
Dividing and Conquering a BlackBox to a Mixture of Interpretable Models: Route, Interpret, Repeat

Anonymous Authors¹

Abstract

ML model design either starts with an interpretable model or a Blackbox and explains it post hoc. Blackbox models are flexible but difficult to explain, while interpretable models are inherently explainable. Yet, interpretable models require extensive ML knowledge and tend to be less flexible, potentially underperforming than their Blackbox equivalents. This paper aims to blur the distinction between a post hoc explanation of a Blackbox and constructing interpretable models. Beginning with a Blackbox, we iteratively *carve out* a mixture of interpretable models and a *residual network*. The interpretable models identify a subset of samples and explain them using First Order Logic (FOL), providing basic reasoning on concepts from the Blackbox. We route the remaining samples through a flexible residual. We repeat the method on the residual network until all the interpretable models explain the desired proportion of data. Our extensive experiments show that our *route, interpret, and repeat* approach (1) identifies a richer diverse set of instance-specific concepts with high concept completeness via interpretable models by specializing in various subsets of data without compromising in performance, (2) identifies the relatively “harder” samples to explain via residuals, (3) outperforms the interpretable by-design models by significant margins during test-time interventions, (4) can be used to fix the shortcut learned by the original Blackbox.

1. Introduction

Model explainability is essential in high-stakes applications of AI, *e.g.*, healthcare. While Blackbox models (*e.g.*, Deep Learning) offer flexibility and modular design, post hoc ex-

planation is prone to confirmation bias (Wan et al., 2022), lack of fidelity to the original model (Adebayo et al., 2018), and insufficient mechanistic explanation of the decision-making process (Rudin, 2019). Interpretable-by-design models overcome those issues but tend to be less flexible than Blackbox models and demand substantial expertise to design. Currently, sing a post hoc explanation or adopting an inherently interpretable model is a mutually exclusive decision to be made at the initial phase of AI model design. This paper blurs the line on that dichotomous model design.

The literature on post hoc explanations is extensive. This includes model attributions ((Simonyan et al., 2013; Selvaraju et al., 2017)), counterfactual approaches (Abid et al., 2021; Singla et al., 2019), and distillation methods (Alharbi et al., 2021; Cheng et al., 2020). Those methods either identify key input features that contribute the most to the network’s output (Shrikumar et al., 2016), generate input perturbation to flip the network’s output (Samek et al., 2016; Montavon et al., 2018), or estimate simpler functions to approximate the network output locally. Post hoc methods preserve the flexibility and performance of the Blackbox, but suffer from lack of fidelity and mechanistic explanation of the network output (Rudin, 2019). Without a mechanistic explanation, recourse to a model’s undesirable behavior is unclear. Interpretable models are alternative designs to the Blackbox without many such drawbacks. For example, modern interpretable methods highlight human understandable *concepts* that contribute to the downstream prediction.

Several families of interpretable models exist for a long time, such as the rule-based approach and generalized additive models (Hastie & Tibshirani, 1987; Letham et al., 2015; Breiman et al., 1984). They primarily focus on tabular data. Such models for high-dimensional data (*e.g.*, images) primarily rely on projecting to a lower dimensional human understandable *concept* or *symbolic* space (Koh et al., 2020) and predicting the output with an interpretable classifier. Despite their utility, the current State-Of-The-Art (SOTA) are limited in design; for example, they do not model the interaction between the concepts except few exceptions (Ciravegna et al., 2021; Barbiero et al., 2022), offering limited reasoning capabilities and robustness. Furthermore, if a portion of the samples does not fit the template design of the inter-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

55 pretable model, they do not offer any flexibility, potentially
 56 compromising performance.

