

Beyond Concept Bottleneck Models: How to Make Black Boxes Intervenable?

Sonia Laguna*, Ri^{*}cards Marcinkevi^{*}cs*, Moritz Vandenhirtz, Julia E. Vogt

L193 Explainable AI, Qianqi Liu

Explanation Types





Concept Bottleneck Models (CBMs) Recap



- Key advantage: user can intervene on concepts
- Key shortcoming: requires concept labels for the entire training set



Motivations

Widespread Deployment of Black-Box Models

Post-hoc intervenability on *any* pretrained black box

• Desire for Instance-Specific Interventions

Formal definition of *intervenability* as a performance gain from concept corrections

Challenge of Full Concept Labeling

Small concept-labeled validation set suffices

Balancing Performance and Intervenability

Fine-tuning black boxes to *increase* intervenability



Method





Method





Method





Intervenability

CBMs

$$\mathbb{E}_{(\boldsymbol{x},\boldsymbol{c},\boldsymbol{y})\sim D}\left[\mathbb{E}_{\boldsymbol{c}'\sim\pi}\Big[\mathcal{L}^{\boldsymbol{y}}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}),\,\boldsymbol{y}\big) - \mathcal{L}^{\boldsymbol{y}}\big(\boldsymbol{g}_{\psi}(\boldsymbol{c}'),\,\boldsymbol{y}\big)\Big]\right]$$

Black-box models

$$\mathbb{E}_{(x,c,y)\sim D,\ c'\sim\pi}\Big[\mathcal{L}^{y}ig(f_{ heta}(x),yig) \ - \ \mathcal{L}^{y}ig(g_{\psi}(z'),yig)\Big],$$

where $z'\in argmin_{\widetilde{z}}\lambda\,\mathcal{L}^{c}ig(q_{\xi}(\widetilde{z}),\,c'ig) \ + \ dig(z,\,\widetilde{z}ig).$



Fine-tuning for Intervenability

Intervenability

$$\begin{split} \min_{\psi, z'} & \mathbb{E}_{(x, c, y) \sim D, c' \sim \pi} \Big[\mathcal{L}^{y} \big(g_{\psi}(z'), y \big) \Big], \\ \text{s.t.} \quad z' \in \operatorname{argmin}_{\tilde{z}} \lambda \, \mathcal{L}^{c} \big(q_{\xi}(\tilde{z}), \, c' \big) \, + \, d(z, \, \tilde{z}). \end{split}$$

Weight

$$\begin{split} \min_{\varphi,\psi,z'} \ \mathbb{E}_{(x,c,y)\sim D,\ c'\sim\pi} \Big[(1-\beta) \, \mathcal{L}^{y} \big(g_{\psi} \big(h_{\varphi}(x) \big), \, y \big) \, + \, \beta \, \mathcal{L}^{y} \big(g_{\psi}(z'), \, y \big) \Big], \\ \text{s.t.} \quad z' \in argmin_{\tilde{z}} \lambda \, \mathcal{L}^{c} \big(q_{\xi}(\tilde{z}), \, c' \big) \, + \, d\big(z, \, \tilde{z}\big). \end{split}$$



Datasets

- Synthetic
 - To test performance under controlled bottleneck vs. incomplete concept scenarios.
- AwA2: 85 binary attributes and species labels
 - Clearly labeled animal attributes. Demonstrates concept-based classification in images.
- CUB
 - Fine-grained bird attributes to test model detail sensitivity.
- CIFAR-10 & ImageNet: 143 and 100 attributes respectively
 - Standard vision benchmark + demonstrating concept extraction with large vision-language models.
- MIMIC-CXR & CheXpert: 14 binary attributes
 - Real-world medical imaging with partial concept labels from radiology reports.

Evaluation metrics: AUROC, AUPR, Brier scores, calibration curves.



Baselines

- **Black-box**: standard neural network f_{θ} trained only on (x, y)
- Concept Bottleneck Models (CBMs): predicts concepts first, then uses them as a "bottleneck" to predict y
- **Post-hoc CBM**: converts a pretrained black-box into a CBM-like architecture by training a probe for concepts and then re-training a final layer.
- Fine-tuned (I): The paper's proposed approach
- Fine-tuned (A): concatenates concept labels with hidden features and fine-tune
- Fine-tuned (MT): fine-tunes the black box with a multi-task loss for label and concept prediction



Results





Results



Figure 4: Intervention results on the (a) synthetic *incomplete*, (b) AwA2, (c) CIFAR-10, and (d) MIMIC-CXR datasets w.r.t. target AUROC (*top*) and AUPR (*bottom*) across ten seeds.



