was less than that of our Sina Weibo dataset. Comparatively, suicide users in this new dataset had more supportive information.

(1) They made more posts (average 83.2 in their last year) than the suicidal ideation users in our Weibo dataset (average 69.2 in the latest one year).

(2) They experienced more stressful periods (average 6.0), showed more negative expressions (23% posts with negative emotions given by SnowNLP) in the last year than the suicidal ideation users (average 69.2 posts, 2.1 stressful periods, 12% posts with negative emotion in the latest one year).

With the help of the supportive information, our knowledge graph based method could deliver a very high performance,

TABLE XV  
PERFORMANCE ON THE NEW DATASET.

Method	Acc.	F1.	Prec.	Rec.
CNN	87.56%	87.65%	87.61%	87.69%
LSTM+Attention	88.86%	87.99%	88.23%	87.77%
SDM	91.78%	90.87%	90.28%	91.47%
Text+History+Graph	92.01%	91.62%	91.76%	91.48%
KG-based method	<b>94.53%</b>	<b>94.31%</b>	<b>94.45%</b>	<b>94.18%</b>

TABLE XVI  
RESULTS OF THE ABLATION STUDY ON THE SINA WEIBO DATASET.

Method	Acc.	F1.	Prec.	Rec.
KG-based method	<b>93.74%</b>	<b>93.69%</b>	<b>93.75%</b>	<b>93.64%</b>
Without-KG method	89.59%	89.21%	89.68%	88.74%

94.53% in accuracy and 94.33% in F1-measure, as shown in Table XV.

To sum up, the proposed knowledge graph based suicidal ideation detection model outperforms the other baseline methods, when facing challenging few-posts and anti-real-posts on microblog.

#### G. An Ablation Study

As the introduction of the user knowledge graph is essential in this study, we conducted an ablation study to see the performance with only the user's post behavior (Without-KG method). As shown in Table XVI, XVII there are huge declines about 4.15%, 7.58% in accuracy and 4.48%, 6.91% in F1-measure on both datasets. The mediocre performance of the Without-KG method introduced the effectiveness of personal suicide-oriented knowledge graph.

TABLE XVII  
RESULTS OF THE ABLATION STUDY ON THE REDDIT DATASET.

Method	Acc.	macro.F1.	Prec.	Rec.
KG-based method	<b>65.92%</b>	<b>65.93%</b>	<b>65.35%</b>	<b>66.21%</b>
Without-KG method	58.34%	59.02%	59.41%	58.63%

#### H. Training and Inference Time Costs

Table XVIII shows the training time costs of the five methods until convergence (getting the best result on the validation set) on a computer server with 4 NVIDIA GTX 1080Ti GPU. As our KG-based method utilized the most modalities (including text, image, social graph, and knowledge graph), its training time is the most. CNN requires the least training time, as it does not use the time-consuming sequential module. The inference time cost per batch (batch size is 16) of all the five methods are close except for CNN.

### V. DISCUSSIONS

#### A. Ethical Considerations

Keeping ethical considerations in mind is essential for the task of suicidal ideation detection. All the data (including

TABLE XVIII  
TRAINING AND INFERENCE TIME COSTS OF THE  
METHODS ON SINA WEIBO DATASET.

Method	Training Time	Inference Time(batch)
CNN (text)	56.7 mins	37.9 seconds
LSTM+Attention (text)	282.4 mins	49.6 seconds
SDM (text+image)	362.3 mins	55.4 seconds
Text+History+Graph (text+graph)	286.8 mins	50.8 seconds
KG-based method (text+image+graph+KG)	386.5 mins	56.4 seconds

name, age, posts, etc.) was crawled from the public social media, and was only used for research. We anonymized the data before labeling. There was no interaction or intervention with the subjects.

### B. Limitations

1) *Insufficiency of Microblog Data*: While the detection result is promising, the construction of users' personal suicide-oriented knowledge graphs, particularly personality-related factors, is still very preliminary in the current study, leaving a few important factors out of the study due to resource limitations. For instance, according to the psychological study [81], individual's *parenting rearing style* is a persistent factor that affects an individual's mental health. It is widely believed that individuals who are poorly educated are more likely to develop suicidal ideation than individuals who are actively educated. Parental rearing style is of great significance to the formation and development of individual's personality. Positive parental rearing style like emotional warmth and understanding, rather than negative excessive interference and excessive protection, contributes to one's healthy growth [82]. Completely relying on social media to analyze such personal knowledge is limited, and some other channels need to be explored to identify whether one was positively or negatively brought up.

