

XAI-CLASS: Explanation-Enhanced Text Classification with Extremely Weak Supervision

Anonymous ACL submission

Abstract

Text classification aims to effectively categorize documents into pre-defined categories. Traditional methods for text classification often rely on large amounts of manually annotated training data, making the process time-consuming and labor-intensive. To address this issue, recent studies have focused on weakly-supervised and extremely weakly-supervised settings, which require minimal or no human annotation, respectively. In previous methods of weakly supervised text classification, pseudo-training data is generated by assigning pseudo-labels to documents based on their alignment (e.g., keyword matching) with specific classes. However, these methods ignore the importance of incorporating the explanations of the generated pseudo-labels, or *saliency* of individual words, as additional guidance during the text classification training process. To address this limitation, we propose XAI-CLASS, a novel explanation-enhanced extremely weakly-supervised text classification method that incorporates word saliency prediction as an auxiliary task. XAI-CLASS begins by employing a multi-round question-answering process to generate pseudo-training data that promotes the mutual enhancement of class labels and corresponding explanation word generation. This pseudo-training data is then used to train a multi-task framework that simultaneously learns both text classification and word saliency prediction. Extensive experiments on several weakly-supervised text classification datasets show that XAI-CLASS outperforms other weakly-supervised text classification methods significantly. Moreover, experiments demonstrate that XAI-CLASS enhances both model performance and explainability.

1 Introduction

Text classification is a fundamental task in natural language processing (NLP), aiming to effectively categorize documents (e.g., news reports) into pre-defined categories (e.g., politics, sports, and busi-

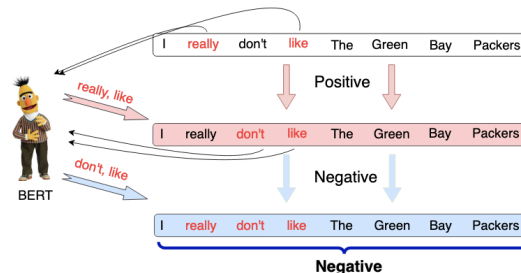


Figure 1: Previous weakly-supervised text classification methods do not model salient words, potentially leading to uncertain predictions. On the other hand, XAI-CLASS generates pseudo-text classification and pseudo-saliency labels by querying two pre-trained language models (PLMs) and updating pseudo-saliency labels by using previously generated pseudo-text classification labels and vice-versa.

ness). It has various downstream applications such as information extraction (Zhang et al., 2022), sentiment analysis (Tang et al., 2015), and question answering (Rajpurkar et al., 2016).

Traditional methods for text classification (Yang et al., 2016, 2019; Zhang et al., 2015) often rely on large amounts of manually annotated training data, making the process time-consuming and labor-intensive. To address this issue, recent studies have focused on weakly-supervised (Chang et al., 2008; Song and Roth, 2014; Gabilovich and Markovitch, 2007; Badene et al., 2019; Ratner et al., 2017; Meng et al., 2018; Mekala and Shang, 2020; Agichtein and Gravano, 2000; Shu et al., 2020; Tao et al., 2018) and extremely weakly-supervised (Meng et al., 2020b; Mekala and Shang, 2020; Wang et al., 2021; Zeng et al., 2022; Zhang et al., 2021) settings, which require minimal or no human annotation, respectively. In this study, we focus on the extremely weakly-supervised setting that utilizes only the class names as supervision. Importantly, we do not assume that the class names need to have appeared in the input documents.

067	Previous methods for extremely weakly-	high confidence in both the predicted class labels	118
068	supervised text classification usually start with	and the saliency words. The resulting pseudo-	119
069	finding initial keywords for each class to construct	training data incorporates both the class labels and	120
070	a keyword vocabulary. This vocabulary is then	the associated explanation words. This pseudo-	121
071	employed to assign pseudo-labels to documents,	training data is then used to train a multi-task frame-	122
072	followed by training the model using traditional	work that simultaneously learns both text classifi-	123
073	supervised learning techniques. For example, LOT-	cation and word saliency prediction. By jointly	124
074	Class (Meng et al., 2020b) leverages a pre-trained	optimizing both tasks, the model can effectively	125
075	masked language model to predict keywords that	enhance both the performance and explainability	126
076	can replace label words. However, this method	of the text classification model. Our contributions	127
077	assumes that the class names must appear in the	are summarized as follows:	128
078	input document, which may not be feasible in		
079	many real-world scenarios. Recent advancements	• We propose XAI-CLASS, a novel extremely	129
080	have relaxed this constraint and do not assume that	weakly-supervised text classification method that	130
081	the class names need to have appeared in the input	leverages multiple-round question answering to	131
082	documents. For example, X-Class (Wang et al.,	promote mutual enhancement between text clas-	132
083	2021) obtains the word and document representa-	sification and word saliency prediction pseudo-	133
084	tions and employs clustering methods for keyword	training data generation.	134
085	grouping and label assignment, while WDDC		
086	(Zeng et al., 2022) applies cloze-style prompting	• We propose a novel explanation-enhanced text	135
087	to identify keywords and assigns pseudo-labels	classification method that trains a multi-task	136
088	based on the representation similarity between the	framework to simultaneously learn both text clas-	137
089	keywords and the documents. However, previous	sification and word saliency prediction.	138
090	methods ignore the importance of incorporating		
091	the explanations of the generated pseudo-labels,	• Experiments on several datasets demonstrate the	139
092	or <i>saliency</i> (Simonyan et al., 2014) of individual	superiority of XAI-CLASS over previous weakly-	140
093	words, as additional guidance during the text	supervised text classification methods for both	141
094	classification training process (Figure 1). This	performance and explainability.	142
095	oversight has limited the potential of these methods		
096	to fully exploit the valuable insights provided by	We will open-source our code and results as a base-	143
097	explanations and word saliency that can greatly	line to facilitate future studies.	144
098	enhance the effectiveness and explainability of the		
099	text classification methods.	2 Related Work	145
		2.1 Text Classification Methods	146
100	To address this limitation, we propose XAI-	Traditional methods for text classification (Yang	147
101	CLASS, a novel explanation-enhanced extremely	et al., 2016, 2019; Zhang et al., 2015) often rely on	148
102	weakly-supervised text classification method that	large amounts of manually annotated training data,	149
103	incorporates word saliency prediction as an aux-	making the process time-consuming and labor-	150
104	iliary task. XAI-CLASS begins by employing a	intensive. To address this issue, recent work has	151
105	multi-round question-answering process to gener-	been proposed for text classification with minimal	152
106	ate pseudo-training data that promotes the mutu-	human annotation.	153
107	al enhancement of class labels and correspond-		
108	ing explanation word generation. Specifically, we	Weakly-Supervised Text Classification To ad-	154
109	first leverage a pre-trained multi-choice question-	dress the above issue of manual annotation, recent	155
110	answering model (Chung et al., 2022) to query the	studies have focused on the weakly-supervised set-	156
111	predicted class labels for given documents. Using	ting that requires minimal human annotation. For	157
112	the predicted class labels as input, we then query	example, Snowball (Agichtein and Gravano, 2000)	158
113	a pre-trained extractive question-answering model	combines pattern-based and distant supervision	159
114	(Devlin et al., 2018) to identify the tokens in the	techniques to extract relations. It uses patterns	160
115	document that were most influential in predicting	based on syntactic dependencies and entity men-	161
116	the class labels. This iterative process continues	tions to identify potential relations in sentences.	162
117	until the predictions remain consistent, indicating	However, this pattern-based approach may struggle	163
		with complex relations involving multiple entities	164

