



TGML Lab | SUT | Winter 1403

---

# An Exploration of Concept-Based Interpretability Methods

Comparison of 5 Concept-Based Papers from NeurIPS, ICLR, ACL

Presentation by: **Maryam Rezaee**

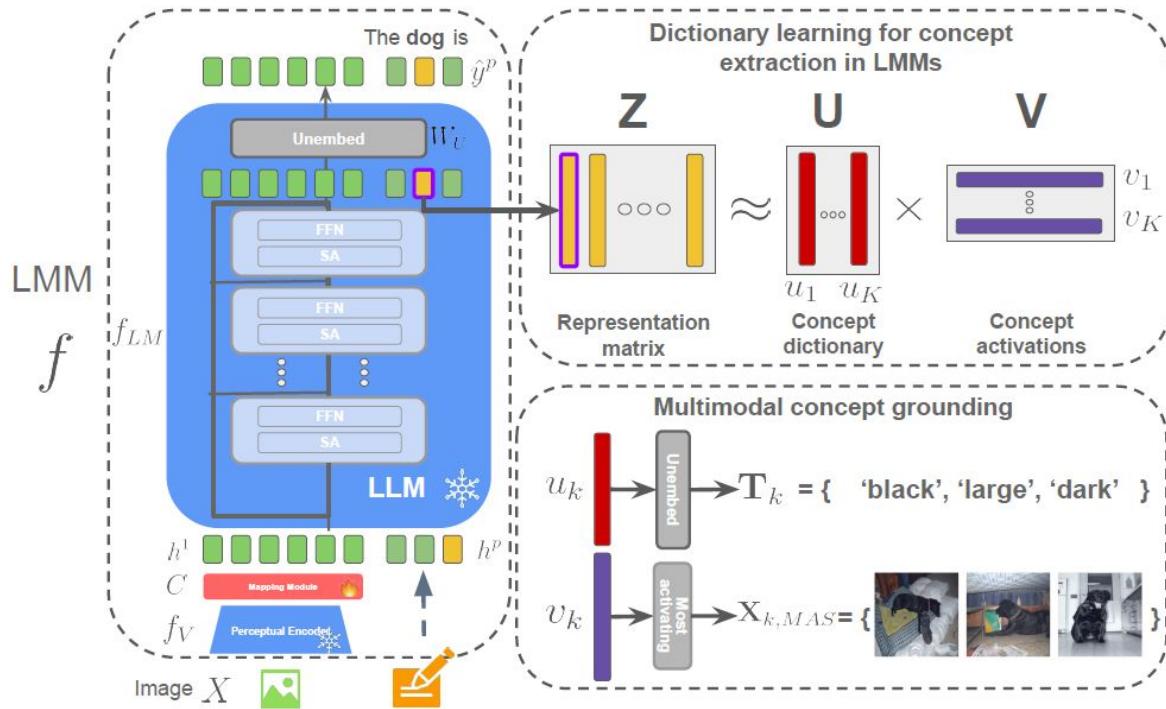
# 01

NeurIPS 2024

---

## A Concept-Based Explainability Framework for Large Multimodal Models (**Cox-LMM**)

# 1.1 Overview



## 1.2 Preliminary

### Model Architecture

- Contains: Visual encoder  $f_V$ ; Trainable connector  $C$ ; LLM  $f_{LM}$
- Internal representation of token at layer  $l$  and position  $p$  inside  $f_{LM}$  is denoted as  $h^p_{(l)}$ , with  $h^p_{(0)} = h^p$
- Input sequence for  $f_{LM}$  is concatenation of:
  - $N_V$  visual tokens provided by the visual encoder  $f_V$  operating on an image  $X$ , followed by the connector  $C$
  - linearly embedded textual tokens previously predicted by  $f_{LM}$

$$\hat{y}^p = f_{LM}(h^1, h^2, \dots, h^{N_V}, \dots, h^p)$$

Where:  $h^1, \dots, h^{N_V} = C(f_V(X))$ ,  $h^{N_V+1}$  is [SOS] and  $h^p = \text{Emb}(\hat{y}^{p-1})$

## 1.2 Preliminary

### ***Training***

- $f_V$  and  $f_{LM}$  are frozen,  $C$  is trained

### ***Transformer Representations View***

- Residual stream view of decoder-only transformers:

$$h_{(l+1)}^p = h_{(l)}^p + a_{(l)}^p + m_{(l)}^p$$



## 1.3 Method

### Overview

- Given a pretrained LMM  $f$  and a token of interest  $t \in Y$ :
  - Select a subset of images  $X$  from dataset  $S$ , relevant for target token  $t$ .
  - Extract representations by processing samples in  $X$  through  $f$  and collect extracted representations of dimension  $B$  in a matrix  $Z \in R^{B \times M}$ , where  $M$  is number of samples in  $X$ .
  - Linearly decompose  $Z \approx UV$ , which includes a dictionary of learnt concepts  $U \in R^{B \times K}$  of size  $K$  and coefficient/activation matrix  $V \in R^{K \times M}$ .
  - Semantically ground the learnt “multimodal concepts”, contained in dictionary  $U$  in both visual and textual modalities.

## 1.3 Method

### Overview

- The main objective is to understand internal representations of an LMM and not to interpret the output of the model, which can be accomplished by combining this pipeline with some concept importance estimation method.

*Thomas Fel, Victor Boutin, Louis Béthune, Rémi Cadène, Mazda Moayeri, Léo Andéol, Mathieu Chalvidal, and Thomas Serre. A **holistic approach to unifying automatic concept extraction and concept importance estimation**. Advances in Neural Information Processing Systems, 36, 2023.*

[https://proceedings.neurips.cc/paper\\_files/paper/2023/hash/abf3682c9cf9245a0294a4bebe4544ff-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2023/hash/abf3682c9cf9245a0294a4bebe4544ff-Abstract-Conference.html)

## 1.3 Method

### ***Representation Extraction***

- Extract samples:

$$\mathbf{X} = \{X_i \mid t \in f(X_i) \text{ and } t \in y_i \text{ and } (X_i, y_i) \in \mathcal{S}\}$$

- Given any  $X \in \mathbf{X}$ , extract the residual stream representation  $h_{(L)}^p$  from a deep layer  $L$ , at the first position in the predicted caption  $p > N_v$ , such that  $\hat{y}^p = t$ .
- Stacked representations  $z_j \in R^B$  of each sample  $X_j \in \mathbf{X}$  as columns of the matrix  $Z = [z_1, \dots, z_M] \in R^{B \times M}$

## 1.3 Method

### ***Representation Decomposition***

- The representation matrix  $Z \approx UV$ , is decomposed as product of two low-rank matrices  $U \in R^{B \times K}$ ,  $V \in R^{K \times M}$
- The columns of  $U = [u_1, \dots, u_k]$  are the basis vectors aka concept vectors/concepts. The rows of  $V$  or columns of  $V^T = [v_1, \dots, v_k]$ ,  $v_i \in R^M$  denote the activations of  $u_i$  for each sample.
- The optimization problem to decompose  $Z$  via Semi-NMF:

$$\mathbf{U}^*, \mathbf{V}^* = \arg \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{Z} - \mathbf{UV}\|_F^2 + \lambda \|\mathbf{V}\|_1 \quad s.t. \quad \mathbf{V} \geq 0, \text{ and } \|u_k\|_2 \leq 1 \quad \forall k \in \{1, \dots, K\}$$

## 1.3 Method

### *Representation Decomposition*

(cont.)

- **How to derive activations for each image?** Given a sample  $X$ ,  $z_X$  is projected on  $U^*$  to compute  $v(X)$ :

$$v(X) = \arg \min_{v \geq 0} \|z_X - \mathbf{U}^* v\|_2^2 + \lambda \|v\|_1$$

- **How to find textual equivalent of concepts?** Since the concept vectors are defined in the token representation space of  $f_{LM'}$  we can leverage the insights from “Lens” based methods that attempt to understand LLM representations using the model’s unembedding layer  $W_U$ .
- **Method?** Map using  $W_U u_k$  and extract tokens with highest probability, then filter for english, non-stop-word words with at least 3 characters. The final set of words is denoted as  $T_k$ .

## 1.3 Method

### ***Using Concept Dictionary for Interpretation***

- To find images from  $\mathbf{X}$  that maximally activate concept  $u_k$  (MAS):

$$\mathbf{X}_{k,MAS} = \operatorname{argmax}_{\hat{X} \subset \mathbf{X}, |\hat{X}|=N_{MAS}} \sum_{X \in \hat{X}} |v_k(X)|$$

- To find most activating concepts for image X:

$$\tilde{u}(X) = \{u_k \in U^* | k \in \arg \max_{k \in \{1, \dots, K\}, |S|=r} |v_k(X)|\}$$

This step could be further combined with concept importance estimation techniques to interpret the model's prediction for token  $t$ .

## 1.4 Results

Token	Simple	PCA	KMeans	Semi-NMF
Dog	0.371	<b>0.004</b>	0.501	<u>0.086</u>
Bus	0.622	<b>0.002</b>	0.487	<u>0.177</u>
Train	0.619	<b>0.015</b>	0.367	<u>0.107</u>
Cat	0.452	<b>0.000</b>	0.500	<u>0.146</u>

Table 2: Overlap between learnt concepts. Lower is better. Best score in **bold**, second best underlined.

