

Interpreting Pretrained Language Models via Concept Bottlenecks

Zhen Tan

Arizona State University
ztan36@asu.edu

Lu Cheng

University of Illinois Chicago
lucheng@uic.edu

Song Wang

University of Virginia
sw3wv@virginia.edu

Yuan Bo

Zhejiang University
byuan@zju.edu.cn

Jundong Li

University of Virginia
jundong@virginia.edu

Huan Liu

Arizona State University
huanliu@asu.edu

Abstract

Pretrained language models (PLMs) have made significant strides in various natural language processing tasks. However, the lack of interpretability due to their “black-box” nature poses challenges for responsible implementation. Although previous studies have attempted to improve interpretability by using, e.g., attention weights in self-attention layers, these weights often lack clarity, readability, and intuitiveness. In this research, we propose a novel approach to interpreting PLMs by employing high-level, meaningful concepts that are easily understandable for humans. For example, we learn the concept of “Food” and investigate how it influences the prediction of a model’s sentiment towards a restaurant review. We introduce C³M, which combines human-annotated and machine-generated concepts to extract hidden neurons designed to encapsulate semantically meaningful and task-specific concepts. Through empirical evaluations on real-world datasets, we manifest that our approach offers valuable insights to interpret PLM behavior, helps diagnose model failures, and enhances model robustness amidst noisy concept labels.

1 Introduction

Although Pretrained Language Models (PLMs) like BERT (Devlin et al., 2018) have achieved remarkable success in various NLP tasks (Zhu et al., 2020; Liu and Lapata, 2019), they are frequently regarded as black boxes, posing significant obstacles to their responsible deployment in real-world scenarios, particularly in critical domains such as healthcare (Koh et al., 2020). Therefore, enabling PLMs’ interpretability is crucial to achieve socially responsible AI (Cheng et al., 2021). To date, many existing works (Belinkov and Glass, 2019; Madsen et al., 2022) leverage attention weights extracted from the self-attention layers to provide token-level or phrase-level importance. These low-level explanations are found unfaithful (Yin and Neubig,

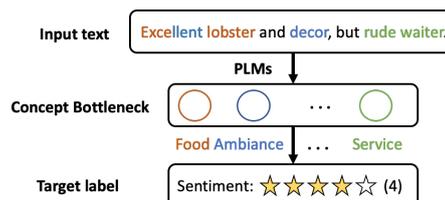


Figure 1: The illustration of CBE-PLMs. Via PLMs, the original texts x is first map into an intermediate layer consisting of a set of human-comprehensible concepts c , which are then used to predict the target label y .

2022) and lack readability and intuitiveness (Losch et al., 2019), leading to unstable or even unreasonable explanations. To address these limitations, we seek to explain via human-comprehensible *concepts* that use more abstract features (e.g., general notions) as opposed to raw input features at the token level (Zarlenga et al., 2022; Liao and Vaughan, 2023). The foundation of this work is the Concept Bottleneck Models (CBMs) (Koh et al., 2020) that interprets deep models (e.g., ResNet (He et al., 2016)) for image classification tasks using high-level concepts (e.g., shape). For NLP tasks such as sentiment analysis, concepts can be Food, Ambiance, and Service as shown in Figure 1, where each concept corresponds to a neuron in the concept bottleneck layer. The final decision layer is then a linear function of these concepts. Using concepts greatly improves the readability and intuitiveness of the explanations compared to low-level features such as “lobster”.

We propose to study *Concept-Bottleneck-Enabled Pretrained Language Models* (CBE-PLMs). There are three key challenges: First, CBMs cannot be directly adapted since PLMs are pre-trained and fine-tuned on separate corpora while CBMs work on the same end-to-end image classification tasks during training and testing. Therefore, the corpora used for pre-training PLMs may contain useful text-concept correlations that are unseen in the downstream task. An investigation of the adaptability of CBMs to CBE-PLMs is needed. Second, the majority of existing

CBMs (Koh et al., 2020; Zarlenga et al., 2022) require human-annotated concepts. This can be challenging for natural language since the annotator may need to read through the entire text to understand the context and label one concept (Németh et al., 2020). This limits the practical usage and scalability of CBE-PLMs. Third, many studies have identified the tradeoff between interpretability and task accuracy using CBMs since the predetermined concepts may leave out important information for target task prediction (Zarlenga et al., 2022). Therefore, it is crucial to improve both interpretability and task performance to achieve optimal interpretability-utility tradeoff.

To tackle the first challenge, we adapt standard training strategies in CBMs (Koh et al., 2020) to learning CBE-PLMs and conduct comprehensive analyses to identify the best way to adapt CBMs to interpret PLMs. For the second challenge of concept discovery and labeling, we propose leveraging Large Language Models (LLMs) trained on extensive human-generated corpora and feedbacks, such as ChatGPT (OpenAI, 2023), to identify novel concepts in text and generate pseudo-labels (via prompting) for unlabeled concepts. Recent studies (Bommasani et al., 2022; OpenAI, 2023) exhibit that these LLMs encapsulate significant amounts of human common sense knowledge. By augmenting the small set of human-specified concepts with machine-generated concepts, we increase concept diversity and useful information for prediction. In addition, generated pseudo-labels offer us a large set of instances with noisy concept labels, complementing the smaller set of instances with clean labels. To further improve interpretability-utility tradeoff (third challenge), we propose to learn from noisy concept labels and incorporate a concept-level MixUp mechanism (Zhang et al., 2017) that allows CBE-PLMs to cooperatively learn from both noisy and clean concept sets. We name our framework for training CBE-PLMs as *ChatGPT-guided Concept augmentation with Concept-level Mixup* (C³M). In summary, our contributions include:

- We provide the first comprehensive investigation of standard training strategies of CBMs for interpreting PLMs and benchmark CBE-PLMs.
- We propose C³M, which leverages LLMs and MixUp to help PLMs learn from human-annotated and machine-generated concepts. C³M liberates CBMs from predefined concepts and enhances the interpretability-utility tradeoff.

- We demonstrate the effectiveness and robustness of test-time concept intervention for the learned CBE-PLMs for common text classification tasks.

2 Related Work

2.1 Interpreting Pretrained Language Models

PLMs such as Word2Vec (Mikolov et al., 2013), BERT (Devlin et al., 2018), and the more recent GPT series (Radford et al., 2019; Brown et al., 2020; OpenAI, 2023) have demonstrated impressive performance in various NLP tasks. However, their opaque nature poses a challenge in comprehending how PLMs work internally (Diao et al., 2022). In order to improve the interpretability and transparency of PLMs, researchers have explored different approaches, such as visualizing attention weights (Galassi et al., 2020), probing feature representations (Mishra et al., 2017; Lundberg and Lee, 2017; Bills et al., 2023), and using counterfactuals (Wu et al., 2021; Ross et al., 2021), among others, to provide explanations at the local token-level, instance-level, or neuron-level. However, these methods often lack faithfulness and intuitiveness, and are of poor readability, which undermines their trustworthiness (Madsen et al., 2022).

Recently, researchers have turned to global concept-level explanations that are naturally understandable to humans. Although this level of interpretability has been less explored in NLP compared to computer vision (Goyal et al., 2019; Kim et al., 2018; Mu and Andreas, 2020), it has gained attention. For instance, a study (Vig et al., 2020) investigates gender classification bias by examining the association of occupation words such as ‘nurse’ with gender. In addition, the CBMs (Koh et al., 2020; Zarlenga et al., 2022) have emerged as novel frameworks for achieving concept-level interpretability in lightweight image classification systems. CBMs typically involve a layer preceding the final fully connected classifier, where each neuron corresponds to a concept that can be interpreted by humans. CBMs also show advantages in improving accuracy through human intervention during testing. Yet, the application of CBMs to larger-scale PLMs interpretation is under-explored. Implementing CBMs necessitates human involvement in defining the concept set and annotating the concept labels. Such requirements are challenging for natural language as humans may need to read through the entire text to understand the context and label one concept (Németh et al., 2020).

