

2022

## **KSAT: Knowledge-infused Self Attention Transformer - Integrating Multiple Domain-Specific Contexts**

Kaushik Roy

Yuxin Zi

Vignesh Narayanan

Manas Gaur

Amit P. Sheth

Follow this and additional works at: [https://scholarcommons.sc.edu/aii\\_fac\\_pub](https://scholarcommons.sc.edu/aii_fac_pub)



Part of the [Computer Engineering Commons](#), and the [Electrical and Computer Engineering Commons](#)

---

---

# KSAT: Knowledge-infused Self Attention Transformer - Integrating Multiple Domain-Specific Contexts

---

**Kaushik Roy**

Artificial Intelligence Institute South Carolina  
University of South Carolina  
kaushikr@email.sc.edu

**Yuxin Zi**

Artificial Intelligence Institute South Carolina  
University of South Carolina  
yzi@email.sc.edu

**Vignesh Narayanan**

Artificial Intelligence Institute South Carolina  
University of South Carolina  
vignar@sc.edu

**Manas Gaur**

Dept of Computer Science and Engineering  
University of Maryland, Baltimore County  
manas@umbdc.edu

**Amit Sheth**

Artificial Intelligence Institute South Carolina  
University of South Carolina  
amit@sc.edu

## Abstract

1 Domain-specific language understanding requires integrating multiple pieces of  
2 relevant contextual information. For example, we see both suicide and depression-  
3 related behavior (multiple contexts) in the text “*I have a gun and feel pretty bad*  
4 *about my life, and it wouldn’t be the worst thing if I didn’t wake up tomorrow*”.  
5 Domain specificity in self-attention architectures is handled by fine-tuning on  
6 excerpts from relevant domain specific resources (datasets and external knowl-  
7 edge - medical textbook chapters on mental health diagnosis related to suicide  
8 and depression). We propose a modified self-attention architecture *Knowledge-*  
9 *infused Self Attention Transformer* (KSAT) that achieves the integration of multiple  
10 domain-specific contexts through the use of external knowledge sources. KSAT  
11 introduces knowledge-guided biases in dedicated self-attention layers for each  
12 knowledge source to accomplish this. In addition, KSAT provides mechanics for  
13 controlling the trade-off between learning from data and learning from knowledge.  
14 Our quantitative and qualitative evaluations show that (1) the KSAT architecture  
15 provides novel human-understandable ways to precisely measure and visualize the  
16 contributions of the infused domain contexts, and (2) KSAT performs competitively  
17 with other knowledge-infused baselines and significantly outperforms baselines  
18 that use fine-tuning for domain-specific tasks.

## 19 1 Motivation

20 Solving domain-specific tasks such as mental health diagnosis (MHD), and triaging, requires in-  
21 tegrating relevant contextual information from data and knowledge sources. Self-Attention based  
22 Language Models (SAMs) capture an aggregated broader context from domain-agnostic, voluminous  
23 training corpora [1]. Fine-tuning SAMs on domain-specific corpora achieves domain-specific context  
24 capture [2, 3]. However, SAM architectures are black-box in nature [4]. Consequently, fine-tuned

25 SAM architectures do not lend themselves to the robust evaluation of the open research aims: **(R1)**  
26 Relevant domain-specific context coverage, and **(R2)** The influence of knowledge context traded-off  
27 against the data context in downstream tasks [5, 6]. We propose a modified self-attention architecture  
28 *Knowledge-infused Self Attention Transformer* (KSAT) to address these aims. KSAT performs well  
29 on select domain-specific tasks (see Task Description, Data, and External Knowledge Sources) while  
30 lending itself to a robust human-understandable evaluation of **R1** and **R2**. Thus KSAT provides a  
31 substantial step towards fostering AI-user trust, and satisfaction [7, 8].