# 1.4 Results



Figure 4: Visual/textual grounding for 8 out of 20 concepts for 'Dog' token (layer 31). For each concept we illustrate the set of 5 most activating samples and 5 most probable decoded words.

# 1.4 Results

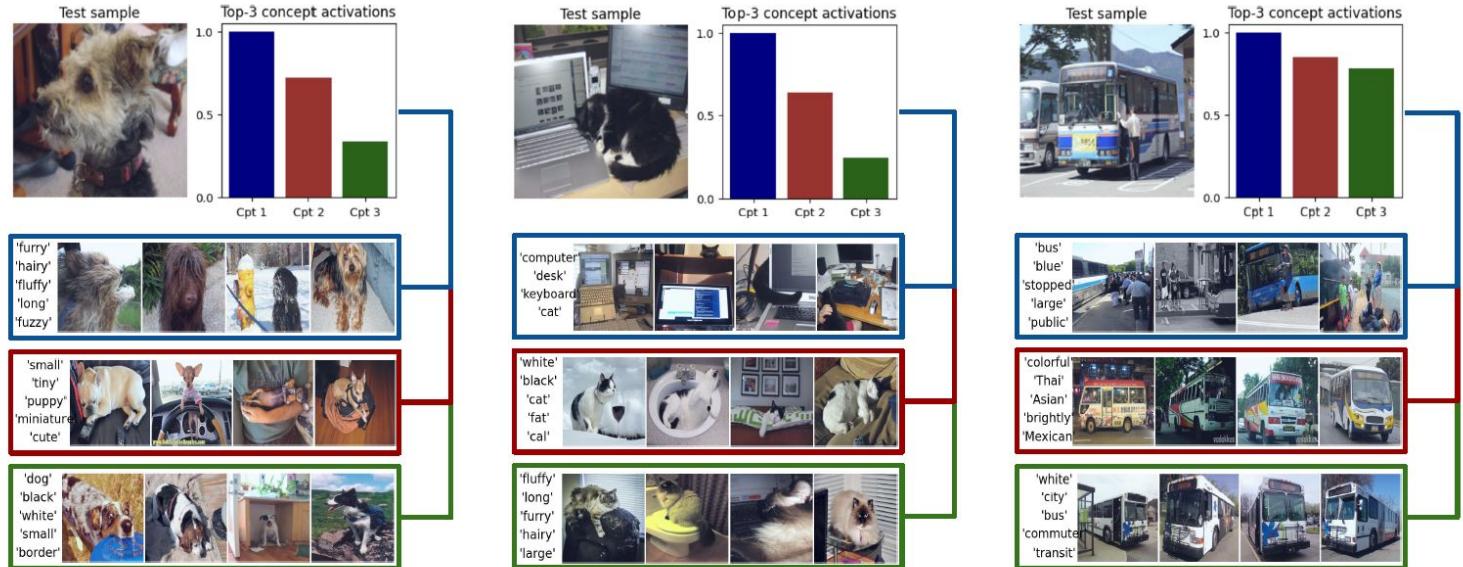


Figure 5: Local interpretations for test samples for different tokens ('Dog', 'Cat', 'Bus') with Semi-NMF (layer 31). Visual/text grounding for three highest concept activations (normalized) is shown.

## 1.4 Results

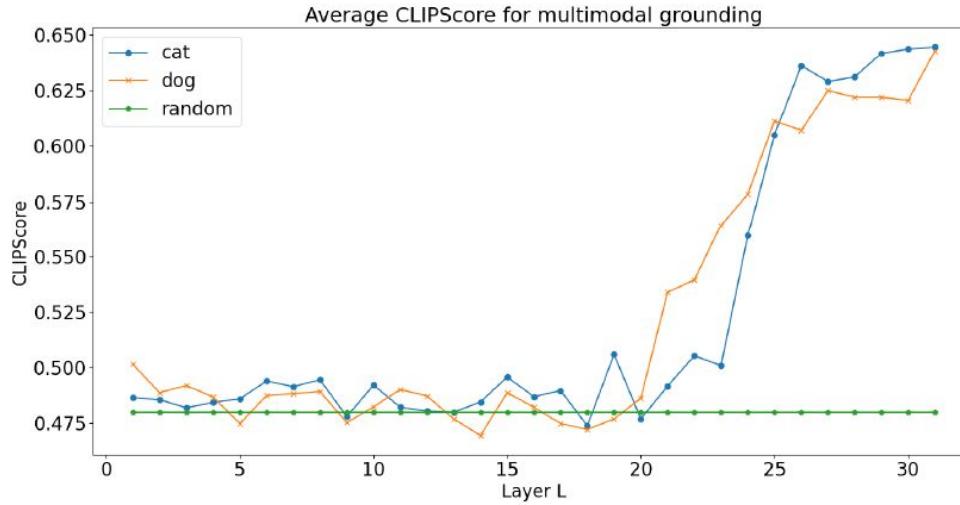


Figure 6: Mean CLIPScore between visual/text grounding  $\mathbf{X}_{k,MAS}, \mathbf{T}_k$  for all concepts (Semi-NMF), across different layers  $L$ . Results are for tokens ‘Dog’ and ‘Cat’.

02

ICML 2023

---

# Probabilistic Concept Bottleneck Models

## 2.1 Overview

**Class:** Green Jay

**Concepts:**

forehead color: blue  
throat color: black  
belly color: yellow  
tail pattern: solid



**Diverse visual contexts**



ambiguity in tail



ambiguity in belly



ambiguity in color

*Figure 1.* Examples of ambiguous cases in the existence of the concept. Images have diverse visual contexts, where partial concepts may become invisible and unclear.

## 2.1 Overview

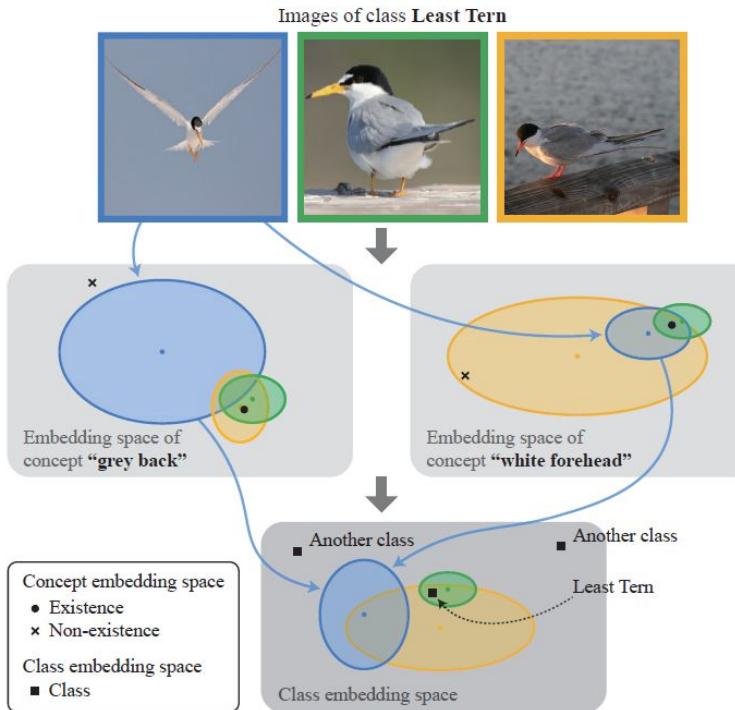


Figure 2. Probabilistic embeddings in ProbCBM. Images are

## 2.1 Overview

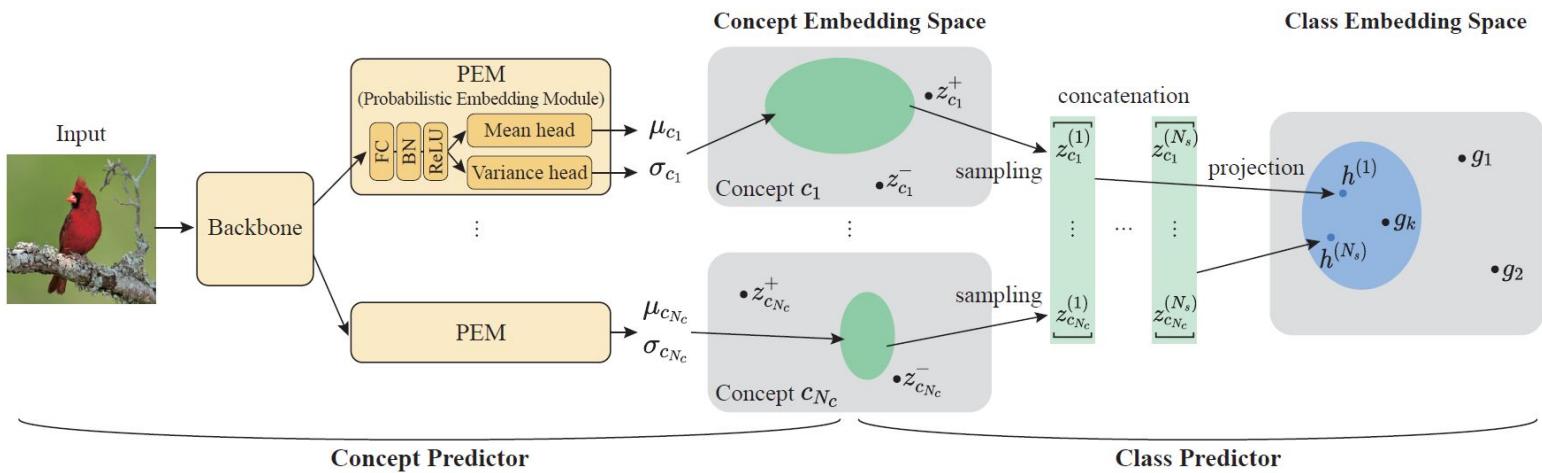


Figure 3. Prediction flow of ProbCBM. Ellipses represent probabilistic embeddings of a given input.