57 **Our contributions** We propose an interpretable method,
 58 aiming to achieve the best of both worlds: not sacrificing
 59 Blackbox performance similar to post hoc explainability
 60 while still providing actionable interpretation. We hypothe-
 61 size that a Blackbox encodes several interpretable models,
 62 each applicable to different portion of data. Thus, a single in-
 63 terpretable model may be insufficient to explain all samples.
 64 We construct a hybrid neuro-symbolic model by progres-
 65 sively *carving out* a mixture of interpretable models and a
 66 *residual network* from the given Blackbox. We coin the term
 67 *expert* for each interpretable model, as they specialize over
 68 a subset of data. Together, all the interpretable models are
 69 termed a “Mixture of Interpretable Experts” (MoIE). Our
 70 design identifies a subset of samples and *routes* them through
 71 the interpretable models to explain the samples with FOL,
 72 providing basic reasoning on concepts from the Blackbox.
 73 The remaining samples are routed through a flexible residual
 74 network. On the residual network, we repeat the method
 75 until MoIE explains the desired proportion of data. We
 76 quantify the sufficiency of the identified concepts to explain
 77 the Blackbox’s prediction using the concept completeness
 78 score (Yeh et al., 2019). Using FOL for interpretable models
 79 offers recourse when undesirable behavior is detected in the
 80 model. We provide an example of fixing a shortcut learning
 81 by modifying the FOL. Potentially, FOL can be used in
 82 human-model interaction (not explored in this paper). Our
 83 method is the divide-and-conquer approach, where the in-
 84 stances covered by the residual network need progressively
 85 more complicated interpretable models. Such insight can
 86 be used to inspect the data and the model further. Finally,
 87 our model allows *unexplainable* category of data, which is
 88 currently not allowed in the interpretable models.

2. Method

92 **Notation:** Assume we have a dataset $\{\mathcal{X}, \mathcal{Y}, \mathcal{C}\}$, where
 93 \mathcal{X} , \mathcal{Y} , and \mathcal{C} are the input images, class labels, and human
 94 interpretable attributes, respectively. $f^0 : \mathcal{X} \rightarrow \mathcal{Y}$, is our
 95 pre-trained initial Blackbox model. We assume that f^0 is
 96 a composition $h^0 \circ \Phi$, where $\Phi : \mathcal{X} \rightarrow \mathbb{R}^l$ is the image
 97 embeddings and $h^0 : \mathbb{R}^l \rightarrow \mathcal{Y}$ is a transformation from the
 98 embeddings, Φ , to the class labels. We denote the learnable
 99 function $t : \mathbb{R}^l \rightarrow \mathcal{C}$, projecting the image embeddings to
 100 the concept space. The concept space is the space spanned
 101 by the attributes \mathcal{C} . Thus, function t outputs a scalar value
 102 representing a concept for each input image.

104 **Method Overview:** Figure 1 summarizes our approach.
 105 We iteratively carve out an interpretable model from the
 106 given Blackbox. Each iteration yields an interpretable
 107 model (the downward grey paths in Figure 1) and a residual
 108 (the straightforward black paths in Figure 1). We start

109 with the initial Blackbox f^0 . At iteration k , we distill the
 110 Blackbox from the previous iteration f^{k-1} into a neuro-
 111 symbolic interpretable model, $g^k : \mathcal{C} \rightarrow \mathcal{Y}$. Our g is
 112 flexible enough to be any interpretable models (Yukse-
 113 gonul et al., 2022; Koh et al., 2020; Barbiero et al., 2022).
 114 The *residual* $r^k = f^{k-1} - g^k$ emphasizes the portion of
 115 f^{k-1} that g^k cannot explain. We then approximate r^k with
 116 $f^k = h^k \circ \Phi$. f^k will be the Blackbox for the subsequent iteration
 117 and be explained by the respective interpretable model.
 118 A learnable gating mechanism, denoted by $\pi^k : \mathcal{C} \rightarrow \{0, 1\}$
 119 (shown as the *selector* in Figure 1) routes an input sample
 120 towards either g^k or r^k . The thickness of the lines in Figure
 121 1 represents the samples covered by the interpretable
 122 models (grey line) and the residuals (black line). With every
 123 iteration, the cumulative coverage of the interpretable
 124 models increases, but the residual decreases. We name our
 125 method *route*, *interpret* and *repeat*.