2.2 Learning from Noisy Labels

Addressing inaccurately labeled or misclassified data in real-world scenarios is the goal of learning from noisy labels, with techniques including noise transition matrix estimation (Liu et al., 2022), robust risk minimization (Engleson and Azizpour, 2021), and more. Recently, the resilience of semi-supervised learning methods like MixMatch (Berthelot et al., 2019) and FixMatch (Sohn et al., 2020) to label noise has been discovered by using pseudo-labels for unlabeled data. Inspired by them, we propose to utilize an LLM (ChatGPT) as a fixed-label guesser, generating noisy intermediate concept labels to potentially predict task labels.

Notably, CBMs specialize in the interpretation and interactability of deep models for general classification tasks. While *Multi-Aspect Sentiment Analysis* (Zhang et al., 2022) (MASA) shares similar goals when using aspects as concepts, it differs as concepts are not confined to fine-grained aspectual features and can be abstract ideas or broader notions throughout entire contexts. Aspect labels in MASA, primarily used for prediction accuracy, are not always mandatory. To summarize, this study pioneers the comprehensive exploration of utilizing concepts for interpreting large-scale PLMs, and provides a robust framework for harnessing the noisy signals from LLMs to achieve interpretable outcomes from lighter-weight PLMs, which can be easily understood by users.

3 Enable Concept Bottlenecks for Pretrained Language Models

3.1 Problem Setup

We focus on interpreting the predictions of fine-tuned PLMs for text classification tasks. Given data $\mathcal{D} = \{(x^{(i)}, y^{(i)}, c^{(i)})_{i=1}^n\}$, where $x \in \mathbb{R}^d$ is the original text input, $y \in \mathbb{R}$ is the target label to predict, and $c \in \mathbb{R}^k$ is a vector of k concepts from the concept set \mathcal{C} with $|\mathcal{C}| = k$, we consider a PLM f_θ parameterized by θ that encodes an input text $x \in \mathbb{R}^d$ into its latent representation $z \in \mathbb{R}^e$. Vanilla fine-tuning strategy, concretely defined in Appendix A, can be abstracted as $x \rightarrow z \rightarrow y$.

Concept-Bottleneck-Enabled Pretrained Language Models. The original concept bottlenecks in CBMs (Koh et al., 2020) come from resizing one of the layers in the CNN encoder to match the number of concepts. However, since PLM encoders typically provide text representations with much higher dimensions than the number of concepts,

directly reducing the neurons in the layer would significantly impact the quality of learned text representation. To address this issue, we instead add a linear layer with the sigmoid activation, denoted as p_ψ , that projects the learned latent representation $z \in \mathbb{R}^e$ into the concept space $c \in \mathbb{R}^k$. This process can be represented as $x \rightarrow z \rightarrow c \rightarrow y$. Note that, unlike the previous works for image classification, each concept here does not need to be binary (i.e., present or not). We allow multi-class concepts, e.g., the concept ‘‘Food’’ in a restaurant review can be positive, negative, or unknown. We refer to the PLM and the projector (f_θ, p_ψ) together as the *concept encoder* and the complete model $(f_\theta, p_\psi, g_\phi)$ as *Concept-Bottleneck-Enabled Pretrained Language Models* (CBE-PLMs).

During training, CBE-PLMs seek to achieve two goals: (1) align concept prediction $\hat{c} = p_\psi(f_\theta(x))$ to x ’s ground-truth concept labels c and (2) align label prediction $\hat{y} = g_\phi(p_\psi(f_\theta(x)))$ to ground-truth task labels y . We accordingly adapt the three conventional strategies, *independent* training, *sequential* training, and *joint* training, proposed in (Koh et al., 2020) to learn the CBE-PLM. Their detailed formulations are given in Appendix A.

3.2 Benchmarking CBE-PLMs

We propose to benchmark the performance of the vanilla fine-tuning and the three training strategies for CBE-PLMs using two text classification datasets: CEBaB (Abraham et al., 2022) and IMDB (Maas et al., 2011). Both datasets contain human-labeled concepts. We consider four typical PLMs following Abraham et al. (2022). Descriptions of the PLM backbones, datasets, and concept labels are detailed in Section 5.1, Section 5.2, and Appendix G. We consider the target task scores and concept prediction scores as the evaluation metrics for utility and interpretability, respectively.

CBM for CBE-PLMs. In this experiment, we aim to identify the optimal training strategy for CBE-PLMs. The results depicted in Figure 2 confirm that standard-PLMs typically yield the highest task scores, demonstrating that the implementation of a concept bottleneck can indeed impact target task performance negatively. However, without considering the concept labels, standard-PLMs lack interpretability. In contrast, CBE-PLMs trained jointly exhibit higher task scores and superior concept prediction scores compared to their counterparts. This divergence from CBMs in the image

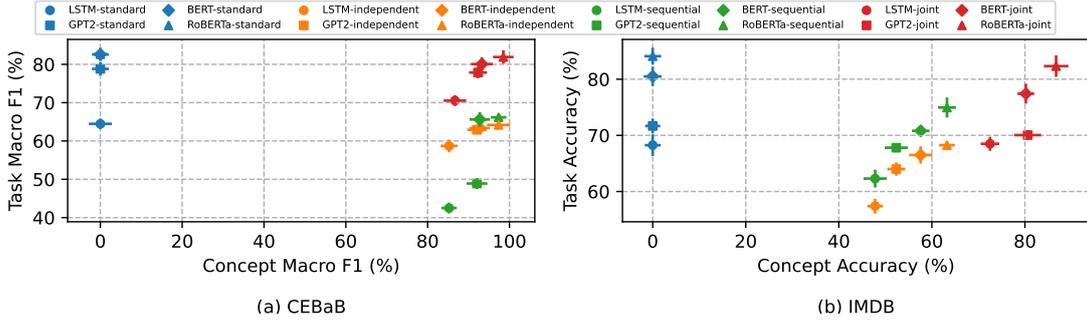


Figure 2: Illustration of the interpretability-accuracy trade-off using various backbones. The top-right indicates a more favorable trade-off, i.e., a better interpretability-utility Pareto front. We also show the confidence intervals for both dimensions.

domain, where all three strategies display similar performance (Koh et al., 2020), is notable. We attribute this to PLMs’ extensive pretraining on numerous human-generated corpora and larger parameter numbers than the studied vision encoders such as ResNets (He et al., 2016). Unlike independent or sequential training where the PLM encoder is fixed after training on the concept labels, joint training allows PLMs to utilize their capacity to learn concepts and target labels jointly, making the learned concept activations from the bottleneck layer better aligned with the task labels. Given this advantage of joint training, we adopt it as the default strategy for training CBE-PLMs in the subsequent sections.

While initial findings from applying vanilla CBM (Koh et al., 2020) for interpreting PLMs appear encouraging, they require human-annotated concepts during training. This proves to be impractical in real-world situations due to the vast number of potential concepts and the time-intensive annotation process (Németh et al., 2020). Often, only a limited number of texts come with manually labeled concepts. Moreover, as humans continuously acquire new concepts, it is desirable for the training framework to discover and incorporate new concepts automatically. Thus, we aim to design a general framework for training CBE-PLMs.

4 C³M: A General Framework for Learning CBE-PLMs

We define the following data portions according to the real-world scenarios. We refer to a dataset with human-annotated concepts as the *source concept dataset*, denoted as $\mathcal{D}_s = \{(x^{(i)}, y^{(i)}, c_s^{(i)})_{i=1}^{n_s}\}$, where n_s denotes the size and $c_s \in \mathbb{R}^{k_s}$ is a vector of k_s concepts from the pre-defined source concept set \mathcal{C}_s . We also consider another dataset without concept labels, referred to as the *unlabeled concept dataset*, denoted as $\mathcal{D}_u = \{(x^{(i)}, y^{(i)})_{i=1}^{n_u}\}$.