## 2.2 Preliminary

### ***Background***

- CBM (Koh et al., 2020) consists of two predictors: a concept predictor and a class predictor.
- Given an input  $x$ , the concept predictor  $g$  maps it to a concept space to predict a set of the existing concepts  $C$  ( $\hat{C} = g(x)$ ). The class predictor  $f$  estimates the class  $y$  from  $\hat{C}$ , i.e., the concepts predicted by the concept predictor ( $\hat{y} = f(\hat{C})$ ).
- To learn the mapping  $g$ , CBM requires the dataset  $\{x(i), C(i), y(i)\}$  while the conventional classifier that directly maps  $x$  to  $y$  require the dataset  $\{x(i), y(i)\}$

## 2.3 Method

### ***Overview***

- ProbCBM: an interpretable model that exploits probabilistic embedding in prediction. It provides concept prediction and concept uncertainty as explanations.
- It adds an interpretable layer into the architecture!



## 2.3 Method

### ***Probabilistic Concept Embedding***

- Given an input  $x$ , the concept predictor makes probabilistic concept embedding for each concept  $c \in C$ , which is formulated as a normal distribution with a mean vector and a diagonal covariance matrix:

$$p(z_c|x) \sim \mathcal{N}(\mu_c, \text{diag}(\sigma_c))$$

- $\mu_c$  and  $\sigma_c$  are predicted by a probabilistic embedding module (PEM) from a feature map extracted from the backbone. Self-attention-based mean and variance head submodules are used in PEM. To reduce the feature dimension, an FC layer was added followed by batch norm and ReLU activation before the mean and variance head modules.

## 2.3 Method

### ***Concept Prediction***

- With  $N_s$  representations sampled from  $p(z_c|x)$ , the probability of the existence of concept  $c$  ( $c = 1$ ) is obtained via Monte-Carlo estimation:

$$\begin{aligned} p(c = 1|x) &\approx \frac{1}{N_s} \sum_{n=1}^{N_s} p(c = 1|z_c^{(n)}), \\ p(c = 1|z_c^{(n)}) &= s \left( a \left( \|z_c^{(n)} - z_c^- \|_2 - \|z_c^{(n)} - z_c^+ \|_2 \right) \right), \end{aligned}$$

- Trainable anchors  $z_c^+$  and  $z_c^-$  represent existence and non-existence of concept  $c$ . For a sampled representation  $z_{c'}^{(n)}$ , the probability that the concept exists is the Euclidean distance with  $z_c^+$  and  $z_c^-$ .

## 2.3 Method

### ***Class Embedding***

- To generate a class embedding from concepts, concatenate sampled concept representations from all concepts, then project with FC:

$$h^{(n)} = \mathbf{w}^T \left( \left[ z_{c_1}^{(n)}, z_{c_2}^{(n)}, \dots, z_{c_{N_c}}^{(n)} \right] \right) + \mathbf{b}$$

### ***Class Prediction***

- Classification probability is obtained via Monte-Carlo estimation, based on Euclidean distance between the class embedding and a trainable anchor point for class  $k$ :

$$p(y_k = 1|x) \approx \frac{1}{N_s} \sum_n^{N_s} \frac{\exp(-d\|h^{(n)} - g_k\|_2)}{\sum_{k'} \exp(-d\|h^{(n)} - g_{k'}\|_2)}$$

## 2.3 Method

### ***Training***

- Add a KL divergence loss between the predicted concept embedding distributions and the standard normal distribution:

$$\mathcal{L}_{\text{concept}} = \mathcal{L}_{\text{BCE}} + \lambda_{\text{KL}} \mathcal{L}_{\text{KL}} \quad \mathcal{L}_{\text{KL}}(c) = \text{KL}(\mathcal{N}(\mu_c, \text{diag}(\sigma_c)) || \mathcal{N}(0, I))$$

- Train the concept predictor and class predictor separately. First, the concept predictor is trained with  $\mathcal{L}_{\text{concept}}$ . Then, the class predictor is trained with  $\mathcal{L}_{\text{class}}$  using the concept embeddings predicted by the concept predictor.

### ***Inference***

- Done via MC sampling or using  $\mu_c$  as  $z_c$  without sampling.

## 2.4 Results

Image	Class: 5 P: 1.000 U: 0.238	Class: 0 P: 0.000 U: 0.241	Class: 0 P: 0.002 U: <b>0.325</b>
Prediction on Concept 1 Classification	Class: 5 P: 0.445 U: 0.013	Class: 0 P: 0.947 U: 0.012	Class 0 P: 0.574 U: 0.014

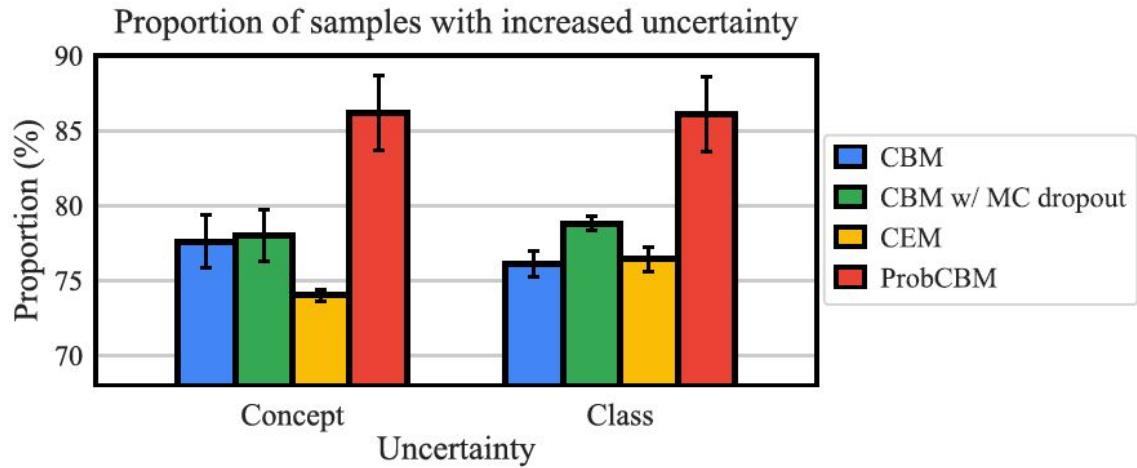
*Figure 4.* Examples of ProbCBM's prediction in diverse visual contexts. P and U represent the probability and uncertainty, respectively. The highest value of concept uncertainty is bolded.

## 2.4 Results

Table 1. Prediction accuracy for the CUB and AwA2 datasets. Results include mean values with standard deviation from three experiment repetitions. “w/o prob.” denotes ProbCBM with deterministic embeddings. “w/o sampling” denotes inference without sampling.

DATA	METHOD	CONCEPT	CLASS
IMAGE: 224 × 224			
CUB	CBM	0.950±0.001	0.670±0.006
	PROBCBM	0.949±0.001	0.680±0.004
	W/O PROB.	0.950±0.001	0.677±0.004
	W/O SAMPLING	0.949±0.001	0.679±0.003
	IMAGE: 299 × 299		
AwA2	CBM	0.956±0.001	0.708±0.006
	CEM	0.954±0.001	0.759±0.002
	PROBCBM	0.956±0.001	0.718±0.005
	CBM	0.975±0.001	0.877±0.004
	PROBCBM	0.975±0.000	0.880±0.002
	W/O PROB.	0.975±0.000	0.880±0.002
	W/O SAMPLING	0.975±0.000	0.880±0.001

## 2.4 Results



*Figure 6.* Proportion (%) of samples whose uncertainty increases after occlusion for the CUB dataset. The experiments are conducted with an image size of  $299 \times 299$ . Results include mean values with standard deviation from three experiment repetitions.

