LINEARLY INTERPRETABLE CONCEPT EMBEDDING MODEL FOR TEXT CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite their success, Large-Language Models (LLMs) still face criticism due to their lack of interpretability. Traditional post-hoc interpretation methods, based on attention and gradient-based analysis, offer limited insight as they only approximate the model's decision-making processes and have been proved to be unreliable. For this reason, Concept-Bottleneck Models (CBMs) have been lately proposed in the textual field to provide interpretable predictions based on human-understandable concepts. However, CBMs still face several criticisms for their architectural constraints limiting their expressivity, for the absence of task-interpretability when employing non-linear task predictors and for requiring extensive annotations that are impractical for real-world text data. In this paper we address these challenges by proposing a novel Linearly Interpretable Concept Embedding Model (LICEM) going beyond the current accuracy-interpretability trade-off. LICEM classification accuracy is better than existing interpretable models and matches black-box models. The provided explanations are more plausible and useful with respect to existing solutions, as attested in a user study. Finally, we show our model can be trained without requiring any concept supervision, as concepts can be automatically predicted by the same LLM backbone.

026 027 028

029

025

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

In recent years, Large-Language Models (LLMs) have revolutionized the way we approach text 031 interpretation, generation, and classification (Devlin et al., 2018; Radford et al., 2018; Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023). Despite their success, LLMs' reliability 033 is insufficient, due to the occurrence of hallucinations (Bang et al., 2023; Huang et al., 2023) and 034 the inconsistency of self-provided explanations that often do not reflect the actual decision-making process (Ye & Durrett, 2022; Madsen et al., 2024; Turpin et al., 2024). Furthermore, traditional explainability methods mainly rely on the attention mechanism (Jain & Wallace, 2019; Wiegreffe 037 & Pinter, 2019) and gradient-based analysis (Chefer et al., 2021b), both of which have been shown 038 to provide limited interpretability as they are often unreliable (Adebayo et al., 2018; Taimeskhanov et al., 2024) and only show where the model looks, but not what it sees in a given input (Rudin, 2019; Fel et al., 2023; Poeta et al., 2023). 040

041 For this reason, Concept-Bottleneck Models (CBMs) (Koh et al., 2020) have been recently proposed 042 in the textual field to improve the interpretability of LLM predictions (Tan et al., 2024b;a). In CBMs, 043 an intermediate layer outputs a set of human-understandable symbols, commonly referred to as 044 concepts, before providing the final classification. Furthermore, CBMs allows concept interventions, i.e., counterfactual predictions based on slight modifications of the predicted concepts. However, they still present several limitations: i) the concept bottleneck architecture prevents to achieve high 046 classification accuracy, particularly in real-world text scenarios where complete concept represen-047 tations are difficult to obtain; ii) when employing non-linear task predictors, standard CBMs are 048 not task-interpretable, i.e., the decision process from the concepts to the final classification is noninterpretable; iii) CBMs concept annotation is expensive and existing generative concept annotation approaches require the employment of multiple modules. 051

This paper addresses these challenges, by proposing a novel Linearly-Interpretable Concept Em bedding Model (LICEM) providing the final prediction in terms of an interpretable linear equation working over concept embeddings. In LICEM both the weights and the concept scores of the linear



Figure 1: Left, LICEM predicting the sentiment of a drug review (Gräßer et al., 2018). LICEM provides accurate predictions and reveals its decision-making process. Middle, LICEM provides the best accuracy/interpretability trade-off. Right, models' concept F1 scores, when increasing the number of concept annotations. Self-LICEM achieves high scores without requiring concept labels.

equation are predicted for each sample. As shown in Figure 1, left, in the context of a drug review 072 classification, LICEM not only allows identifying the important concepts in the text, such as 'Effective' or 'Side Effects', but also to interpret by-design its local decision process. In the experiment, 074 we positively answer all our research questions. In particular, we show that i) LICEM achieves 075 higher accuracy than existing task-interpretable models while matching or surpassing black-box 076 methods (Figure 1, middle); ii) LICEM explanations are more plausible and useful with respect to 077 existing solutions by means of a user study; iii) LICEM can be trained without any concept annotation (Self-LICEM), as concepts can be automatically predicted by its LLM backbone, providing higher 078 concept accuracy than an existing method (Figure 1, right). 079

081

082

066

067

068

069

071

073

BACKGROUND 2

083 084 **CBMs.** CBMs (Koh et al., 2020; Tan et al., 2024a) are interpretable models that break the standard 085 end-to-end learning paradigm into the training of two neural modules $f \circ q$. The concept encoder $g: X \to C$ maps raw features $x \in X \subset \mathbb{R}^d$ into m higher-level abstractions $c \in C \subset [0, 1]^m$ (i.e., the concepts); the task encoder $f: C \to Y$ predicts n downstream classes based on the learned 087 concepts $\hat{y} = f(g(x)), y \in Y \subset [0,1]^n$. CEMs (Espinosa Zarlenga et al., 2022; Kim et al., 2023) 880 decompose the concept encoder into two functions $g = q \circ h$. The inner function $h: X \to H \subset \mathbb{R}^b$ 089 provides a representation of an input sample, while $q: H \to \mathbb{C}$ maps this representation into m k-dimensional concept embeddings $\mathbf{c} \in \mathbb{C} \subset \mathbb{R}^{m,k}$. The concept prediction \hat{c}_j is then given by a 090 091 neural function over the concept embeddings $\hat{c}_i = s(\mathbf{c}_i)$, where s is shared among the m concepts. A 092 probabilistic formalization can be found in Appendix A.1. However, the interpretability of the CEM task predictor $f(\mathbf{c})$ is limited, as the individual dimensions of concept embeddings lack meaningful 094 interpretation. Additionally, adapting CEM to text scenarios remains an open question. 095

LLM-based Textual Encoders. When considering transformer models, there exist several methods 096 for implementing a text encoder h(x). An immediate choice is to employ an encoder-only architecture, such as BERT (Devlin et al., 2018), and extracting the embedding associated to the [CLS] token. 098 However, as recently shown in Jiang et al. (2023b), one can also exploit the remarkable performance 099 of existing decoder-only LLMs. The architecture of an LLM can be conceptualized as comprising two 100 distinct components: the stacked decoder blocks which are responsible for generating a contextualized 101 representation e, and a classification head that processes this representation to predict the next token. 102 The first can be interpreted as sampling a representation by the distribution p_h , where h is the 103 pre-trained LLM (without the classification head). This distribution can be conditioned toward the 104 generation of specific embeddings by using a prompt t, i.e., $e \sim p_h(e|t, x)$. To induce the LLM to 105 generate an embedding which is representative of a sentence, Jiang et al. (2023b) proposed to use the prompt "this sentence: '[sentence]' means in one word: " and substituting '[sentence]' with 106 the sequence of tokens x. In order to obtain a rich representation of the x sequence, exploiting the 107 available knowledge in the pretrained LLM, we thus use the embedding $e \sim p_h(e|t, x)$.



Figure 2: LICEMs visualization. Using a pretrained LLM model, we i) require it to provide an encoding of the input text following; ii) prompt the LLM to generate the concepts predictions \hat{c} (e.g., Side Effects = 0.8) in Self-LICEM, while in LICEM they are provided by a concept embedding layer; iii) make the final prediction in an interpretable way by first predicting the equation weights w_{ij} (e.g., $w_{\text{Side Effects}} = -1.8$) for predicting the *i*-th class, then executing the resulting linear equation.

3 Method

133 134 135

127

128

129

130

131 132

In this paper, we aim to develop an interpretable concept-based model for text classification. To 136 achieve this, we need to rely on rich text and concept representations. We create an LLM-based 137 CEM by first using an LLM to model the text encoder and extract an embedding e (as proposed 138 in Jiang et al. (2023b) and discussed in the previous section), which is then fed into a concept 139 embedding layer (Espinosa Zarlenga et al., 2022). Using an LLM as the text encoder allows us to create a powerful task predictor (LICEM, Section 3.1) without the need for fine-tuning the text 140 encoder. Additionally, leveraging pretrained LLMs enables the self-generation of concept predictions 141 (Self-LICEM, Section 3.2), extending the scalability of CBMs to scenarios without available concept 142 annotations. A description of the overall pipeline is provided in Figure 2. 143

144 145

3.1 LINEARLY-INTERPRETABLE CONCEPT EMBEDDING MODEL (LICEM)

To create an interpretable predictor, it is essential to utilize both an interpretable data representation and an interpretable model (Ribeiro et al., 2016). Concept-based models allow employing an interpretable data representation within the network. However, to prevent a loss in generalization, CEMs provide the task predictions over concept embeddings whose single dimensions are noninterpretable. Thus, even when using an interpretable task predictor (e.g., a linear layer), CEM does not allow providing an interpretable prediction.

To address this issue, in this work, we propose to *predict a linear equation* that can be executed over the concept predictions and that outputs the final classification as an interpretable aggregation of the most important concepts. In this approach, the neural network's output is modeled as a linear equation where the independent variables are the concepts, and their weights reflect their importance in the task prediction, followed by a bias term. We employ two neural modules to predict *for each sample* the weights and the bias of a linear equation that is executed over the concepts (that are also predicted). Formally:

159 160

LICEM :
$$\hat{y}_i = \sigma \left(\sum_j \hat{w}_{ij} \hat{c}_j + \hat{b}_i \right)$$
 $\hat{w}_{ij} = \rho_i(\mathbf{c_j}), \ \hat{b}_i = \beta_i(\mathbf{c}), \ \hat{c}_j = s(\mathbf{c_j})$ (1)

162 where, as in common logistic regressions, \hat{w}_{ij} is the weight for the *j*-th concept in predicting the 163 *i*-th task, b_i is the bias for the *i*-th task, while \hat{c}_i and c_i are, respectively, the prediction and the 164 embedding of the *j*-th concept provided by CEM, and σ represents the activation function. For a 165 single concept j, the weights for all classes \hat{w}_j are predicted by a neural module $\rho : \mathbf{C}_{\mathbf{i}} \to \mathbb{R}^n$ 166 working on the corresponding concept embedding c_i . As commonly, $\hat{w}_{ij} < 0$ indicates a negatively important concept, $\hat{w}_{ij} > 0$ a positively important one, and $\hat{w}_{ij} \sim 0$ a non-important concept. To 167 improve readability, we aim for sparse weights, where few concepts have $\hat{w}_i \neq 0$. We achieve this by 168 adding L_1 regularization to the training loss. The bias term b is predicted over all concept embeddings 169 by a function $\beta : \mathbf{C} \to \mathbb{R}$, representing the overall bias for each class. This term is optional, but it 170 allows for positive predictions even when no concept is positively predicted. Indeed, when $\hat{c}_i = 0$ 171 for all $j \in \{1, ..., m\}$, the prediction would be $\hat{y}_i = 0$ regardless of \hat{w}_{ij} . To prevent over-reliance on 172 the bias term, we add L_2 regularization to encourage small bias values, minimizing its influence on 173 task prediction. Finally, we use a sigmoid activation function σ for binary classification tasks and a 174 softmax for multi-class classification tasks. To understand the contribution of a concept to the final 175 prediction of a class, we propose considering the combined contribution $\hat{w}_{ij}\hat{c}_j$ and plotting them in a 176 LIME-like feature importance plot, as shown in the output of Figure 2. 177

Training. LICEM is trained similarly to any supervised concept-based model with a cross-entropy
 H loss over both the predicted concepts and the tasks:

180

 $\mathcal{L}_{sup} = H(c, \hat{c}) + \lambda_y H(y, \hat{y}) + \lambda_w ||w||_1 + \lambda_b ||b||_2$ (2)

where we indicate the loss over the concept predictions as $\mathcal{L}_c = H(\hat{c}, c)$, the loss over the task predictions as $\mathcal{L}_t = H(\hat{y}, y)$, with ||w|| and ||b|| the regularization terms over the weights and biases and with λ_y, λ_w and λ_b the optimization weights for each term. In the rest of the paper, we will refer to this strategy as *supervised*.