165	or deeper semantic understanding. Dataless (Chang et al., 2008) proposes a classification method using semantic representation. It leverages external knowledge sources to capture the semantic information in the text. However, the limitation is its dependence on the availability and quality of external knowledge sources. Doc2cube (Tao et al., 2018) clusters similar documents and assigns them to text cubes. It leverages the inherent structure and patterns within the collection for guidance. However, the effectiveness of Doc2Cube depends on the quality of document similarity measures used for clustering. Inaccurate or inadequate similarity metrics can impact document allocation accuracy.	
179	Extremely Weakly-Supervised Text Classification Compared with weakly-supervised text classification, extremely weakly supervised text classification goes a step further by using even weaker supervision or no labeled data during training. For example, LOTClass (Meng et al., 2020b) consists of three steps: substituting label names to enable the model to understand the meaning of each label, identifying category-relevant words for word-level classification, and finally conducting generalized self-training. Conwea (Mekala and Shang, 2020) utilizes contextualized word representations generated by PLMs to capture the rich semantic information of words in context for label assignment. XClass (Wang et al., 2021) expands label words and generates document representations based on BERT (Devlin et al., 2018) for clustering and the best documents are selected to train the classifier. WDDC (Zeng et al., 2022) uses cloze-style completion to generate summary text words, which serve as supervised signals for training the document classifier. However, these methods all have high requirements for the frequency of occurrence of labels and their closely related words in the text. ClassKG (Zhang et al., 2021) constructs a keyword graph by extracting important keywords from the documents, which serves as a representation of the document collection. Then ClassKG utilizes the connectivity and similarity of keywords in the graph to train the model. However, the efficiency and scalability of the method can be a concern when dealing with large-scale datasets.	
211	2.2 Explainable Text Classification	
212	Explainable text classification methods can be decomposed into two categories: post-hoc explainability and intrinsic explainability.	
	Post-hoc Explainability Post-hoc explainability explain inputs <i>after</i> a model has already been trained. This category consists of perturbation methods, such as LIME (Ribeiro et al., 2016), which learns an interpretable model of points in the neighborhood of a given input. Post-hoc explainability techniques can also be categorized by backpropagation-based methods. For example, Simonyan et al. attempts to explain instances by introducing the concept of saliency maps, which calculate gradients of inputs with respect to the inputs’ features. Kindermans et al. extends this idea by computing the partial derivatives of the prediction with respect to the input and multiplies them with the input (Ancona et al., 2017).	215 216 217 218 219 220 221 222 223 224 225 226 227 228 229
	Intrinsic Explainability In contrast to post-hoc explainability, intrinsic explainability methods attempt to create models that offer explanations. This has been accomplished through a handful of measures, one of which being constraining features (Freitas, 2014) to be sparse and by measuring feature sensitivity (Simonyan et al., 2014). XAI-CLASS aligns with this class of explainable text classification, as we generate and inject saliency information in our framework directly.	230 231 232 233 234 235 236 237 238 239
	3 Methodology	240
	We propose XAI-CLASS, an explanation-enhanced extremely weakly-supervised text classification method. The XAI-CLASS framework (Figure 2) consists of two major steps: (1) iterative pseudo-label generation, and (2) explainable multi-task learning. In this section, we describe XAI-CLASS framework in detail.	241 242 243 244 245 246 247
	3.1 Preliminaries	248
	Problem Formulation Our framework operates under the extremely weakly supervised text classification scenario, whose goal is to predict the correct class of a document with only its contents and the possible classes it could be categorized into. Mathematically, we represent a corpus as \mathcal{X} which contains documents $\mathcal{D} = \{t_i \forall i \in [1, \mathcal{D}]\}$ made up of tokens t_i . The set of all labels is denoted by $\mathcal{Y} = \{y_i \forall i \in [1, \mathcal{Y}]\}$.	249 250 251 252 253 254 255 256 257
	Saliency Representation XAI-CLASS employs salient tokens of a given document to identify which parts of the input should be attended to. We represent the set of all salient tokens of an input	258 259 260 261

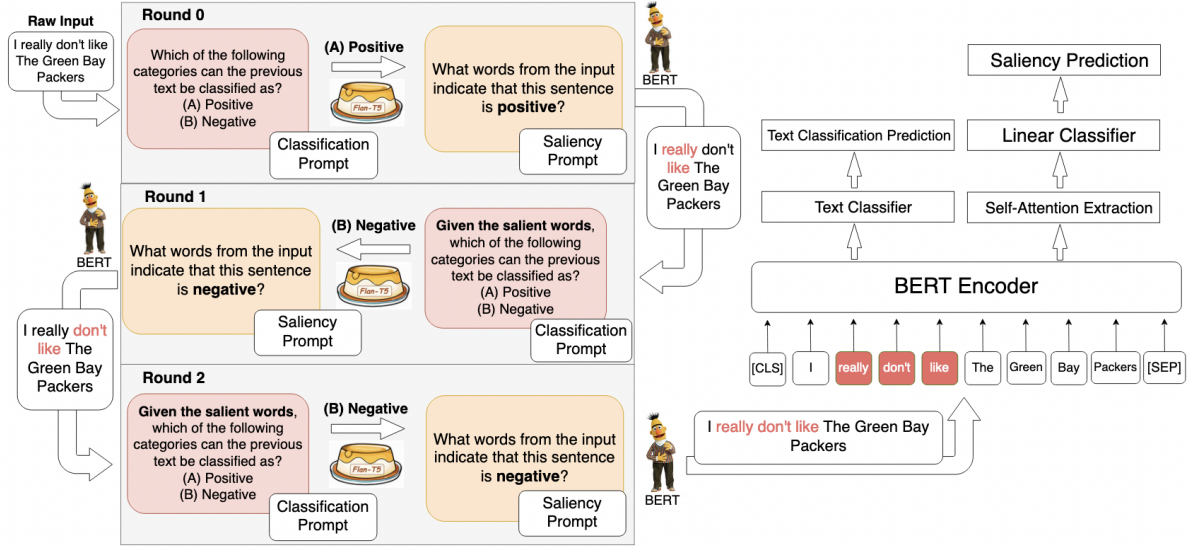


Figure 2: XAI-CLASS architecture. (Left) Given an input document \mathcal{D} ("I really don't like The Green Bay packers"), we first query the class prediction from a PLM \mathcal{T}^C (FLAN-T5) and then query the indicative words (highlighted in red) from another PLMs \mathcal{T}^E (BERT), forming our initial setup. We introduce the notion of a *round*, where we once again query \mathcal{T}^C using the queried indicative words and use this more confident prediction to query the salient words from \mathcal{T}^E once more. We repeat this operation until a variable number of rounds. (Right) We then tokenize \mathcal{D} and feed this along with the salient tokens into our BERT-based multi-task learning model, learning to predict both text classification and saliency labels using the contextualized representations.

document as $\mathcal{E} = \{t_i | \forall i \in [1, |\mathcal{E}]]\}$ (Simonyan et al., 2014), where token t_i is salient.