## 2.4 Results

(a) CUB

Target class:  
American Goldfinch



**Concept Prediction**

	P	U
wing color: white	1.000	0.178
upper tail color: black	1.000	0.119
breast color: yellow	0.999	0.164
belly color: yellow	1.000	0.093

	P	U
	0.988	0.258
	0.986	<b>0.311</b>
	0.998	0.170
	1.000	0.145

	P	U
	0.967	<b>0.301</b>
	0.995	0.257
	0.306	<b>0.264</b>
	0.016	<b>0.303</b>

**Class Prediction**

American Goldfinch  
P: 0.972 U: 2.436

American Goldfinch  
P: 0.954 U: 3.019

American Goldfinch  
P: 0.791 U: 3.663

(b) AwA2

Target class:  
Deer



**Concept Prediction**

	P	U
furry	0.999	<b>0.517</b>
spots	1.000	0.033
horns	1.000	0.129
forest	1.000	0.059

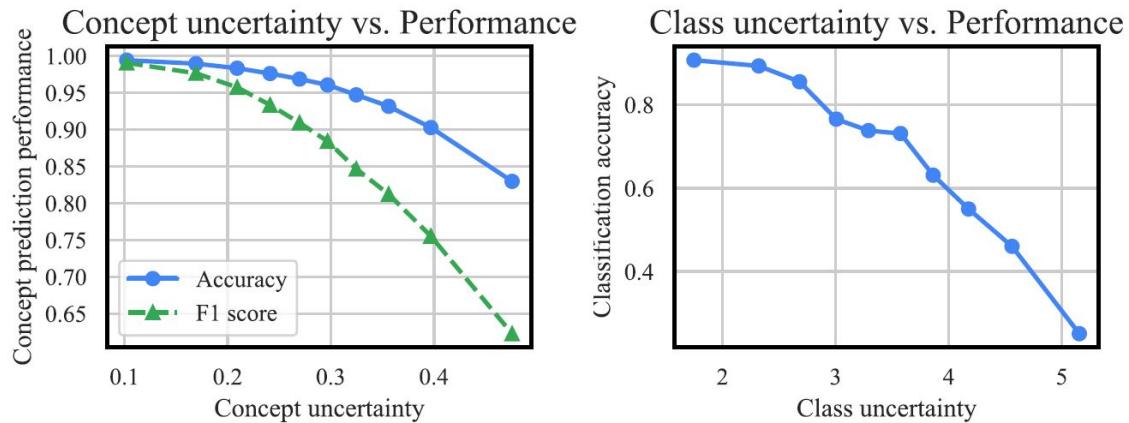
**Class Prediction**

Deer  
P: 1.000 U: 0.188

Deer  
P: 0.999 U: 0.233

Figure 7. Examples of ProbCBM’s prediction in diverse visual contexts for the (a) CUB and (b) AwA2 datasets. P and U represent the probability and uncertainty, respectively. The highest value of uncertainty in each concept is bolded.

## 2.4 Results



*Figure 10.* Correlation between performance and uncertainty for the CUB dataset. A point represents one group of samples with similar uncertainties and the x-axis indicates the median uncertainty for each group.

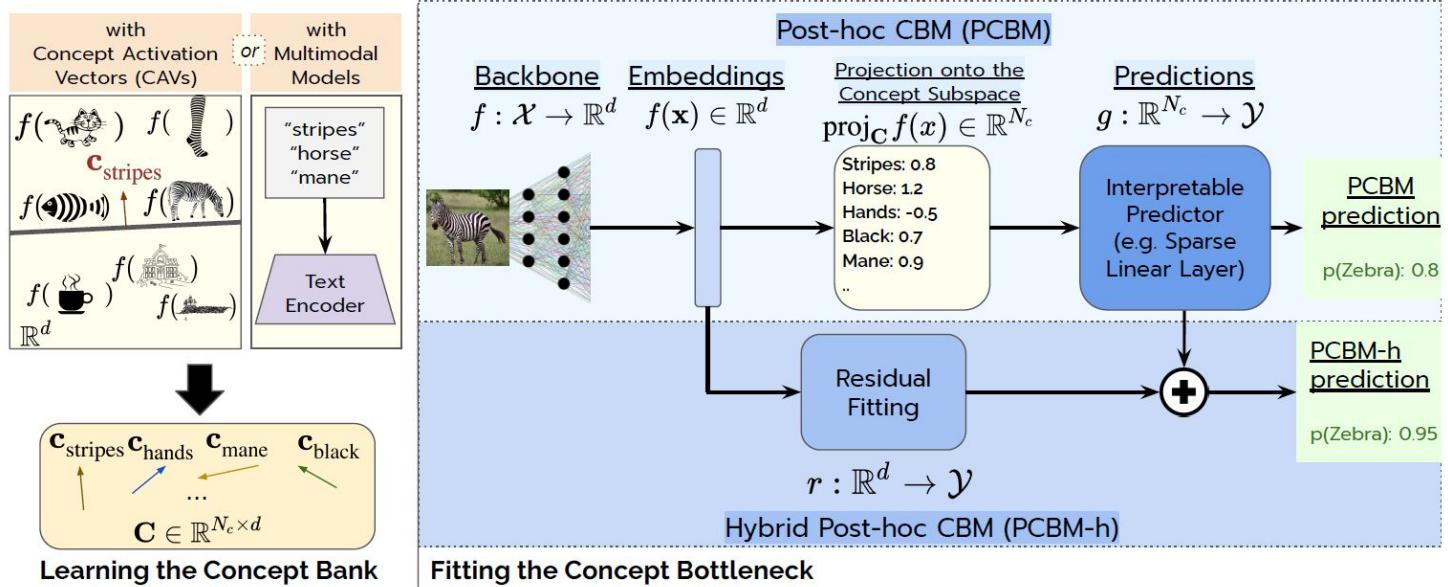
# 03

ICLR 2023

---

## **Post-Hoc Concept Bottleneck Models**

# 3.1 Overview



## 3.2 Preliminary

### **Model Architecture**

- Any pretrained backbone model defined as  $f : \mathcal{X} \rightarrow \mathbb{R}^d$   
 $d$  is the size of embedding space and  $X$  is the input space
- **Examples?**  $f$  can be:
  - the image encoder of CLIP
  - the model layers up to the penultimate layer of a ResNet



## 3.3 Method

### ***Learning the Concept Subspace (c)***

- First define a concept library  $I = \{i_1, i_2, \dots, i_{N_c}\}$

*Concepts can be selected by a domain expert or learned automatically from the data (Ghorbani et al., 2019; Yeh et al., 2020)*

- For each concept  $i$ , collect embeddings for:

- positive examples  $P_i = \{f(x_{p1}), \dots, f(x_{pNp})\}$

- negative examples  $N_i = \{f(x_{n1}), \dots, f(x_{pNn})\}$

*Unlike normal CBMc, these samples can be different from the data used to train the backbone model*



## 3.3 Method

### ***Learning the Concept Subspace ( $\mathbf{c}$ )***

**(cont.)**

- Train a linear SVM using  $P_i$  and  $N_i$  to learn the corresponding CAV (the vector normal to the linear classification boundary)
- Denote the CAV for concept  $i$  as  $c_i$
- Denote  $C \in R^{N_c \times d}$  as matrix of CAVs (rows are  $c_i$ )
- Given an input, project the embedding of the input onto the subspace spanned by CAVs (the concept subspace):

$$f_C(x) = \text{proj}_C f(x) \in \mathbb{R}^{N_c} \quad f_C^{(i)}(x) = \frac{\langle f(x), c_i \rangle}{\|c_i\|_2^2} \in \mathbb{R}$$

## 3.3 Method

### **Leveraging Multimodal Models to Learn Concepts**

- No need for annotation and labeled data!
- Multimodal models such as CLIP have a text encoder  $f_{\text{text}}$  along with an image encoder, which maps a description to the shared embedding space.
- Therefore use:  
 $c_{\text{stripes}}^{\text{text}} = f_{\text{text}}(\text{"stripes"})$
- Collect the text embeddings and construct our multimodal concept bank  $C^{\text{text}}$  as the concept subspace.

*For a given classification task, we can use ConceptNet to obtain concepts relevant to these classes. ConceptNet is an open knowledge graph, where we can find concepts that have particular relations to a query concept.*

## 3.3 Method

### ***Learning the Interpretable Predictor***

- An interpretable sparse predictor that maps the concept subspace to the model prediction:  
$$g : \mathbb{R}^{N_c} \rightarrow \mathcal{Y}$$
- To learn the PCBM, we solve:

$$\min_g \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(g(f_C(x)), y)] + \frac{\lambda}{N_c K} \Omega(g)$$

$$g(f_C(x)) = \mathbf{w}^T f_C(x) + b \quad \Omega(g) = \alpha \|\mathbf{w}\|_1 + (1 - \alpha) \|\mathbf{w}\|_2^2$$

*However, even when we have a rich concept subspace, concepts may not be enough to solve a task of interest.*

## 3.3 Method

### ***Recovering Original Model Performance with Residual Modeling***

- We introduce Hybrid Post-Hoc CBMs (PCBM-h).
- After fixing the concept bottleneck and the interpretable predictor, re-introduce the embeddings to ‘fit the residuals’:

$$\min_r \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}(g(f_C(x)) + r(f(x)), y)]$$

$r : \mathbb{R}^d \rightarrow \mathcal{Y}$  is the residual predictor  $r(f(x)) = w_r^T f(x) + b_r$

- $r$  can be discarded to observe normal performance.