The complete dataset is then the combination of these two datasets: $\mathcal{D} = \{\mathcal{D}_s, \mathcal{D}_u\}$. n_s and k_s are typically small, limiting the effectiveness of CBE-PLMs. Specifically, small n_s leads to sparse concept labels in \mathcal{D} , and vanilla CBM cannot be trained on datasets with unlabeled concepts \mathcal{D}_u . Additionally, small k_s indicates that we may not have sufficient information for model prediction. To address these limitations, we propose *ChatGPT-guided Concept augmentation with Concept-level Mixup* (C³M), a novel framework for training CBE-PLMs effectively. As illustrated in Figure 3, at the high level, we augment the concept set \mathcal{C}_s and annotate pseudo concept labels for the unlabeled concept dataset using ChatGPT. Since these pseudo labels are noisy, we propose a novel concept-level MixUp to train the CBE-PLMs effectively on the augmented dataset with noisy concept labels.

4.1 ChatGPT-guided Concept Augmentation

In this section, we detail how to leverage ChatGPT (GPT4) to automatically (1) augment the concept set, and (2) annotate missing concept labels.

4.1.1 Concept Set Augmentation

The goal of concept set augmentation is to automatically generate high-quality concepts using human-specified concepts \mathcal{C}_s as references. These generated concepts should be semantic meaningful and useful for target task prediction. Inspired by LFCBM (Oikarinen et al., 2023), we query ChatGPT with appropriate prompts to generate additional concepts. Our prompts are designed using “in-context learning” (Brown et al., 2020; Min et al., 2022; Xie et al., 2022), and include examples from human annotations. Below is an example of a ChatGPT prompt designed for a sentiment classification task using the IMDB dataset (Maas et al., 2011):

Besides $\{Acting, Storyline, Emotional Arousal, Cinematography\}$, what are the additional important features to judge if a $\{movie\}$ is good or not?

Parentheses represent fields that can be customized for different tasks. The concepts *Acting*, *Storyline*, *Emotional Arousal*, and *Cinematography* are from the source concept set \mathcal{C}_s with labels manually annotated following procedures in Appendix B. Different from LFCBM which generates concepts merely relying on GPT3 (Brown et al., 2020), we further include a small set of human-specified concepts in the prompt to improve the quality of generated concepts. This additional information can help effectively filter out undesired output without additional operations (e.g., deletions). Rarely seen concepts are discarded using a predefined threshold and the remaining generated concepts are referred to as *augmented concepts set* \mathcal{C}_a with size k_a . Results are given in Table 5 in Appendix G.

4.1.2 Noisy Concept Label Annotation

The next step is to automatically annotate unlabeled concepts using *noisy* labels. We again leverage the power of ChatGPT which has been shown to encapsulate significant amounts of human common sense knowledge (Bommasani et al., 2022; OpenAI, 2023; Singh et al., 2023) and show strong performance for some text annotation tasks (Gilardi et al., 2023). As we will also show here, LLMs are surprisingly proficient at identifying language concepts when suitably prompted. Using the same example of movie reviews, the prompt for this step is designed as follows:

- According to the review " $\{text_1\}$ ", the " $\{concept_1\}$ " of the movie is "positive".
- According to the review " $\{text_2\}$ ", the " $\{concept_2\}$ " of the movie is "negative".
- According to the review " $\{text_3\}$ ", the " $\{concept_3\}$ " of the movie is "unknown".
- According to the review " $\{text_i\}$ ", how is the " $\{concept_i\}$ " of the movie? Please answer with one option in "positive, negative, or unknown".

Following a similar "in-context learning" strategy described in Section 4.1.1, Prompts a-c are three human-annotated examples randomly selected to represent positive, negative, and unknown concept labels, respectively. Prompt d is the query instance. The goal is to obtain noisy labels for any given $\{text_i\}$ and $\{concept_i\}$. There are three types of noisy concept annotations:

- Noisy labels for human-specified concepts in \mathcal{D}_s . The resulting dataset $\tilde{\mathcal{D}}_s$ is used to validate

the quality of labels generated by ChatGPT only (See Table 4 in Appendix F).

- Noisy labels for ChatGPT-generated concepts in \mathcal{D}_s . The augmented concept set is denoted as $c_{sa} = (c_s || c_a) \in \mathbb{R}^{k_s+k_a}$, where $||$ refers to the concatenation operator and $c_a \in \mathbb{R}^{k_a}$ stands for the generated concepts. For example, we identify new important concepts such as *Soundtrack* using ChatGPT for the IMDB movie reviews.
- Noisy labels for both human-specified and ChatGPT-generated concepts in unlabeled concept datasets \mathcal{D}_u . The augmented concept set is denoted as $\tilde{c}_{sa} = (\tilde{c}_s || \tilde{c}_a) \in \mathbb{R}^{k_s+k_a}$ and $\tilde{c}_s \in \mathbb{R}^{k_s}, \tilde{c}_a \in \mathbb{R}^{k_a}$ stand for the generated concept labels for human-specified and ChatGPT-generated concepts, respectively.

In summary, we transform the original dataset with sparse concept labels into an augmented dataset with new concepts and noisy labels: $\mathcal{D} = \{\mathcal{D}_s, \mathcal{D}_u\} \rightarrow \tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_{sa}, \tilde{\mathcal{D}}_u\}$. Examples of these two types of queries are illustrated in Appendix J.

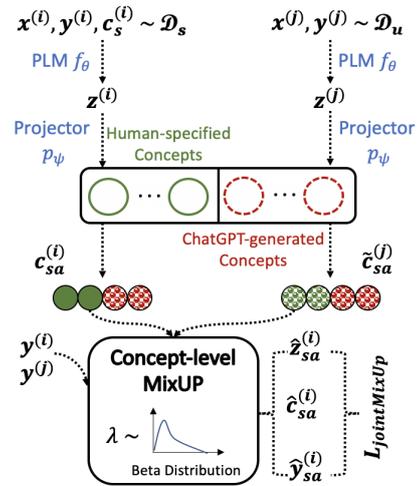


Figure 3: Illustration of the proposed framework C³M.

4.2 Learning from Noisy Concept Labels

While directly training CBE-PLMs on the transformed dataset $\tilde{\mathcal{D}}$ is straightforward, this method's drawback is its equal treatment of human annotations and ChatGPT-generated noisy labels, potentially leading to prediction and interpretation inaccuracies. To improve interpretability and accuracy, we introduce a novel *Concept-level MixUp* (CM) approach. It advocates for a convex behavior of PLMs between human-annotated and ChatGPT-generated concepts, thereby enhancing its robustness against noisy concept labels.

4.2.1 Concept-level MixUp

To better utilize the noisy concept labels, CM first linearly interpolates the texts and concept labels between human-annotated concepts ($\tilde{\mathcal{D}}_{sa}$) and ChatGPT-generated concepts ($\tilde{\mathcal{D}}_u$). Specifically, we interpolate any two text-concept-label ternaries $(x^{(i)}, c^{(i)}, y^{(i)})$, $(x^{(j)}, c^{(j)}, y^{(j)})$ for both their latent representation $(z^{(i)}, z^{(j)})$, concepts $(c^{(i)}, c^{(j)})$, and the task labels $(y^{(i)}, y^{(j)})$ using the MixUp (\cdot) defined as follows:

$$\begin{aligned} \lambda &\sim \text{Beta}(\alpha, \alpha); \quad \hat{\lambda} = \max(\lambda, 1 - \lambda); \\ z^{(i)} &= f_{\theta}(x^{(i)}); \quad z^{(j)} = f_{\theta}(x^{(j)}); \\ \hat{z}^{(i,j)} &= \hat{\lambda}z^{(i)} + (1 - \hat{\lambda})z^{(j)}; \\ \hat{c}^{(i,j)} &= \hat{\lambda}c^{(i)} + (1 - \hat{\lambda})c^{(j)}; \\ \hat{y}^{(i,j)} &= \hat{\lambda}y^{(i)} + (1 - \hat{\lambda})y^{(j)}, \end{aligned} \quad (1)$$

where α is a hyperparameter for the Beta distribution. Notably, $\hat{\lambda} \geq 0.5$ preserves the order of human-annotated concepts and ChatGPT-generated concepts for computing individual loss components in Eq. (4) appropriately. Then, we combine and shuffle human-annotated and ChatGPT-annotated data in the transformed dataset $\tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_{sa}, \tilde{\mathcal{D}}_u\}$:

$$\mathcal{W} = \text{Shuffle}(\tilde{\mathcal{D}}) = \text{Shuffle}(\tilde{\mathcal{D}}_{sa} || \tilde{\mathcal{D}}_u), \quad (2)$$

where $||$ indicates the concatenation of two portions of datasets. Next, we perform MixUp (\cdot) for the i th instance as follows:

$$\begin{aligned} (\hat{z}_{sa}^{(i)}, \hat{c}_{sa}^{(i)}, \hat{y}_{sa}^{(i)}) &= \text{MixUp}(\tilde{\mathcal{D}}_{sa}^{(i)}, \mathcal{W}^{(i)}), \\ (\hat{z}_u^{(i)}, \hat{c}_u^{(i)}, \hat{y}_u^{(i)}) &= \text{MixUp}(\tilde{\mathcal{D}}_u^{(i)}, \mathcal{W}^{(i)}). \end{aligned} \quad (3)$$

Through these steps, we can generate a "mixed version" for each instance in $\tilde{\mathcal{D}}_{sa}$ and $\tilde{\mathcal{D}}_u$, while preserving a larger portion of the original instance.