The XAI-CLASS framework is depicted in Figure 2, which incorporates both input text and saliency representations to learn contextualized mappings that are mapped to both text and saliency classifiers.

3.2 Iterative Pseudo-Label Generation

Pseudo-Text Classification Label Generation

Using a PLM \mathcal{T}^C , we first derive pseudo-text classification labels automatically using only input text. For example, given the sentence "I really don't like The Green Bay packers" in Figure 2, we feed this sentence through \mathcal{T}^C to determine the appropriate classification label (in this case, negative sentiment). We formally define this query process using \mathcal{D} as the input document to generate a pseudo-text classification label \hat{y}^T below:

$$\hat{y}^C = \mathcal{T}^C(\mathcal{D}). \quad (1)$$

Pseudo-Explanation Label Generation It is possible that \mathcal{T}^C may not produce confident predictions. For instance, \mathcal{T}^C may classify the example sentence in Figure 2 as positive sentiment because of the words "really" and "like", disregarding the phrase "don't like". To further enhance these

pseudo-text classification label predictions, we utilize another PLMs \mathcal{T}^E that captures the reasoning of \mathcal{T}^C ; namely, identifying the salient tokens in the input that were responsible for the pseudo-text classification label.

Formally, for a given input document \mathcal{D} and previously generated pseudo-text classification label \hat{y}^C , we query \mathcal{T}^E to determine the salient tokens based on the predicted label:

$$\hat{y}^E = \mathcal{T}^E(\mathcal{D}, \hat{y}^C), \quad (2)$$

where \hat{y}_i^E is a binary vector with cardinality $|\mathcal{D}|$ that's formulated based on the following equation:

$$\begin{cases} \mathcal{D}_i \text{ is salient,} & \hat{y}_i^E = 1 \\ \mathcal{D}_i \text{ is not salient,} & \hat{y}_i^E = 0. \end{cases} \quad (3)$$

The generation of pseudo-label text classification and explanation labels, respectively, form one *round*.

Iterative Mutual Enhancement Using the pseudo-text classification and explanation labels generated, we once again query \mathcal{T}^C , but now we additionally provide the pseudo-explanation labels as input. For example, the sentence in round 1 of Figure 2 and the salient tokens (highlighted in

red) are used as input to the classification prompt, which is fed into \mathcal{T}^C . This extension of equation 1 is defined below:

$$\hat{y}^C = \mathcal{T}^C(\mathcal{D}, \hat{y}^E). \quad (4)$$

We repeat equations 4 and 2, respectively, to ensure high confidence in both \mathcal{T}^C and \mathcal{T}^E predictions, i.e., the predictions from both PLMs do not further change after one round.

3.3 Explainable Multi-Task Architecture

Once \mathcal{T}^C and \mathcal{T}^E have generated confident labels, we then input both of these into a multi-task text classification model. In Figure 2 for example, we take the "negative" text classification label and the "really don't like" salient labels as input.

Specifically, we first tokenize the input document \mathcal{D} using a BERT-based tokenizer. We then pass this tokenized document into our BERT-based multi-task model and extract the following information from the model:

$$l^C, \mathbf{A} = \mathcal{T}(\mathcal{D}), \quad (5)$$

where l^C is the loss of the text classification task and $\mathbf{A} \in \mathbb{R}^{L \times H \times |\mathcal{D}| \times |\mathcal{D}|}$ is the multi-head attention tensor. L is the number of layers, and H is the number of attention heads in \mathbf{A} from the BERT-based model. We extract the attention matrix $\tilde{\mathbf{A}} \in \mathbb{R}^{|\mathcal{D}| \times |\mathcal{D}|}$ from the last layer and the last attention head of \mathbf{A} . We then apply a linear classifier $\mathbf{W} \in \mathbb{R}^{|\mathcal{D}| \times 1}$ to this attention matrix $\tilde{\mathbf{A}}$:

$$\hat{y} = \tilde{\mathbf{A}}\mathbf{W} + \mathbf{b} \quad (6)$$

where $\mathbf{b} \in \mathbb{R}^{|\mathcal{D}| \times 1}$ is the bias vector. We apply a sigmoid layer $\sigma(\cdot)$ on top of a binary cross-entropy loss function to get the attention-based loss l^E of the saliency word prediction task:

$$l^E = -w[y \cdot \log\sigma(\hat{y}) + (1 - y) \cdot \log(1 - \sigma(\hat{y}))], \quad (7)$$

Our multi-task loss function is thus a linear combination of the aforementioned loss as well as the loss l^C from the text classification task:

$$l = l^C + \lambda l^E, \quad (8)$$

where $\lambda \in [0, 1]$ is a hyper-parameter controlling the performance balance between the text classification and saliency word prediction.

4 Experiments

4.1 Experimental Setup

Datasets We conducted experiments across 10 datasets. Dataset statistics and statistics are shown in Table 1 and listed below, respectively.

- **AGNews** (Zhang et al., 2015) consists of news articles collected from the AG's online news corpus, with articles from four different categories.
- **20News** (Lang, 1995) consists of documents from 20 different news groups, covering a wide range of topics.
- **UCINews** (Gasparetti, 2016) has a substantial number of news articles covering four categories: entertainment, technology, business, and health.
- **NYT-Topic** (Meng et al., 2020a) is a collection of New York Times articles whose labels correspond to an article's topic.
- **NYT-Location** (Meng et al., 2020a) uses the same articles as NYT-Topic but the label space corresponds to locations.
- **Yelp** (Zhang et al., 2015) is a sentiment analysis dataset consisting of reviews on restaurants, bars, and other businesses.
- **Books** (Wan and McAuley, 2018) is a corpus of book titles and their descriptions, originating from Goodreads¹, which is used for book genre classification.
- **IMDB** (Zaidan et al., 2007) contains movie reviews from IMDB, where each review is considered to be either of positive or negative sentiment.
- **Twitter**² is a collection of tweets that have been labeled or annotated with sentiment labels, indicating whether the sentiment expressed in the tweet is positive, negative, or neutral.
- **MIMIC-III** (Johnson et al., 2018) is a public electronic health record (EHR) database with patient discharge summaries as text and diagnostic-related group (DRG) codes as class labels used in our experiments.

Baselines Our baselines include both fully supervised and weakly supervised text classification methods below.

¹<https://www.goodreads.com/>

²<https://www.kaggle.com/competitions/tweet-sentiment-extraction>

Table 1: Dataset statistics, depicting the sizes of the training, testing, and development set as well as the total number of classes.