2.1. Neuro-Symbolic Knowledge Distillation

Knowledge distillation in our method involves 3 parts: (1) a series of trainable selectors, *routing* each sample through the interpretable models and the residual networks, (2) a sequence of learnable neuro-symbolic interpretable models, each providing FOL explanations to *interpret* the Blackbox, and (3) *repeating* with Residuals for the samples that cannot be explained with their interpretable counterparts. We detail each component below.

2.1.1. THE SELECTOR FUNCTION

As the first step of our method, the selector π^k *routes* the j^{th} sample through the interpretable model g^k or residual r^k with probability $\pi^k(c_j)$ and $1 - \pi^k(c_j)$ respectively, where $k \in [0, K]$, with K being the number of iterations. We define the empirical coverage of the k^{th} iteration as $\zeta(\pi^k) = \frac{1}{m} \sum_{j=1}^m \pi^k(c_j)$, the empirical mean of the samples selected by the selector for the associated interpretable model g^k , with m being the total number of samples in the training set. Thus, the entire selective risk is:

$$\mathcal{R}^k(\pi^k, g^k) = \frac{\frac{1}{m} \sum_{j=1}^m \mathcal{L}_{(g^k, \pi^k)}^k(x_j, c_j)}{\zeta(\pi^k)}, \quad (1)$$

where $\mathcal{L}_{(g^k, \pi^k)}^k$ is the optimization loss used to learn g^k and π^k together, discussed in Section 2.1.2. For a given coverage of $\tau^k \in (0, 1]$, we solve the following optimization problem:

$$\begin{aligned} \theta_{s^k}^*, \theta_{g^k}^* &= \arg \min_{\theta_{s^k}, \theta_{g^k}} \mathcal{R}^k(\pi^k(\cdot; \theta_{s^k}), g^k(\cdot; \theta_{g^k})) \\ \text{s.t. } \zeta(\pi^k(\cdot; \theta_{s^k})) &\geq \tau^k, \end{aligned} \quad (2)$$

where $\theta_{s^k}^*$, $\theta_{g^k}^*$ are the optimal parameters at iteration k for the selector π^k and the interpretable model g^k respectively. In this work, π^k s of different iterations are neural networks

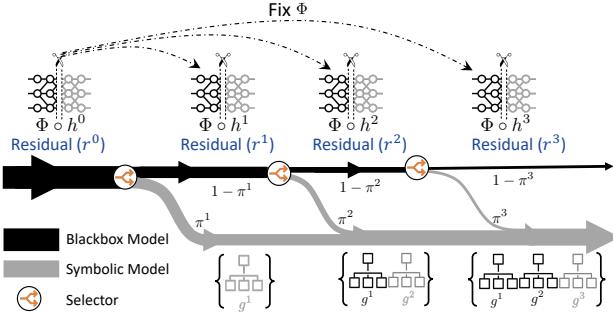


Figure 1. Schematic view of *route*, *interpret* and *repeat*. At iteration k , the selector *routes* each sample either towards the interpretable model g^k (to *interpret*) with probability $\pi^k(\cdot)$ or the residual $r^k = f^{k-1} - g^k$ with probability $1 - \pi^k(\cdot)$ (to *repeat* in the further iterations). f^{k-1} is the Blackbox of the $(k-1)^{th}$ iteration. g^k generates FOL-based explanations for the samples it covers. Otherwise, the selector routes through the next step until it either goes through a subsequent interpretable model or reaches the last residual. Components in black and grey indicate the fixed and trainable modules in our model respectively.

with sigmoid activation. At inference time, the selector routes the j^{th} sample with concept vector c_j to g^k if and only if $\pi^k(c_j) \geq 0.5$ for $k \in [0, K]$.