## 3.4 Results

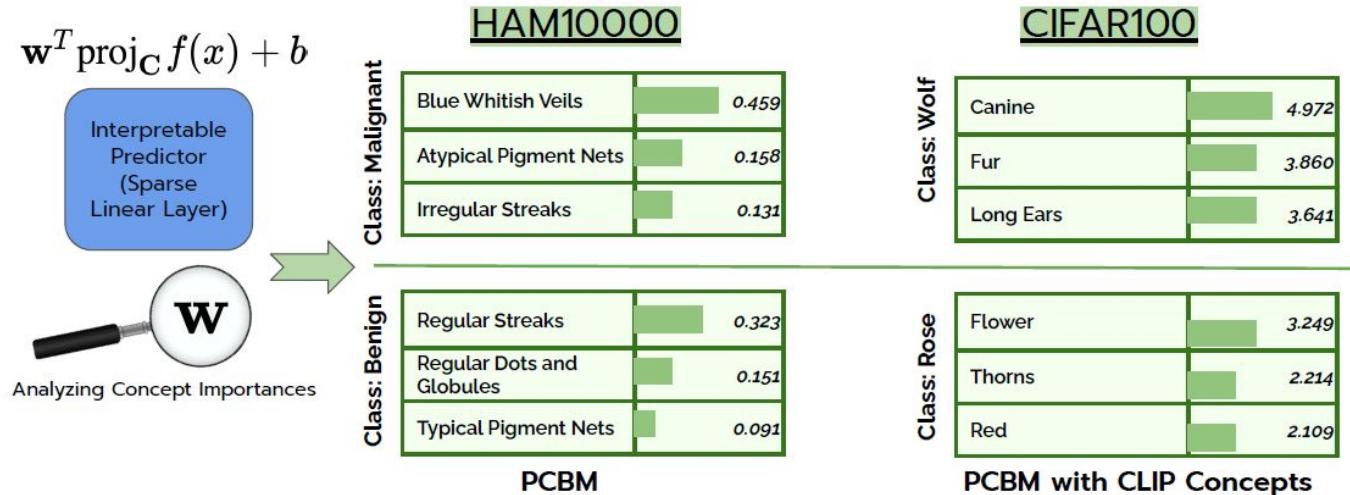


Figure 2: **Explaining Post-hoc CBMs.** We report the top 3 largest weights in the linear layer for the shown classes. For instance, *Blue Whitish Veils*, *Atypical Pigment Networks*, and *Irregular Streaks* have large weights for classifying whether a skin lesion is malignant. These are consistent with dermatologists' findings (Menzies et al., 1996).

## 3.4 Results

Table 1: **PCBMs achieve comparable performance to the original model.** We report performance over different scenarios for the original model and PCBMs with concept datasets. In CIFAR100, PCBm performs poorly since the concept bank is not expressive enough to solve a finer-grained task; however, PCBm-h recovers the original model’s accuracy. Strikingly, PCBMs match the performance of the original model in HAM10000 and ISIC, *with as few as 8 human-interpretable concepts*. Original CBMs cannot be trained on CIFAR/HAM10000/ISIC/COCO-Stuff, as they do not have concept labels in the training dataset. The mean and standard errors are reported over 10 random seeds. We report AUROC for HAM10000 and ISIC, mAP for COCO-Stuff, and accuracy for CIFAR and CUB.

	CIFAR10	CIFAR100	COCO-Stuff	CUB	HAM10000	ISIC
Original Model	0.888	0.701	0.770	0.612	0.963	0.821
PCBM	$0.777 \pm 0.003$	$0.520 \pm 0.005$	$0.741 \pm 0.002$	$0.588 \pm 0.008$	$0.947 \pm 0.001$	$0.736 \pm 0.012$
PCBM-h	$0.871 \pm 0.001$	$0.680 \pm 0.001$	$0.768 \pm 0.01$	$0.610 \pm 0.010$	$0.962 \pm 0.002$	$0.801 \pm 0.056$

## 3.4 Results

Table 2: **Concept Bottlenecks with CLIP concepts.** When a concept bank is not available or is insufficient, we can use natural language descriptions of concepts with CLIP to implement CBMs.

	<b>CIFAR10</b>	<b>CIFAR100</b>	<b>COCO-Stuff</b>
Original Model (CLIP)	0.888	0.701	0.770
PCBM & labeled concepts	$0.777 \pm 0.003$	$0.520 \pm 0.005$	$0.741 \pm 0.002$
PCBM & CLIP concepts	$0.833 \pm 0.003$	$0.600 \pm 0.003$	$0.755 \pm 0.001$
PCBM-h & CLIP concepts	$0.874 \pm 0.001$	$0.691 \pm 0.006$	$0.769 \pm 0.001$



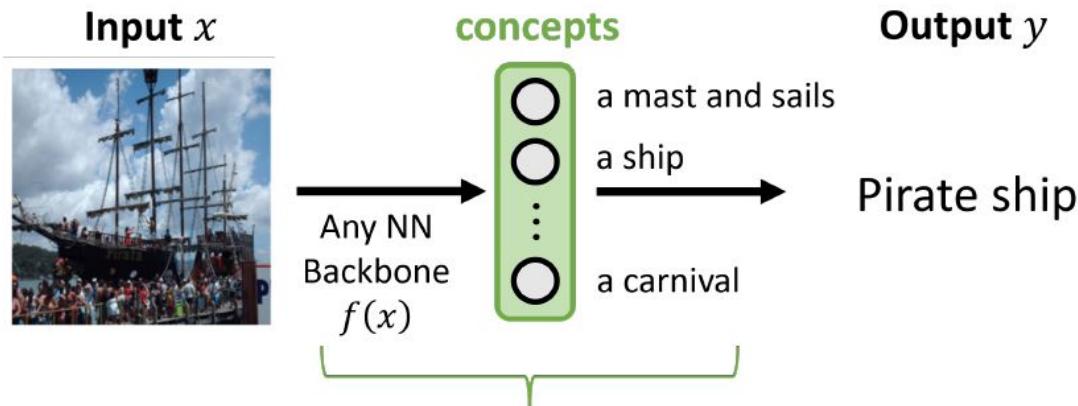
04

ICLR 2023

---

# **Label-Free Concept Bottleneck Models**

## 4.1 Overview



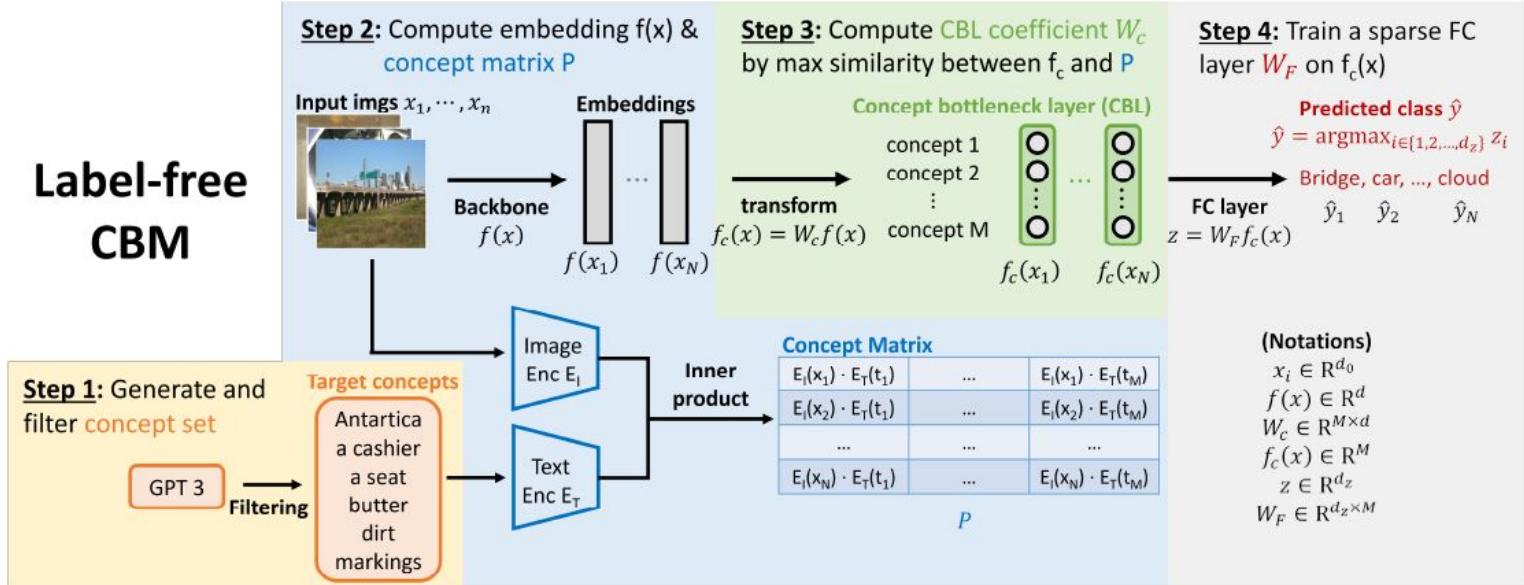
### Label-free CBM:

automated, scalable, efficient, no concept labels required



# 4.1 Overview

## Label-free CBM



## 4.2 Preliminary

### ***Background***

- CBMs create a layer before the last fully connected layer where each neuron corresponds to a human interpretable concept.
- Post-Hoc CBM only retrains the last FC layer to solve this.
- The problems of requiring annotated concept data for TCAV remains and their use of CLIP model can only be applied to if the NN backbone is the CLIP image encoder.
- The performance of PCBMs is uninterpretable due to res con.