Datasets	# Train	# Dev	# Test	# Class
AGNews	108,000	12,000	7,600	4
20News	14,609	1,825	1,825	6
UCINews	26,008	2,560	27,556	4
NYT-Topic	19,197	6,400	6,400	9
NYT-Location	19,197	6,400	6,400	10
Yelp	22,800	7,600	7,600	2
Books	20,165	6,719	6,719	8
IMDB	1,600	200	200	2
Twitter	21,983	2,747	2,748	3
MIMIC-III	20,266	2,252	2,252	369

- **BERT** (Devlin et al., 2018) is a fully supervised baseline that trains a transformer model using labeled data.
- **Clinical-BERT** (Alsentzer et al., 2019) is a supervised baseline that trains the BERT model on the clinical text.
- **ConWea**³ (Mekala and Shang, 2020) expands the keyword vocabulary based on contextual representations of the labels and the corpus.
- **LOTClass**⁴ (Meng et al., 2020b) Constructs a keyword vocabulary for pseudo-label generation.
- **X-Class**⁵ (Wang et al., 2021) uses clustering to choose the representative documents for each class.
- **ClassKG**⁶ (Zhang et al., 2021) iteratively constructs keyword sub-graphs consisting of keywords across data points and derives pseudo-labels by annotating the corresponding sub-graphs.
- **WDDC-MLM**⁷ (Zeng et al., 2022) employs a masked language model to generate signal words. They combine the generated words with category names and utilize them for training.
- **NPPrompt**⁸ (Zhao et al., 2022) is a zero-shot technique that identifies similar words via non-

³<https://github.com/dheeraj7596/ConWea>

⁴<https://github.com/yumeng5/LOTClass>

⁵<https://github.com/ZihanWangKi/XClass>

⁶<https://github.com/zhanglu-cst/ClassKG>

⁷<https://github.com/HKUST-KnowComp/WDDC>

⁸<https://github.com/XuandongZhao/NPPrompt>

parametric prompts and uses them as pseudo-labels.

- **MEGClass**⁹ (Kargupta et al., 2023) generates pseudo-training labels by iteratively estimating class distribution and contextualized document embeddings.

Evaluation Metrics We use micro- F_1 and macro- F_1 as the evaluation metrics to compare the performance of the text classification methods. More details can be found in Appendix A.

Parameter Settings For each baseline method, we use the default parameter settings as reported in the original papers. More details about the parameter settings of XAI-CLASS can be found in Appendix C.

4.2 Main Results

Our main results are displayed in Table 2. XAI-CLASS outperforms all other baselines on the Yelp, NYT-Topic, Books, and UCINews datasets while providing comparable results on AGNews. We hypothesize our SOTA performance on Yelp is primarily due to its sentimental nature (as it is a polarity dataset) and the label space being distinct (positive or negative sentiment), allowing for there to be more salient words XAI-CLASS can identify compared to other types of datasets used. We provide results on two other polarity datasets, IMDB and Twitter, in Table 3. Our hypothesis is validated by XAI-CLASS outperforming baselines on Yelp and IMDB but not on the Twitter dataset, due to the introduction of the "neutral" class in Twitter.

XAI-CLASS's performance on the Books dataset drastically outperforms all other baselines. We believe this is the result of the indicative and sentiment words that often appear in the description of many books. For example, words commonly found in book descriptions such as "seduce", "murder", and "paranormal" clearly indicate the genres are "romance", "thriller", and "fantasy", respectively.

We believe much of the performance drop-off in 20News is due to labels not being completely disjoint (Zeng et al., 2022). For example, the "electronics" fine-grained class is categorized under the "science" class, although one could argue it would be more appropriate to classify instances of type "electronics" in the "computer" class (Lang, 1995).

⁹<https://github.com/pkargupta/MEGClass>

Table 2: Micro/macro F_1 scores of baseline methods compared with XAI-CLASS. XAI-CLASS results are based on the optimal number of rounds associated with each dataset. Bolded results correspond to the best-performing model.

Model	Yelp	20News	NYT-Topic	NYT-Loc	Books	AGNews	UCINews
BERT (Supervised)	95.70/95.70	96.60/96.60	95.98/95.01	96.00/95.00	81.00/81.00	93.05/93.06	93.13/93.15
ConWea	71.40/71.20	75.73/73.26	81.67/71.54	85.31/83.81	52.30/52.60	74.43/74.01	32.93/32.69
LOTClass	87.40/87.20	73.78/72.53	67.11/43.38	58.49/58.96	19.90/16.10	86.59/86.56	73.20/72.36
X-Class	86.80/86.80	73.17/73.07	79.01/68.62	89.51/89.68	53.60/54.20	85.74/85.66	68.85/69.62
ClassKG	91.20/91.20	81.00/82.00	72.06/65.76	86.84/83.35	55.00/54.70	88.80/88.80	N/A
WDDC-MLM	81.20/81.10	81.21/68.82	81.50/69.20	88.84/86.91	53.86/53.75	88.26/88.25	81.50/81.34
NPPrompt	81.20/81.10	68.90/68.80	64.60/64.20	53.90/53.80	49.60/49.70	85.20/85.20	N/A
MEGClass	87.41/87.41	81.72/80.63	85.42/68.03	93.06/91.93	56.35/55.71	N/A	N/A
XAI-CLASS	95.45/95.45	75.29/71.30	88.39/80.35	82.50/86.52	70.56/70.67	88.20/88.15	83.95/83.87

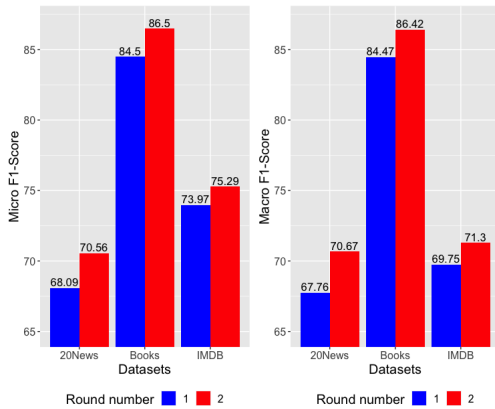


Figure 3: Micro F_1 and macro F_1 scores of two rounds of XAI-CLASS on 20News, Books, and IMDB test sets.

4.3 Ablation Study

Iterative Mutual Enhancement Effectiveness

To determine the effectiveness of iterative mutual enhancement, we identify the performance of datasets across multiple rounds. Figure 3 shows these results, clearly indicating that the performance increases when iterating up to a specified number of rounds. It should be noted that the optimal number of rounds is dependent on the dataset, with datasets that have high performance without many rounds most likely requiring fewer rounds than otherwise.

Analyzing Salient Token Utility To analyze the utility of incorporating salient tokens in XAI-CLASS, we conduct experiments on the IMDB and Twitter datasets (Table 3) as they have ground truth salient labels available. Results on both datasets indicate the XAI-CLASS-FS, a variant of XAI-CLASS that includes ground truth saliency labels during training, outperforms XAI-CLASS. This performance increase when utilizing ground truth saliency tokens justifies utilizing salient tokens as it shows that the gold-standard ground-truth saliency

Table 3: F_1 scores of BERT baseline against XAI-CLASS variants. XAI-CLASS-FS is the fully supervised version (with respect to saliency labels) of XAI-CLASS, consisting of ground truth salient labels.