## 32 **2 Background**

### 33 **2.1 Related Work**

34 Prior approaches that are relevant to **R1** and **R2** and incorporate multiple knowledge contexts can be  
35 broadly categorized based on the knowledge-infusion technique as **(1)** knowledge modulated SAMs  
36 and **(2)** knowledge infused input embedding-based SAMs [9, 10]. The former uses knowledge to guide  
37 the self-attention mechanism in SAMs, and the latter embeds the knowledge into a vector space before  
38 passing the inputs into SAMs. Here, we briefly summarize their contributions towards **R1** and **R2**.  
39 Both Category **(1)**, and Category **(2)** methods’ domain coverage is evaluated through performance  
40 on domain-specific task descriptions (**R1**). These methods’ ablations highlight contributions of  
41 knowledge context (**R2**). However, inspecting the numerical outputs from the model components  
42 (projection matrices and vectors) does not easily lend themselves to human-understandable scrutiny.  
43 Explainable AI techniques (post-processing of the numerical outputs that transform them into human-  
44 understandable information) are required to confirm the author(s) perspectives [11]. Post-processing-  
45 based explanations are local approximations of the SAM reasoning for particular inputs and therefore  
46 do not present the global picture, casting doubts on the SAM evaluation validity. KSAT presents  
47 a SAM architecture whose numerical outputs lend themselves to robust human-understandable  
48 evaluations of **R1** and **R2**.

### 49 **2.2 Task Description, Data, and External Knowledge Sources**

50 Although the KSAT architecture broadly applies to any domain-specific task, we choose the specific  
51 task of Mental Health Diagnostic Assistance for Suicidal Tendencies by Gaur et al. [12]. We  
52 denote this dataset as MHDA. The data contains high-quality expert annotations on Reddit posts  
53 from suicide-related subreddits. The annotation method ensures minimal noise from measurement  
54 artifacts and high agreement among the expert annotators. We use the clinically established diagnostic  
55 process information contained in the *Columbia Suicide Severity Rating Scale* (CSSRS) for knowledge  
56 contexts. Figure 1 (a, b) illustrate the various contexts (each tree path represents a context) under  
57 which suicidal patterns can arise. The task is to predict the suicidal patterns, namely - *indication*,  
58 *ideation1*, *ideation2*, and *behavior or attempt*. Figure 1 (c) shows examples from the MHDA dataset,  
59 augmented with knowledge context annotations. We denote the augmented dataset as k-MHDA.  
60 k-MHDA contains knowledge context annotations at the post and sentence level (see Figure 1 (c)).  
61 We defer construction details of k-MHDA from the CSSRS knowledge and the MHDA data to the  
62 appendix Section Constructing k-MHDA as it is not the main focus of the paper<sup>1</sup>.

## 63 **3 KSAT - Proposed Architecture**

64 For an input  $x$  in the vanilla SAM layer  $l$ , a  $CLS$  token is introduced that encodes the information  
65 in the data for downstream tasks. For example, for classification,  $Z_{CLS}$ , the representation of the  
66  $CLS$  token from the final layer  $L$  is passed through a softmax layer which outputs class probabilities.  
67 In the KSAT layer, we introduce a  $Z_{KCLS}$  token that encodes the knowledge for that layer. The  
68  $Z_{KCLS}$  values are determined by the graph context between the input  $x$  sentences and the concepts

---

<sup>1</sup>We will release the k-MHDA dataset, along with code to construct it along with the KSAT code for reproducibility of results.