## 4.2 Preliminary

Method:	(I) Flexibility		(II) Interpretability		(III) Performance	
	Without labeled concept data	Any network architecture	Sparse final layer	All features interpretable	Preserves accuracy	Extends to ImageNet scale
CBM	No	Yes	No	Yes	No	No
IBD	No	Yes	No	No	Yes	No
P-CBM	No	Yes	Yes	Yes	No	No
P-CBM (CLIP)	Yes	No	Yes	Yes	No	Maybe
P-CBM-h	No	Yes	Yes	No	Yes	No
P-CBM-h (CLIP)	Yes	No	Yes	No	Yes	Maybe
Label-free CBM <b>(This work)</b>	Yes	Yes	Yes	Yes	Yes	Yes

Table 1: Comparison of our method against existing methods for creating Concept Bottleneck models, CBM (Koh et al., 2020), IBD (Zhou et al., 2018) and 4 versions of P-CBM (Yuksekgonul et al., 2022), where ‘-h’ indicates the hybrid model that uses uninterpretable residual term, ‘(CLIP)’ means models using CLIP concepts. We used maybe to indicate models that could in theory extend to ImageNet but have not been tested.

## 4.3 Method

### ***Step 1.A: Initial Concept Set Creation***

- Ask GPT-3 to:
  - List the most important features for recognizing something as a {class}:
  - List the things most commonly seen around a {class}:
  - Give superclasses for the word {class}:

*We found using GPT-3 to generate initial concepts to perform better than using the knowledge-graph ConceptNet which was used in PCBM. Reasons for this and a comparison between the two are presented in Appendix A.6.*



## 4.3 Method

### **Step 1.B: Concept Set Filtering**

- Filters to improve quality and reduce the size of concept set:
  - Concept length
  - Remove concepts too similar to classes
  - Remove concepts too similar to each other
  - Remove concepts not present in training data
  - Remove concepts we can't project accurately  
*(last item is performed after step 3)*

## 4.3 Method

### Step 2: Compute Embeddings

- Given concepts  $C = \{t_1, \dots, t_M\}$ , training dataset  $D = \{x_1, \dots, x_N\}$ , calculate and save the CLIP concept activation matrix  $P$  using the CLIP image and text encoders:

$$P_{i,j} = E_I(x_i) \cdot E_T(t_j)$$

- This is independent of the backbone model, thus we need to learn a projection from the backbone model's feature space into a space where axis directions correspond to interpretable concepts.*

## 4.3 Method

### ***Step 3: Learn Projection to Create a Concept Bottleneck Layer***

- Initialize projection weights  $W_c$  as a random  $M \times d_0$  matrix where  $d_0$  is the dimensionality of backbone features  $f(x)$ .
- Define  $f_c(x) = W_c f(x)$ , where  $f_c(x_i) \in R^M$ .
- Use  $k$  to denote a neuron of interest in the projection layer, and its activation pattern is  $q_k$  where  $q_k = [f_{c,k}(x_1), \dots, f_{c,k}(x_N)]^\top$ , with  $q_k \in R^N$  and  $f_{c,k}(x) = [f_c(x)]_k$ .
- Now, how do we train each neuron? Use CLIP-Dissect.

## 4.3 Method

### **Step 3: Learn Projection to Create a CBL**

**(cont.)**

- To make the neurons in the CBL interpretable, we need to enforce the projected neurons to activate in correlation with the target concept, which we do by optimizing  $W_c$  to maximize the CLIP-Dissect similarity between the neuron's activation pattern and the target concept.
- To improve this similarity, optimize:

$$L(W_c) = \sum_{i=1}^M -\text{sim}(t_i, q_i) := \sum_{i=1}^M -\frac{\bar{q}_i^3 \cdot \bar{P}_{:,i}^3}{\|\bar{q}_i^3\|_2 \|\bar{P}_{:,i}^3\|_2}$$

## 4.3 Method

### **Step 4: Learning the Sparse Final Layer**

- Learn the predictor with the fully connected layer  $W_F \in R^{dz \times M}$  where  $dz$  is the number of output classes.
- Keep  $W_F$  sparse, since sparse layers have been demonstrated to be more interpretable.
- Solve elastic net objective:

$$\min_{W_F, b_F} \sum_{i=1}^N L_{ce}(W_F f_c(x_i) + b_F, y_i) + \lambda R_\alpha(W_F)$$

$$R_\alpha(W_F) = (1 - \alpha) \frac{1}{2} \|W_F\|_F + \alpha \|W_F\|_{1,1}$$

## 4.4 Results

Model	Sparse final layer	Dataset				
		CIFAR10	CIFAR100	CUB200	Places365	ImageNet
Standard	No	88.80%	70.10%	76.70%	48.56%	76.13%
Standard (sparse)	Yes	82.96%	58.34%	<b>75.96%</b>	38.46%	<b>74.35%</b>
P-CBM*	Yes	70.50%	43.20%	59.60%	N/A	N/A
P-CBM (CLIP)*	Yes	84.50%	56.00%	N/A	N/A	N/A
Label-free CBM <b>(Ours)</b>	Yes	<b>86.40%</b> ± 0.06%	<b>65.13%</b> ± 0.12%	74.31% ± 0.29%	<b>43.68%</b> ± 0.10%	71.95% ± 0.05%

Table 2: Accuracy comparison, best performing sparse model bolded. We can see our method outperforms Posthoc-CBM and performs similarly to a sparse standard model. The results for our method are mean and standard deviation over three training runs. \*Indicates reported accuracy.

## 4.4 Results

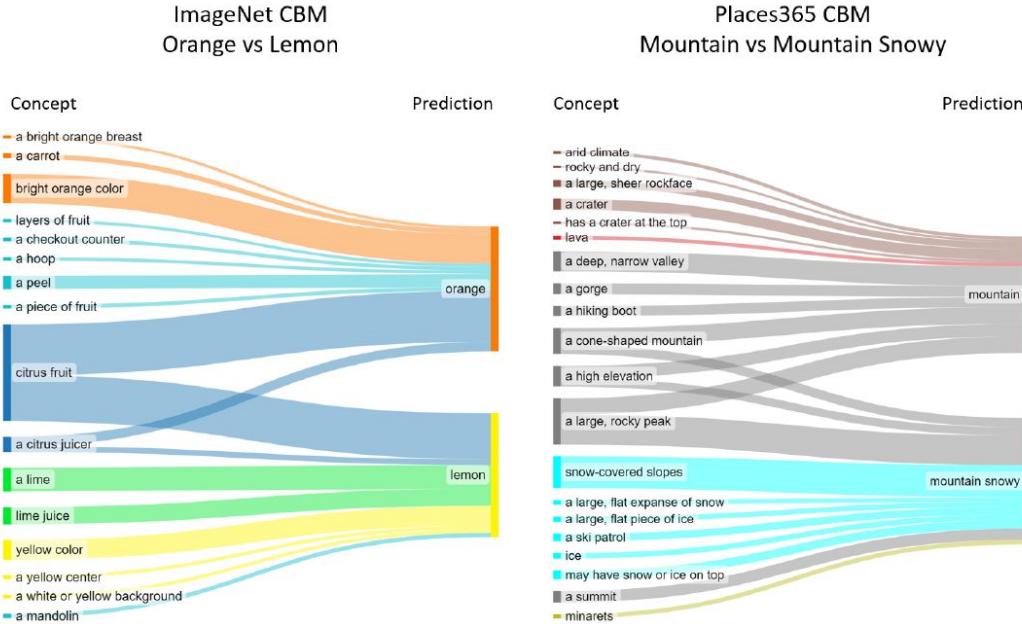


Figure 3: Visualization of the final layer weights of our Label-free CBM. Showcasing how our models differentiate between two similar classes.

# 05

A C L 2 0 2 4

---

## **Explaining Language Model Predictions with High-Impact Concepts**

## 5.1 Overview

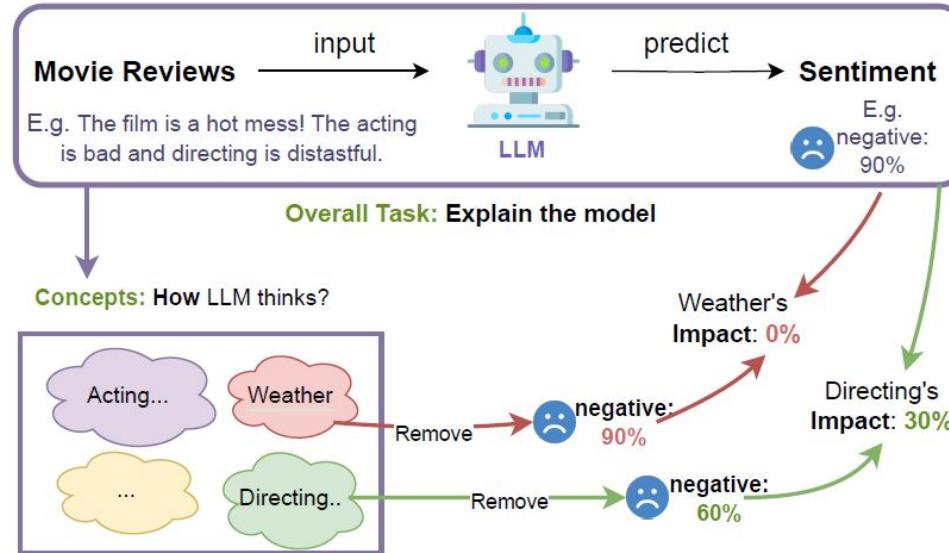


Figure 1: Illustration of concept-based explanations that result in high impact (green line) or not (red line) when explaining the LLMs in a sentiment classification task.