2) *Noise of Microblog data*: Due to the less restrictive and free-style nature of social media like microblog, it is quite possible that someone may have a few more accounts, or may not want to fill in real information on the social media, such as Weibo. In this study, we empirically examined the impact of data noise on suicidal ideation detection performance. Two types of data noise were considered: (1) users didn't fill in real information of themselves on the social media; and (2) users made only a few posts on the social media. We conducted experiments on a few derived datasets and another new dataset, and the results were expected. That is, the lower data noise exhibits, the higher detection performance can be achieved. This raised a very interesting question: "*can we firstly take a look at the data first before feeding them to train a detection model?*" In addition, user identification and inference of users real thoughts and feelings will also be quite desirable, deserving further investigation.

## VI. CONCLUSION

In this paper, we built and used a personal suicide-oriented knowledge graph for suicidal ideation detection on social

media. A two-layered attention mechanism was deployed to explicitly reason and establish key risk factors to individuals' suicidal ideation. The performance study on 7,329 microblog users (3,652 with suicidal ideation and 3,677 without suicidal ideation) show that: 1) with the constructed personal knowledge graph, the social media-based suicidal ideation detection can achieve over 93% accuracy, outperforming the state-of-art approach; and 2) among the six categories of personal factors, *post*, *personality*, and *experience* were the top-3 key indicators. Under these categories, *posted text*, *stress level*, *stress duration*, *posted image*, and *ruminant thinking* contribute the most to one's suicidal ideation detection.

While the paper shows some promising results of the proposed method, leveraging social media to accurately identify personal properties is still limited. More reliable data sources and domain-specific expert knowledge need to be used and integrated into suicidal ideation detection. Dynamic maintenance of user-centric knowledge graph also deserves deep investigation for building really handy solutions.

## ACKNOWLEDGMENT

We thank all the anonymous reviewers' constructive comments, enabling us to improve the manuscript.

The work is supported by the National Natural Science Foundation of China (61872214, 61532015, 61521002).

## REFERENCES

- [1] M. Nock, G. Borges, E. Bromet, and et al., "Cross-national prevalence and risk factors for suicidal ideation, plans and attempts," *British Journal of Psychiatry*, vol. 192, no. 2, pp. 98–105, 2008.
- [2] N. Jacob, J. Scourfield, and R. Evans, "Suicide prevention via the Internet," *J. of Crisis*, 2014.
- [3] C. Bagge and A. Osman, "The suicide probability scale: Norms and factor structure," *Psychological reports*, vol. 83, no. 2, pp. 637–638, 1998.
- [4] K.-w. Fu, K. Y. Liu, and P. S. Yip, "Predictive validity of the chinese version of the adult suicidal ideation questionnaire: Psychometric properties and its short version," *Psychological Assessment*, vol. 19, no. 4, p. 422, 2007.
- [5] K. M. Harris, J.-J. Syu, O. D. Lello, Y. E. Chew, C. H. Willcox, and R. H. Ho, "The abc's of suicide risk assessment: Applying a tripartite approach to individual evaluations," *PLoS One*, vol. 10, no. 6, p. e0127442, 2015.
- [6] H. Christensen, P. Batterham, and B. O'Dea, "E-health interventions for suicide prevention," *International journal of environmental research and public health*, vol. 11, no. 8, pp. 8193–8212, 2014.
- [7] C. A. Essau, "Frequency and patterns of mental health services utilization among adolescents with anxiety and depressive disorders," *Depression and anxiety*, vol. 22, no. 3, pp. 130–137, 2005.
- [8] D. J. Rickwood, F. P. Deane, and C. J. Wilson, "When and how do young people seek professional help for mental health problems?" *Medical journal of Australia*, vol. 187, no. S7, pp. S35–S39, 2007.
- [9] H. D. Zachrisson, K. Rödje, and A. Mykletun, "Utilization of health services in relation to mental health problems in adolescents: a population based survey," *BMC public health*, vol. 6, no. 1, p. 34, 2006.
- [10] S. Ji, S. Pan, X. Li, E. Cambria, G. Long, and Z. Huang, "Suicidal ideation detection: A review of machine learning methods and applications," *IEEE Transactions on Computational Social Systems*, 2020.
- [11] A. Alambo, M. Gaur, U. Lokala, U. Kursuncu, K. Thirunarayan, A. Gyrrard, A. Sheth, R. S. Welton, and J. Pathak, "Question answering for suicide risk assessment using reddit," in *2019 IEEE 13th International Conference on Semantic Computing (ICSC)*. IEEE, 2019, pp. 468–473.
- [12] Q. Cheng, T. M. Li, C. L. Kwok, T. Zhu, and P. S. Yip, "Assessing suicide risk and emotional distress in chinese social media: a text mining and machine learning study," *Journal of Medical Internet Research*, vol. 19, no. 7, p. e243, 2017.
- [13] J. Du, Y. Zhang, J. Luo, Y. Jia, Q. Wei, C. Tao, and H. Xu, "Extracting psychiatric stressors for suicide from social media using deep learning,"