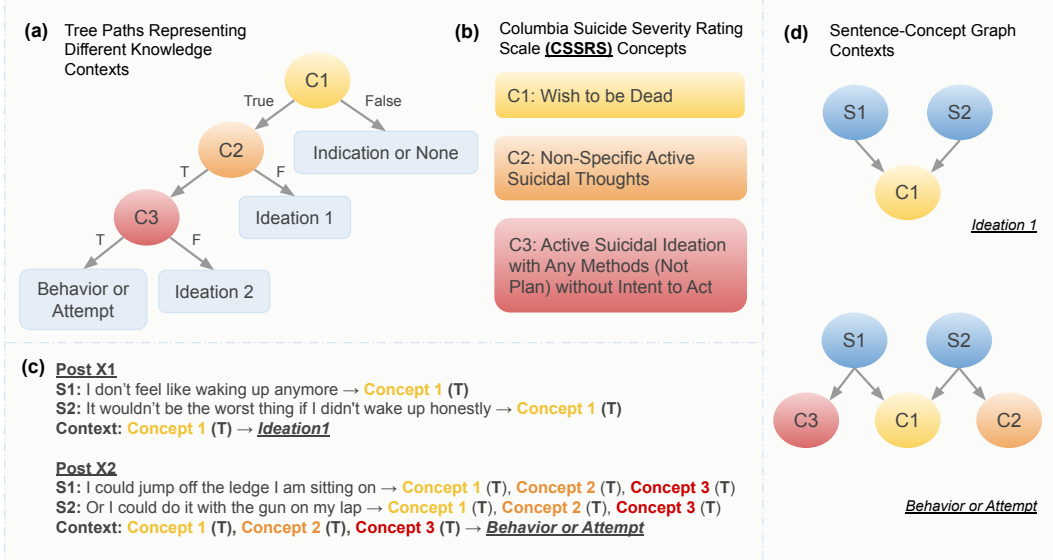


Figure 1: The knowledge source we utilize for our experiments is the *Columbia Suicide Severity Rating Scale (CSSRS)* and the data are posts from Reddit subreddits related to suicide. (a) Shows the knowledge contexts (tree-path represents a context) present in the CSSRS leading to varying suicidal patterns. (b) Shows the different concepts listed in the CSSRS. (c) Shows the posts and sentences from the data annotated with relevant knowledge concepts and contexts from the CSSRS. (d) Shows how sentences from the input posts (**Post X1** and **Post X2**) and concepts from the CSSRS are graphically connected in different contexts (In this case, the Ideation1 and the Behavior or Attempt contexts).

69 in the knowledge context annotations for input  $x$ . As an example, consider the inputs  $x$  to be like  
70 the posts **Post X1** and **Post X2** shown in Figure 1(c), then the graphs in Figure 1 (d), illustrate the  
71 graph contexts that show sentence-concept connections. One would expect that the  $Z_{KCLS}$  encoding  
72 for the sentences **S1** and **S2** in both posts would be similar as they have similar graph contexts. The  
73  $\mathbf{KG}_{\text{bias}}$  term ensures that the values in  $Z_{KCLS}$  captures this behavior (see Figure 2 (b)).

### 74 3.1 The KSAT layer

75 In contrast with the vanilla SAM layer, every KSAT layer has the following key differences:

#### 76 3.1.1 Layer Parameters and Outputs

77 Let  $Y$  denote the set of suicidal outcomes pertaining to each context:  $\{\text{Indication or None}, \text{Ideation1},$   
78  $\text{Ideation2}, \text{Behavior or Attempt}\}$  (see Figure 1 (a)). Let  $x$  denote an input post (see Figure 1 (c) for  
79 example posts). For each input  $x$ , every KSAT layer  $l$  outputs a vector of probabilities for every  
80 outcome  $y \in Y$  given by **Equation 1**:

$$P_l^{\text{ksat}}(y|x) = \sigma(W^T(\alpha_l Z_{KCLS}(x) + (1 - \alpha_l)Z_{CLS}(x)) + \mathbf{KG}_{\text{bias}}(x)) \quad (1)$$