185 186

187

3.2 EXPLOITING LLMS TO AVOID CONCEPT ANNOTATION: A SELF-GENERATIVE APPROACH

To alleviate human annotators from the burden of providing concept supervision, a few works are 188 starting to exploit the knowledge already available in pre-trained LLMs, both in the image (Yang 189 et al., 2023; Oikarinen et al., 2023) and in the textual domains (Ludan et al., 2023). First, an LLM is 190 asked to provide several attributes that describe each class. Each attribute is considered a concept for 191 that class, possibly shared with other classes. E.g., a *parrot* may be described as being a bird, with 192 bright feathers and of medium size. Then another LLM is required to predict whether the concept is 193 present in the input samples. The LLM, in this case, is formally represented by the distribution p_{θ} , 194 where θ denotes the parameters of a pre-trained LLM with classification head. When conditioned on a prompt t, the model generates the token "yes" if a specific concept is identified in the input 195 196 text sequence x, and "no" otherwise. Thus, the predicted concept is sampled as $c' \sim p_{\theta}(c'|t, x)$. In Appendix A.2 we report some examples of prompts. 197

Generative approach. In Ludan et al. (2023), these concept predictions c' are used as labels to train a textual concept encoder. Formally, $\mathcal{L}_{gen} = \mathcal{L}_{c'} + \lambda \mathcal{L}_t = H(c', \hat{c}) + \lambda H(y, \hat{y})$. We will refer to this strategy as *generative*, as a generative model provides concept annotations.

201 **Self-generative approach.** While the generative approach reduces human annotation efforts, it 202 requires training an additional concept encoder to learn the LLM-provided labels. In this paper, since 203 we already employ an LLM as a text encoder, we propose using the same LLM to directly make the 204 concept predictions. More precisely, we prompt the LLM to provide both a representation e for each 205 sample x and the concept predictions, i.e., $\hat{c} = c' \sim p_{\theta}(c'|t, x)$. This results in a modification of both 206 CEM and LICEM as the concept predictions are self-generated by the same LLM, as shown in Figure 207 2. We will refer to this approach as *self-generative*, as the same model directly provides the concept predictions. This method eliminates the need for concept annotations, but also reduces the number of 208 parameters to train and improves concept performance if compared to the generative method. Indeed, 209 the concept accuracy of the self-generative method represents an optimum for the generative one. 210 In the former, the concepts c' provided by the LLM are directly used as concept predictions, while 211 in the latter, they serve as training labels for an external text encoder, which aims to replicate c'. 212 Self-LICEM is obtained by substituting the concept predictions \hat{c} with c' from Equation 1: 213

- 214
- 215

Self-LICEM
$$\hat{y}_i = \sigma \left(\sum_j \hat{w}_{ij} c'_j + \hat{b}_i \right).$$
 (3)

The concept embedding encoder q and the neural modules ρ and β producing the interpretable linear equation are trained as in Equation 2, but minimizing, this time, only the loss over the task:

$$\mathcal{L}_{selfgen} = H(y, \hat{y}) + \lambda_w ||w||_1 + \lambda_b ||b||_2, \tag{4}$$

This approach is not limited to LICEM; it can also be extended to CBM-based and CEM-based models. In these cases, the LLM provides the concept predictions (CBM) or both the predictions and the embedding (CEM). In both cases, the optimization strategy involves minimizing only the crossentropy on the task predictions $H(y, \hat{y})$, as shown in Eq. 4. This allows converting any pre-trained LLM into a concept-based model without the need for concept annotations.

4 EXPERIMENTS

219 220

221

222

223

224

225 226

227 228

229 230

231

232

233

234

235

236

237 238

239

In this section, we want to answer the following research questions:

- **Generalization.** Does LICEM achieve superior performance in text classification compared to other interpretable models, and is it on par with non-interpretable ones? (Section 4.2)
- **Concept Efficiency.** How many concept supervisions are required to match Self-LICEM accuracy? Does the self-generative strategy outperform the generative one in concept accuracy? (Section 4.3)
- **Interpretability.** Are LICEM explanations more interpretable than those of other methods? Can we effectively interact with LICEM? (Section 4.4)

4.1 Setup

We test LICEM performance over different datasets (both with and without concept-supervisions), comparing against several models and for different metrics. For all experiments, we report the average and standard deviation across three repetitions. The models were trained on a dedicated server equipped with an AMD EPYC 7543 32-Core processor and one NVIDIA A100 GPU. Our code is publicly available at www.example.com¹

Dataset. We evaluated LICEM performance on three text-classification datasets for which concept annotation is available: CEBaB (Abraham et al., 2022), MultiEmotions-IT (Sprugnoli et al., 2020), and Drug review (Gräßer et al., 2018). Additionally, we tested the generative and self-generative approaches on the Depression dataset (Yates et al., 2017), where concept annotations are unavailable, but where an LLM (Jiang et al., 2024) identified six depression-related concepts which are: 'Self-deprecation', 'Loss of Interest', 'Hopelessness', 'Sleep Disturbances', 'Appetite Changes', and 'Fatigue'. Further information regarding the datasets is reported in Appendix A.3.

252 **Baselines.** We compare LICEM against several baselines, including black-box and concept-based 253 models, both task-interpretable and non-interpretable approaches. For all models, we use a non 254 fine-tuned Mixtral 8x7B (Jiang et al., 2024) encoder h(x), following the encoding strategy proposed 255 in Jiang et al. (2023b). In Appendix A.4 we also report all results based on a fine-tuned BERT 256 encoder (Devlin et al., 2018) as backbone. The results show that the decoder-only LLM achieves 257 similar performance without fine-tuning the whole LLM. Besides, it enables the self-generative approach: in Appendix A.5 we report a comparison of the concept annotation performance when 258 using different LLMs. For black-box models (E2E), we evaluate an end-to-end model directly 259 classifying the task with a Mixtral encoder h(x) and few layers as classification head (MLP), and 260 the same Mixtral used in Zero-shot and Few-shot prompting. CBM+LL and CBM+MLP are the 261 two CBMs originally proposed in (Koh et al., 2020) and recently adapted to text in (Tan et al., 262 2024b). They employ a concept bottleneck layer followed, the first one, by an interpretable linear 263 layer, while the second by a non-interpretable multi-layer perceptron. CBM+DT and CBM+XG are 264 respectively two CBM variants proposed in (Barbiero et al., 2023), using an interpretable decision 265 tree and a non-interpretable XGBoost classifier (Chen & Guestrin, 2016) on top of the concept 266 bottleneck layer, respectively. CBM+DT is task-interpretable, as one can extract a decision rule based 267 on concepts, whereas the second variant CBM+XG is non-interpretable. As described in Section 2, 268 CEM (Espinosa Zarlenga et al., 2022) employs embeddings to represent concepts and enhance CBM

²⁶⁹

¹We will release the code upon paper acceptance.

Table 1: Task accuracy (%) of the compared models. We report in **bold** the best result among the same type of models (e.g., supervised, interpretable ones) considering models equally best if their standard deviations overlap. We use \checkmark to indicate models requiring concept supervision (C. Sup.) or having a task-interpretable predictor (T. Inter.). We highlight in light gray the models we propose in this work. The Self-Generative approach extends the scalability of concept-based models to datasets without concept annotations, where supervised models cannot be applied (-).

277	Туре	Method	C. Sup.	T. Inter.	CEBaB	Multiemo-It	Drug	Depression
78 79 80	E2E	Mixtral–MLP Mixtral–Zero-shot Mixtral–Few-shot	× × ×	× × ×	$\begin{array}{c} \textbf{88.80} \pm 0.75 \\ 86.80 \pm 0.31 \\ 84.79 \pm 0.42 \end{array}$	$\begin{array}{c} 80.01 \pm 0.63 \\ 80.06 \pm 0.66 \\ \textbf{84.17} \pm 0.67 \end{array}$	$\begin{array}{c} \textbf{63.66} \pm 1.20 \\ 60.81 \pm 0.28 \\ \textbf{62.16} \pm 0.27 \end{array}$	$\begin{array}{c} \textbf{97.18} \pm 0.03 \\ 73.77 \pm 0.23 \\ 76.38 \pm 0.08 \end{array}$
81 82		CBM+MLP CBM+XG CEM	\$ \$ \$	× × ×	$\begin{array}{c} 78.41 \pm 9.30 \\ 83.01 \pm 0.10 \\ \textbf{89.60} \pm 0.49 \end{array}$	$\begin{array}{c} 45.43 \pm 8.20 \\ 69.01 \pm 0.02 \\ \textbf{83.33} \pm 0.47 \end{array}$	$\begin{array}{c} 45.42 \pm 4.90 \\ 55.00 \pm 0.13 \\ \textbf{66.81} \pm 0.40 \end{array}$	
33 84 85 86	SUP.	CBM+LL CBM+DT DCR LICEM (ours)	\$ \$ \$	\$ \$ \$	$\begin{array}{c} 71.43 \pm 9.71 \\ 77.20 \pm 0.40 \\ 88.05 \pm 0.53 \\ \textbf{89.89} \pm 0.77 \end{array}$	$\begin{array}{c} 42.67 \pm 7.01 \\ 65.00 \pm 0.02 \\ 82.01 \pm 0.71 \\ \textbf{83.47} \pm 0.49 \end{array}$	$\begin{array}{c} 34.60 \pm {\scriptstyle 10.10} \\ 47.20 \pm {\scriptstyle 0.40} \\ 65.40 \pm {\scriptstyle 0.80} \\ \textbf{66.80} \pm {\scriptstyle 0.29} \end{array}$	
37 38 39	SELF	Self-CBM+MLP Self-CBM+XG Self-CEM	× × ×	× × ×	$\begin{array}{c} 82.71 \pm 0.01 \\ 82.70 \pm < 0.01 \\ \textbf{89.14} \pm 0.38 \end{array}$	$\begin{array}{c} 75.42 \pm 4.42 \\ 79.09 \pm < 0.01 \\ \textbf{84.06} \pm 0.09 \end{array}$	$\begin{array}{c} 47.59 \pm 0.33 \\ 53.28 \pm < 0.01 \\ \textbf{65.20} \pm 0.73 \end{array}$	$\begin{array}{c} 82.31 \pm 0.04 \\ 82.28 \pm < 0.01 \\ \textbf{97.16} \pm 0.08 \end{array}$
)0)1)2	GEN. (OURS)	Self-CBM+LL Self-CBM+DT Self-DCR Self-LICEM	× × ×	\ \ \ \	$\begin{array}{c} 82.71 \pm 1.23 \\ 83.95 \pm < 0.01 \\ 87.72 \pm 0.66 \\ \textbf{89.56} \pm 0.29 \end{array}$	$\begin{array}{c} 77.15 \pm 0.96 \\ 78.44 \pm \! <\!\! 0.01 \\ 83.47 \pm 0.43 \\ \textbf{84.49} \pm 0.25 \end{array}$	$\begin{array}{c} 47.35 \pm 0.29 \\ 53.28 \pm < 0.01 \\ 63.29 \pm 0.36 \\ \textbf{65.89} \pm 0.39 \end{array}$	$\begin{array}{c} 82.12 \pm 0.15 \\ 82.28 \pm \! <\!\! 0.01 \\ 97.11 \pm 0.03 \\ \textbf{97.23} \pm 0.02 \end{array}$