- BMC medical informatics and decision making*, vol. 18, no. 2, p. 43, 2018.
- [14] R. Sawhney, P. Manchanda, P. Mathur, R. Shah, and R. Singh, "Exploring and learning suicidal ideation connotations on social media with deep learning," in *Proc. 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, 2018, pp. 167–175.
  - [15] G. Coppersmith, R. Leary, P. Crutchley, and A. Fine, "Natural language processing of social media as screening for suicide risk," *Biomedical informatics insights*, vol. 10, p. 1178222618792860, 2018.
  - [16] M. J. Vioulés, B. Moulahi, A. J., and S. Bringay, "Detection of suicide-related posts in twitter data streams," *IBM Journal of Research & Development*, vol. 62, no. 1, pp. 7:1–7:12, 2018.
  - [17] J. Robinson, M. Rodrigues, S. Fisher, E. Bailey, and H. Herrman, "Social media and suicide prevention: findings from a stakeholder survey," *Shanghai Arch. Psychiatry*, vol. 27, no. 9, pp. 27–35, 2015.
  - [18] M. Gould, P. Jamieson, and D. Romer, "Media contagion and suicide among the young," *American Behavioral Scientist*, vol. 46, no. 9, pp. 1269–1284, 2003.
  - [19] M. Gaur, A. Alambo, J. P. Sain, U. Kursuncu, K. Thirunarayan, R. Kavuluru, A. Sheth, R. Welton, and J. Pathak, "Knowledge-aware assessment of severity of suicide risk for early intervention," in *The World Wide Web Conference*. ACM, 2019, pp. 514–525.
  - [20] L. Brisset, Y. Leanza, E. Rosenberg, B. Vissandjée, L. J. Kirmayer, G. Muckle, S. Xenocostas, and H. Laforce, "Language barriers in mental health care: A survey of primary care practitioners," *Journal of immigrant and minority health*, vol. 16, no. 6, pp. 1238–1246, 2014.
  - [21] L. Nie, Y.-L. Zhao, M. Akbari, J. Shen, and T.-S. Chua, "Bridging the vocabulary gap between health seekers and healthcare knowledge," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 2, pp. 396–409, 2015.
  - [22] G. Ning, Z. Zhang, and Z. He, "Knowledge-guided deep fractal neural networks for human pose estimation," *IEEE Transactions on Multimedia*, vol. 20, no. 5, pp. 1246–1259, 2017.
  - [23] F. Xue, R. Hong, X. He, J. Wang, S. Qian, and C. Xu, "Knowledge based topic model for multi-modal social event analysis," *IEEE Transactions on Multimedia*, 2019.
  - [24] C. Chaudhary, P. Goyal, D. N. Prasad, and Y.-P. P. Chen, "Enhancing the quality of image tagging using a visio-textual knowledge base," *IEEE Transactions on Multimedia*, vol. 22, no. 4, pp. 897–911, 2019.
  - [25] L. Nut and C. Nigel, "Normalising medical concepts in social media texts by learning semantic representation," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 2016, pp. 1014–1023.
  - [26] P. P. Sinha, R. Mishra, R. Sawhney, D. Mahata, R. R. Shah, and H. Liu, "#suicidal - a multipronged approach to identify and explore suicidal ideation in Twitter," in *Proc. 28th ACM International Conference on Information and Knowledge Management*, 2019.
  - [27] A. Mbarek, S. Jamoussi, A. Charfi, and A. B. Hamadou, "Suicidal profiles detection in Twitter," in *Proc. of the 15th Intl. Conf. on Web Information Systems and Technologies (WEBIST)*, 2019, pp. 289–296.
  - [28] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. [Online]. Available: <https://openreview.net/forum?id=rJXMpikCZ>