4.2.2 Loss Function

The loss function $L_{\text{jointMixUp}}$ for training CBE-PLMs with the MixUped dataset is defined below:

$$\begin{aligned} L_{sa} &= L_{\text{joint}}(\hat{z}_{sa}^{(i)}, \hat{c}_{sa}^{(i)}, \hat{y}_{sa}^{(i)}); \\ L_u &= L_{\text{joint}}(\hat{z}_u^{(i)}, \hat{c}_u^{(i)}, \hat{y}_u^{(i)}); \\ L_{\text{jointMixUp}} &= L_{sa} + \tau L_u, \end{aligned} \quad (4)$$

where τ is a hyperparameter and L_{joint} is the joint training loss used in vanilla CBM formulated in Appendix A. In this way, We backpropagate gradients of the mixed noisy concept labels and gold concept labels to update the parameters in CBE-PLMs.

5 Experiments

5.1 Datasets

In this section, we give detailed descriptions of the experimented datasets. Each of the datasets has two components: source concept dataset and unlabeled concept dataset ($\mathcal{D} = \{\mathcal{D}_s, \mathcal{D}_u\}$). Existing datasets with human-annotated concept labels are very limited. One source concept dataset is CEBaB (Abraham et al., 2022; Wu et al., 2022), a common sentiment classification dataset for restaurant reviews. Its corresponding \mathcal{D}_u is the restaurant reviews from the Yelp Dataset¹. We also curate another dataset for movie reviews. Specifically, we randomly sample two portions of reviews from the IMDB datasets (Maas et al., 2011) to represent \mathcal{D}_s and \mathcal{D}_u , respectively. Following a previous NLP work (Cai et al., 2021), we manually annotate the concept labels for \mathcal{D}_s in the movie reviews. More annotation details are included in Appendix B. For convenience, we still refer to these two new datasets as CEBaB and IMDB. Each concept contains three values, i.e., Negative, Positive, and Unknown. As described in Section 4.1, each dataset \mathcal{D} is then transformed into $\tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_{sa}, \tilde{\mathcal{D}}_u\}$. The basic statistics of the transformed datasets and their human-annotated concepts are given in Table 3 in Appendix E and Table 4 in Appendix F, respectively. Note that the last column in Table 4 indicates the accuracy of ChatGPT-labeled concepts in \mathcal{D}_s , as described in Section 4.1. Table 5 in Appendix G provides statistics about augmented concepts. Both the human-annotated and ChatGPT-generated data, along with the framework implementation are released².

5.2 PLM Backbones

We experiment with the same PLM backbones as in the CEBaB paper (Abraham et al., 2022): GPT2 (Radford et al., 2019), BERT (Devlin et al., 2018), RoBERTa (Liu et al., 2019), and BiLSTM (Hochreiter and Schmidhuber, 1997) with CBOW (Mikolov et al., 2013). For better performance, we obtain the representations of the input texts by pooling the embedding of all tokens. Reported scores are the averages of six independent runs, each taking 5 to 40 minutes. More implementation details and parameter values are included in Appendix C and Table 2 in Appendix D.

¹<https://www.kaggle.com/datasets/omkarsabnis/yelp-reviews-dataset>

²https://github.com/Zhen-Tan-dmml/CBM_NLP.git

Table 1: Comparisons of task accuracy and interpretability using CEBaB and IMDB datasets. Metrics for both task and concept labels are written as **Accuracy/Macro F1**. Scores are reported in %. Scores in **bold** indicate that the CBE-PLM under the current setting outperforms its standard PLM counterpart. CM denotes Concept-level MixUp.

Dataset		CEBaB				IMDB			
Model		\mathcal{D}		$\tilde{\mathcal{D}}$		\mathcal{D}		$\tilde{\mathcal{D}}$	
		Task	Concept	Task	Concept	Task	Concept	Task	Concept
PLMs	LSTM	40.57/60.67	-	43.34/64.47	-	68.25/53.37	-	90.5/90.46	-
	GPT2	66.69/77.25	-	67.26/78.81	-	71.67/67.53	-	97.64/97.55	-
	BERT	68.75/78.71	-	71.81/82.58	-	80.5/78.4	-	98.89/98.68	-
	RoBERTa	71.36/80.17	-	73.12/82.64	-	84.1/82.5	-	99.13/99.12	-
CBE-PLMs	LSTM	56.47/67.82	86.46/85.24	54.54/65.84	83.46/84.74	68.5/55.4	72.5/77.5	93.02/91.53	76.92/75.41
	GPT2	64.04/77.75	92.14/92.05	63.57/74.71	90.17/90.13	70.05/69.53	80.6/82.5	96.85/96.81	86.14/88.06
	BERT	67.27/79.24	93.65/92.75	68.23/78.13	89.64/90.45	77.42/74.57	80.2/83.7	97.62/97.58	92.57/92.05
	RoBERTa	70.98/79.89	96.12/95.34	69.85/79.29	91.45/92.23	82.33/80.13	86.7/85.3	98.45/98.12	93.99/94.28
CBE-PLMs-CM	LSTM	-	-	59.67/70.53	88.75/86.67	-	-	94.35/92.32	83.83/84.52
	GPT2	-	-	65.54/77.87	93.58/92.32	-	-	97.89/97.88	89.64/88.25
	BERT	-	-	70.58/80.07	94.43/93.26	-	-	98.18/98.06	94.87/94.32
	RoBERTa	-	-	72.88/81.91	96.3/98.5	-	-	99.69/99.66	96.35/96.36

5.3 Task Accuracy vs Interpretability

Table 1 presents the results for the two original datasets (\mathcal{D}) and their transformed versions ($\tilde{\mathcal{D}}$). We have the following observations:

CBE-PLMs offer interpretability and competitive task prediction performance. Compared to standard PLMs (trained solely with task labels), CBE-PLMs provide concept-level interpretability with only a minor decrease in task prediction. Interestingly, a smaller PLM, i.e., LSTM with CBOW embeddings, achieves improved task accuracy when learning from concept labels. This suggests that the accuracy-interpretability tradeoff in concept learning is not necessary, as opposed to the prevailing view. Concepts can help guide PLMs trained on smaller corpora with fewer parameters towards better prediction performance.

Noisy concept labels can facilitate the training of CBE-PLMs on small datasets. The extremely limited size of the IMDB source concept dataset (deliberately set to 100) yields unsurprisingly low test scores. Transforming \mathcal{D} into $\tilde{\mathcal{D}}$ using ChatGPT for noisy labeled concept instances leads to significant improvements in both concept and task predictions for CBE-PLMs-CM.

Uncritical learning from noisy concept labels can impair performance. Results for CEBaB in Table 1 demonstrate that, learning from the transformed dataset $\tilde{\mathcal{D}}$ directly leads to inferior performance for CBE-PLMs. Unlike IMDB, the source concept dataset in CEBaB contains sufficient training instances, therefore, enforcing CBE-PLMs to learn from noisy concept labels will undesirably mislead the model, exacerbating both the concept and task prediction performance.