generalization performance, but at the cost of losing task interpretability. Finally, DCR (Barbiero et al., 2023) is a neuro-symbolic approach designed to improve the interpretability of CEM. It generates propositional rules executed by a fuzzy system on top of concept predictions. We adapt CEM and DCR to work in the text classification scenario, and we compare their performance against the proposed model. For the training details regarding each model, please refer to Appendix A.3.

300 **Metrics.** We evaluate LICEM using various metrics. To assess generalization performance, we 301 compute the task accuracy and the macro-averaged concept F1 score (as concept classes are highly 302 imbalanced); for self-generative models, the macro-averaged F1 score evaluates the concept predic-303 tions directly provided by the LLM (Section 3.2). To measure efficiency, we examine the concept F1 304 score of all models when increasing the number of concept annotations. For interpretability, we first 305 evaluate LICEM explanations through a user study, comparing their plausibility and usefulness to 306 that of DCR; secondly we evaluate the effectiveness of concept interventions over LICEM to enhance 307 classification accuracy (Espinosa Zarlenga et al., 2024); third we measure the Causal-Concept Effect 308 (CaCE) (Goyal et al., 2019), which assesses the causal relevance of concepts for task predictions.

309 310

311

295

296

297

298

299

4.2 LICEM GENERALIZATION (TABLE 1)

312 LICEM matches black-box task performance and outperforms all task-interpretable models. 313 The initial finding from analyzing Table 1 is that LICEM consistently delivers task performances 314 that are comparable or even better than black-box and non task-interpretable models. Interestingly, 315 although with overlapping standard deviation, E2Es are never the best performing models, which is three times a LICEM and once a CEM. When it comes to interpretable models, LICEM invariably 316 emerge as the best interpretable predictor among the compared, with an improvement of at least 317 7-19% over CBM interpretable variants. With respect to DCR, we believe the improvement is due 318 to the way LICEM provide the final classification: the parameters of a linear equation are easier to 319 predict than constructing a logic rule, and to optimize, as they do not require passing through a fuzzy 320 system. 321

Self-generative approach increasing CBMs scalability while maintaining task performance.
 Self-generative CBMs maintain the task accuracy of supervised CBMs while increasing their scalability, as they can be applied to scenarios where concept annotations are not available, such as



Figure 3: Concepts prediction performance vs number of concept labels used during training. To increase plot readability, we only included the CBM+LL and the average F1 score for the generative approaches (Gen). Self-Gen. and Gen. approaches are reported with a straight line, as they do not require concept annotation.

342 the Depression dataset. We reported the performance of all concept-based baselines (not only Self-LICEM) when trained along the self generative approach to show that it enables all CBMs to work 343 344 on top of a pretrained LLM without concept annotations. When comparing the model performance in the two approaches, we can generally notice that the confidence intervals are overlapping. In a few 345 cases, such as CBM+LL, we can notice a stable improvement over all the datasets when using the 346 self-generative approach up to +30% on the Multiemo-It dataset. This is likely due to the trade-off 347 posed when training a concept-bottleneck layer, which has to favor either the task or the concept 348 performance: when directly working over good concept predictions, CBM performance improves. 349 In Appendix A.6, we also report the task accuracy of models trained along the standard generative 350 approach, showing similar results. 351

4.3 LICEMs CONCEPT EFFICIENCY (FIGURE 3)

354 Self-Generative approach strongly reduces the human annotation effort. In Figure 3, we report 355 the concept prediction performance of the compared methods when increasing the number of concept 356 labels used for training. Self-generative and generative approaches are reported with a straight line 357 since they do not require any concept supervision². Generative and self-generative models achieve a 358 concept macro-averaged F1 score that is higher or close to that of supervised models when using all available annotations, and significantly higher otherwise. When considering the CEBaB and Drug 359 datasets, supervised models do not surpass Self-Gen even when using all concept annotations, with 360 the latter achieving the highest concept accuracy. Likely, the amount of concept annotations required 361 to match the accuracy of the self-generative approach exceeds what is available in these datasets. 362

The self-generative concept accuracy exceeds that of the generative approach. The concepts prediction performance of the generative approach tends to be lower than that of the self-generative approach, with a reduction ranging from 2% to 7% in F1 macro score. This is because the concepts predicted by generative models are approximations of the self-generated concepts c' used in the self-generative approach. These self-generated concepts serve as the labels for training the concept encoders in the generative learning process. Detailed concepts prediction performance is presented in Appendix A.6, Table 6 for all models across all datasets, when provided with full concept annotations.

370 371

372

373

337

338

339

340 341

352

4.4 LICEM INTERPRETABILITY

LICEM explanations are more plausible and more useful than DCR (Fig. 4). To evaluate the interpretability of LICEM explanations, we conducted a user study comprising 21 ques-

 ³⁷⁵ ²Generative approaches results are reported with variance because the concepts are still learnt and thus the performance vary across models. For the self-generative approach, instead, the result does not vary because the concepts are predicted equally by the LLM for all models since we set the LLM's temperature to zero, which results in a deterministic annotation.

tions and involving 46 participants, consisting of both machine learning experts and non-experts.
It is structured as follows. First, participants are asked to choose the most plausible explanation (Rajagopal et al., 2021) from three options: the LICEM explanation, the DCR explanation, or neither. This process is repeated across three datasets (we excluded Multiemo-It, as it contains only Italian comments). Examples of the questions are shown in Figure 6, 7, and 8. In a second task, we assess explanation usefulness by computing how much participants can guess the model predictions based on the provided explanation (Fel et al., 2023).

This experiment is carried out for 385 both LICEM and DCR explana-386 tions and is repeated across the 387 same datasets as in the previous 388 step. In both cases, the samples 389 have been randomly drawn from 390 each dataset. A complete charac-391 terization of the user study is re-392 ported in Appendix A.7. The left image of Fig. 4 presents the re-393 sults related to explanation plau-394 sibility. It is evident that the 395 LICEM explanation is consis-396 tently considered more plausible 397 over the rule-based DCR expla-398 nation by both expert and non-399 expert users. Contrary to our 400 expectations, LICEM was espe-



Figure 4: Averaged survey results for the two user groups. On the left, we report the explanation plausibility; on the right, users' accuracy in guessing the model prediction based on its explanation.

cially favored by expert users, with nearly 80% of them appreciating its explanations. The right image
of Fig. 4 illustrates the accuracy achieved by users when tasked with selecting a class label based on
a given explanation. Both groups of users demonstrated good accuracy when making classifications
using the LICEM explanations. Expert people, in particular, nearly double the accuracy when using
LICEM compared to when using DCR explanations.

406 LICEM is responsive to concept interventions (Figure 5). To assess the possibility to interact 407 with LICEM, we evaluated the effect of concept interventions, i.e., modifications at test time of 408 the predicted concepts with a concept provided by a human expert. Figure 5 shows the test task 409 accuracy gain with increasing intervention probability on the CEBaB dataset, demonstrating LICEM's responsiveness and significant performance improvement. A similar behaviour can also be observed 410 for CBMs, even though they were starting from a lower task accuracy and a higher increase could 411 have also been expected. Results for all datasets are reported in Appendix A.8, showing similar 412 results, with LICEM always improving its task accuracy through interactions. For comparison, we 413 also report the E2E model with a flat line, since it does not offer this possibility. 414

LICEM predictions are caused by most important concepts (Table 2). We assess the responsive ness of concept-based models to *do-interventions* over concepts (Pearl et al., 2016), by computing
 the causal concept effect (CaCE) (Goyal et al., 2019). CaCE measures the impact of modifying input



418 419

420

421

422

423

424

425

426

427

428 429 Table 2: Causal Concept Effect (CaCE) for different methods. A high (absolute) value implies a strong responsiveness of a model to modifications to the concept.

Model	Food	Amb.	Service	Noise
CBM+LL	-0.02	0.01	0.01	-0.01
CEM	0.29	0.08	0.13	-0.05
DCR	0.33	0.02	0.20	-0.02
LICEM	0.62	0.18	0.37	0.15
Self-LICEM	0.63	0.20	0.35	0.15

Figure 5: Concept interventions on the CE-BaB dataset. We report the task accuracy gain when varying the probability of intervention.