## 5.1 Overview

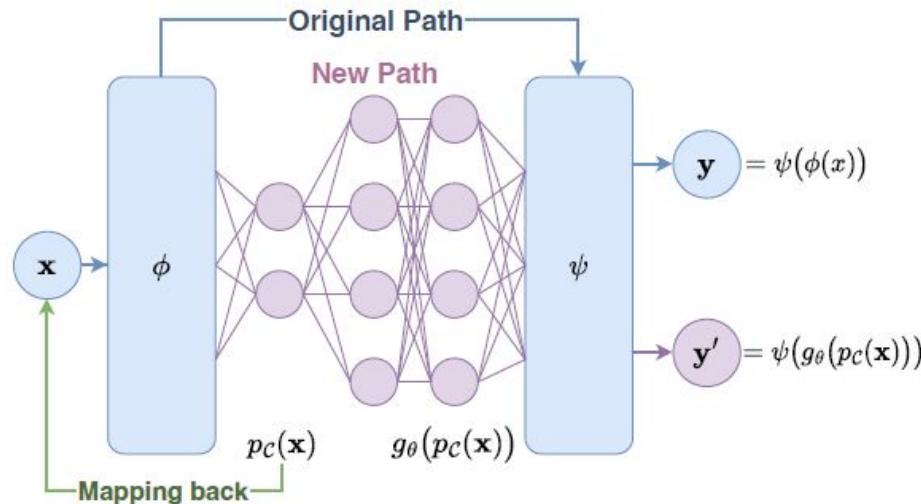


Figure 2: The overall concept generation process of a concept bottleneck model.

## 5.1 Overview

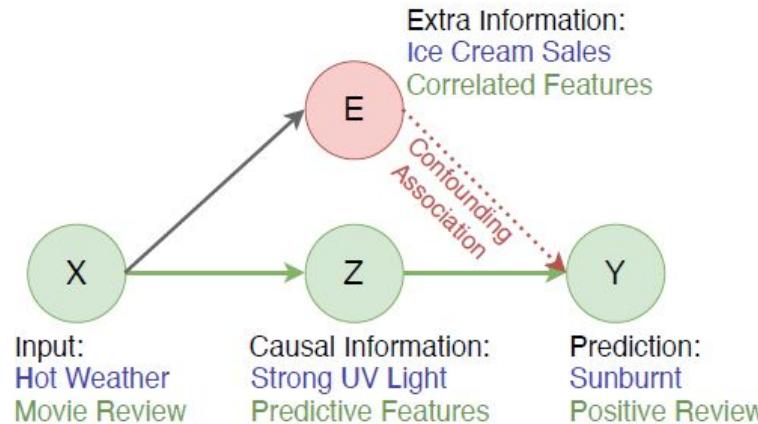


Figure 3: Illustration of the causal graph indicating the confounding association in explanation models. Blue is a real-life example. Green is the correspondence in a movie review classification task.

## 5.2 Preliminary

### *Background*

- Traditional CBMs that do not consider impact may falsely use the correlational feature to explain causation.
- Previous causality evaluations focused on assessing the causal effect via local (i.e., instance-level) change and removal intervention, leading to potential bias.
- **Alternative?** High-Impact Concepts (HI-concept)
  - Concept explanation framework with causal impact optimization.
  - Can be used on a pretrained whitebox LLM.
  - Employs a global (i.e., model-level) accuracy change metric and insertion operation to diagnose the causality measurement.

## 5.2 Preliminary

### CBM Losses

- **Reconstruction loss:** To recover the original model's predictions after inserting a step in the middle of two sections of model (hidden layers and final prediction).

$$\begin{aligned}\mathcal{L}_{\text{rec}}(\theta, \mathcal{C}) &= \text{CE}\left(\psi(\phi(\mathbf{x})), \psi(g_{\theta}(p_{\mathcal{C}}(\mathbf{x})))\right) \\ &= - \sum_{b \in \mathcal{B}} \psi(\phi(\mathbf{x}))_b \log (\psi(g_{\theta}(p_{\mathcal{C}}(\mathbf{x})))_b)\end{aligned}$$

- **Regularization loss:** To make each concept vector correspond to actual examples and concepts to be distinct from each other.

$$\begin{aligned}\mathcal{L}_{\text{reg}}(\mathcal{C}) &= -\lambda_1 \frac{\sum_{i=1}^n \sum_{\mathbf{x}_t \in \mathcal{T}_{\mathbf{c}_i}} \mathbf{c}_i^\top \phi(\mathbf{x}_t)}{nN} \\ &\quad + \lambda_2 \frac{\sum_{i_1 \neq i_2} \mathbf{c}_{i_1}^\top \mathbf{c}_{i_2}}{n(n-1)}.\end{aligned}$$

## 5.2 Preliminary

### ***Causal Analysis***

- To define impact in our method, we utilize two important definitions in causal analysis: Individual Treatment Effect (ITE) and Average Treatment Effect (ATE), which measure the effect of interventions in randomized experiments.

$$\begin{aligned} \text{ITE}(x) &:= \mathbb{E}[y|\mathbf{X} = x, \text{do}(T = 1)] \\ &\quad - \mathbb{E}[y|\mathbf{X} = x, \text{do}(T = 0)]; \\ \text{ATE} &:= \mathbb{E}[\text{ITE}(x)]. \end{aligned}$$

## 5.3 Method

### ***Defining Impact***

- To separate correlation from causation, enforce explanations to be predictive by considering their “impact”.
- A concept  $c_i$  is a cause if its removal changes the outcome, and corresponds to a direction in latent space. As  $f$  is fixed, its prediction process is deterministic and reproducible, allowing us to conduct experiments.
- Thus, we define the do-operation of ITE as concept removal and define impact of the concept as:

$$I(\mathbf{c}_i, \mathbf{x}) = \mathbb{E}[y | \mathbf{X} = \mathbf{x}, \mathbf{c}_i = \mathbf{0}] - \mathbb{E}[y | \mathbf{X} = \mathbf{x}, \mathbf{c}_i = \mathbf{c}_i]$$

## 5.3 Method

### ***Optimizing for Impact***

- To incorporate consideration for impact into the concept discovery process, we introduce two new losses to the original framework.
- **Auto-encoding loss:** Reconstructs the hidden representations.

$$\begin{aligned}\mathcal{L}_{\text{enc}}(\theta, \mathcal{C}) &= \text{MSE}\left(\phi(\mathbf{x}), g_{\theta}(pc(\mathbf{x}))\right) \\ &= \frac{1}{d} \|\phi(\mathbf{x}) - g_{\theta}(pc(\mathbf{x}))\|_2^2.\end{aligned}$$

- **Causality loss:** Disentangles concept directions that have a greater impact by optimizing the following loss:

$$\begin{aligned}\mathcal{L}_{\text{cau}}(\theta, \mathcal{C}) &= - \sum_{\mathbf{c}_i \in \mathcal{S}} \sum_{\mathbf{x}_j \in \mathcal{D}} \left| \psi\left(g_{\theta}(pc(\mathbf{x}_j) | \mathbf{c}_i = \mathbf{0})\right) \right. \\ &\quad \left. - \psi\left(g_{\theta}(pc(\mathbf{x}_j) | \mathbf{c}_i = \mathbf{c}_i)\right) \right| \approx -|I_{\text{avg}}(\mathcal{C})|.\end{aligned}$$

## 5.3 Method

### ***Final Method Steps***

- Extract concepts from the model's hidden representations.
  - Run the dataset to collect activations of a layer from all samples.
- Discover human-understandable concepts.
  - Assume there is a lower-dimensional concept space  $p_c(x)$  that can approximate the model's hidden space  $\phi(x)$ .
  - Initialize a random set of concept vectors  $C$  in the hidden space.
  - Optimize  $c_i$  using the losses:
$$\begin{aligned}\mathcal{L}(\theta, \mathcal{C}) = & \mathcal{L}_{\text{rec}}(\theta, \mathcal{C}) + \mathcal{L}_{\text{reg}}(\mathcal{C}) \\ & + \lambda_e \mathcal{L}_{\text{enc}}(\theta, \mathcal{C}) + \lambda_c \mathcal{L}_{\text{cau}}(\theta, \mathcal{C})\end{aligned}$$
- Apply concepts to explain the model.
  - Project activations of a sample  $x$  onto concept space using  $C^T \phi(x)$ .
  - This  $P_c(x)$  is prob of concept importance, use to highlight words.