2.1.2. NEURO-SYMBOLIC INTERPRETABLE MODELS

In this stage, we design interpretable model g^k of k^{th} iteration to *interpret* the Blackbox f^{k-1} from the previous $(k-1)^{th}$ iteration by optimizing the following loss function:

$$\mathcal{L}_{(g^k, \pi^k)}^k(x_j, c_j) = \ell\left(f^{k-1}(x_j), g^k(c_j)\right) \underbrace{\pi^k(c_j)}_{\text{trainable component for current iteration } k} \underbrace{\prod_{i=1}^{k-1} (1 - \pi^i(c_j))}_{\text{fixed component trained in the previous iterations}}, \quad (3)$$

where the term $\pi^k(c_j) \prod_{i=1}^{k-1} (1 - \pi^i(c_j))$ denotes the probability of j^{th} sample being routed through the interpretable model g^k . It is the probability of the sample going through the residuals for all the previous iterations from 1 through $k-1$ (*i.e.*, $\prod_{i=1}^{k-1} (1 - \pi^i(c_j))$) times the probability of going through the interpretable model at iteration k (*i.e.*, $\pi^k(c_j)$). Refer to Figure 1 for an illustration. We learn π^1, \dots, π^{k-1} in the prior iterations and are not trainable at iteration k . As each interpretable model g^k specializes in explaining a specific subset of samples (denoted by coverage τ), we refer to it as an *expert*. We use SelectiveNet’s (Geifman & El-Yaniv, 2019) optimization method to optimize Equation (2) since selectors need a rejection mechanism to route samples through residuals. Appendix A.4 details the optimization procedure in Equation (3). We refer to the interpretable experts of all the iterations as a “Mixture of Interpretable Experts” (MoIE) cumulatively after training. Furthermore, we utilize Entropy-based linear layer neural

Table 1. Datasets and Blackboxes.

DATASET	BLACKBOX	# EXPERTS
CUB-200 (Wah et al., 2011)	RESNET101 (He et al., 2016)	6
CUB-200 (Wah et al., 2011)	VIT (Wang et al., 2021)	6
AWA2 (Xian et al., 2018)	RESNET101 (He et al., 2016)	4
AWA2 (Xian et al., 2018)	VIT (Wang et al., 2021)	6
HAM1000 (Tschandl et al., 2018)	INCEPTION (Szegedy et al., 2015)	6
SIIM-ISIC (Rotemberg et al., 2021)	INCEPTION (Szegedy et al., 2015)	6
EFFUSION IN MIMIC-CXR (Johnson et al.)	DENSENET121 (Huang et al., 2017)	3

network (ELL) (Barbiero et al., 2022) as the interpretable symbolic model g to construct First Order Logic (FOL) explanations of a given prediction.

2.1.3. THE RESIDUALS

The last step is to *repeat* with the residual r^k , as $r^k(x_j, c_j) = f^{k-1}(x_j) - g^k(c_j)$. We train $f^k = h^k(\Phi(\cdot))$ to approximate the residual r^k , creating a new Blackbox f^k for the next iteration ($k+1$). This step is necessary to specialize f^k over samples not covered by g^k . Optimizing the following loss function yields f^k for the k^{th} iteration:

$$\mathcal{L}_f^k(x_j, c_j) = \ell(r^k(x_j, c_j), f^k(x_j)) \underbrace{\prod_{i=1}^k (1 - \pi^i(c_j))}_{\text{trainable component for iteration } k} \underbrace{\prod_{i=k+1}^k (1 - \pi^i(c_j))}_{\text{non-trainable component for iteration } k} \quad (4)$$

Notice that we fix the embedding $\Phi(\cdot)$ for all the iterations. Due to computational overhead, we only finetune the last few layers of the Blackbox (h^k) to train f^k . At the final iteration K , our method produces a MoIE and a Residual, explaining the interpretable and uninterpretable components of the initial Blackbox f^0 , respectively. Appendix A.5 describes the training procedure of our model, the extraction of FOL and the architecture of our model at inference.

Selecting number of iterations K : We follow two principles to select the number of iterations K as a stopping criterion: 1) Each expert should have enough data to be trained reliably (coverage ζ^k). If an expert covers insufficient samples, we stop the process. 2) If the final residual (r^K) underperforms a threshold, it is not reliable to distill from the Blackbox. We stop the procedure to ensure the overall accuracy is maintained. Appendix A.11.1 shows the comparison of the computational time of MoIE and the Blackbox.