432 samples on model predictions. For concept-based models, interventions can be made at the concept 433 level (Dominici et al., 2024). In the evaluated dataset, several concepts are globally relevant for task 434 classification (positively or negatively), thus we expect models to exhibit high absolute CaCE values. 435 In Table 2, we report the results for the CEBaB dataset: both LICEMs demonstrate high CaCE values, 436 particularly for 'Food' and 'Service' which are crucial concepts. These values are higher than CEM and DCR, suggesting a stronger reliance on the prediction over these concepts. Conversely, CBMs 437 report low CaCE values that may indicate concept leakage issues (Marconato et al., 2022), possibly 438 due to the constraints of the concept-bottleneck representation. Results for all datasets are reported in 439 Appendix A.9, showing consistent findings. 440

441 442

443

5 RELATED WORK

444 LLM interpretability. Recent studies have highlighted the unreliability of LLMs, as they often 445 occur hallucinations (Ji et al., 2023), and when prompted for explanations, their responses frequently 446 do not reflect the actual decision-making process (Ye & Durrett, 2022; Madsen et al., 2024; Turpin 447 et al., 2024). Although the attention mechanism in transformer models offers some interpretability, it has been criticized for its lack of clarity and consistency (Jain & Wallace, 2019; Wiegreffe & Pinter, 448 2019). To improve LLM explainability, various standard XAI techniques, such as LIME (Ribeiro 449 et al., 2016) and Shapley values (Lundberg & Lee, 2017), along with newer methods (Kokalj et al., 450 2021; Heyen et al., 2024; Chefer et al., 2021b;a), have been employed. However, these standard 451 techniques have limitations (Kindermans et al., 2019; Ghorbani et al., 2019; Adebayo et al., 2018; 452 Taimeskhanov et al., 2024), primarily because they explain predictions in terms of input features that 453 often lack meaningful interpretations for non-experts (Poursabzi-Sangdeh et al., 2021). Consequently, 454 researchers are now exploring interpretable-by-design models also in the textual domain (Rajagopal 455 et al., 2021; Jain et al., 2022; Tan et al., 2024b;a). 456

Concept-based models. Concept-based models (Alvarez Melis & Jaakkola, 2018; Koh et al., 457 2020; Ciravegna et al., 2023; Kim et al., 2023) are transparent and interactive models that utilize an 458 intermediate layer to represent concepts. To increase the representation capability of the concept layer, 459 Espinosa Zarlenga et al. (2022) proposed using concept embeddings. However, the interpretability 460 of CEM task predictor is limited, as individual embedding dimensions lack clear meaning. In 461 this work, we demonstrate how to create an interpretable task predictor over these embeddings. 462 A recent neurosymbolic method (DCR, Barbiero et al. (2023)) based on fuzzy logic attempted to 463 tackle this issue. We extend CEM and DCR applicability to the textual domain, while showing that 464 LICEM achieves superior predictive performance than DCR and higher interpretability than both. 465 Additionally, supervised concept-based models (Koh et al., 2020; Espinosa Zarlenga et al., 2022) 466 often require extensive concept annotations, which are frequently unavailable, particularly in text. We enhance a recent generative approach (Yang et al., 2023; Oikarinen et al., 2023; Ludan et al., 2023) 467 by using the same LLM for self-generated concept predictions and sample representations. 468

400

6 CONCLUSION

470 471 472

In this paper, we propose LICEM, a novel linearly interpretable concept-based model for text classification. The experimental results show this model matches black-box models performance, is interpretable and can be trained without concept supervision (Self-LICEM). Besides a technological impact, we believe this work can also positively impact the society by enhancing LLM transparency and interpretability, thus facilitating their employment in several fields such as Healthcare, Finance, Legal Systems and Autonomous Vehicles.

478 Future work. In this analysis, we focus on binary or ternary sentiment analysis for the ease of 479 identifying concepts, and to texts composed of a few sentences. In future work, we will extend our 480 analysis to other NLP tasks and to longer texts, to ensure the scalability of this approach. Specifically, 481 we plan to extend the capability of this model to work in language modelling tasks, similarly to Ismail 482 et al. (2023) employing CBMs to solve generative tasks in computer vision. Furthermore, other 483 interpretable functions could be generated and used to provide an interpretable prediction, besides linear equations. As an example, we could also generate a text describing how each concept has been 484 predicted and its role in the final prediction, together with the indication of the task prediction. We 485 leave these investigations for future research.



Eldar D Abraham, Karel D'Oosterlinck, Amir Feder, Yair Gat, Atticus Geiger, Christopher Potts, Roi
 Reichart, and Zhengxuan Wu. Cebab: Estimating the causal effects of real-world concepts on nlp
 model behavior. Advances in Neural Information Processing Systems, 35:17582–17596, 2022.

540 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 541 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. 542 arXiv preprint arXiv:2303.08774, 2023. 543 Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity 544 checks for saliency maps. Advances in neural information processing systems, 31, 2018. 546 David Alvarez Melis and Tommi Jaakkola. Towards robust interpretability with self-explaining neural 547 networks. Advances in neural information processing systems, 31, 2018. 548 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, 549 Ziwei Ji, Tiezheng Yu, Willy Chung, et al. A multitask, multilingual, multimodal evaluation of 550 chatgpt on reasoning, hallucination, and interactivity. arXiv preprint arXiv:2302.04023, 2023. 551 552 Pietro Barbiero, Gabriele Ciravegna, Francesco Giannini, Mateo Espinosa Zarlenga, Lucie Charlotte 553 Magister, Alberto Tonda, Pietro Lio, Frederic Precioso, Mateja Jamnik, and Giuseppe Marra. Inter-554 pretable neural-symbolic concept reasoning. In Proceedings of the 40th International Conference 555 on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 1801–1825. 556 PMLR, 23-29 Jul 2023. Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 558 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are 559 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020. 560 561 Hila Chefer, Shir Gur, and Lior Wolf. Generic attention-model explainability for interpreting bi-modal 562 and encoder-decoder transformers. In Proceedings of the IEEE/CVF International Conference on 563 Computer Vision, pp. 397-406, 2021a. Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. In 565 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 782–791, 566 2021b. 567 568 Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In Proceedings of the 569 22nd acm sigkdd international conference on knowledge discovery and data mining, pp. 785–794, 570 2016. 571 Gabriele Ciravegna, Pietro Barbiero, Francesco Giannini, Marco Gori, Pietro Lió, Marco Maggini, 572 and Stefano Melacci. Logic explained networks. Artificial Intelligence, 314:103822, 2023. 573 574 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep 575 bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018. 576 Gabriele Dominici, Pietro Barbiero, Mateo Espinosa Zarlenga, Alberto Termine, Martin Gjoreski, and 577 Marc Langheinrich. Causal concept embedding models: Beyond causal opacity in deep learning. 578 arXiv preprint arXiv:2405.16507, 2024. 579 580 Mateo Espinosa Zarlenga, Pietro Barbiero, Gabriele Ciravegna, Giuseppe Marra, Francesco Gian-581 nini, Michelangelo Diligenti, Zohreh Shams, Frederic Precioso, Stefano Melacci, Adrian Weller, 582 Pietro Lió, and Mateja Jamnik. Concept embedding models: Beyond the accuracy-explainability 583 trade-off. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), 584 Advances in Neural Information Processing Systems, volume 35, pp. 21400–21413. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/ 585 2022/file/867c06823281e506e8059f5c13a57f75-Paper-Conference.pdf. 586 Mateo Espinosa Zarlenga, Katie Collins, Krishnamurthy Dvijotham, Adrian Weller, Zohreh Shams, 588 and Mateja Jamnik. Learning to receive help: Intervention-aware concept embedding models. 589 Advances in Neural Information Processing Systems, 36, 2024. 590 591 Thomas Fel, Agustin Picard, Louis Bethune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. Craft: Concept recursive activation factorization for explainability. 592 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2711-2721, 2023.

594 Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In 595 Proceedings of the AAAI conference on artificial intelligence, volume 33, pp. 3681–3688, 2019. 596 Yash Goyal, Amir Feder, Uri Shalit, and Been Kim. Explaining classifiers with causal concept effect 597 (cace). arXiv preprint arXiv:1907.07165, 2019. 598 Felix Gräßer, Surya Kallumadi, Hagen Malberg, and Sebastian Zaunseder. Aspect-based sentiment 600 analysis of drug reviews applying cross-domain and cross-data learning. In Proceedings of the 601 2018 international conference on digital health, pp. 121–125, 2018. 602 Henning Heyen, Amy Widdicombe, Noah Yamamoto Siegel, Philip Colin Treleaven, and Maria 603 Perez-Ortiz. The effect of model size on llm post-hoc explainability via lime. In ICLR 2024 604 Workshop on Secure and Trustworthy Large Language Models, 2024. 605 606 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong 607 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language 608 models: Principles, taxonomy, challenges, and open questions. arXiv preprint arXiv:2311.05232, 2023. 609 610 Aya Abdelsalam Ismail, Julius Adebayo, Hector Corrada Bravo, Stephen Ra, and Kyunghyun Cho. 611 Concept bottleneck generative models. In The Twelfth International Conference on Learning 612 Representations, 2023. 613 614 Rishabh Jain, Gabriele Ciravegna, Pietro Barbiero, Francesco Giannini, Davide Buffelli, and Pietro Lio. Extending logic explained networks to text classification. In Proceedings of the 2022 615 Conference on Empirical Methods in Natural Language Processing, pp. 8838–8857. Association 616 for Computational Linguistics, 2022. 617 618 Sarthak Jain and Byron C Wallace. Attention is not explanation. arXiv preprint arXiv:1902.10186, 619 2019. 620 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, 621 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM 622 Computing Surveys, 55(12):1–38, 2023. 623 624 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 625 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 626 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023a. 627 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris 628 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 629 Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024. 630 631 Ting Jiang, Shaohan Huang, Zhongzhi Luan, Deqing Wang, and Fuzhen Zhuang. Scaling sentence 632 embeddings with large language models, 2023b. 633 Eunji Kim, Dahuin Jung, Sangha Park, Siwon Kim, and Sungroh Yoon. Probabilistic concept 634 bottleneck models. In Proceedings of the 40th International Conference on Machine Learning, 635 ICML'23. JMLR.org, 2023. 636 637 Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven 638 Dähne, Dumitru Erhan, and Been Kim. The (un) reliability of saliency methods. Explainable AI: 639 Interpreting, explaining and visualizing deep learning, pp. 267–280, 2019. 640 Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and 641 Percy Liang. Concept bottleneck models. In International conference on machine learning, pp. 642 5338-5348. PMLR, 2020. 643 644 Enja Kokalj, Blaž Škrlj, Nada Lavrač, Senja Pollak, and Marko Robnik-Šikonja. BERT meets 645 shapley: Extending SHAP explanations to transformer-based classifiers. In Hannu Toivonen and Michele Boggia (eds.), Proceedings of the EACL Hackashop on News Media Content Analysis 646 and Automated Report Generation, pp. 16–21, Online, April 2021. Association for Computational 647 Linguistics. URL https://aclanthology.org/2021.hackashop-1.3.