81  $W$ ,  $Z_{CLS}(x)$ , and  $Z_{KCLS}(x)$  are of dimension size  $384 \times 1$  as we use the sentence-transformer  
82 published by Reimers et al. fine-tuned on the MHDA corpus, for embedding inputs [13]. Recall that  
83 the term  $\mathbf{KG}_{\text{bias}}$  is used to ensure that  $Z_{CLS}(x)$  encodes the knowledge context represented as a  
84 graph (details in Encoding Graph Contexts in Layers). The dimension size of  $\mathbf{KG}_{\text{bias}}$  is  $1 \times 1$ . The  
85 scalar term  $\alpha_l$  is used as a trade-off factor modulating the contributions between data and knowledge  
86 contexts using a convex combination (see Equation 1). We contrast KSAT layer outputs against the  
87 vanilla SAM layer, where a similar probability vector is output only in the final layer and not at  
88 every layer (see Figure 2). The parameters within the vanilla SAM layer (the query, key, and value  
89 projection matrices) used to compute the self-attention matrix are retained in the KSAT layer [1].  
90 Unlike in the vanilla SAM layer, where the parameters may or may not be shared across layers, the

91 parameters are not shared across different KSAT layers as different layers encode different knowledge  
 92 contexts.

### 93 3.1.2 Encoding Graph Contexts in Layers

94 Even though every KSAT layer  $l$  outputs a vector of probabilities for every outcome  $y \in \{\textit{Indication}$   
 95  $\textit{or None, Ideation1, Ideation2, Behavior or Attempt}\}$ , representing each context, it encodes only a  
 96 single context in  $Z_{KCLS}$ . We show how to compute the  $\mathbf{KG}_{\text{bias}}$  term by referencing the posts **Post**  
 97 **X1** and **Post X2** and the graph contexts *Ideation1* and *Behavior or Attempt* from Figure 1 (c,d).

98 **Computing Sentence-Concept Connection Vectors:** We first compute a sentence-concept connec-  
 99 tion vector for each sentence in the posts. For **Post X1**, both sentence **S1** and **S2** are connected to the  
 100 concept C1 in the concept set {C1,C2,C3} (see Figure 1 (b)). Therefore both their sentence-concept  
 101 connection vectors are computed as: [1, 0, 0]. Similarly for **Post X2**, the sentence-concept connection  
 102 vectors for both **S1** and **S2** are computed as: [1, 1, 1].

103 **Computing Graph Context Distances:** Recall that KSAT layers take as input a single post  
 104 ( $x$  in Equation 1). Denoting the sentence-concept vectors for a sentence **S1** as  $C_{S1}$ , to compute  
 105  $\mathbf{KG}_{\text{bias}}(x)$ , we first need to compute graph context distances for the posts **Post X1** and **Post X2**:  
 106  $d(C_{S1}, C_{S2})$ . If the posts had more than two sentences the graph context distances would include  
 107 all pairs (Eg:  $d(C_{S1}, C_{S2}), d(C_{S2}, C_{S3}), d(C_{S1}, C_{S3})$  for three sentences). The term  $d(C_{S_i}, C_{S_j})$   
 108 for a pair of sentences  $(S_i, S_j)$  in post  $x$  captures the graph context-based distance between the  
 109 sentences. Intuitively, sentences that have equivalent graph contexts should have  $d(S_i, S_j) = 0$ . We  
 110 use hamming distance in our experiments.

111 The  $\mathbf{KG}_{\text{bias}}$  term for an input  $x$  is thus given by Equation 2 as:

$$\mathbf{KG}_{\text{bias}}(x) = - \sum_{(S_i, S_j) \in x} \frac{(Z_{KCLS}(x)[S_i] - Z_{KCLS}(x)[S_j])^2}{d(S_i, S_j) + \epsilon} \quad (2)$$

112 The formulation for the  $\mathbf{KG}_{\text{bias}}(x)$  term in Equation 2 encourages the  $Z_{KCLS}$  representations of  
 113 sentences that are in the same graph context to be similar. The  $\epsilon$  term is to prevent dividing by zero  
 114 errors.