  - [29] J. P. Pestian, M. Sorter, B. Connolly, K. Bretonnel Cohen, C. McCullumsmith, J. T. Gee, L.-P. Morency, S. Scherer, L. Rohlf, and S. R. Group, "A machine learning approach to identifying the thought markers of suicidal subjects: a prospective multicenter trial," *Suicide and Life-Threatening Behavior*, vol. 47, no. 1, pp. 112–121, 2017.
  - [30] J. R. Crawford and J. D. Henry, "The depression anxiety stress scales (dass): Normative data and latent structure in a large non-clinical sample," *British journal of clinical psychology*, vol. 42, no. 2, pp. 111–131, 2003.
  - [31] J. D. Henry and J. R. Crawford, "The short-form version of the depression anxiety stress scales (dass-21): Construct validity and normative data in a large non-clinical sample," *British journal of clinical psychology*, vol. 44, no. 2, pp. 227–239, 2005.
  - [32] M. A. Just, L. Pan, V. L. Cherkassky, D. L. McMakin, C. Cha, M. K. Nock, and D. Brent, "Machine learning of neural representations of suicide and emotion concepts identifies suicidal youth," *Nature human behaviour*, vol. 1, no. 12, p. 911, 2017.
  - [33] K. M. Harris and M. T.-T. Goh, "Is suicide assessment harmful to participants? findings from a randomized controlled trial," *International journal of mental health nursing*, vol. 26, no. 2, pp. 181–190, 2017.
  - [34] S. R. Braithwaite, C. Giraud-Carrier, J. West, M. D. Barnes, and C. L. Hanson, "Validating machine learning algorithms for twitter data against established measures of suicidality," *JMIR mental health*, vol. 3, no. 2, p. e21, 2016.
  - [35] K. M. Harris, J. P. McLean, and J. Sheffield, "Suicidal and online: How do online behaviors inform us of this high-risk population?" *Death studies*, vol. 38, no. 6, pp. 387–394, 2014.
  - [36] J. Pestian, H. Nasrallah, P. Matykievicz, A. Bennett, and A. Leenaars, "Suicide note classification using natural language processing: A content analysis," *Biomedical informatics insights*, vol. 3, pp. BII–S4706, 2010.
  - [37] Y.-P. Huang, T. Goh, and C. L. Liew, "Hunting suicide notes in web 2.0-preliminary findings," in *Ninth IEEE International Symposium on Multimedia Workshops (ISMW 2007)*. IEEE, 2007, pp. 517–521.
  - [38] T. M. Li, B. C. Ng, M. Chau, P. W. Wong, and P. S. Yip, "Collective intelligence for suicide surveillance in web forums," in *Pacific-Asia Workshop on Intelligence and Security Informatics*. Springer, 2013, pp. 29–37.
  - [39] N. Masuda, I. Kurahashi, and H. Onari, "Suicide ideation of individuals in online social networks," *PloS one*, vol. 8, no. 4, p. e62262, 2013.
  - [40] M. De Choudhury, E. Kiciman, M. Dredze, G. Coppersmith, and M. Kumar, "Discovering shifts to suicidal ideation from mental health content in social media," in *Proceedings of the 2016 CHI conference on human factors in computing systems*. ACM, 2016, pp. 2098–2110.
  - [41] H.-C. Shing, P. Resnik, and D. W. Oard, "A prioritization model for suicidality risk assessment," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 8124–8137.
  - [42] N. Jones, N. Jaques, P. Pataranutaporn, A. Ghandeharioun, and R. Picard, "Analysis of online suicide risk with document embeddings and latent dirichlet allocation," in *2019 8th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW)*. IEEE, 2019, pp. 1–5.