CBE-PLMs-CM trained via the proposed C^3M framework consistently deliver superior interpretability-utility trade-offs. By encouraging the CBE-PLMs to linearly interpolate between examples with gold-labeled concepts and those with ChatGPT-generated concepts, the model is able to extract useful semantic knowledge meanwhile becoming robust to noisy concept labels. The result is promising: We achieve the best concept-level prediction (interpretability measure) without sacrificing the task prediction performance, and in some cases, CBE-PLMs trained through C^3M can even outperform their standard PLM counterparts.

5.4 Explainable Predictions

A unique advantage of CBMs is that its decision rules can be interpreted as a linear combination of comprehensible variables (Koh et al., 2020). Inheriting this strength, our proposed CBE-PLMs can deliver intuitive concept-level explanations for predictions by assessing the activations of each concept. We measure concept contribution using the product of activation and the corresponding weight in the linear label predictor g_ϕ (Oikarinen et al., 2023). Concepts with negative activation are designated as ‘‘Neg Concept’’. We highlight the concepts contributing the most in our visualizations. Visualization results are demonstrated in Figure 4 for a toy example, while real-world CEBaB and IMDB case studies can be found in Appendix H. These visualizations provide new intriguing insights into real-world applications. For instance, negative concepts (e.g., Service) contribute more to the final prediction of positive sentiment in Figure 4, making the predicted sentiment second highest ($Y = 4$) rather than the highest ($Y = 5$). Moreover, inter-

pretability results such as Figure 6 in Appendix H imply that concepts such as “Food” and “Ambiance” weigh more heavily in customers’ restaurant evaluations compared to “Noise” and “Menu Variety”.

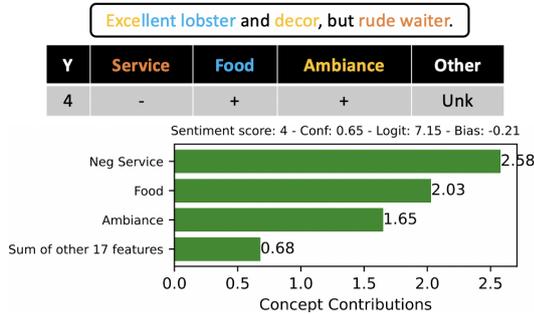


Figure 4: Illustration of the explainable prediction for a toy example in restaurant review sentiment analysis.

5.5 Test-time Intervention

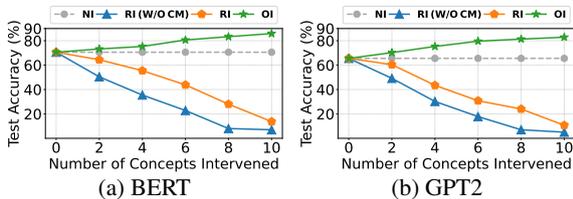


Figure 5: The results of Test-time Intervention. "NI" denotes "no intervention", "RI (W/O CM)" denotes "random intervention on CBE-PLMs without the concept-level MixUp", "RI" denotes "random intervention on CBE-PLMs", and "OI" denotes "oracle intervention".

Another strength of CBE-PLMs is that they allow test-time concept intervention (inherited from CBMs), facilitating deeper, user-friendly interactions. To assess this strength, we follow Koh et al. (2020) to intervene in the predicted concepts and investigate the impact of such interventions on test-time prediction accuracy. Concept mispredictions arise from ChatGPT’s incorrect labels or inaccurate concept activation. Recall that the input of the task label predictor is the predicted concept activations $\hat{a} = p_\phi(f_\theta(x))$ rather than the predicted ternary concepts \hat{c} . In a concept-level intervention I , the activation \hat{a}_j of the j th concept with a target concept c_j is set to the 5th, 95th, or 50th percentile of \hat{a}_j over the training distribution for Negative, Positive, or Unknown c_j respectively. Multiple concepts can be intervened upon by replacing all related predicted concept activations and updating the prediction. Experiments were conducted on the transformed version \tilde{D} of the CEBaB dataset. Figure 5 exhibits results for CBE-PLMs using BERT and GPT2 as the PLM backbones (with similar observations for LSTM and

RoBERTa). A case study is further illustrated in Appendix I. The results reveal that task accuracy improves substantially when more concepts are corrected by the oracle. Additionally, while the performance of CBE-PLMs declines as more concepts are intervened upon incorrectly (randomly), the proposed concept-level MixUp effectively mitigates this impact. Notably, the decline in performance is marginal when only two concepts are erroneously intervened upon. These findings underscore the pronounced advantages of test-time intervention for CBE-PLMs trained through C³M. First, domain experts can interact with the model to rectify any inaccurately predicted concept values. Second, in reality, even experts might inadvertently implement incorrect interventions. Yet, despite this susceptibility, our proposed concept-level MixUp strategy effectively curbs performance degradation, particularly when inaccuracies affect only a small subset of the intervention. This attests to the robustness of the proposed framework.

6 Conclusion

Our analysis began with an exhaustive examination of three training strategies, identifying joint training as the most efficacious. Further, we proposed the C³M framework, designed to streamline the training process of CBE-PLMs in the presence of incomplete concept labels. Moreover, we showcased the interpretability of our models in their decision-making process and elucidated how this comprehensibility can be harnessed to boost test accuracy via concept intervention.

Outlook: Our research lays the groundwork for future studies focused on enhancing the transparency and robustness of PLMs. We foresee that CBE-PLMs could potentially show more resilience to data biases compared to standard PLMs, which have been known to display biased performance due to spurious correlations between sensitive attributes (e.g., gender) and task labels (Wang and Culotta, 2021; Udomcharoenchaikit et al., 2022). For instance, a biased PLM might wrongly infer patterns like female users writing more extreme reviews while male users tend towards moderate ones. CBE-PLMs, by focusing on concept labels and relying solely on these concepts for classifications, might reduce such biases. If the concepts are not associated with sensitive attributes and their relationship with task labels is consistent, CBE-PLMs could offer enhanced fairness.

Limitations

While our approach presents a significant step towards more interpretable pretrained language models, several limitations warrant further exploration. First, our approach relies heavily on the accuracy of the predefined concepts. Despite the promising results, this dependency raises the issue of potential bias present in the concept selection process (for both human-specified and ChatGPT-generated concepts). If a concept is not well-defined or if important concepts are missing, this could lead to incomplete or skewed interpretations. Second, the methodology proposed in this paper has not been experimented on very large language models, such as Bloom (Scao et al., 2022). The core idea of this framework is to utilize large language models (LLMs) to provide explanations for comparatively lighter-weight pretrained language models (PLMs). Nevertheless, the proposed framework is of a universal nature and should be compatible with any PLMs. Investigations utilizing larger PLMs are reserved for future research endeavors. Third, the process of prompting large language models to generate concept labels remains somewhat of an art. While we have proposed a systematic method for constructing desired prompts, the performance of the model may still be sensitive to the quality and structure of these prompts. Lastly, while our proposed method shows promising results in English language tasks, it has not been tested extensively on other languages. This restricts its applicability in a multilingual setting. Future work should extend this method to other languages and conduct cross-lingual analysis. We hope future research will build upon our work to address these limitations, moving us closer to truly interpretable, responsible, and universally applicable language models.

Ethics Statement

In conducting this research, we strictly adhered to the [ACL Ethics Policy](#). All data used in our work were either publicly available or anonymized, ensuring no personally identifiable information was involved. The work presented in this paper significantly contributes to the field of natural language processing and machine learning. By improving the interpretability of pre-trained language models, we are contributing to the creation of more transparent and trustworthy AI systems. This advancement is expected to have broad-ranging impacts across numerous domains that increasingly rely on AI,

including healthcare, education, business, and finance, enhancing decision-making processes and user interaction with AI systems. However, the increased efficacy of these models could also raise potential societal concerns if not used responsibly. The misuse of these advanced NLP technologies could lead to privacy breaches, the propagation of misinformation, or the amplification of existing biases in data. As with any powerful technology, it is essential to consider its ethical implications and manage its deployment with care to ensure it's used for the betterment of society. Our work also underscores the need for continual research into strategies that mitigate potential bias in AI systems and protect user privacy. As researchers, we are committed to working towards these goals and urge those employing this technology to adhere to the same principles.