### 115 3.2 Aggregating KSAT Layer Outputs

116 Combining KSAT layer probabilities given by Equation 1 is application domain dependent. For  
 117 the suicidal outcomes in the set:  $\{\textit{Indication or None, Ideation1, Ideation2, Behavior or Attempt}\}$ ,  
 118 it is reasonable to expect that suicidal ideation (both *Ideation1* and *Ideation2*) precedes the act  
 119 of attempting suicide (*Behavior or Attempt*). In other application domains, the contexts could be  
 120 independent of each other.

121 In our experiments, we stack four KSAT layers corresponding to the outcomes *Indication or None*,  
 122 *Ideation1, Ideation2, Behavior or Attempt*, in that order (the order is derived from the tree structure  
 123 in Figure 1 (a)). Typically, for dependent probability outcomes,  $X$ , and  $Y$ , where  $X$  precedes  $Y$ , the  
 124 probability  $P(X, Y)$  would be modeled as  $P(Y | X)P(X)$ . However, since we stack KSAT layers  
 125 in a particular order, we use a product approximation ( $P(X, Y) = P(X)P(Y)$ ) as information is  
 126 propagated upwards through the KSAT layers. Thus, the final layer probabilities from the KSAT  
 127 layers is computed using Equation 3 as:

$$P_{\text{final}}^{\text{ksat}}(y|x) = \prod_l P_l^{\text{ksat}}(y|x), \quad (3)$$

128 where  $P_l^{\text{ksat}}(y|x)$  is given by Equation 1

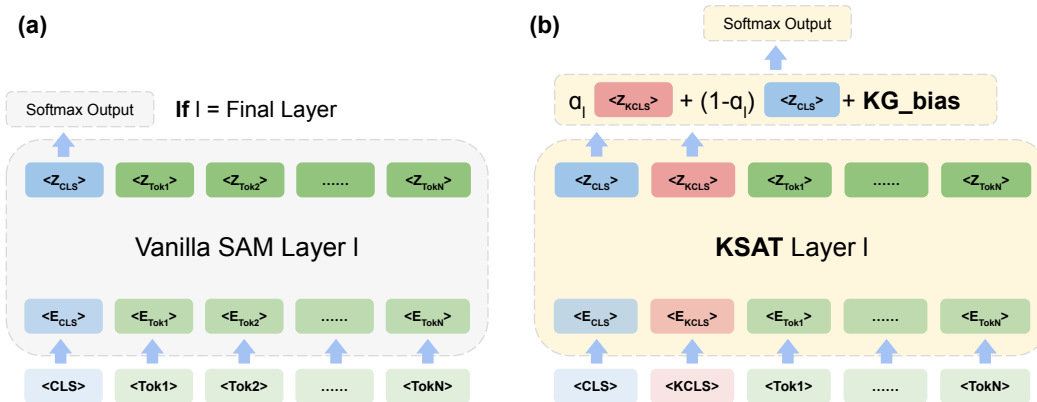


Figure 2: The comparative architecture of KSAT with existing vanilla SAMs. (a) Shows the vanilla SAM layer where the  $Z_{CLS}$  token from the final layer is passed as input to the softmax function for prediction. (b) Shows the introduction of the additional  $Z_{KCLS}$  token in a KSAT layer. The convex combination added with the knowledge bias as input to the softmax function is used to balance the use of data and knowledge contexts selectively.

## 129 4 KSAT - Results and Analysis

130 Recall the research aims that KSAT addresses when integrating multiple contexts - **(R1)** Relevant  
 131 domain-specific context coverage, and **(R2)** The influence of knowledge context traded-off against  
 132 the data context in downstream tasks (see **Section Motivation**).

### 133 4.1 KSAT - Quantitative Results (Addresses R1)

134 Table 1 shows the accuracy / AUC-ROC scores (rounded-of) for KSAT vs two best-performing  
 135 fine-tuned transformer models, and knowledge infused baseline models **K-type(i)** and **K-type(ii)** (see  
 136 section Related Work). Due to space concerns we describe details of the models **K-type(i)** and **K-**  
 137 **type(ii)** in the Appendix Section Construction of the **K-type(i)** and **K-type(ii)** baseline models. **KSAT**  
 138 outperforms the fine-tuned transformer models and performs comparably with models **K-type(i)** and  
**K-type(ii)** on the k-MHDA dataset (see section Constructing k-MHDA).