Results

Dataset	Model	Concepts			Target		
		AUROC	AUPR	Brier	AUROC	AUPR	Brier
Synthetic	BLACK BOX CBM Post hoc CBM Fine-tuned, A Fine-tuned, MT Fine-tuned, I	0.716±0.018 0.837±0.008 0.714±0.017 0.784±0.013 0.716±0.018	$\begin{array}{c} 0.710 \pm 0.017 \\ \textbf{0.835} \pm \textbf{0.008} \\ 0.707 \pm 0.018 \\ \hline \\ 0.780 \pm 0.014 \\ 0.710 \pm 0.017 \end{array}$	0.208±0.006 0.196±0.006 0.207±0.009 0.186±0.006 0.208±0.006	$\begin{array}{c} 0.686 {\pm} 0.043 \\ \textbf{0.713} {\pm} \textbf{0.040} \\ 0.707 {\pm} 0.049 \\ 0.682 {\pm} 0.047 \\ 0.687 {\pm} 0.046 \\ 0.695 {\pm} 0.051 \end{array}$	$\begin{array}{c} 0.675 {\pm} 0.046 \\ \textbf{0.700} {\pm} \textbf{0.038} \\ 0.698 {\pm} 0.048 \\ 0.668 {\pm} 0.046 \\ 0.668 {\pm} 0.043 \\ 0.685 {\pm} 0.051 \end{array}$	0.460±0.003 0.410±0.012 0.285±0.015 0.470±0.004 0.471±0.003 0.285±0.014
AwA2	BLACK BOX CBM Post hoc CBM Fine-tuned, A Fine-tuned, MT Fine-tuned, I	0.991±0.002 0.993±0.001 0.992±0.002 0.994±0.002 0.991±0.002	$\begin{array}{c} 0.979 \pm 0.006 \\ 0.979 \pm 0.002 \\ 0.976 \pm 0.005 \\ \hline \\ 0.985 \pm 0.004 \\ 0.979 \pm 0.005 \end{array}$	$\begin{array}{c} 0.027 \pm 0.006 \\ 0.025 \pm 0.001 \\ 0.025 \pm 0.005 \\ \hline \\ 0.022 \pm 0.005 \\ 0.027 \pm 0.006 \end{array}$	0.996±0.001 0.988±0.001 0.996±0.001 0.996±0.001 0.997±0.001 0.996±0.001	0.926±0.020 0.892±0.005 0.929±0.018 0.938±0.016 0.938±0.017 0.925±0.020	$\begin{array}{c} 0.199 {\pm} 0.038 \\ 0.234 {\pm} 0.009 \\ \textbf{0.170} {\pm} \textbf{0.033} \\ \textbf{0.170} {\pm} \textbf{0.036} \\ 0.178 {\pm} 0.038 \\ 0.195 {\pm} 0.040 \end{array}$
CIFAR-10	BLACK BOX CBM Post hoc CBM Fine-tuned, A Fine-tuned, MT <u>Fine-tuned, I</u>	0.713±0.002 	$\begin{array}{c} 0.802 \pm 0.001 \\$	0.110±0.000 — 0.125±0.004 — 0.109±0.000 0.110±0.000	0.879±0.001 0.888±0.001 0.876±0.002 0.870±0.004 0.873±0.003	0.504±0.004 	0.920 ± 0.006 0.624 \pm 0.003 0.896 \pm 0.005 0.890 \pm 0.014 0.902 \pm 0.021
MIMIC-CXR	BLACK BOX CBM Post hoc CBM Fine-tuned, A Fine-tuned, MT Fine-tuned, I	0.743±0.006 0.744±0.006 0.707±0.006 0.748±0.008 0.744±0.005	$\begin{array}{c} 0.170 \pm 0.004 \\ \textbf{0.224} \pm \textbf{0.003} \\ 0.154 \pm 0.006 \\ \hline \\ 0.187 \pm 0.003 \\ 0.172 \pm 0.005 \end{array}$	$\begin{array}{c} 0.046 \pm 0.001 \\ 0.053 \pm 0.001 \\ 0.046 \pm 0.001 \\ \hline \\ 0.045 \pm 0.001 \\ 0.046 \pm 0.001 \end{array}$	0.789±0.006 0.765±0.007 0.801±0.006 0.773±0.009 0.785±0.006 0.808±0.007	0.706±0.009 0.699±0.006 0.727±0.008 0.665±0.013 0.696±0.009 0.733±0.009	$\begin{array}{c} 0.444 {\pm} 0.003 \\ 0.427 {\pm} 0.003 \\ \textbf{0.301} {\pm} \textbf{0.005} \\ 0.459 {\pm} 0.004 \\ 0.450 {\pm} 0.008 \\ 0.314 {\pm} 0.015 \end{array}$



Pros and Cons

Pros

- •Can be applied post hoc to any pre-trained network.
- •Small validation set
- •Effective even with incomplete concepts
- •Architecture-agnostic
- Instance-specific interventions

Cons

- •Still needs some concept labels
- Added complexity
- •Highly depends on probe's quality
- •Potential misalignment of concepts
- •User burden



Thanks and Q&A?