## 5.3 Method

### **Evaluating Impact of Concepts**

- Recovering Accuracy Change ( $\Delta\text{Acc}$ ) is introduced to measure whether the model's ability to faithfully recover predictions after concept omission is disrupted:

$$\Delta\text{Acc}(\mathcal{C}) = \frac{1}{|\mathcal{C}|} \sum_{\mathbf{c}_i \in \mathcal{C}} |\text{Acc}(\mathcal{C}) - \text{Acc}(\mathcal{C} \setminus \{\mathbf{c}_i\})|$$

- We follow previous work to use Causal Concept Effect (CACE) to evaluate the causal effect of the set of concepts  $C$ :

$$\text{CACE}(\mathbf{c}_i) := \sum_{\mathbf{x}_j \in \mathcal{D}_{\text{test}}} |\psi(g_{\theta}(p_{\mathcal{C}}(\mathbf{x}_j))) - \psi(g_{\theta}(p_{\mathcal{C} \setminus \{i\}}(\mathbf{x}_j)))|;$$

$$\text{CACE}(\mathcal{C}) = \frac{1}{|\mathcal{C}|} \sum_{\mathbf{c}_i \in \mathcal{C}} \text{CACE}(\mathbf{c}_i)$$

## 5.4 Results

$p_{\text{cor}}$	Cls.Acc	Method	Acc	CACE	$\Delta\text{Acc}$
0.50	95.4%	ConceptSHAP	97.6%	0.070	6.1%
		<i>HI-concept</i>	<b>98.4%</b>	<b>0.102</b>	<b>9.4% (+3.3%)</b>
0.65	99.0%	ConceptSHAP	<b>99.7%</b>	0.038	3.5%
		<i>HI-concept</i>	99.3%	<b>0.084</b>	<b>6.8% (+3.4%)</b>
0.75	96.1%	ConceptSHAP	98.3%	0.069	6.0%
		<i>HI-concept</i>	<b>98.9%</b>	<b>0.123</b>	<b>12.2% (+6.2%)</b>

Table 1: Faithfulness (Acc) and Causality (CACE,  $\Delta\text{Acc}$ ) evaluation on the toy dataset. Cls.Acc denotes the original classification model’s accuracy.

## 5.4 Results

Dataset	Model	Method	Faithfulness					Causality	
			Acc	Precision	Recall	F1	Completeness	CACE	$\Delta$ Acc
IMDB	Transformer	$\beta$ -TCVAE (Chen et al., 2018)	43.53%	50.23	50.03	33.08	27.36	0.037	1.50%
		K-means (Likas et al., 2003)	83.64%	84.74	85.05	83.63	61.87	<u>0.047</u>	<u>2.59%</u>
		PCA (F.R.S., 1901)	<u>85.18%</u>	<u>85.56</u>	<u>86.20</u>	<u>85.15</u>	<b>62.36</b>	0.001	0.01%
		ConceptSHAP (Yeh et al., 2020)	84.36%	85.04	85.56	84.34	<u>62.05</u>	0.031	1.30%
		<i>HI-concept</i>	<b>88.78%</b>	<b>90.07</b>	<b>87.50</b>	<b>88.24</b>	58.10	<b>0.150</b>	<b>11.06%</b>
IMDB	BERT	$\beta$ -TCVAE (Chen et al., 2018)	93.86%	94.31	93.43	93.68	10.71	<u>0.057</u>	<u>4.05%</u>
		K-means (Likas et al., 2003)	<b>98.69%</b>	<b>96.16</b>	<b>96.23</b>	<b>96.19</b>	15.69	0.037	0.97%
		PCA (F.R.S., 1901)	<u>96.68%</u>	<b>96.65</b>	<b>96.68</b>	<b>96.67</b>	15.33	0.002	0.02%
		ConceptSHAP (Yeh et al., 2020)	95.84%	95.78	95.96	95.83	<u>17.16</u>	0.050	0.06%
		<i>HI-concept</i>	92.97%	93.25	93.34	92.97	<b>21.04</b>	<b>0.099</b>	<b>8.99%</b>
Llama	Llama	$\beta$ -TCVAE (Chen et al., 2018)	20.56%	33.41	33.36	13.30	-14.29	0.001	0.15%
		K-means (Likas et al., 2003)	15.31%	5.10	33.33	8.85	-21.82	<u>0.019</u>	0.00%
		PCA (F.R.S., 1901)	<b>95.15%</b>	<b>67.97</b>	<b>77.66</b>	<b>69.80</b>	<b>64.19</b>	0.001	0.03%
		ConceptSHAP (Yeh et al., 2020)	18.83%	42.83	34.95	14.88	-1.78	0.005	<u>1.60%</u>
		<i>HI-concept</i>	<u>87.87%</u>	<u>53.27</u>	<u>68.60</u>	<u>55.29</u>	<u>59.83</u>	<b>0.042</b>	<b>28.69%</b>
AG-News	Transformer	$\beta$ -TCVAE (Chen et al., 2018)	98.91%	98.94	98.93	<b>66.73</b>	<u>0.049</u>	<u>6.62%</u>	
		K-means (Likas et al., 2003)	98.16%	98.32	98.11	98.18	65.99	0.044	0.07%
		PCA (F.R.S., 1901)	<b>99.99%</b>	<b>99.99</b>	<b>99.99</b>	<b>99.99</b>	66.66	0.000	0.03%
		ConceptSHAP (Yeh et al., 2020)	73.01%	59.36	74.34	64.88	47.07	0.000	0.00%
		<i>HI-concept</i>	<u>99.50%</u>	<u>99.50</u>	<u>99.51</u>	<u>99.50</u>	<u>66.70</u>	<b>0.046</b>	<b>7.12%</b>
AG-News	BERT	$\beta$ -TCVAE (Chen et al., 2018)	92.30%	94.93	91.89	92.91	57.25	<u>0.044</u>	5.32%
		K-means (Likas et al., 2003)	86.83%	92.74	85.42	87.53	52.62	0.028	<u>7.15%</u>
		PCA (F.R.S., 1901)	<b>99.79%</b>	<b>99.82</b>	<b>99.77</b>	<b>99.79</b>	61.04	0.001	0.01%
		ConceptSHAP (Yeh et al., 2020)	93.46%	93.70	94.62	93.66	<b>62.69</b>	0.025	4.44%
		<i>HI-concept</i>	<b>99.90%</b>	<b>99.89</b>	<b>99.90</b>	<b>99.89</b>	<u>61.12</u>	<b>0.058</b>	<b>10.54%</b>
Llama	Llama	$\beta$ -TCVAE (Chen et al., 2018)	1.27%	0.25	20.00	0.50	-23.89	0.000	0.01%
		K-means (Likas et al., 2003)	37.00%	7.40	20.00	10.80	1.09	<u>0.007</u>	0.02%
		PCA (F.R.S., 1901)	<b>85.41%</b>	<b>65.78</b>	<b>67.98</b>	<b>66.73</b>	<b>51.46</b>	0.000	0.03%
		ConceptSHAP (Yeh et al., 2020)	17.01%	35.37	35.20	15.87	-7.73	0.002	<u>2.83%</u>
		<i>HI-concept</i>	<u>81.52%</u>	<u>48.59</u>	<u>55.99</u>	<u>51.53</u>	<u>43.07</u>	<b>0.039</b>	<b>27.79%</b>

Table 2: Faithfulness (Acc, Precision, Recall, F1, Completeness) and causality (CACE,  $\Delta$ Acc) evaluation of different text classification methods. The best result is bolded, and the second-best result is underlined.

## 5.4 Results

Method	CACE	Keywords
CS	0.134	apple, NASA, Microsoft, new, sun, red, super, game
CS	0.000	one, two, gt, new, cl, lt, first, world, mo, last
HI-C	0.130	us, bush, u, eu, new, peoples, china, high, gt, world
HI-C	0.003	us, update, new, mo, two, first, knicks, last, one, hen

Table 3: Generated concepts with Average Impact (CACE) from AG-News dataset, BERT model. CS is ConceptSHAP, HI-C is *HI-concept*. Each line is one concept, represented by keywords, which are ordered by descending importance.

## 5.4 Results

Method	Visualization
ConceptSHAP	dream team leads spain 44 - 42 at halftime athens, greece - as expected, the u.s. men's basketball team had its hands full in a quarterfinal game against spain on thursday...
<i>HI-concept</i>	dream team leads spain 44 - 42 at halftime athens, greece - as expected, the u.s. men's basketball team had its hands full in a quarterfinal game against spain on thursday ...

Figure 4: Qualitative comparison from AG-News: “World” news misclassified as “Sports” by BERT.

## 5.4 Results

	Accuracy	Confidence	Time Spent
Plain	72.5%	3.2	10.7
ConceptSHAP	68.5%	2.7	10.6
Polyjuice	73.5%	2.6	7.6
<i>HI-concept</i>	<b>80.5%</b>	<b>3.5</b>	<b>9.3</b>

Table 4: Human study for explainability evaluation.



## 5.4 Results

Method	Acc	CACE	$\Delta$ Acc
Without Auto-Encoding Loss	93.46%	0.028	6.11%
Without Prediction Loss	68.00%	0.035	<b>17.41%</b>
Without Regularizer Loss	95.76%	0.041	6.23%
Without Causality Loss	<b>99.92%</b>	0.029	2.95%
<i>HI-concept</i>	99.90%	<b>0.058</b>	10.54%

Table 5: Ablation on BERT for IMDB with faithfulness (Acc) and impact (CACE,  $\Delta$ Acc) evaluation.

# THANKS!

**Any questions?**

---



**Maryam Rezaee**

TGML Lab | Sharif University of Technology

*Under the supervision of*

Dr. Fatemeh SeyyedSalehi