Dataset	KSAT	K-type(i)	K-type(ii)	XLNET	RoBERTa
k-MHDA	83% / %78	84% / 71%	84% / 72%	68% / 57%	68% / 63%

Table 1: Shows the accuracy / AUC-ROC scores (rounded-of) for KSAT vs fine-tuned transformer models, and knowledge infused baseline models **K-type(i)** and **K-type(ii)**. **KSAT** outperforms all fine-tuned transformer models and performs comparably with models **K-type(i)** and **K-type(ii)** on the k-MHDA dataset.

139

### 140 4.2 KSAT-Qualitative Results (Addresses R2)

141 Figure 3 illustrates the final KSAT layer representations (the  $Z_{KCLS}$  vectors) and  $\alpha_l$  values of sample  
 142 test posts to visualize data and knowledge contexts (see section **Layer Parameters and Outputs**).

## 143 5 Conclusion

144 We proposed KSAT that integrates multiple contexts and data and shows its utility in the domain-  
 145 specific use-case of MHDA. Although we test KSAT on MHDA, KSAT applies to other domain-  
 146 specific tasks that require contextualization from multiple knowledge sources. The architecture  
 147 of KSAT allows for precise measurement and visualization of the contributions from the different  
 148 knowledge contexts and the data, thus addressing the research aims **R1** and **R2**. In future work, we will  
 149 apply KSAT to other downstream domain-specific knowledge-intensive tasks such as conversational

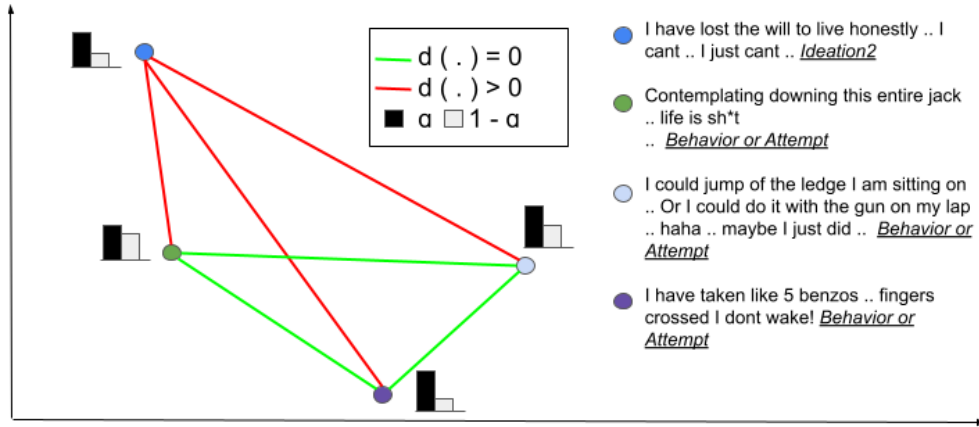


Figure 3: Shows that the first two *Behavior or Attempt* posts in the list of posts are far off based on the  $Z_{CLS}$  representations (*false*), where as KSAT representations depict them as closer i.e  $d(\cdot) < \epsilon$  (*true*). The  $\alpha$  shows the data and knowledge contributions for the samples.

150 question answering that requires drawing context from multiple knowledge sources. We will replicate  
 151 KSAT’s qualitative evaluation in MHDA for robust human-understandable evaluation on future tasks.