3. Related work

Post hoc explanations: Post hoc explanations retain the flexibility and performance of the Blackbox. The post hoc explanation has many categories, including feature attribution (Simonyan et al., 2013; Smilkov et al., 2017; Binder et al., 2016) and counterfactual approaches (Singla et al.,

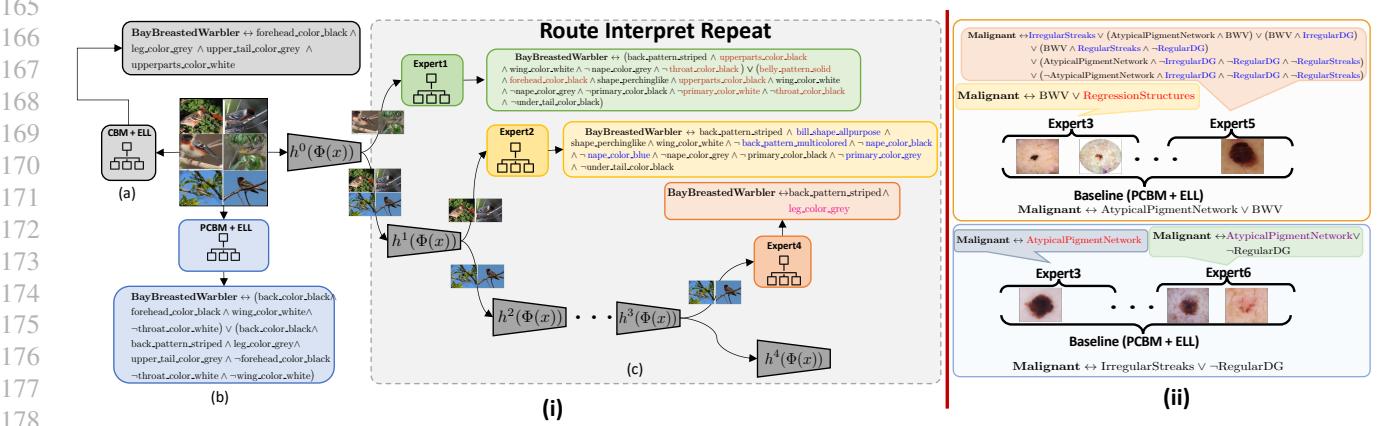


Figure 2. MoIE identifies diverse concepts for specific subsets of a class, unlike the generic ones by the baselines. (i) We construct the FOL explanations of the samples of, ‘‘Bay breasted warbler’’ in the CUB-200 dataset for VIT-based (a) CBM + ELL as an interpretable-by-design baseline, (b) PCBM + ELL as a posthoc baseline, (c) experts in MoIE at inference. We highlight the unique concepts for experts 1,2, and 3 in red, blue, and magenta, respectively. (ii) Comparison of FOL explanations by MoIE with the PCBM + ELL baselines for HAM10000 (top) and ISIC (down) to classify Malignant lesion. We highlight unique concepts for experts 3, 5, and 6 in red, blue, and violet, respectively. For brevity, we combine FOLs for each expert for the samples covered by them.

2019; Abid et al., 2021). For example, feature attribution methods associate a measure of importance to features (e.g., pixels) that is proportional to the feature’s contribution to BlackBox’s predicted output. Many methods were proposed to estimate the importance measure, including gradient-based methods (Selvaraju et al., 2017; Sundararajan et al., 2017), game-theoretic approach (Lundberg & Lee, 2017), and many more. The post hoc approaches suffer from lack of fidelity to input (Adebayo et al., 2018) and ambiguity in explanation due to a lack of correspondence to human-understandable concepts. Recently, Posthoc Concept Bottleneck models (PCBMs) (Yuksekgonul et al., 2022) learn the concepts from a trained Blackbox embedding and use an interpretable classifier for classification. Also, they fit a residual in their hybrid variant (PCBM-h) to mimic the performance of the Blackbox. We will compare against the performance of the PCBMs method. Another major shortcoming is that, due to a lack of mechanistic explanation, post hoc explanations do not provide a recourse when an undesirable property of a Blackbox is identified. Interpretable-by-design provides a remedy to those issues (Rudin, 2019).