  - [43] M. M. Tadesse, H. Lin, B. Xu, and L. Yang, "Detection of suicide ideation in social media forums using deep learning," *Algorithms*, vol. 13, no. 1, p. 7, 2020.
  - [44] J. Jashinsky, S. H. Burton, C. L. Hanson, J. West, C. Giraud-Carrier, M. D. Barnes, and T. Argyle, "Tracking suicide risk factors through twitter in the US," *Crisis*, 2014.
  - [45] H. Sueki, "The association of suicide-related twitter use with suicidal behaviour: a cross-sectional study of young internet users in Japan," *Journal of affective disorders*, vol. 170, pp. 155–160, 2015.
  - [46] F. Ren, X. Kang, and C. Quan, "Examining accumulated emotional traits in suicide blogs with an emotion topic model," *J. of biomedical and health informatics*, vol. 20, no. 5, pp. 1384–1396, 2015.
  - [47] L. Guan, B. Hao, and T. Zhu, "How did the suicide act and speak differently online? behavioral and linguistic features of china's suicide microblog users," *arXiv preprint arXiv:1407.0466*, 2014.
  - [48] L. Zhang, X. Huang, T. Liu, A. Li, Z. Chen, and T. Zhu, "Using linguistic features to estimate suicide probability of chinese microblog users," in *Proceedings of International Conference on Human Centered Computing*, 2015, pp. 549–559.
  - [49] X. Huang, X. Li, T. Liu, D. Chiu, T. Zhu, and L. Zhang, "Topic model for identifying suicidal ideation in chinese microblog," in *Proceedings of the 29th Pacific Asia Conference on Language, Information and Computation*, 2015, pp. 553–562.
  - [50] L. Guan, B. Hao, Q. Cheng, P. S. Yip, and T. Zhu, "Identifying chinese microblog users with high suicide probability using internet-based profile and linguistic features: classification model," *JMIR mental health*, vol. 2, no. 2, p. e17, 2015.
  - [51] B. O'Dea, S. Wan, P. J. Batterham, A. L. Calear, C. Paris, and H. Christensen, "Detecting suicidality on twitter," *Internet Interventions*, vol. 2, no. 2, pp. 183–188, 2015.
  - [52] G. Coppersmith, R. Leary, E. Whyne, and T. Wood, "Quantifying suicidal ideation via language usage on social media," in *Joint Statistics Meetings Proceedings, Statistical Computing Section, JSM*, 2015.
  - [53] J. Parraga-Alava, R. A. Caicedo, J. M. Gómez, and M. Inostroza-Ponta, "An unsupervised learning approach for automatically to categorize potential suicide messages in social media," in *2019 38th International Conference of the Chilean Computer Science Society*. IEEE, 2019, pp. 1–8.
  - [54] S. Fodeh, T. Li, K. Menczynski, T. Burgette, A. Harris, G. Ilita, S. Rao, J. Gemmell, and D. Raicu, "Using machine learning algorithms to detect suicide risk factors on twitter," in *2019 International Conference on Data Mining Workshops (ICDMW)*. IEEE, 2019, pp. 941–948.
  - [55] S. Ji, X. Li, Z. Huang, and E. Cambria, "Suicidal ideation and mental disorder detection with attentive relation networks," *arXiv preprint arXiv:2004.07601*, 2020.
  - [56] A. Malhotra and R. Jindal, "Multimodal deep learning based framework for detecting depression and suicidal behaviour by affective analysis of social media posts," *EAI Endorsed Trans. Pervasive*