## 152 References

- 153 [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of  
 154 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,  
 155 2018.
- 156 [2] Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. How to fine-tune bert for text classifi-  
 157 cation? In *China national conference on Chinese computational linguistics*, pages 194–206.  
 158 Springer, 2019.
- 159 [3] Laila Rasmy, Yang Xiang, Ziqian Xie, Cui Tao, and Degui Zhi. Med-bert: pretrained contex-  
 160 tualized embeddings on large-scale structured electronic health records for disease prediction.  
 161 *NPJ digital medicine*, 4(1):1–13, 2021.
- 162 [4] Manas Gaur, Keyur Faldu, and Amit Sheth. Semantics of the black-box: Can knowledge graphs  
 163 help make deep learning systems more interpretable and explainable? *IEEE Internet Computing*,  
 164 25(1):51–59, 2021.
- 165 [5] Jerome R Bellegarda. Statistical language model adaptation: review and perspectives. *Speech*  
 166 *communication*, 42(1):93–108, 2004.
- 167 [6] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,  
 168 and Noah A Smith. Don’t stop pretraining: adapt language models to domains and tasks. *arXiv*  
 169 *preprint arXiv:2004.10964*, 2020.
- 170 [7] Amit Sheth, Manas Gaur, Kaushik Roy, and Keyur Faldu. Knowledge-intensive language  
 171 understanding for explainable ai. *IEEE Internet Computing*, 25(5):19–24, 2021.
- 172 [8] Amit Sheth, Manas Gaur, Kaushik Roy, Revathy Venkataraman, and Vedant Khandelwal.  
 173 Process knowledge-infused ai: Toward user-level explainability, interpretability, and safety.  
 174 *IEEE Internet Computing*, 26(5):76–84, 2022.
- 175 [9] Matthew E Peters, Mark Neumann, Robert L Logan IV, Roy Schwartz, Vidur Joshi, Sameer  
 176 Singh, and Noah A Smith. Knowledge enhanced contextual word representations. *arXiv preprint*  
 177 *arXiv:1909.04164*, 2019.

- 178 [10] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Guihong Cao, Daxin  
179 Jiang, Ming Zhou, et al. K-adapter: Infusing knowledge into pre-trained models with adapters.  
180 *arXiv preprint arXiv:2002.01808*, 2020.
- 181 [11] Feiyu Xu, Hans Uszkoreit, Yangzhou Du, Wei Fan, Dongyan Zhao, and Jun Zhu. Explainable  
182 ai: A brief survey on history, research areas, approaches and challenges. In *CCF international  
183 conference on natural language processing and Chinese computing*, pages 563–574. Springer,  
184 2019.
- 185 [12] Manas Gaur, Amanuel Alambo, Joy Prakash Sain, Ugur Kursuncu, Krishnaprasad Thirunarayan,  
186 Ramakanth Kavuluru, Amit Sheth, Randy Welton, and Jyotishman Pathak. Knowledge-aware  
187 assessment of severity of suicide risk for early intervention. In *The world wide web conference*,  
188 pages 514–525, 2019.
- 189 [13] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-  
190 networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language  
191 Processing*. Association for Computational Linguistics, 11 2019.
- 192 [14] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko.  
193 Translating embeddings for modeling multi-relational data. *Advances in neural information  
194 processing systems*, 26, 2013.
- 195 [15] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
196 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
197 approach. *arXiv preprint arXiv:1907.11692*, 2019.

## 198 A Constructing k-MHDA

199 There are 500 Reddit posts in the MHDA dataset. The knowledge contexts in the CSSRS can be  
200 illustrated as a tree (see Figure 1 (d)). We can construct a probabilistic decision tree that takes input  
201 post  $x$  and outputs an outcome  $y$  from among the leaves. We can write the tree in algebraic form as  
202 shown in **Equation 4**.

$$P(y | x, \{\theta_i\}) = \sum_{y \in \text{Outcomes}} p_y \prod_{i=1}^3 \sum_{x_{sub} \in x} \left( \cos\_sim \left( x_{sub}^R, q_i^R \right) \geq \theta_i \right) \geq 0.5 \quad (4)$$

203  $p_y$  is the ground truth probability for each outcome. Index  $i$  iterates through the 3 concepts in  
204 Figure 1 (d).  $x_{sub}$  denotes a sub-fragment of the input post (1 sentence, 2 sentence, etc.).  $q_i$  denotes  
205 the concept texts from the 3 concepts that  $i$  indexes.  $x_{sub}^R$  and  $q_i^R$  are representations of the post  
206 sub-fragment and the concept texts using the sentence-transformer published by Reimers et al. [13].