Concept-based interpretable models: Our approach falls into the category of concept based interpretable models. Such methods provide a mechanistically interpretable prediction that is a function of human-understandable concepts. The concepts are usually extracted from the activation of the middle layers of the Neural Network (bottleneck). Examples include Concept Bottleneck models (CBMs) (Koh et al., 2020), antehoc concept decoder (Sarkar et al., 2021), and a high-dimensional Concept Embedding model (CEMs) (Zarlenga et al., 2022) that uses high di-

mensional concept embeddings to allow extra supervised learning capacity and achieves SOTA performance in the interpretable-by-design class. Most concept-based interpretable models do not model interaction between concepts and cannot be used for reasoning. An exception is ELL (Barbiero et al., 2022) which uses an entropy-based approach to derive explanations in terms of FOL using the concepts. The underlying assumption of those methods is that one interpretable function can explain the entire set of data, which can limit flexibility and consequently hurt the performance of the models. Our approach relaxes that assumption by allowing multiple interpretable functions and a residual. Each function is appropriate for a portion of the data, and a small portion of the data is allowed to be uninterpretable by the model (*i.e.*, residual). We will compare our method with CBMs, CEMs and their ELL-enhanced variants.

Application in fixing the shortcut learning: Shortcuts are spurious features that correlate with both input and the label on the training dataset but fail to generalize in more challenging real-world scenarios. Explainable AI (X-AI) aims to identify and fix such an undesirable property. Related work in X-AI includes LIME (Ribeiro et al., 2016), utilized to detect spurious background as a shortcut to classify an animal. Recently interpretable model (Rosenzweig et al., 2021), involving local image patches, are used as a proxy to the Blackbox to identify shortcuts. However, both methods operate in pixel space, not concept space. Also, both approaches are post hoc and do not provide a way to eliminate the shortcut learning problem. Our MoIE discovers shortcuts using the high-level concepts in the FOL explanation of the Blackbox’s prediction and eliminates them via metadata

220 normalization (MDN) (Lu et al., 2021).
 221

222 4. Experiments

223 We perform experiments on a variety of vision and medical
 224 imaging datasets to show that 1) MoIE captures a diverse set
 225 of concepts, 2) the performance of the residuals degrades
 226 over successive iterations as they cover “harder” instances,
 227 3) MoIE does not compromise the performance of the Black-
 228 box, 4) MoIE achieves superior performances during test
 229 time interventions, and 5) MoIE can fix the shortcuts using
 230 the Waterbirds dataset (Sagawa et al., 2019). We repeat
 231 our method until MoIE covers at least 90% of samples or
 232 the final residual’s accuracy falls below 70%. Refer to Ta-
 233 ble 1 for the datasets and Blackboxes experimented with.
 234 For ResNets and Inception, we flatten the feature maps
 235 from the last convolutional block to extract the concepts.
 236 For VITs, we use the image embeddings from the trans-
 237 former encoder to perform the same. We use SIIM-ISIC
 238 as a real-world transfer learning setting, with the Blackbox
 239 trained on HAM10000 and evaluated on a subset of the
 240 SIIM-ISIC Melanoma Classification dataset (Yuksekgonul
 241 et al., 2022). Appendix A.6 and Appendix A.7 expand on
 242 the datasets and hyperparameters.
 243