References

- Eldar D Abraham, Karel D'Oosterlinck, Amir Feder, Yair Gat, Atticus Geiger, Christopher Potts, Roi Reichart, and Zhengxuan Wu. 2022. Cebab: Estimating the causal effects of real-world concepts on nlp model behavior. *Advances in Neural Information Processing Systems*, 35:17582–17596.
- Yonatan Belinkov and James Glass. 2019. Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 7:49–72.
- David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. 2019. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. 2023. Language models can explain neurons in language models. <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny

- Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2022. [On the opportunities and risks of foundation models](#).
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Hongjie Cai, Rui Xia, and Jianfei Yu. 2021. Aspect-category-opinion-sentiment quadruple extraction with implicit aspects and opinions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 340–350.
- Lu Cheng, Kush R Varshney, and Huan Liu. 2021. Socially responsible ai algorithms: Issues, purposes, and challenges. *Journal of Artificial Intelligence Research*, 71:1137–1181.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang. 2022. Black-box prompt learning for pre-trained language models. *arXiv preprint arXiv:2201.08531*.
- Erik Engleson and Hossein Azizpour. 2021. Generalized jensen-shannon divergence loss for learning with noisy labels. *Advances in Neural Information Processing Systems*, 34:30284–30297.
- Andrea Galassi, Marco Lippi, and Paolo Torrioni. 2020. Attention in natural language processing. *IEEE transactions on neural networks and learning systems*, 32(10):4291–4308.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks. *arXiv preprint arXiv:2303.15056*.
- Yash Goyal, Amir Feder, Uri Shalit, and Been Kim. 2019. Explaining classifiers with causal concept effect (cace). *arXiv preprint arXiv:1907.07165*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, pages 2668–2677. PMLR.
- Evgeny Kim and Roman Klinger. 2018. Who feels what and why? annotation of a literature corpus with semantic roles of emotions. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1345–1359.
- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. 2020. Concept bottleneck models. In *International Conference on Machine Learning*, pages 5338–5348. PMLR.
- Q Vera Liao and Jennifer Wortman Vaughan. 2023. Ai transparency in the age of llms: A human-centered research roadmap. *arXiv preprint arXiv:2306.01941*.
- Yang Liu, Hao Cheng, and Kun Zhang. 2022. Identifiability of label noise transition matrix. *arXiv preprint arXiv:2202.02016*.
- Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Max Losch, Mario Fritz, and Bernt Schiele. 2019. Interpretability beyond classification output: Semantic bottleneck networks. *arXiv preprint arXiv:1907.10882*.
- Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.

- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150.
- Andreas Madsen, Siva Reddy, and Sarath Chandar. 2022. Post-hoc interpretability for neural nlp: A survey. *ACM Computing Surveys*, 55(8):1–42.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*.
- Saumitra Mishra, Bob L Sturm, and Simon Dixon. 2017. Local interpretable model-agnostic explanations for music content analysis. In *ISMIR*, volume 53, pages 537–543.
- Jesse Mu and Jacob Andreas. 2020. Compositional explanations of neurons. *Advances in Neural Information Processing Systems*, 33:17153–17163.
- Renáta Németh, Domonkos Sik, and Fanni Máté. 2020. Machine learning of concepts hard even for humans: The case of online depression forums. *International Journal of Qualitative Methods*, 19:1609406920949338.
- Tuomas Oikarinen, Subhro Das, Lam M Nguyen, and Tsui-Wei Weng. 2023. Label-free concept bottleneck models. In *The Eleventh International Conference on Learning Representations*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch. In *NeurIPS*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Alexis Ross, Ana Marasović, and Matthew E Peters. 2021. Explaining nlp models via minimal contrastive editing (mice). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3840–3852.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Manmeet Singh, Vaisakh SB, Neetiraj Malviya, et al. 2023. Mind meets machine: Unravelling gpt-4’s cognitive psychology. *arXiv preprint arXiv:2303.11436*.
- Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. 2020. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608.
- Can Udomcharoenchaikit, Wuttikorn Ponwitayarat, Patomporn Payoungkhamdee, Kanruethai Masuk, Weerayut Buaphet, Ekapol Chuangsuwanich, and Sarana Nutanong. 2022. Mitigating spurious correlation in natural language understanding with counterfactual inference. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11308–11321.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. *Advances in neural information processing systems*, 33:12388–12401.
- Zhao Wang and Aron Culotta. 2021. Robustness to spurious correlations in text classification via automatically generated counterfactuals. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14024–14031.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Huggingface’s transformers: State-of-the-art natural language processing](#).
- T Wu, M Tulio Ribeiro, J Heer, and D Weld. 2021. Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models. In *Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)*.
- Zhengxuan Wu, Karel D’Oosterlinck, Atticus Geiger, Amir Zur, and Christopher Potts. 2022. Causal proxy models for concept-based model explanations. *arXiv preprint arXiv:2209.14279*.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2022. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*.
- Jie Yang, Yue Zhang, Linwei Li, and Xingxuan Li. 2017. Yedda: A lightweight collaborative text span annotation tool. *arXiv preprint arXiv:1711.03759*.

Kayo Yin and Graham Neubig. 2022. Interpreting language models with contrastive explanations. *arXiv preprint arXiv:2202.10419*.

Mateo Espinosa Zarlenga, Pietro Barbiero, Gabriele Ciravegna, Giuseppe Marra, Francesco Giannini, Michelangelo Diligenti, Frederic Precioso, Stefano Melacci, Adrian Weller, Pietro Lio, et al. 2022. Concept embedding models. In *NeurIPS 2022-36th Conference on Neural Information Processing Systems*.

Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.

Wenxuan Zhang, Xin Li, Yang Deng, Lidong Bing, and Wai Lam. 2022. A survey on aspect-based sentiment analysis: tasks, methods, and challenges. *IEEE Transactions on Knowledge and Data Engineering*.

Jinhua Zhu, Yingce Xia, Lijun Wu, Di He, Tao Qin, Wengang Zhou, Houqiang Li, and Tiejun Liu. 2020. Incorporating bert into neural machine translation. In *International Conference on Learning Representations*.

A Definitions of Training Strategies

Given a text input $x \in \mathbb{R}^d$, concepts $c \in \mathbb{R}^k$ and its label y , the strategies for fine-tuning the text encoder f_θ , the projector p_ψ and the label predictor g_ϕ are defined as follows:

i) Vanilla fine-tuning a PLM: The concept labels are ignored, and then the text encoder f_θ and the label predictor g_ϕ are fine-tuned either as follows:

$$\theta, \phi = \operatorname{argmin}_{\theta, \phi} L_{CE}(g_\phi(f_\theta(x)), y),$$

or as follows (frozen text encoder f_θ):

$$\phi = \operatorname{argmin}_{\phi} L_{CE}(g_\phi(f_\theta(x)), y),$$

where L_{CE} indicates the cross-entropy loss. In this work we only consider the former option for its significant better performance.

ii) Independently training PLM with the concept and task labels: The text encoder f_θ , the projector p_ψ and the label predictor g_ϕ are trained separately with ground truth concepts labels and task labels as follows:

$$\begin{aligned} \theta, \psi &= \operatorname{argmin}_{\theta, \psi} L_{CE}(p_\psi(f_\theta(x)), c), \\ \phi &= \operatorname{argmin}_{\phi} L_{CE}(g_\phi(c), y). \end{aligned}$$

During inference, the label predictor will use the output from the projector rather than the ground-truth concepts.

iii) Sequentially training PLM with the concept and task labels: We first learn the concept encoder as the independent training strategy above, and then use its output to train the label predictor:

$$\phi = \operatorname{argmin}_{\phi} L_{CE}(g_\phi(p_\psi(f_\theta(x)), y)).$$

iv) Jointly training PLM with the concept and task labels: Learn the concept encoder and label predictor via a weighted sum L_{joint} of the two objectives described above:

$$\begin{aligned} \theta, \psi, \phi &= \operatorname{argmin}_{\theta, \psi, \phi} L_{joint}(x, c, y) \\ &= \operatorname{argmin}_{\theta, \psi, \phi} [L_{CE}(g_\phi(p_\psi(f_\theta(x)), y) \\ &\quad + \gamma L_{CE}(p_\psi(f_\theta(x)), c)]. \end{aligned}$$

It’s worth noting that the CBE-PLMs trained jointly are sensitive to the loss weight γ . We report the most effective results here, tested value for γ are given in Table 2 in Appendix D.