648 649	Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. <i>CoRR</i> , abs/1711.05101, 2017. URL http://arxiv.org/abs/1711.05101.
651 652 653	Josh Magnus Ludan, Qing Lyu, Yue Yang, Liam Dugan, Mark Yatskar, and Chris Callison-Burch. Interpretable-by-design text classification with iteratively generated concept bottleneck. <i>arXiv</i> preprint arXiv:2310.19660, 2023.
654 655	Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. <i>Advances in neural information processing systems</i> , 30, 2017.
657 658 659	Andreas Madsen, Sarath Chandar, and Siva Reddy. Are self-explanations from large language models faithful? In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pp. 295–337, 2024.
660 661	Anita Mahinpei, Justin Clark, Isaac Lage, Finale Doshi-Velez, and Weiwei Pan. Promises and pitfalls of black-box concept learning models. <i>arXiv preprint arXiv:2106.13314</i> , 2021.
662 663 664 665	Emanuele Marconato, Andrea Passerini, and Stefano Teso. Glancenets: Interpretable, leak-proof concept-based models. <i>Advances in Neural Information Processing Systems</i> , 35:21212–21227, 2022.
666 667 668	Tuomas Oikarinen, Subhro Das, Lam M. Nguyen, and Tsui-Wei Weng. Label-free concept bottleneck models. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=FlCg47MNvBA.
669 670 671	Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. <i>Causal inference in statistics: A primer</i> . John Wiley & Sons, 2016.
672 673	Eleonora Poeta, Gabriele Ciravegna, Eliana Pastor, Tania Cerquitelli, and Elena Baralis. Concept- based explainable artificial intelligence: A survey. <i>arXiv preprint arXiv:2312.12936</i> , 2023.
674 675 676	Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Wort- man Vaughan, and Hanna Wallach. Manipulating and measuring model interpretability. In <i>Proceedings of the 2021 CHI conference on human factors in computing systems</i> , pp. 1–52, 2021.
678 679	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
680 681 682	Dheeraj Rajagopal, Vidhisha Balachandran, Eduard H Hovy, and Yulia Tsvetkov. Selfexplain: A self-explaining architecture for neural text classifiers. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pp. 836–850, 2021.
684 685 686	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pp. 1135–1144, 2016.
687 688	Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. <i>Nature machine intelligence</i> , 1(5):206–215, 2019.
689 690 691 692	Rachele Sprugnoli et al. Multiemotions-it: a new dataset for opinion polarity and emotion analysis for italian. In <i>Proceedings of the Seventh Italian Conference on Computational Linguistics (CLiC-it 2020)</i> , pp. 402–408. Accademia University Press, 2020.
693 694 695	Magamed Taimeskhanov, Ronan Sicre, and Damien Garreau. Cam-based methods can see through walls. In <i>Joint European Conference on Machine Learning and Knowledge Discovery in Databases</i> , pp. 332–348. Springer, 2024.
696 697 698 699	Zhen Tan, Tianlong Chen, Zhenyu Zhang, and Huan Liu. Sparsity-guided holistic explanation for llms with interpretable inference-time intervention. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 21619–21627, 2024a.
700 701	Zhen Tan, Lu Cheng, Song Wang, Bo Yuan, Jundong Li, and Huan Liu. Interpreting pretrained language models via concept bottlenecks. In <i>Pacific-Asia Conference on Knowledge Discovery and Data Mining</i> , pp. 56–74. Springer, 2024b.

702 703 704	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, and Lukas Blecher et al. Llama 2: Open foundation and fine-tuned chat models, 2023.
705	
706	Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. Language models don't always
707 708	say what they think: unfaithful explanations in chain-of-thought prompting. Advances in Neural Information Processing Systems, 36, 2024.
709	Sarah Wiegreffe and Yuval Pinter. Attention is not not explanation. arXiv preprint arXiv:1908.04626,
710	2019.
711	Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin, Chris Callison-Burch, and Mark
712	Yatskar, Language in a bottle: Language model guided concept bottlenecks for interpretable image
713 714	classification. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19187–19197, 2023.
715	
716 717	Andrew Yates, Arman Cohan, and Nazli Goharian. Depression and self-harm risk assessment in online forums. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language</i>
718	<i>Processing</i> , pp. 2968–2978, 2017.
719 720	Seonghyeon Ye, Hyeonbin Hwang, Sohee Yang, Hyeongu Yun, Yireun Kim, and Minjoon Seo. Investigating the effectiveness of task-agnostic prefix prompt for instruction following, 2023.
721	
722	Xi Ye and Greg Durrett. The unreliability of explanations in few-shot prompting for textual reasoning.
723	Advances in neural information processing systems, 35:30378–30392, 2022.
724	
725	
726	
727	
728	
729	
730	
730	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
750	
753	
754	
755	

A APPENDIX / SUPPLEMENTAL MATERIAL

A.1 PROBABILISTIC FORMALIZATION

CBMs. As stated in Section 2, CBMs (Kim et al., 2023) provide explanation operating abstract human-understandable concepts. Let $x \sim p(x)$ represent the random variable drawn from the data distribution, and $c_j \sim p(c_j|x)$ denote the concept *j* derived from sample *x*, where $c_j \in [0, 1]$. The prediction of the categorical variable *y* is then formulated as:

$$y \sim p(y|x) = p(y|c_1, \dots c_m) \prod_{\substack{j=1\\p(c_1, \dots, c_m|x)}}^m p(c_j|x)$$
(5)

where $p(y|c_1,...,c_m)$ is a categorical distribution usually modeled using a fully connected layer, and $\prod_{j=1}^{m} p(c_j|x)$ is parameterized by a neural model. For simplicity we define $p(y|c_1,...,c_m) = p(y|c)$ and $p(c|x) = p(c_1,...,c_m|x)$. The different components are optimized by maximizing the following loss function:

$$\mathcal{L} = E_{x,c \sim p(x,c)}[-\log p(c|x)] + \lambda_y E_{x,y \sim p(x,y)}[-\log p(y|x)]$$
(6)

where $\lambda_y \in [0, 1]$ is the coefficient used to prioritize the concept learning relative the task learning.

CEMs. CEMs (Espinosa Zarlenga et al., 2022) addresses the low task performance of CBMs, 778 which is attributed to the bottleneck created by the intermediate concept layer, generating a concept 779 embedding for each concept. Initially, a vector representation of the raw data is generated, denoted 780 as $e \sim p(e|x)$, where $e \in \mathbb{R}^{l}$. Subsequently, both the active and inactive concept states, represented 781 as $\mathbf{c}_i^+, \mathbf{c}_i^- \in \mathbb{R}^k$, are derived from the two conditional distributions $p(\mathbf{c}_i^+|e)$ and $p(\mathbf{c}_i^-|e)$. At this 782 stage the concept score $\hat{c}_j \in [0, 1]$ is sampled as $\hat{c}_j \sim p(\hat{c}_j | \mathbf{c}_j^+, \mathbf{c}_j^-)$. The representational embedding 783 of concept j, denoted as $\mathbf{c}_j \in \mathbb{R}^k$, is computed as a convex combination of the active and inactive 784 states, given by $\mathbf{c}_j = \hat{c}_j \cdot \mathbf{c}_j^+ + (1 - \hat{c}_j) \cdot \mathbf{c}_j^-$. Finally, all concept embeddings are utilized to condition 785 the generation of the target variable, expressed as $y \sim p(y|\mathbf{c}_1,...,\mathbf{c}_m)$. This process can then be 786 formalized by 787

788

758

770

771

772

773 774 775

791 792 $y \sim p(y|x) = \underbrace{p(y|\mathbf{c}_1, ..., \mathbf{c}_m)}_{\text{CLASSIFIER}} \underbrace{\prod_{j=1}^m \left[p(\hat{c}_j|\mathbf{c}_j^+, \mathbf{c}_j^-) p(\mathbf{c}_j^+|e) p(\mathbf{c}_j^-|e) \right]}_{\text{CEM}} \underbrace{p(e|x)}_{\text{Encoder}}$ (7)

The classifier, which operates on the concatenation of concept embeddings, is typically structured as a deep neural network to enable end-to-end optimization. The loss function to optimize is analogous to 6.

LICEMs. LICEM builds over the CEM's output, modeling the distribution associated to the classifier 797 in 7. It utilizes both concept embeddings and concept scores to generate explanations, representing 798 them as a linear combination of concepts, were the weight associated to each concept j regarding 799 class $i, \hat{w}_{ij} \sim p_{\rho_i}(\hat{w}_{ij} | \mathbf{c_j})$, changes according to the concept embedding. Additionally, a dynamic 800 bias is sampled using all the concept embeddings $\hat{b}_i \sim p_{\beta_i}(\hat{b}_i | \mathbf{c}_1, ..., \mathbf{c}_m)$. The logit corresponding 801 to class i is calculated as $l_i = \sum_j \hat{w}_{ij} \hat{c}_j + \hat{b}_i$. For a multiclass classification task, the softmax 802 function is applied to the computed logits. The final probability associated with class i is given 803 by $\pi_i = Softmax(l_i)$. The predicted class label \hat{y} is subsequently sampled from a categorical 804 distribution defined by $\hat{y} \sim Cat(\hat{y}|\pi_1,...,\pi_n)$, where *n* denotes the total number of classes. 805

Self-generative. With the self-generative approach CEM is modified in order to allow external concept scores injection. This traduces into eliminating the neural module which models $p(\hat{c}_j | \mathbf{c}_j^+, \mathbf{c}_j^-)$, and using the LLM generated scores $\hat{c}_j \sim p(c_j | x, t)$, where *t* represents the prompt, to select the state of the concept embedding $\mathbf{c}_j = \hat{c}_j \cdot \mathbf{c}_j^+ + (1 - \hat{c}_j) \cdot \mathbf{c}_j^-$. Using $\mathbf{c} = (\mathbf{c}_1, ..., \mathbf{c}_m)$ and $\hat{c} = (\hat{c}_1, ..., \hat{c}_m)$ to simplify the notation, the label prediction process can be formalized as:

$$y \sim p(y|x) = \underbrace{p(y|\mathbf{c}, \hat{c})}_{\text{LICEM}} \underbrace{\prod_{j=1}^{m} \left[p(\mathbf{c}_j|e, \hat{c}_j) \right]}_{\text{MODIFIED CEM}} \underbrace{p(\hat{c}|t, x)}_{\text{LLM}} \underbrace{p(e|x)}_{\text{Encoder}} \tag{8}$$

A.2 PROMPTS FOR ANNOTATION

Here we report the prompts used to instruct *Mistral 7B* and *Mixtral 8x7B* to perform the annotations on the 4 different datasets used in this work. We adopted the in-context instruction learning prompting strategy (Ye et al., 2023).