- Health Technol.*, vol. 6, no. 21, p. e1, 2020. [Online]. Available: <https://doi.org/10.4108/eai.13-7-2018.164259>
- [57] H. A. Bouarara, "Detection and prevention of twitter users with suicidal self-harm behavior," *IJKBO*, vol. 10, no. 1, pp. 49–61, 2020. [Online]. Available: <https://doi.org/10.4018/IJKBO.2020010103>
- [58] R. Sawhney, P. Manchanda, R. Singh, and S. Aggarwal, "A computational approach to feature extraction for identification of suicidal ideation in tweets," in *Proceedings of the ACL Student Research Workshop*, 2018, pp. 91–98.
- [59] L. Cao, H. Zhang, L. Feng, Z. Wei, X. Wang, N. Li, and X. He, "Latent suicide risk detection on microblog via suicide-oriented word embeddings and layered attention," in *Proceedings of the 2019 conference on empirical methods in natural language processing*, 2019.
- [60] R. Mishra, P. P. Sinha, R. Sawhney, D. Mahata, P. Mathur, and R. R. Shah, "Snap-batnet: Cascading author profiling and social network graphs for suicide ideation detection on social media," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, 2019, pp. 147–156.
- [61] N. A. Cooperman and J. M. Simoni, "Suicidal ideation and attempted suicide among women living with hiv/aids," *Journal of behavioral medicine*, vol. 28, no. 2, pp. 149–156, 2005.
- [62] C. Steele, "The psychology of self-affirmation: sustaining the integrity of the self," *Advances in experimental social psychology*, vol. 21, no. 2, pp. 261–302, 1988.
- [63] A. R. Rich and R. L. Bonner, "Concurrent validity of a stress-vulnerability model of suicidal ideation and behavior: A follow-up study," *Suicide and Life - Threatening Behavior*, vol. 17, no. 4, pp. 265–270, 1987.
- [64] V. Blüml, N. D. Kapusta, S. Doering, E. Brähler, B. Wagner, and A. Kersting, "Personality factors and suicide risk in a representative sample of the german general population," *PLoS one*, vol. 8, no. 10, p. e76646, 2013.
- [65] T. S. Greenspon, "Is there an antidote to perfectionism?" *Psychology in the Schools*, vol. 51, no. 9, pp. 986–998, 2014.
- [66] M. Lv, A. Li, T. Liu, and T. Zhu, "Creating a chinese suicide dictionary for identifying suicide risk on social media," *PeerJ*, p. <https://peerj.com/articles/1455/>, 2015.
- [67] L. Mandelli, F. Nearchou, C. Vaiopoulos, C. Stefanis, S. Vitoratou, A. Serretti, and N. Stefanis, "Neuroticism, social network, stressful life events: Association with mood disorders, depressive symptoms and suicidal ideation in a community sample of women," *J. of Psychiatry Research*, pp. 38–44, 2014.
- [68] Q. Li, Y. Xue, L. Zhao, J. Jia, and L. Feng, "Analyzing and identifying teens stressful periods and stressor events from a microblog," *IEEE Journal of Biomedical and Health Informatics (J-BHI)*, vol. 21, no. 5, pp. 1–15, 2017.
- [69] G. Coppersmith, M. Dredze, and C. Harman, "Quantifying mental health signals in twitter," in *Proceedings of the workshop on computational linguistics and clinical psychology: From linguistic signal to clinical reality*, 2014, pp. 51–60.
- [70] H. Xiao, "bert-as-service," <https://github.com/hanxiao/bert-as-service>, 2018.
- [71] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [72] C. Du, H. Sun, J. Wang, Q. Qi, and J. Liao, "Adversarial and domain-aware bert for cross-domain sentiment analysis," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 4019–4028.
- [73] Y.-C. Chen, Z. Gan, Y. Cheng, J. Liu, and J. Liu, "Distilling knowledge learned in bert for text generation," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7893–7905.
- [74] Z.-J. Zha, J. Liu, D. Chen, and F. Wu, "Adversarial attribute-text embedding for person search with natural language query," *IEEE Transactions on Multimedia*, 2020.
- [75] P. Dai, H. Zhang, and X. Cao, "Deep multi-scale context aware feature aggregation for curved scene text detection," *IEEE Transactions on Multimedia*, vol. 22, no. 8, pp. 1969–1984, 2019.
- [76] J. Sabbath, "The suicidal adolescent: The expendable child," *J. of the American Academy of Child Psychiatry*, pp. 272–285, 1969.
- [77] M. Lv, A. Li, T. Liu, and T. Zhu, "Creating a chinese suicide dictionary for identifying suicide risk on social media," *PeerJ*, vol. 3, p. e1455, 2015.
- [78] C.-L. Huang, C. K. Chung, N. Hui, Y.-C. Lin, Y.-T. Seih, B. C. Lam, W.-C. Chen, M. H. Bond, and J. W. Pennebaker, "The development of the chinese linguistic inquiry and word count dictionary," *Chinese Journal of Psychology*, 2012.
- [79] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=SJU4ayYgl>
- [80] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, 2017, pp. 1024–1034. [Online]. Available: <http://papers.nips.cc/paper/6703-inductive-representation-learning-on-large-graphs>
- [81] N. Darling and L. Steinberg, "Parenting style as context: an integrative model," *Psychological Bulletin*, vol. 113, no. 3, pp. 487–496, 1993.
- [82] K. Lai and C. McBride-Chang, "Suicidal ideation, parenting style, and family climate among hong kong adolescents," *Int J. Psychology*, vol. 36, no. 2, pp. 81–87, 2001.



**Lei Cao** is a PhD candidate in the Department of Computer Science and Technology, Tsinghua University, Beijing, China. His research interests include computational psychology and sentiment analysis.



**Huijun Zhang** is a PhD candidate in the Department of Computer Science and Technology, Tsinghua University, Beijing, China. Her research interests include computational psychology and sentiment analysis.



**Ling Feng** is a professor of computer science and technology with Tsinghua University, China. Her research interests include computational mental healthcare, context-aware data management and services toward ambient intelligence, data mining and warehousing, and distributed object-oriented database systems.