B Details of the Manual Concept Annotation for the IMDB Dataset

Our annotation policy is following a previous work (Cai et al., 2021) for NLP datasets annotating. For the IMDB dataset, we annotate the four concepts (Acting, Storyline, Emotional Arousal, Cinematography) manually. Even though the concepts are naturally understandable by humans, two Master students familiar with sentiment analysis are selected as annotators for independent annotation with the annotation tool introduced by Yang et al. (2017). The strict quadruple matching F1 score between two annotators is 85.74%, which indicates a consistent agreement between the two annotators (Kim and Klinger, 2018). In case of disagreement, a third expert will be asked to make the final decision.

C Implementation Detail

In this section, we provide more details on the implementation settings of our experiments. Specifically, we implement our framework with PyTorch (Paszke et al., 2017) and HuggingFace (Wolf et al., 2020) and train our framework on a single 80GB Nvidia A100 GPU. We follow a prior work (Abraham et al., 2022) for backbone implementation. All backbone models have a maximum token number of 512 and a batch size of 8. We use the Adam optimizer to update the backbone, projector, and label predictor according to Section 3.1.

The values of other hyperparameters (Table 2 in Appendix D) for each specific PLM type are determined through grid search. We run all the experiments on an Nvidia A100 GPU with 80GB RAM.

D Parameters and Notations

In this section, we provide used notations in this paper along with their descriptions for comprehensive understanding. We also list their experimented values and optimal ones, as shown in Table 2.

E Statistics of Data Splits

The Statistics and split policies of the experimented datasets, including the source concept dataset \mathcal{D}_s , the unlabeled concept dataset \mathcal{D}_u , and their augmented versions. The specific details are presented in Table 3.

F Statistics of Human-Annotated Concepts

The Statistics of Human-Annotated Concepts in both CEBaB and IMDB datasets. We also include the accuracy of ChatGPT’s concept prediction here. The specific details are presented in Table 4.

G Statistics of Concepts in Transformed Datasets

The Statistics and split policies of the transformed datasets of experimented datasets are presented in Table 5.

H More Results on Explainable Predictions

Case studies on explainable predictions for both CEBaB and IMDB datasets are given in Figure 6 and Figure 7 respectively.

I A case study on Test-time Intervention

We present a case study of Test-time Intervention using an example from the transformed unlabeled concept data $\hat{\mathcal{D}}_u$ of the CEBaB dataset, as shown in Figure 8. The first row displays the target concept labels generated by ChatGPT. The second row shows the predictions from the trained CBE-PLM model, which mispredicts two concepts ("Waiting time" and "Waiting area"). The third row demonstrates test-time intervention using ChatGPT as the oracle, which corrects the predicted task labels.

Finally, the fourth row implements test-time intervention with a human oracle, rectifying the concept that ChatGPT originally mislabeled.

J Examples of Querying ChatGPT

In this paper, we query ChatGPT for 1) augmenting the concept set, and 2) annotate missing concept labels. Note that in practice, we query ChatGPT (GPT4) via [OpenAI API](#). Here we demonstrate examples from the ChatGPT (GPT4) [GUI](#) for better illustration. The illustrations are given in Figure 9 and Figure 10.

Table 2: Key parameters in this paper with their annotations and evaluated values. Note that **bold** values indicate the optimal ones.

Notations	Specification	Definitions or Descriptions	Values
max_len	-	maximum token number of input	128 / 256 / 512
batch_size	-	batch size	8
plm_epoch	-	maximum training epochs for PLM and Projector	20
clf_epoch	-	maximum training epochs for the linear classifier	20
hidden_dim	-	hidden dimension size	128
emb_dim	LSTM	embedding dimension for LSTM	300
lr	LSTM	learning rate when the backbone is LSTM	1e-1 / 1e-2 / 5e-2 / 1e-3 / 1e-4
	GPT2	learning rate when the backbone is GPT2	1e-3 / 5e-3 / 1e-4 / 5e-4 / 1e-5
	BERT	learning rate when the backbone is BERT	1e-4 / 5e-4 / 1e-5 / 3e-5 / 5e-5
	RoBERTa	learning rate when the backbone is RoBERTa	1e-4 / 5e-4 / 1e-5 / 3e-5 / 5e-5
\gamma	-	loss weight in the joint loss L_{joint}	0.1 / 0.3 / 0.5 / 0.7 / 1.0
\tau	-	loss weight in the joint-MixUp loss $L_{jointMixUp}$	0.1 / 0.5 / 1.0 / 1.5 / 2.0

Table 3: Statistics of experimented datasets. k denotes the number of concepts.

Dataset	\mathcal{D}_s		\mathcal{D}_u		$\tilde{\mathcal{D}}_{sa}$		$\tilde{\mathcal{D}}_u$		Task
	Train/Dev/Test	k	Train/Dev/Test	k	Train/Dev/Test	k	Train/Dev/Test	k	
CEBaB	1755/1673/1685	4	2000/500/500	0	1755/1673/1685	10	2000/500/500	10	5-way classification
IMDB	100/50/50	4	1000/1000/1000	0	100/50/50	8	1000/1000/1000	8	2-way classification

Table 4: Statistics of human-specified concepts in \mathcal{D}_s and the accuracy of ChatGPT’s concept prediction.

Dataset (\mathcal{D}_s)	Concept	Negative	Positive	Unknown	Total	ChatGPT Acc.
CEBaB	Food	1693 (33.1%)	2087 (40.8%)	1333 (26.1%)	5113	77.9%
	Ambiance	787 (15.4%)	994 (19.4%)	3332 (65.2%)	5113	69.2%
	Service	1249 (24.4%)	1397 (27.3%)	2467 (48.2%)	5113	78.7%
	Noise	645 (12.6%)	442 (8.6%)	4026 (78.7%)	5113	77.7%
IMDB	Acting	76 (38%)	66 (33%)	58 (29%)	200	73.0%
	Storyline	80 (40%)	77 (38.5%)	43 (21.5%)	200	64.0%
	Emotional Arousal	74 (37%)	73 (36.5%)	53 (26.5%)	200	60.5%
	Cinematography	118 (59%)	43 (21.5%)	39 (19.4%)	200	66.5%

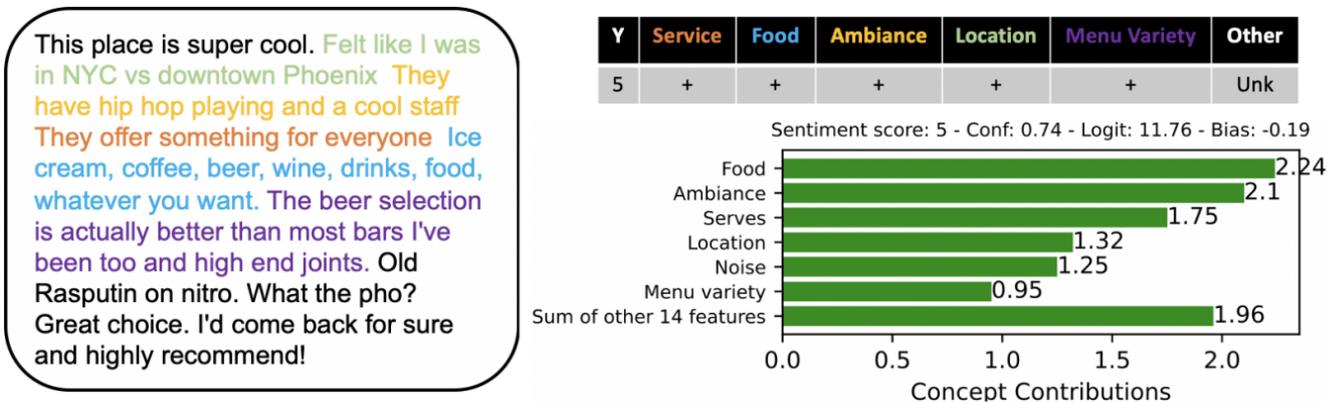


Figure 6: Illustration of the explainable prediction for an example from the CEBaB dataset.