CEBAB

822	
823	In a dataset of restaurant reviews there are 4 possible concepts: Good Food, Good Ambiance, Good Service and Good Noise. Given a certain review you have to detect if those concepts are present or not in the review
824	book house, siven a certain review, you have to accel if those concepts are present of not in the review.
825	Answer format: Good Food:score, Good Ambiance:score, Good Service:score, Good Noise:score.
826	Do not add any text other than that specified by the answer format.
827	The score should be equal to 1 if the concept is present or zero otherwise, no other values are accepted.
828	The following are examples:
829	Review: "The food was delicious and the service fantastic".
830	Answer: Good Food:1, Good Ambiance:0, Good Service:1, Good Noise:0
831	Review. "The staff was very rough but the restaurant decorations were great. Other than that there was a very
832	relaxing background music".
833	Answer: Good Food:0, Good Ambiance:1, Good Service:0, Good Noise:1
834	Now it's your turn:
835	Raview (review)
836	Answer:
837	

Drug

In a dataset of drug reviews there are 2 possible concepts:
- Effectiveness: 1 if the drug was highly effective and 0 if it was marginally or not effective,
- Side effects: 1 if the drug gave side effects and 0 otherwise.
Given a certain review, you have to detect if those concepts are present or not in the review.
Answer format: Effectveness:score, Side effects:score.
Do not add any text other than that specified by the answer format. The score should be equal to 1 if the concept is present or zero otherwise, no other values are accepted
The following are examples:
Review: "The medicine worked wonders for me. However, I did experience some side effects. Despite this,
I still found it easy to use and incredibly effective". Answer: Effectiveness:1. Side effects:1
Review: "Not only it did fail to alleviate my symptoms, but it also led to unpleasant side effects".
Now it's your turn:
Review: <review></review>
Answer:

Multiemo-it

859	In a dataset containing comments in Italian, you need to identify the following concepts:
860	
861	-Joy: the user who wrote the comment expresses joy, -Trust: the user who wrote the comment expresses trust,
862	-Sadness: the user who wrote the comment expresses sadness,
863	-Surprise: the user who wrote the comment is surprised.

Response format: Joy:score, Trust:score, Sadness:score, Surprise:score.

864 865 866

873

890

The score must be equal to 1 if the concept is present and 0 otherwise; other values are not accepted. The following is an example: 867 Comment: "Mi piace la rivisitazione di questa canzone, dolce, raffinata, elegante, bellissima!" Answer: Joy:1, Trust:1, Sadness:0, Surprise:1 868 869 Now it's your turn: Comment: <comment> 870 Answer: 871 872

DEPRESSION

874 You have to identify the presence or absence of 6 concepts in a given text. The concepts to be identified are: 875 - Self-Deprecation: the text exhibits self-critical or self-deprecating language, expressing feelings of guilt, shame, or inadequacy. 877 - Loss of Interest: diminished pleasure or motivation in the writer's descriptions of hobbies or pursuits. 878 - Hopelessness: the writer express feelings of futility or a lack of optimism about their prospects. - Sleep Disturbances: the writer mentions insomnia, oversleeping, or disrupted sleep as part of their 879 experience - Appetite Changes: there are references to changes in eating habits. - Fatigue: there are references to exhaustion or lethargy. 882 Answer format: Self-Deprecation:score, Loss of Interest:score, Hopelessness:score, Sleep Disturbances:score, Appetite Changes:score, Fatigue:score. 883 884 The score has to be 1 if the concept is detected and 0 otherwise. Do not add any other text besides the one specified in the answer format. 885 Text: <text> Answer: 887 888 889

EXPERIMENTAL DETAILS A.3

891 **Dataset** To check the performance of LICEM, we first selected three text-classification datasets 892 for which concept annotations are provided or in which attribute annotations can be employed. The 893 first dataset is CEBaB (Abraham et al., 2022), a dataset designed to study the causal effects of 894 real-world concepts on NLP models. It includes short restaurant reviews annotated with sentiment ratings at both overall-review level and for four dining experience aspects (food quality, noise 895 level, ambiance, and service). The second dataset is MultiEmotions-IT (Sprugnoli et al., 2020), 896 a dataset designed for opinion polarity and emotion analysis and containing comments related to 897 videos and advertisements posted on social media platforms. These comments have been manually annotated according to different aspects, among which we choose two dimensions: opinion polarity, 899 describing the overall sentiment expressed by users (that we employed as task labels), and basic 900 emotions from which we selected joy, trust, sadness, and surprise (concept labels). The third dataset 901 is Drug review (Gräßer et al., 2018), a dataset that provides patient reviews on specific drugs. The 902 reviews are annotated with the overall satisfaction of the users (which we discretize to a binary 903 representation) and drug experience annotations as effectiveness and side effects. Furthermore, to 904 test the generalization capability of self-supervised methods in a scenario where concept annotations are not actually provided, we chose the Depression dataset (Yates et al., 2017)³ which consists of 905 Reddit posts for users who claimed to have been diagnosed with depression and control users. The 906 set of concepts utilized for the Depression dataset was generated by the same LLM employed for 907 the annotations, *Mixtral 8x7B* (Jiang et al., 2024). Upon prompting the model to identify concepts 908 relevant to depression-related comments, it returned the following six key concepts: self-deprecation, 909 loss of interest, hopelessness, sleep disturbances, appetite changes, and fatigue. 910

911 **Evaluation** We evaluate LICEM against the baselines according to different metrics, each one 912 analysing a different characteristic of the models. First, to check LICEM generalization performance, 913 we compute the task accuracy and the macro-averaged F1 score for concepts prediction. For 914 GENERATIVE and SELF-SUP methods, we train the model without employing the actual concept 915 annotations but by prompting an LLM as described in Section 3.2. To test the efficiency of the 916 models, we report the concepts prediction performance of the models when increasing the number of 917

³For the Depression dataset, we employed the cleaned version available on Kaggle.

918 concept annotations (provided by humans). Finally, to test the interpretability of the model, we first 919 conducted a user study involving 30 participants, consisting of both machine learning experts and non-920 experts to evaluate LICEM explanations. Secondly, we checked whether it is possible to intervene on 921 the predicted concepts (Espinosa Zarlenga et al., 2024) and improve the classification accuracy even 922 when using an interpretable predictor. Thirdly, we checked the Causal-Concept Effect (CaCE) (Goyal et al., 2019), a measure introduced to assess the causality of a model with respect to a given concept. 923 Concept-based models, indeed, are generally required to make task predictions according to the 924 predicted concepts. However, the employment of vectorial concept representations (Mahinpei et al., 925 2021) may lead to model ignoring the predicted concepts. We see in the results that this is not the 926 case for LICEM. 927

928

Experimental settings For the E2E, CBMs, CEM, DCR and LICEM models, the training process involved utilizing an AdamW optimizer (Loshchilov & Hutter, 2017). The λ_y coefficient (2) was set to 0.5 to emphasize concept learning over task loss while $\lambda_w = 1 \times 10^{-6}$ and $\lambda_b = 10^{-6}$. Moreover, a scheduler was implemented with a gamma of 0.1 and a step size of 10 epochs throughout the training period of 100 epochs. After every hidden layer we have used a ReLU activation function. Here are further insights into the methodologies' architectures, with the number of output neurons indicated within brackets.

936 937

938

939

940 941

942

943

944 945

- E2E: layer 1 (100), layer 2 (number of classes);
- CEM: concept embedding size of 768, layer 1 (10), layer 2 (number of classes);
- CBMs, concept prediction: layer 1 (10), layer 2 (number of concepts);
 - LL, task prediciton: layer (number of classes);
 - MLP, task prediction: layer 1 (3 · number of concepts), layer 2 (number of classes).
- DCR: the temperature parameter is set to 0.1.

The text's embedding size varies depending on the chosen backbone. When employing BERT, it
remains at 768, whereas adopting the LLMs approach (Jiang et al., 2023b) it increases to 4096.
For Dtree and XGBoost, we employed the default hyperparameter settings. The DTree model was
implemented using the sklearn library, while the XGBoost model was implemented using the xgboost
library⁴. We conducted five experiments for each methodology. The training time for the different
experiments averages around 10 minutes using the setup specified in Section 4.1.

The CEBaB dataset (Abraham et al., 2022) does not necessitate any splitting procedure as it inherently offers training, validation, and test sets. In the training set, modifications include counterfactual examples, while both the validation and test sets exclusively contain original reviews. For the remaining datasets, we partitioned the data into training, validation, and test sets using stratified sampling based on the task labels. The proportions allocated are 0.7 for training, 0.1 for validation, and 0.2 for testing. Each experiment was conducted with a different seed.

957 958 959

A.4 ENCODER COMPARISON

This section presents all the results obtained using a fine-tuned BERT backbone as the encoder h(x). In the remainder of the paper, we consistently reported results when utilizing *Mixtral 8x7B* (Jiang et al., 2024) as the backbone model. In this section, we instead provide the performance of all models in terms of task accuracy (see Table 3) and of concept macro-averaged F1 score (refer to Table 4) when employing BERT as the backbone (Devlin et al., 2018), which is an encoder-only model.

Both tables show that there is no great difference with respect to Tables 1, 6, with BERT providing
slightly lower performance on Multiemo-It and on the Drug dataset. This result shows that the
proposed approach can be applied also to other architectures. We chose to employ Mixtral in the
remainder of the paper since it can be also effectively used to provide concept annotations, therefore
having a single model for both encoding the sample and predicting the concept scores.

⁴The xgboost library we used can be found at https://github.com/dmlc/xgboost.

972Table 3: This table presents the performance in terms of task accuracy (%) of different models973utilizing BERT as backbone. We report in **bold** the best result among the same type of models (e.g.,974supervised, interpretable ones) considering models equally best if their standard deviations overlap.975We use \checkmark to indicate models requiring concept supervision (C. Sup.) or having an interpretable task976predictor (T. Inter.). We highlight in light gray the models we propose in this work. We do not report977supervised model results for depression (-) since it does not provide concept annotations.