Model	Dev	Test
Dataset: IMDB		
BERT (Supervised)	85.90	85.60
XAI-CLASS-FS	89.50	87.80
XAI-CLASS	91.50	86.40
Dataset: Twitter		
BERT (Supervised)	77.20	78.10
XAI-CLASS-FS	78.40	79.20
XAI-CLASS	61.20	63.40

labels are being incorporated. The dramatic performance increase when incorporating ground truth salient tokens for Twitter leads us to hypothesize that there’s more of a need for proper pseudo-salient representation for datasets that have labels with limited salient words, as the majority of XAI-CLASS misclassifications on the Twitter dataset are on data points whose ground truth is the "neutral" class, which doesn’t have many indicative salient words.

Backbone Pre-trained Language Models

In our experiments, we compared multiple pre-trained language models and chose FLAN-T5 (Chung et al., 2022) as \mathcal{T}^C for the text classification label generation, and BERT (Devlin et al., 2018) as \mathcal{T}^E for the explanation label generation. More information regarding our justification for our choice of \mathcal{T}^C and \mathcal{T}^E can be found in Appendix B.

4.4 Explainability Study

To evaluate the explainability of XAI-CLASS over baseline methods, we qualitatively assess the explainability of Clinical-BERT and XAI-CLASS using six explanation techniques: **Saliency (Si-**

Table 4: Explainability of Clinical-BERT and XAI-CLASS using six explanation techniques on five explanation evaluation metrics (HA, CI, F, RC, DC) on MIMIC-III. Results are in Clinical-BERT/XAI-CLASS format.

Method	F	HA	DC	RC	CI
Random	38.45/ 38.56	0.21/ 0.24	0.02/ 0.03	0.06/0.06	0.13/0.13
ShapSampl	29.43 /29.28	0.56/ 0.61	0.23/ 0.25	0.21/ 0.23	0.13/ 0.14
LIME	38.00 /37.89	0.31/ 0.33	0.36/ 0.39	0.61/0.61	0.12/ 0.14
Occlusion	23.00/ 25.02	0.55/ 0.56	0.19/ 0.21	0.34/ 0.41	0.12/ 0.14
Saliency _{μ}	51.01 /49.23	0.57/ 0.59	0.34 /0.32	0.26/ 0.36	0.14/ 0.19
Saliency _{L2}	44.30/44.30	0.31/ 0.37	0.33/ 0.39	0.24/ 0.31	0.15 /0.13
InputXGrad _{μ}	20.20/ 28.73	0.53/ 0.57	0.41/ 0.42	0.19/ 0.18	0.15/ 0.17
InputXGrad _{L2}	48.72/ 49.54	0.22/ 0.24	0.41/ 0.43	0.22 /0.21	0.15/ 0.16
GuidedBP _{μ}	36.66 /35.76	0.37 /0.34	0.40/ 0.43	0.02/ 0.04	0.13 /0.12
GuidedBP _{L2}	49.31 /48.38	0.45 /0.43	0.40/ 0.43	0.19/0.19	0.14 /0.11

Table 5: Sample of instances with incorrect/ambiguous ground truths in the 20News dataset.

Input Text	Class Prediction	Ground Truth	Salient Word Prediction
72 Chevelle SS for sale. [...] I need money for college. [...] 1972 chevelle super sport rebuilt 402 [...] \$ 5995.	Sale	Sports	sale, money, sport
[...] key would appear to be cryptographically useless. [...] The same key is used for both encryption and decryption.	Computer	Science	crypto-graphically, encryption, key
What exactly is an IBM 486 SLC processor? Could someone please tell me if the 486 SLC and 486 SLC2 processors IBM is putting in their Thinkpad 700's.	Computer	Science	IBM, processor, 486
Cultural enquiries more like those who use their backs instead of their minds [...] intolerant of anything outside of their group [...] there is no justification for taking away individuals freedom.	Politics	Sports	cultural, freedom

monyan et al., 2014), **InputXGradient** (Kindermans et al., 2016), **Guided Backpropagation** (Springenberg et al., 2014), **Occlusion** (Zeiler and Fergus, 2014), **Shapley Value Sampling** (Castro et al., 2009), and **LIME** (Ribeiro et al., 2016) over five explanation evaluation metrics (Atanasova et al., 2020) **Agreement with Human Rationales** (HA), **Confidence Indication** (CI), **Faithfulness** (F), **Rationale Consistency** (RC), and **Dataset Consistency** (DC) on the MIMIC-III dataset. Details of the above explanation evaluation metrics can be found in a previous study of explanation techniques in text classification (Atanasova et al., 2020). The results in Table 4 demonstrate that XAI-CLASS improved the model explainability by capturing the saliency information during the training process for all explanation evaluation metrics excluding faithfulness. More results on the explainability case study can be found in Appendix D.

4.5 Case Study

We further explore some cases with incorrect/ambiguous ground truths for multiple reasons, depicted in Table 5. The text in the first row of Table 5 is most likely supposed to be assigned to the "sale" class but is instead labeled with the "sports" class as ground truth, most likely because the word

"sport" appears in the text. XAI-CLASS predicted the "sale" class, even though it determined that "sport" was a salient token. This suggests that the model is robust to a small number of words dictating the classification prediction. The second row in Table 5 coincides with the cryptograph example in section 4.2, where one could argue all salient words picked up by the model could be categorized under the term "computer", instead of the ground truth "science". The last two rows of Table 5 appear to be mislabelled, as the third row's text talks exclusively about processors and the fourth example talks only about political issues, yet they are labeled as "science" and "sports", respectively.

5 Conclusion

We propose XAI-CLASS, a novel extremely weakly-supervised text classification method that employs a multi-round question-answering process to generate pseudo-training data and trains a multi-task framework that simultaneously learns both text classification and word saliency prediction. XAI-CLASS has superior performance over baselines for both model performance and explainability. Future work includes extending XAI-CLASS to the multi-label setting.

561 Limitations

562 XAI-CLASS, although effective, operates under
563 the assumption of a disjoint label space and is
564 not specifically tailored for fine-grained or multi-
565 label text classification tasks. As a result, it may
566 not perform optimally on datasets like 20News,
567 where there are instances where ground truth labels
568 have some degree of overlap. However, exploring
569 weakly-supervised methods for fine-grained, multi-
570 label text classification is an intriguing direction
571 for future research. Furthermore, it's important
572 to note that XAI-CLASS requires careful consid-
573 eration when selecting the number of rounds of
574 question answering. It is not designed to scale to
575 a large number of rounds, and typically, no more
576 than three rounds are used. This limitation arises
577 because each round involves two queries for the
578 question answering models: one for generating text
579 classification labels and the other for saliency word
580 generation. This process can be computationally
581 expensive, necessitating a mindful balance between
582 computational resources and desired performance.

583 Ethics Statement

584 Given our current methodology, we do not antic-
585 ipate any significant ethical concerns. We have
586 utilized datasets and models from open-source
587 domains, promoting transparency and accessibil-
588 ity of information. Text classification is a well-
589 established task in natural language processing,
590 widely studied and applied in various domains.
591 However, we acknowledge that our architecture
592 relies on PLMs, which may make decisions based
593 on biases present in the training data. Although
594 our experiments have not revealed any apparent
595 performance issues related to bias, it is important
596 to recognize that this observation may be limited
597 to the datasets we have used. It is crucial to remain
598 vigilant and continue exploring ways to mitigate
599 and address biases that may arise from the use of
600 pre-trained models.