Table 5: Statistics of concepts in transformed datasets (\tilde{D}). Human-specified concepts are underlined. Concepts shown in gray are not used in experiments as the portion of the "Unknown" label is too large.

Dataset	Concept	Negative	Positive	Unknown	Total
CEBaB	<u>Food</u>	2043(25.2%)	4382(54.0%)	1688(20.8%)	8113
	<u>Ambiance</u>	868(10.7%)	1659(20.4%)	5586(68.9%)	8113
	<u>Service</u>	1543(19.0%)	2481(30.6%)	4089(50.4%)	8113
	<u>Noise</u>	668(8.2%)	477(5.9%)	6968(85.9%)	8113
	Cleanliness	55(0.7%)	610(7.5%)	7448(91.8%)	8113
	Price	714(8.8%)	527(6.5%)	6872(84.7%)	8113
	Location	303(3.7%)	2598(32.0%)	5212(64.2%)	8113
	Menu Variety	238(2.9%)	2501(30.8%)	5374(66.2%)	8113
	Waiting Time	572(7.1%)	608(7.5%)	6933(85.5%)	8113
	Waiting Area	267(3.3%)	1136(14.0%)	6710(82.7%)	8113
	Parking	53(0.7%)	107(1.3%)	7953(98.0%)	8113
	Wi-Fi	9(0.1%)	39(0.5%)	8065(99.4%)	8113
	Kids-Friendly	15(0.2%)	536(6.6%)	7562(93.2%)	8113
	IMDB	<u>Sentiment</u>	1624(50.7%)	1576(49.2%)	0(0.0%)
<u>Acting</u>		663(20.7%)	1200(37.5%)	1337(41.8%)	3200
<u>Storyline</u>		1287(40.2%)	1223(38.2%)	690(21.6%)	3200
<u>Emotiona Arousal</u>		1109(34.7%)	1136(35.5%)	955(29.8%)	3200
Cinematography		165(5.2%)	481(15.0%)	2554(79.8%)	3200
Soundtrack		107(3.3%)	316(9.9%)	2777(86.8%)	3200
Directing		537(16.8%)	850(26.6%)	1813(56.7%)	3200
Background Setting		288(9.0%)	581(18.2%)	2331(72.8%)	3200
Editing	304(9.5%)	240(7.5%)	2656(83.0%)	3200	

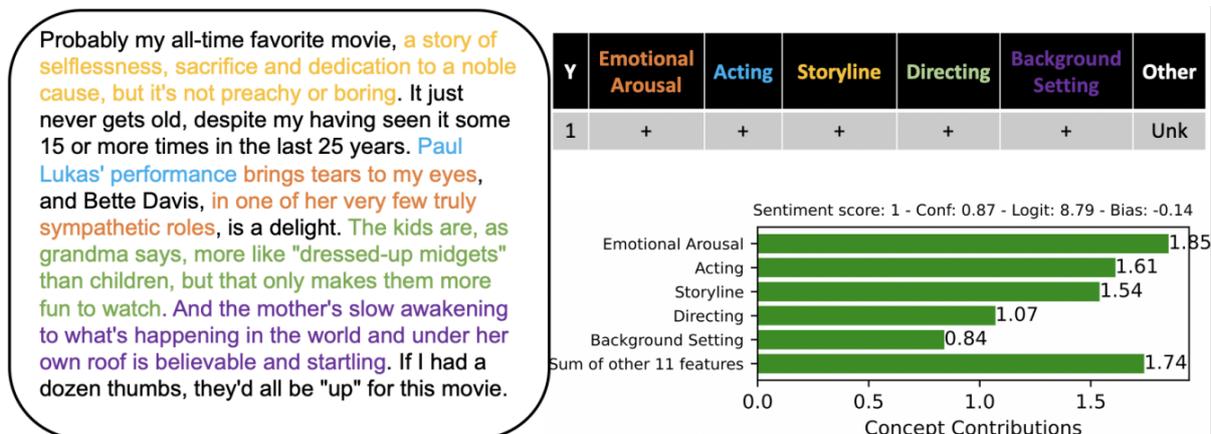


Figure 7: Illustration of the explainable prediction for an example from the IMDB dataset.

Finally 85c bakery in Las Vegas! :D This place was one of my favorite places to go when I lived in Cali, so many asian style bread and sea salt cream drinks. Bread and drinks here are above average taste (at least to me) and also price is reasonable. This time I bought 4 bread and a fresh strawberry milk tea and it was like 12 bucks, which I think it's reasonable in the price wise. The store is also clean, spacious, and modernly decorated! Be aware that there might be a huge line when you come around lunch time in weekends.

Y/P	Cleanliness	Food	Ambiance	Location	Menu Variety	Price	Waiting time	Waiting area	Other
4	+	+	+	+	+	+	-	+	Unk
5	+	+	+	+	+	+	Unk	Unk	Unk
4	+	+	+	+	+	+	-	+	Unk
4	+	+	+	+	+	+	-	Unk	Unk

Figure 8: Illustration of the explainable prediction for an example from the transformed unlabeled concept data \tilde{D}_u of the CEBaB dataset. The brown box with dash lines indicates the test-time intervention on corresponding concepts.

DI Besides Acting, Storyline, Emotional Arousal, Cinematography, what are the additional important features to judge if a movie is good or not?

There are several other elements that contribute to the overall quality and success of a movie:

- Direction:** The director's vision and the ability to effectively bring the script to life are critical.
- Screenplay:** The quality of the script, including dialogue and plot development, plays a significant role.
- Editing:** Good editing is essential for creating a smooth narrative flow and pacing.
- Sound Design and Score:** Effective use of sound and music can significantly enhance the mood and emotions conveyed in a film.
- Production Design:** This involves the sets, costumes, and props that help to create a believable and engaging world.
- Performance of the Supporting Cast:** The performances of secondary characters can greatly influence the overall quality of a film.
- Special Effects:** For certain genres, the quality of special effects can greatly impact the viewing experience.
- Originality:** Fresh and unique concepts or perspectives can make a movie stand out.
- Themes and Messages:** A movie can be judged on its exploration of themes and the messages it conveys.
- Cultural Impact:** Films that spark significant cultural or social conversation can be seen as successful.
- Audience and Critical Reception:** The responses and reviews from both audiences and critics often play a role in determining a film's perceived quality.

Figure 9: The illustration of querying ChatGPT for additional concepts for the IMDB dataset.

DI a. According to the review "If you like original gut wrenching laughter you will like this movie. If you are young or old then you will love this movie, hell even my mom liked it.

Great Camp!!!", the "Emotional Arousal" of the movie is "positive".

b. According to the review "Besides being boring, the scenes were oppressive and dark. The movie tried to portray some kind of moral, but fell flat with its message. What were the redeeming qualities?? On top of that, I don't think it could make librarians look any more unglamorous than it did.", the "cinematography" of the movie is "negative".

c. According to the review "If you like original gut wrenching laughter you will like this movie. If you are young or old then you will love this movie, hell even my mom liked it.

Great Camp!!!", the "acting" of the movie is "unknown".

d. According to the review "This a fantastic movie of three prisoners who become famous. One of the actors is george clooney and I'm not a fan but this roll is not bad. Another good thing about the movie is the soundtrack (The man of constant sorrow). I recommend this movie to everybody. Greetings Bart", how is the "Background Setting" of the movie? Please answer with one option in "positive, negative, or unknown" without any explanation.

The "Background Setting" of the movie based on the given review is "unknown".

Figure 10: The illustration of querying ChatGPT for annotating a missing concept label for the IMDB dataset.