Туре	Method	C. Sup.	T. Inter.	CEBaB	Multiemo-It	Drug	Depression
e2e	MLP	X	×	$\textbf{90.68} \pm 0.47$	$\textbf{75.67} {\scriptstyle \pm 0.47}$	$\textbf{59.33} \pm 0.56$	$\textbf{97.80} \pm 0.23$
	CBM+MLP	1	X	$78.01 \pm \textbf{6.51}$	$54.10 \pm \textbf{4.51}$	$36.67 \pm \textbf{6.24}$	_
	CBM+XG	1	×	$80.00 \pm $	69.02 ± 0.64	$51.00 \pm \textbf{0.28}$	_
	CEM	1	×	$\textbf{90.67} \pm 0.47$	$\textbf{77.00} \pm 0.82$	$\textbf{58.33} \pm 1.70$	_
SUP.	CBM+LL	1	✓	$61.00 \pm {\scriptstyle 12.02}$	$49.67 \pm \scriptscriptstyle 5.46$	$34.33 \pm \textbf{7.38}$	_
	CBM+DT	1	1	75.67 ± 0.47	65.02 ± 0.34	$46.23 \pm \textbf{0.78}$	_
	DCR	1	1	86.55 ± 0.58	$74.01 \pm \textbf{0.24}$	$\textbf{59.75} \pm 0.45$	-
	LICEM (ours)	1	1	$\boldsymbol{87.89} \pm 0.38$	$\textbf{75.31} \pm 0.15$	$\textbf{60.14} \pm 0.44$	_
	CBM+MLP	X	X	$73.93 \pm \textbf{5.67}$	44.19 ± 2.07	$35.16 \pm \textbf{4.3}$	$83.20 \pm \textbf{2.18}$
	CBM+XG	X	×	83.29 ± 0.43	69.85 ± 1.55	$34.94 \pm \textbf{0.91}$	87.00 ± 1.01
GEN.	CEM	×	×	$\textbf{85.88} \pm 0.95$	$\textbf{73.15} \pm 0.67$	$\textbf{56.95} \pm 0.36$	96.12 ± 0.50
	CBM+LL	X	1	$58.81 \pm \textbf{7.16}$	$58.35 \pm \textbf{1.59}$	$36.84 \pm {\scriptstyle 11.52}$	$51.48 \pm \textbf{2.16}$
	CBM+DT	X	1	79.28 ± 0.52	62.61 ± 2.08	34.17 ± 0.11	80.55 ± 0.03
	DCR	X	1	$\textbf{85.63} \pm 0.81$	70.02 ± 2.70	57.46 ± 0.02	$95.98 \pm \textbf{0.27}$
	LICEM (ours)	X	1	$\textbf{86.22} \pm 0.66$	$\textbf{74.45} \pm 0.57$	$\textbf{60.23} \pm 0.58$	96.87 ± 0.20

Table 4: This table presents the performance in terms of concept prediction of the models that utilize
BERT as backbone. Concept prediction (%) of the compared models for datasets equipped with
concept annotations is measured using the macro-averaged F1 score. We report in **bold** the best result
among the same type of models (e.g., supervised, interpretable ones) considering models equally best
if their standard deviations overlap. We highlight in light gray the models we propose in this work.
The methods using the self-generative have the same macro-averaged F1 score, therefore we use - to
represent all methods.

Туре	Method	CEBaB	Multiemo-It	Drug
E2E	MLP	$\textbf{79.92} {\scriptstyle \pm 1.77}$	$\textbf{63.25} \pm 1.09$	$\textbf{79.01} \pm 2.9$
	CBM+MLP	$75.17{\scriptstyle\pm}_{\scriptstyle3.11}$	$64.08 {\scriptstyle \pm 1.22}$	$74.26 {\scriptstyle \pm 0.9}$
	CEM	79.97 ± 1.29	64.42 ± 1.21	77.32 ± 1.2
	CBM+XG	$\textbf{79.92} {\scriptstyle \pm 1.77}$	$\textbf{63.25} \pm 1.09$	$\textbf{79.01} \pm 0.9$
SUP.	CBM+LL	$74.25{\scriptstyle\pm}{\scriptstyle 4.55}$	$62.08 \pm \textbf{0.88}$	$73.11 {\scriptstyle \pm 1.7}$
	CBM+DT	$79.92 {\scriptstyle \pm 1.77}$	63.25 ± 1.09	79.01 ± 2.9
	DCR	$82.06 {\scriptstyle \pm 0.40}$	$64.29 {\scriptstyle \pm 0.42}$	$80.10 {\scriptstyle \pm 0.2}$
	LICEM (ours)	$\pmb{82.93} \pm 0.13$	$\textbf{65.61} \pm 0.69$	$\boldsymbol{81.59} \pm 0.42$
	CBM+MLP	$75.05 \pm \textbf{8.31}$	$49.59 \pm \scriptscriptstyle 10.01$	$43.58 \pm \scriptscriptstyle 14.99$
	CEM	$\pmb{81.08} \pm 0.44$	$\textbf{58.30} \pm 1.79$	$\textbf{80.99} \pm 0.42$
	CBM+XG	$79.24 \pm \textbf{1.21}$	$\textbf{60.79} \pm 0.71$	64.72 ± 0.45
GEN.	CBM+LL	$\textbf{78.75} \pm 0.59$	$\boldsymbol{61.72} \pm 0.24$	$66.72 \pm \textbf{19.48}$
	CBM+DT	$\textbf{79.24} \pm 1.21$	$\textbf{60.79} \pm 0.70$	64.72 ± 0.45
	DCR	$\textbf{80.25} \pm 1.02$	59.11 ± 0.84	81.47 ± 0.49
	LICEM (ours)	$\textbf{77.79} \pm 2.49$	58.87 ± 0.66	$\pmb{81.18} \pm 0.33$
SELF GEN.	-	$\textbf{84.08} \pm 0.00$	64.27 ± 0.00	$\textbf{83.00} \pm 0.00$



Figure 10: Comparison among concept annotation methods where the annotation quality is measured 1042 in terms of macro-averaged F1 score. On average, Mixtral 8x7B yields the best results. 1043

1044 1045

1026 1027

1028 1029

1030 1031

1032 1033

1034 1035

1036

1039

1040 1041

A.5 LLM-based concept annotation vs Class-level annotation 1046

1047 This section presents a comparison between the usage of two different LLMs, *Mistral 7B* (Jiang 1048 et al., 2023a) and Mixtral 8x7B (Jiang et al., 2024), as concept annotators. In Figure 10 we report the 1049 results in terms of macro-averaged F1 score (as concept classes are highly imbalanced) on the three 1050 datasets for which human concept annotation is available. We also report, as a baseline, a global 1051 (class-level) annotation strategy, providing to all samples belonging to a given class the same concept 1052 annotation. In this case, we label the positive class with positive concepts and negated negative 1053 concepts (e.g. for all samples of the class *Good Drug* we use 'Efficient' and 'Not Side Effects'). 1054 We can observe that between the two LLMs there is not a significant difference in performance, with Mixtral 8x7B providing on average slightly better results. Comparing against the baseline, 1055 instead, we can observe that there is a great improvement in CEBaB and in the Drug dataset, while in 1056 Multiemo-It the improvement is more modest. 1057

1058

1059 A.6 TASK ACCURACY AND CONCEPTS PREDICTION PERFORMANCE

In this section we report the task accuracy and the concepts prediction performance results for all 1061 the different experiments conducted, generative approach included when using Mixtral 8x7B as a 1062 backbone. As shown in Table 5, LICEM outperforms the other task interpretable models, reaching 1063 the highest task accuracy for the CEBaB dataset using the generative approach. 1064

We also report the averaged F1 macro to measure the concepts prediction performance of all models when provided with all the available concept annotations. The results shown in Figure 3 are here confirmed. We again see that Self-supervised strategy is a very good approach since without human 1067 effort it provides better concept macro-averaged F1 score in CEBaB and Drug. Only on Multiemo-It 1068 the performance are significantly lower. This result may be due to the fact that the latter dataset is in 1069 Italian while the other datasets are in English, a language for which the LLMs have certainly seen 1070 more training samples. 1071

1072

A.7 SURVEY CHARACTERIZATION 1073

1074

In this section, we provide further details regarding the conducted survey. A total of 46 participants 1075 with varying levels of experience in machine learning, from complete beginners to experts, were 1076 recruited (see Figure 11). The gender distribution was nearly balanced, with 40% identifying as 1077 female and 60% as male. The majority of participants, 91.3%, were within the 20 - 40 age range, 1078 while only 8.7% were aged over 40. 1079

The survey was structured in the following manner:

1080Table 5: Task accuracy (%) of the compared models. We report in **bold** the best result among the1081same type of models (e.g., supervised, interpretable ones) considering models equally best if their1082standard deviations overlap. We use \checkmark to indicate models requiring concept supervision (C. Sup.) or1083having a task-interpretable predictor (T. Inter.). We highlight in light gray the models we propose in1084this work. The Generative and the Self-generative approaches extend the scalability of concept-based1085models to datasets without concept annotations, where supervised models cannot be applied (-).

Туре	Method	C. Sup.	T. Inter.	CEBaB	Multiemo-It	Drug	Depression
	Mixtral–MLP	X	×	88.80 ± 0.75	80.01 ± 0.63	$\textbf{63.66} \pm 1.20$	$\textbf{97.18} \pm 0.03$
E2E	Mixtral-Zero-shot	X	×	86.80 ± 0.31	80.06 ± 0.66	$60.81 \pm \textbf{0.28}$	73.77 ± 0.23
	Mixtral-Few-shot	X	×	84.79 ± 0.42	84.17 ± 0.67	62.16 ± 0.27	76.38 ± 0.08
	CBM+MLP	1	X	78.41 ± 9.30	45.43 ± 8.20	$45.42 \pm \textbf{4.90}$	_
	CBM+XG	1	X	83.01 ± 0.10	69.01 ± 0.02	$55.00 \pm \textbf{0.13}$	_
	CEM	1	×	89.60 ± 0.49	$\textbf{83.33} \pm 0.47$	66.81 ± 0.40	-
SUP.	CBM+LL	1	1	71.43 ± 9.71	42.67 ± 7.01	34.60 ± 10.10	_
	CBM+DT	1	1	77.20 ± 0.40	65.00 ± 0.02	47.20 ± 0.40	_
	DCR	1	1	88.05 ± 0.53	82.01 ± 0.71	65.40 ± 0.80	-
	LICEM (ours)	1	1	89.89 ± 0.77	$\textbf{83.47} \pm 0.49$	$\textbf{66.80} \pm 0.29$	-
	CEM	X	X	89.97 ± 0.66	$\textbf{82.41} \pm 0.11$	63.80 ± 0.38	$\textbf{97.06} \pm 0.11$
CEN	CBM	X	1	62.07 ± 0.22	$68.66 \pm \textbf{4.20}$	33.14 ± 2.10	50.25 ± 0.39
GEN.	DCR	X	1	88.97 ± 0.18	$80.82 \pm \textbf{0.54}$	63.74 ± 1.16	$95.35 \pm \textbf{0.21}$
	LICEM (ours)	X	1	$\textbf{90.64} \pm 0.38$	$\pmb{81.85} \pm 0.71$	$\textbf{66.15} \pm 0.44$	96.50 ± 0.18
	Self-CBM+MLP	X	X	82.71 ± 0.01	$75.42{\scriptstyle\pm}{\scriptstyle 4.42}$	47.59 ± 0.33	82.31 ± 0.04
	Self-CBM+XG	X	X	$82.70 \pm < 0.01$	$79.09{\scriptstyle~\pm<0.01}$	$53.28 \pm < 0.01$	$82.28 \pm \scriptscriptstyle < 0.01$
SELE	Self-CEM	X	×	$\textbf{89.14} \pm 0.38$	$\textbf{84.06} \pm 0.09$	$\textbf{65.20} \pm 0.73$	$\textbf{97.16} \pm 0.08$
GEN.	Self-CBM+LL	X	1	$82.71 \pm \textbf{1.23}$	77.15 ± 0.96	47.35 ± 0.29	82.12 ± 0.15
(OURS	Self-CBM+DT	X	1	$83.95 \pm < 0.01$	$78.44 \pm < 0.01$	$53.28 \pm < 0.01$	$82.28 \pm \scriptscriptstyle < 0.01$
	Self-DCR	X	1	87.72 ± 0.66	$83.47 \pm \textbf{0.43}$	63.29 ± 0.36	$\textbf{97.11} \pm 0.03$
	Self-LICEM	X	1	89.56 ± 0.29	84.49 ± 0.25	$\textbf{65.89} \pm 0.39$	$\textbf{97.23} \pm 0.21$



Figure 11: Distribution of users by expertise level.