601 References

602 Eugene Agichtein and Luis Gravano. 2000. [Snowball:
603 Extracting relations from large plain-text collections.](#)
604 In *Proceedings of the Fifth ACM Conference on Dig-
605 ital Libraries*, DL '00, page 85–94, New York, NY,
606 USA. Association for Computing Machinery.

607 Emily Alsentzer, John Murphy, William Boag, Wei-
608 Hung Weng, Di Jindi, Tristan Naumann, and

Matthew McDermott. 2019. Publicly available clini- 609
cal bert embeddings. In *Proceedings of the 2nd Clin- 610
ical Natural Language Processing Workshop*, pages 611
72–78. 612

Marco Ancona, Enea Ceolini, A. Cengiz Öztireli, and 613
Markus H. Gross. 2017. [A unified view of gradient-
614 based attribution methods for deep neural networks.](#)
615 *CoRR*, abs/1711.06104. 616

Pepa Atanasova, Jakob Grue Simonsen, Christina Li- 617
oma, and Isabelle Augenstein. 2020. [A diagnostic
618 study of explainability techniques for text classifi-
619 cation.](#) In *Proceedings of the 2020 Conference on
620 Empirical Methods in Natural Language Processing
621 (EMNLP)*, pages 3256–3274, Online. Association for
622 Computational Linguistics. 623

Sonia Badene, Kate Thompson, Jean-Pierre Lorré, and 624
Nicholas Asher. 2019. [Data programming for learn-
625 ing discourse structure.](#) In *Proceedings of the 57th
626 Annual Meeting of the Association for Computational
627 Linguistics*, pages 640–645, Florence, Italy. Associa-
628 tion for Computational Linguistics. 629

Javier Castro, Daniel Gómez, and Juan Tejada. 2009. 630
[Polynomial calculation of the shapley value based
631 on sampling.](#) *Computers & Operations Research*,
632 36(5):1726–1730. 633

Ming Wei Chang, Lev Ratinov, Dan Roth, and Vivek 634
Srikumar. 2008. Importance of semantic represen- 635
tation: Dataless classification. In *AAAI-08/IAAI-08
636 Proceedings - 23rd AAAI Conference on Artificial
637 Intelligence and the 20th Innovative Applications of
638 Artificial Intelligence Conference*, Proceedings of the
639 National Conference on Artificial Intelligence, pages
640 830–835. 23rd AAAI Conference on Artificial Intel-
641 ligence and the 20th Innovative Applications of Arti-
642 ficial Intelligence Conference, AAAI-08/IAAI-08 ;
643 Conference date: 13-07-2008 Through 17-07-2008. 644

Hyung Won Chung, Le Hou, Shayne Longpre, Barret 645
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi 646
Wang, Mostafa Dehghani, Siddhartha Brahma, Al- 647
bert Webson, Shixiang Shane Gu, Zhuyun Dai, 648
Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh- 649
ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, 650
Dasha Valter, Sharan Narang, Gaurav Mishra, Adams 651
Yu, Vincent Zhao, Yanping Huang, Andrew Dai, 652
Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja- 653
cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, 654
and Jason Wei. 2022. [Scaling instruction-finetuned
655 language models.](#) 656

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 657
Kristina Toutanova. 2018. [BERT: pre-training of
658 deep bidirectional transformers for language under-
659 standing.](#) *CoRR*, abs/1810.04805. 660

Alex A Freitas. 2014. Comprehensible classification 661
models: a position paper. *ACM SIGKDD explo- 662
663 rations newsletter*, 15(1):1–10.

664	Evgeniy Gabrilovich and Shaul Markovitch. 2007.	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	717
665	Computing semantic relatedness using wikipedia-	Percy Liang. 2016. SQuAD: 100,000+ questions for	718
666	based explicit semantic analysis. In <i>Proceedings</i>	machine comprehension of text . In <i>Proceedings of</i>	719
667	<i>of the 20th International Joint Conference on Artif-</i>	<i>the 2016 Conference on Empirical Methods in Natu-</i>	720
668	<i>icial Intelligence, IJCAI'07</i> , page 1606–1611, San	<i>ral Language Processing</i> , pages 2383–2392, Austin,	721
669	Francisco, CA, USA. Morgan Kaufmann Publishers	Texas. Association for Computational Linguistics.	722
670	Inc.		
671	Fabio Gaspiretti. 2016. News Aggregator.	Alexander Ratner, Christopher De Sa, Sen Wu, Daniel	723
672	UCI Machine Learning Repository. DOI:	Selsam, and Christopher Ré. 2017. Data program-	724
673	https://doi.org/10.24432/C5F61C .	ming: Creating large training sets, quickly .	725
674	Alistair EW Johnson, David J Stone, Leo A Celi, and	Marco Tulio Ribeiro, Sameer Singh, and Carlos	726
675	Tom J Pollard. 2018. The mimic code repository: en-	Guestrin. 2016. Model-agnostic interpretability of	727
676	abling reproducibility in critical care research. <i>Jour-</i>	machine learning. <i>arXiv preprint arXiv:1606.05386</i> .	728
677	<i>nal of the American Medical Informatics Association</i> ,	Kai Shu, Subhabrata Mukherjee, Guoqing Zheng,	729
678	25(1):32–39.	Ahmed Hassan Awadallah, Milad Shokouhi, and Su-	730
679	Priyanka Kargupta, Tanay Komarlu, Susik Yoon, Xuan	san Dumais. 2020. Learning with weak supervision	731
680	Wang, and Jiawei Han. 2023. Megclass: Text	for email intent detection . In <i>Proceedings of the 43rd</i>	732
681	classification with extremely weak supervision via	<i>International ACM SIGIR Conference on Research</i>	733
682	mutually-enhancing text granularities. <i>arXiv preprint</i>	<i>and Development in Information Retrieval, SIGIR</i>	734
683	<i>arXiv:2304.01969</i> .	'20, page 1051–1060, New York, NY, USA. Associa-	735
684	Daniel Khashabi, Sewon Min, Tushar Khot, Ashish	tion for Computing Machinery.	736
685	Sabharwal, Oyvind Tafjord, Peter Clark, and Han-	Karen Simonyan, Andrea Vedaldi, and Andrew Zis-	737
686	naneh Hajishirzi. 2020. Unifiedqa: Crossing format	serman. 2014. Deep inside convolutional networks:	738
687	boundaries with a single qa system. <i>arXiv preprint</i>	Visualising image classification models and saliency	739
688	<i>arXiv:2005.00700</i> .	maps .	740
689	Pieter-Jan Kindermans, Kristof Schütt, Klaus-Robert	Yangqiu Song and Dan Roth. 2014. On dataless hi-	741
690	Müller, and Sven Dähne. 2016. Investigating	erarchical text classification. In <i>Proceedings of the</i>	742
691	the influence of noise and distractors on the in-	<i>Twenty-Eighth AAAI Conference on Artificial Intelli-</i>	743
692	terpretation of neural networks. <i>arXiv preprint</i>	<i>gence, AAAI'14</i> , page 1579–1585. AAAI Press.	744
693	<i>arXiv:1611.07270</i> .	Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas	745
694	Ken Lang. 1995. Newsweeder: Learning to filter net-	Brox, and Martin Riedmiller. 2014. Striving for sim-	746
695	news . In Armand Prieditis and Stuart Russell, editors,	licity: The all convolutional net. <i>arXiv preprint</i>	747
696	<i>Machine Learning Proceedings 1995</i> , pages 331–339.	<i>arXiv:1412.6806</i> .	748
697	Morgan Kaufmann, San Francisco (CA).	Duyu Tang, Bing Qin, and Ting Liu. 2015. Document	749
698	Dheeraj Mekala and Jingbo Shang. 2020. In <i>Proceed-</i>	modeling with gated recurrent neural network for	750
699	<i>ings of the 58th Annual Meeting of the Association for</i>	sentiment classification . In <i>Proceedings of the 2015</i>	751
700	<i>Computational Linguistics</i> , pages 323–333, Online.	<i>Conference on Empirical Methods in Natural Lan-</i>	752
701	Association for Computational Linguistics. [link].	<i>guage Processing</i> , pages 1422–1432, Lisbon, Portu-	753
702	Yu Meng, Jiaxin Huang, Guangyuan Wang, Zihan Wang,	gal. Association for Computational Linguistics.	754
703	Chao Zhang, Yu Zhang, and Jiawei Han. 2020a. Dis-	Fangbo Tao, Chao Zhang, Xiusi Chen, Meng Jiang,	755
704	criminative topic mining via category-name guided	Tim Hanratty, Lance Kaplan, and Jiawei Han. 2018.	756
705	text embedding. In <i>Proceedings of The Web Confer-</i>	Doc2cube: Allocating documents to text cube with-	757
706	<i>ence 2020</i> , pages 2121–2132.	out labeled data . In <i>2018 IEEE International Confer-</i>	758
707	Yu Meng, Jiaming Shen, Chao Zhang, and Jiawei Han.	<i>ence on Data Mining (ICDM)</i> , pages 1260–1265.	759
708	2018. Weakly-supervised neural text classification .	Mengting Wan and Julian McAuley. 2018. Item rec-	760
709	In <i>Proceedings of the 27th ACM International Confe-</i>	ommendation on monotonic behavior chains. In	761
710	<i>rence on Information and Knowledge Management</i> .	<i>Proceedings of the 12th ACM conference on recom-</i>	762
711	ACM.	<i>mender systems</i> , pages 86–94.	763
712	Yu Meng, Yunyi Zhang, Jiaxin Huang, Chenyan	Zihan Wang, Dheeraj Mekala, and Jingbo Shang. 2021.	764
713	Xiong, Heng Ji, Chao Zhang, and Jiawei Han.	X-class: Text classification with extremely weak su-	765
714	2020b. Text classification using label names only:	pervision . In <i>Proceedings of the 2021 Conference of</i>	766
715	A language model self-training approach . <i>CoRR</i> ,	<i>the North American Chapter of the Association for</i>	767
716	abs/2010.07245 .	<i>Computational Linguistics: Human Language Tech-</i>	768
		<i>nologies</i> , pages 3043–3053, Online. Association for	769
		Computational Linguistics.	770