207 **Equation 5** determines the presence or absence of concept  $q_i$  in a post sub-fragment  $x_{sub}$ . First, we  
208 compute the cosine similarity between their sentence-transformer representations  $x_{sub}^R$  and  $q_i^R$ . If the  
209 resulting value is  $\geq \theta_i$ , we determine that the concept  $q_i$  is present in  $x_{sub}$ , else we determine that the  
210 concept  $q_i$  is absent in  $x_{sub}$ .

211  $\sum_{x_{sub} \in x} (\cdot) \geq 0.5$  in Equation 4 is the algebraic form of the  $\vee$  operation as we determine that concept  
212  $q_i$  is present in the post  $x$ , if any of the post fragments  $x_{sub} \in x$  show presence of concept  $q_i$ .

$$\left( \cos\_sim \left( x_{sub}^R, q_i^R \right) \geq \theta_i \right) \quad (5)$$

213 We can then evaluate the Bernoulli Loss  $\mathcal{L}$  given an input, outcome pair  $(x, y)$  and parameters  $\{\theta_i\}$   
214 as:

$$\mathcal{L}(x, y, \{\theta_i\}) = P(y | x, \{\theta_i\}) \log(P(y | x, \{\theta_i\})) + (1 - P(y | x, \{\theta_i\})) \log(1 - P(y | x, \{\theta_i\})) \quad (6)$$



215 We use grid-search to find a configuration of parameters  $\{\theta_i\}$  and post sub-fragment  $x_{sub}$  that has the  
216 maximum value for  $\prod_{(x,y) \in \text{MHDA}} \mathcal{L}(x, y, \{\theta_i\})$ . We vary each individual  $\theta_i$  in the range  $-1$  to  $1$   
217 (the *range of the cosine function*) and  $x_{sub}$  takes values from the set  $\{1, 2, 3\}$ .

218 Inference is carried out as it is in a decision tree classifier with the concept presence or absence at  
219 each branch, evaluated using **Equation 5**.

220 **Knowledge Context Annotation with outputs from grid-search:** The grid-search yielded outputs  
221  $\{\theta_i\} = \{0.3, 0.5, 0.3\}$ , and post sub-fragment size  $|x_{sub}| = 1$  (one sentence). Therefore the post  
222 “I don’t feel like waking up and have a gun. Oh well.” is annotated with the knowledge context:  
223 **(Concept 1 (T)), Concept 2 (T), Concept 3 (T) = Behavior or Attempt**, as evaluation of **Equation**  
224 **5** determines absence of **Concept 1**, **Concept 2**, and **Concept 3** in the post sentence “I don’t feel  
225 like waking up and have a gun”. The evaluation uses the grid-search outputs of  $\{\theta_i\}$ . The second  
226 sentence “Oh well” is not necessary to evaluate as we determine a concept’s presence or absence in  
227 the post if any of the post fragments  $x_{sub}$  (one sentence) show the presence of the concept.

## 228 **B Construction of the K-type(i) and K-type(ii) baseline models**

229 For our baseline implementations we adapt the state-of-the-art models for each of the model types  
230 **K-type(i)** and **K-type(ii)** in Section 2.1 for our task description (see Section 2.2. Specifically, we  
231 use the model KnowBERT and K-Adapter as **K-type(i)** and **K-type(ii)** model instances respectively  
232 [9, 10]. For knowledge graph embeddings in the KnowBERT case, we utilize TransE embeddings.  
233 Similar to their work, we use the RoBERTa model for the task-specific adapter module. We retain all  
234 hyperparameters from the original implementations [14, 15].