• **Introduction to Explanations**: We provided an introduction to the various types of explanations, ensuring that participants had sufficient background information to understand and interpret these explanations.

- Questionnaire: Participants were asked a total of 7 questions for each of the three datasets that contained english text: CEBaB, Drug, and Depression. The questions were divided as follows:
 - The first 3 questions asked participants to select their preferred explanation for a given text. Examples of these questions can be found in Figure 6, 7, 8.
 - The remaining 4 questions asked participants to predict the label of the text based on a provided explanation, with two questions pertaining to DCR and two to LICEM. Examples of these questions are presented in Figure 9, 12, 13.

For both types of questions, we randomly selected samples from the three datasets (CEBaB, Drug, and Depression) where both models (LICEM and DCR) made the correct predictions.

Table 6: This table presents the performance in terms of concept prediction of the models that utilize Mixtral 8x7B as backbone. Concept prediction (%) of the compared models for datasets equipped with concept annotations is measured using the macro-averaged F1 score. We report in **bold** the best result among the same type of models (e.g., supervised, interpretable ones) considering models equally best if their standard deviations overlap. Self-supervised methods are reported with the same concept accuracy with zero standard deviation, since the concept predictions are provided by an LLM with temperature set to zero. The methods using the self-generative have the same macro-averaged F1 score, therefore we use - to represent all methods.

	Туре	Method	CEBaB	Multiemo-It	Drug
	E2E	MLP	$\textbf{75.92} {\scriptstyle \pm 0.77}$	$\textbf{74.25} {\scriptstyle \pm 1.02}$	$\textbf{78.50} \pm 0.23$
		CBM+MLP CEM CBM+XG	$\begin{array}{c} 65.17 \pm 2.35 \\ \textbf{78.83} \pm 0.85 \\ \textbf{75.92} \pm 0.77 \end{array}$	$\begin{array}{c} 61.75 \pm 1.02 \\ \textbf{77.12} \pm 1.38 \\ \textbf{74.25} \pm 1.02 \end{array}$	$\begin{array}{c} 65.33 \pm 2.46 \\ \textbf{80.79} \pm 0.47 \\ 78.50 \pm 0.23 \end{array}$
	SUP.	CBM+LL CBM+DT DCR LICEM (ours)	$\begin{array}{c} 64.25 \pm 2.56 \\ \textbf{75.92} \pm 0.77 \\ \textbf{78.45} \pm 1.92 \\ \textbf{75.45} \pm 0.93 \end{array}$	$\begin{array}{c} 59.12 \pm 2.13 \\ 74.25 \pm 1.02 \\ \textbf{75.67} \pm 1.43 \\ \textbf{76.36} \pm 0.39 \end{array}$	$\begin{array}{c} 64.83 \pm 1.20 \\ 78.50 \pm 0.23 \\ 79.96 \pm 0.43 \\ \textbf{80.83} \pm 0.36 \end{array}$
	GEN.	CBM+MLP CEM CBM+XG	$\begin{array}{c} 71.87 \pm 0.14 \\ \textbf{74.70} \pm 0.98 \\ \textbf{75.02} \pm 0.57 \end{array}$	$\begin{array}{c} 52.60 \pm {\scriptstyle 14.32} \\ \textbf{63.61} \pm {\scriptstyle 0.44} \\ 61.69 \pm {\scriptstyle 0.44} \end{array}$	$\begin{array}{c} 55.68 \pm {\scriptstyle 19.84} \\ \textbf{79.45} \pm {\scriptstyle 0.41} \\ \textbf{79.15} \pm {\scriptstyle 0.30} \end{array}$
		CBM+LL CBM+DT DCR LICEM (ours)	$\begin{array}{c} \textbf{72.15} \pm 0.59 \\ \textbf{75.02} \pm 0.57 \\ \textbf{75.62} \pm 2.59 \\ \textbf{74.44} \pm 0.25 \end{array}$	$\begin{array}{c} \textbf{63.72} \pm 0.84 \\ 61.69 \pm 0.44 \\ 62.79 \pm 0.44 \\ \textbf{63.75} \pm 0.36 \end{array}$	$\begin{array}{c} 66.72 \pm {\scriptstyle 19.48} \\ \textbf{79.04} \pm {\scriptstyle 0.30} \\ \textbf{79.04} \pm {\scriptstyle 0.33} \\ \textbf{79.05} \pm {\scriptstyle 0.58} \end{array}$
	SELF GEN.	-	$\textbf{84.08} \pm 0.00$	64.27 ± 0.00	$\textbf{83.00} \pm 0.00$

According to the explanation, select the correct label. *



Neutral drug review

O Positive drug review

Figure 12: Example of label prediction given LICEM explanation, Drug dataset.

1181
1182A.8CONCEPT INTERVENTIONS

As introduced in Section 4.4, LICEM is sensible to concept interventions. This characteristic is very important since it implies that a human can interact with the model, providing counterfactual predictions when prompted with different concept predictions. In Figure 14, 15, 16 we simulate this situation by correcting mispredicted concepts with the correct concept predictions and check whether the task prediction has been also modified. More in details, we report the improvement in task accuracy when increasing the probability to correct the concepts, demonstrating LICEM's

 1188
 According to the explanation, select the correct label.
 *

 1189
 Explanation: Deprecation ^ - Loss_of_Interest ^ - Sleep_Disturbances ^ - Appetite_Changes

 1191
 Depressed comment

 1193
 Not depressed comment

 1195
 Figure 13: Example of label prediction given DCR explanation, CEBaB dataset.

 1197

responsiveness and significant performance improvement. A similar behaviour can also be observed for CBMs, even though they were starting from a lower task accuracy and a higher increase could have also been expected. For comparison, we also report the E2E model with a flat line, since it does not offer this possibility. As noted in (Espinosa Zarlenga et al., 2022), CEMs (which are not task interpretable) may not respond well to concept interventions, especially without conducting them during training. Thus, we trained all CEM-based models with a 0.5 intervention probability during the forward pass.



Figure 14: Concept interventions on the CEBaB dataset for (left) supervised approaches and (right) self-supervised ones.



Figure 15: Concept interventions on the Multiemo-it dataset for (left) supervised approaches and (right) self-supervised ones.

A.9 CAUSAL CONCEPT EFFECT (CACE)

As anticipated in Section 4.4, Concept-based models predictions must be causally influenced by the predicted concepts. We assess concept-based models' responsiveness to *do-interventions* using the Causal Concept Effect (CaCE) (Goyal et al., 2019), which measures the impact of input modifications



Figure 16: Concept interventions on the Drug dataset for (left) supervised approaches and (right) self-supervised ones.

Table 7: Causal Concept Effect (CaCE) for different methods. A high (absolute) value implies a strong responsiveness of a model to modifications to a certain concept.

	Concept	CBM+LL	CEM	DCR	LICEM	SELF-LICEM
В	Good Food	$\textbf{-0.02} \pm 0.01$	$0.29 \hspace{0.1 in} \pm \hspace{0.1 in} 0.03 \hspace{0.1 in}$	$0.33 \scriptstyle \pm 0.04 $	$0.62 \pm $	0.63 ± 0.01
3A)	Good Amb.	$0.01 \pm 0.05 $	$0.08 \pm $	$0.02 \scriptstyle \pm 0.01 $	0.18 ± 0.03	$0.20 \pm 0.04 $
ЗeЕ	Good Service	$0.01 \pm $	$0.13 \pm $	$0.20{\scriptstyle~\pm 0.08}$	$0.37 \pm $	0.35 ± 0.02
0	Good Noise	$\textbf{-0.01} \pm 0.10$	$-0.05 \hspace{0.1 in} \pm \hspace{0.1 in} 0.01 \hspace{0.1 in}$	$\textbf{-0.02} \hspace{0.1in} \pm 0.01 \hspace{0.1in}$	$0.15 \pm $	$0.15 \hspace{0.1cm} \pm \hspace{0.1cm} 0.03 \hspace{0.1cm}$
no	Joy	$0.04\pm$	$0.18 \pm $	$0.16 \pm $	$0.28 \pm $	$0.27 \pm $
ier	Trust	$0.02 \pm $	0.60 ± 0.04	$0.47 \pm $	$0.62 \pm $	0.63 ± 0.01
ult	Sadness	$\textbf{-0.04} \pm 0.05$	-0.06 ± 0.01	$\textbf{-0.04} \pm 0.02$	$\textbf{-0.04} \pm 0.01$	-0.10 ± 0.02
Σ	Surprise	$\textbf{-0.01} \pm 0.06$	$0.03 \hspace{0.1 in} \pm \hspace{0.1 in} 0.01$	0.06 ± 0.05	$\textbf{-0.02} \hspace{0.1in} \pm \hspace{0.1in} \textbf{0.01}$	$0.01 \pm $
gu	Effectiveness	0.02 ±0.10	$0.43 \scriptstyle \pm 0.02 $	$0.28 \pm 0.02 $	$0.45 \pm 0.04 $	0.46 ± 0.02
Dr	Side Effects	$\textbf{-0.07} \hspace{0.1in} \pm 0.14$	-0.52 ± 0.01	-0.25 ± 0.02	-0.55 ± 0.06	-0.55 ± 0.03

on model predictions. Higher absolute CaCE values indicate stronger conditioning on relevant
concepts. Tables 7 shows that both supervised and self-supervised LICEM have higher CaCE values
compared to CBM, CEM and DCR, suggesting stronger reliance on predicted concepts. This result
is positive since all concepts considered in this work are relevant for the task at hand. We leave for
future work the exploration of tasks where there are confounding concepts and checking whether
LICEM is capable to not consider them.