771 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-
772 bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.
773 Xlnet: Generalized autoregressive pretraining for lan-
774 guage understanding. *Advances in neural informa-*
775 *tion processing systems*, 32.

776 Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He,
777 Alex Smola, and Eduard Hovy. 2016. [Hierarchical](#)
778 [attention networks for document classification](#). In
779 *Proceedings of the 2016 Conference of the North*
780 *American Chapter of the Association for Computa-*
781 *tional Linguistics: Human Language Technologies*,
782 pages 1480–1489, San Diego, California. Associa-
783 tion for Computational Linguistics.

784 Omar Zaidan, Jason Eisner, and Christine Piatko. 2007.
785 Using “annotator rationales” to improve machine
786 learning for text categorization. In *Human Language*
787 *Technologies 2007: The Conference of the North*
788 *American Chapter of the Association for Computa-*
789 *tional Linguistics; Proceedings of the Main Confer-*
790 *ence*, pages 260–267, Rochester, New York. Associa-
791 tion for Computational Linguistics.

792 Matthew D Zeiler and Rob Fergus. 2014. Visualiz-
793 ing and understanding convolutional networks. In
794 *Computer Vision–ECCV 2014: 13th European Con-*
795 *ference, Zurich, Switzerland, September 6–12, 2014,*
796 *Proceedings, Part I 13*, pages 818–833. Springer.

797 Ziqian Zeng, Weimin Ni, Tianqing Fang, Xiang Li,
798 Xinran Zhao, and Yangqiu Song. 2022. [Weakly su-](#)
799 [pervised text classification using supervision signals](#)
800 [from a language model](#). In *Findings of the Associ-*
801 *ation for Computational Linguistics: NAACL 2022*,
802 pages 2295–2305, Seattle, United States. Association
803 for Computational Linguistics.

804 Lu Zhang, Jiandong Ding, Yi Xu, Yingyao Liu, and
805 Shuigeng Zhou. 2021. [Weakly-supervised text classi-](#)
806 [fication based on keyword graph](#). In *Proceedings of*
807 *the 2021 Conference on Empirical Methods in Natu-*
808 *ral Language Processing*, pages 2803–2813, Online
809 and Punta Cana, Dominican Republic. Association
810 for Computational Linguistics.

811 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.
812 Character-level convolutional networks for text classi-
813 fication. *Advances in neural information processing*
814 *systems*, 28.

815 Yunyi Zhang, Fang Guo, Jiaming Shen, and Jiawei Han.
816 2022. [Unsupervised key event detection from mas-](#)
817 [sive text corpora](#). In *Proceedings of the 28th ACM*
818 *SIGKDD Conference on Knowledge Discovery and*
819 *Data Mining*. ACM.

820 Xuandong Zhao, Siqi Ouyang, Zhiguo Yu, Ming Wu,
821 and Lei Li. 2022. Pre-trained language models
822 can be fully zero-shot learners. *arXiv preprint*
823 *arXiv:2212.06950*.

824 A Evaluation Metrics

825 We report performance based on the micro and
826 macro F_1 scores, which are defined below.

$$F_1 = 2 \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$F_1 \text{ macro} = \frac{1}{n} \sum_{i=1}^n F_{1,i}$$

$$F_1 \text{ micro} = \frac{2 \sum_{i=1}^n \text{TP}_i}{2 \sum_{i=1}^n \text{TP}_i + \sum_{i=1}^n \text{FP}_i + \sum_{i=1}^n \text{FN}_i}$$

827 where TP is true positive, FP is false positive, and
828 FN is false negative. We use the sklearn¹⁰ library
829 to obtain these metrics.

830 B Pseudo-Label and Text Classification 831 Backbone Analysis

832 We conduct experiments to identify the most ap-
833 propriate PLMs for \mathcal{T}^C and \mathcal{T}^E . To identify the
834 most appropriate \mathcal{T}^E , we conduct zero-shot text
835 classification on 7 datasets (Table 6). Flan-T5 per-
836 forms better than all other models across the seven
837 datasets, indicating why we chose Flan-T5 as the
838 backbone PLM for \mathcal{T}^C .

839 We perform a similar experiment to identify the
840 most appropriate backbone for \mathcal{T}^E . Concretely,
841 we perform zero-shot salient label prediction using
842 the IMDB and Twitter datasets, as these are the
843 only datasets we’ve experimented with that have
844 ground truth saliency labels (Table 7). The results
845 show that BERT and Unified-QA (Khashabi et al.,
846 2020) should be the \mathcal{T}^E backbone of choice when
847 using the IMDB and Twitter datasets for training,
848 respectively.

849 C Parameter Settings

850 **Runtime Analysis** We conduct all of our experi-
851 ments on an NVIDIA DGX A100 GPU (640GB).
852 The run times for optimal configurations across all
853 datasets can be found in Table 8.

854 **Hyper-parameters** The optimal hyper-
855 parameters for our results in Tables 2 and 3 are
856 listed in Table 9. The possible values each of the
857 hyper-parameters can take are listed below:

- 858 • $\mathcal{T}^C \in \{\text{FLAN-T5-SMALL}, \text{FLAN-T5-BASE},$
859 $\text{FLAN-T5-LARGE}, \text{FLAN-T5-XL}, \text{FLAN-T5-}$
860 $\text{XXL}\}$

¹⁰<https://scikit-learn.org/stable/>

- PLM for psuedo-text classification label
generation

- $\mathcal{T}^E \in \{\text{BERT-BASE}, \text{BERT-LARGE}, \text{UNIFIED-}$
 $\text{QA-LARGE}, \text{UNIFIED-QA-3B}\}$

- PLM for psuedo-saliency label generation

- $\lambda \in \{0.5, 0.7, 0.9\}$

- Hyper-parameter for determining how much
of the saliency loss should be incorporated

- Round # $\in \{0, 1, 2, 3\}$

- Learning Rate $\in \{2e - 04, 2e - 05, 5e - 05\}$

- Dropout $\in \{0.1, 0.2, 0.3, 0.4\}$

- Number of Epochs $\in \{1, 2, 3\}$

We implement the PLMs in Python using the Hug-
gingFace Transformer library¹¹.

875 D Explanability Case Study

876 To further evaluate the explainability of XAI-
877 CLASS over the baseline methods, we qualitatively
878 assess the explainability of Clinical-BERT and
879 XAI-CLASS using attention distribution (heatmap).
880 The results in Figure 4 demonstrate that XAI-
881 CLASS improved the model explainability by cap-
882 turing the saliency information during the training
883 process, particularly in all evaluation metrics ex-
884 cluding faithfulness. The results align well with
885 human-given ICD-9 codes as the explanation for
886 the DRG code prediction.

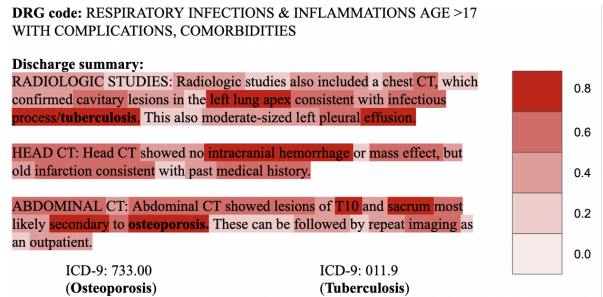
¹¹<https://github.com/huggingface/transformers>

Table 6: Zero-shot text classification label generation micro/macro F_1 -scores across multiple \mathcal{T}^C models. Flan-T5 outperforms all other models across all datasets used, thus serving as our backbone for \mathcal{T}^C .

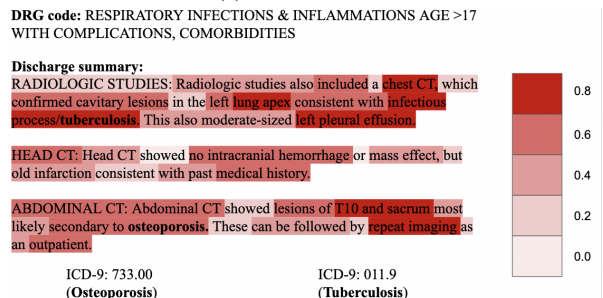
Model	Yelp	20News	NYT-Topic	NYT-Loc	Books	AGNews	UCINews
GPT-2	50.74/42.88	13.97/9.74	7.58/2.72	4.02/2.64	9.88/4.89	26.30/20.16	25.57/11.85
BERT	49.63/39.38	24.99/11.40	31.77/5.58	19.31/4.04	12.56/12.06	26.07/13.18	24.92/11.72
Unified-QA	95.66/95.66	69.75/66.10	76.17/67.30	72.38/75.11	48.13/48.64	86.21/86.12	80.88/80.67
Flan-T5	97.42/97.42	75.34/72.15	87.53/78.87	81.36/85.90	72.03/72.45	88.51/88.48	84.27/84.16

Table 7: Zero-shot salient label generation micro/macro F_1 -scores across multiple \mathcal{T}^E models on the IMDB and Twitter datasets. We report on these datasets as these are the only datasets with salient labels.

Model	IMDB	Twitter
BERT	12.26/10.92	67.56/40.32
Flan-T5	11.26/10.12	69.31/40.94
GPT-2	9.25/8.47	72.16/41.92
Unified-QA	11.57/10.37	73.8/42.47



(a) Clinical-BERT



(b) XAI-CLASS

Table 8: Average run time for each dataset for best hyper-parameter configuration.

Dataset	Runtime (hours)
AGNews	10
20News	4
UCINews	4
IMDB	1
Twitter	3

Figure 4: The attention distribution (heatmap) of of Clinical-BERT and XAI-CLASS. A darker red color indicates that the model assigns higher importance to that particular word for explaining the prediction of the DRG code.

Table 9: Optimal hyper-parameters for XAI-CLASS’s results in Tables 2 and 3.

Dataset	\mathcal{T}^C	\mathcal{T}^E	Round #	λ	Learning Rate	Dropout	# Epochs
Books	FLAN-T5-XXL	BERT-BASE	2	0.5	$2e - 05$	0.3	3
NYT-Topic	FLAN-T5-XXL	BERT-BASE	1	0.5	$2e - 05$	0.3	1
NYT-Location	FLAN-T5-XXL	BERT-BASE	2	0.5	$2e - 05$	0.3	3
Yelp	FLAN-T5-XXL	BERT-BASE	1	0.5	$2e - 05$	0.3	1
AGNews	FLAN-T5-XXL	BERT-BASE	1	0.5	$2e - 05$	0.3	1
20News	FLAN-T5-XL	BERT-BASE	2	0.7	$2e - 05$	0.3	3
UCINews	FLAN-T5-XL	BERT-BASE	1	0.5	$2e - 05$	0.3	1
IMDB	FLAN-T5-XL	BERT-BASE	1	0.9	$2e - 05$	0.4	3
Twitter	FLAN-T5-XL	BERT-BASE	0	0.7	$2e - 05$	0.1	3