244 Baselines: We compare our methods to two concept based
 245 baselines – 1) interpretable-by-design and 2) posthoc. They
 246 consist of two parts: a) a concept predictor $\Phi : \mathcal{X} \rightarrow \mathcal{C}$,
 247 predicting concepts from images; and b) a label predictor
 248 $g : \mathcal{C} \rightarrow \mathcal{Y}$, predicting labels from the concepts. . The
 249 end-to-end CEMs, and sequential CBMs serve as interpretable-
 250 by-design baselines. Similarly, PCBMs and PCBMs-h serve
 251 as posthoc baselines. Convolution-based Φ includes all
 252 layers till the last convolution block. VIT-based Φ consists
 253 of the transformer encoder block. The standard CBM and
 254 PCBMs do not show how the concepts are composed
 255 to make the label prediction. So, we create CBM + ELL,
 256 PCBMs + ELL and PCBMs-h + ELL by using the identical
 257 g of MOIE (shown in Appendix A.7), as a replacement
 258 for the standard classifiers of CBM and PCBMs. We train
 259 the Φ and g in these new baselines to sequentially generate
 260 FOLs (Barbiero et al., 2022). Due to the unavailability
 261 of concept annotations, we extract the concepts from the
 262 Derm7pt dataset (Kawahara et al., 2018) using the pretrained
 263 embeddings of the Blackbox (Yuksekgonul et al., 2022) for
 264 HAM10000. Thus, we do not have interpretable-by-design
 265 baselines for HAM10000 and ISIC.
 266

267 4.1. Results

268 4.1.1. EXPERT DRIVEN EXPLANATIONS BY MOIE

269 First, we show that MoIE captures a rich set of diverse
 270 instance-specific concepts qualitatively. Next, we show
 271 quantitatively that MoIE-identified concepts are faithful to
 272

273 Blackbox’s final prediction using the metric “completeness
 274 score” and zeroing out relevant concepts.

Heterogeneity of Explanations: At each iteration of MoIE,
 the blackbox ($h^k(\Phi(\cdot))$) splits into an interpretable expert
 (g^k) and a residual (r^k). Figure 2i shows this mechanism for
VIT-based MoIE and compares the FOLs with CBM + ELL
and PCBMs + ELL baselines to classify “Bay Breasted Warbler” of CUB-200. The experts of different iterations specialize in specific instances of “Bay Breasted Warbler”. Thus, each expert’s FOL comprises its instance-specific concepts of the same class (Figure 2i-c). For example, the concept, *leg_color_grey* is unique to expert4, but *belly_pattern_solid* and *back_pattern_multicolored* are unique to experts 1 and 2, respectively, to classify the instances of “Bay Breasted Warbler”. Unlike MoIE, the baselines employ a single interpretable model g , resulting in a generic FOL with identical concepts for all the samples of “Bay Breasted Warbler” (Figure 2(a-b)). Thus the baselines fail to capture the heterogeneity of explanations. For additional results of CUB-200, refer to Appendix A.11.2.

Figure 2ii show such diverse explanations for HAM10000 (top) and ISIC (bottom). In Figure 2ii-(top), the baseline-FOL consists of concepts such as *AtypicalPigmentNetwork* and *BlueWhitishVeil (BWV)* to classify “Malignancy” for all the instances for HAM10000. However, expert 3 relies on *RegressionStructures* along with *BWV* to classify the same for the samples it covers while expert 5 utilizes several other concepts e.g., *IrregularStreaks*, *Irregular dots and globules (IrregularDG)* etc. Due to space constraints, Appendix A.11.3 reports similar results for the Awa2 dataset. Also, VIT-based experts compose less concepts per a sample than the ResNet-based experts, shown in Appendix A.11.4.

MoiE-identified concepts attain higher completeness scores. Figure 5(a-b) shows the completeness scores (Yeh et al., 2019) for varying number of concepts. Completeness score is a post hoc measure, signifying the identified concepts as “sufficient statistic” of the predictive capability of the Blackbox. Recall that g utilizes ELL (Barbiero et al., 2022), associating each concept with an attention weight after training. A concept with high attention weight implies its high predictive significance. Iteratively, we select the top relevant concepts based on their attention weights and compute the completeness scores for the top concepts for MoIE and the PCBMs + ELL baseline in Figure 5(a-b) (Appendix A.8 for details). For example, MoIE achieves a completeness score of 0.9 compared to 0.75 of the baseline($\sim 20\% \uparrow$) for the 10 most significant concepts for CUB-200 dataset with VIT as Blackbox.

MoiE identifies more meaningful instance-specific concepts. Figure 5(c-d) reports the drop in accuracy by zeroing out the significant concepts. Any interpretable model (g) supports concept-intervention (Koh et al., 2020). After iden-

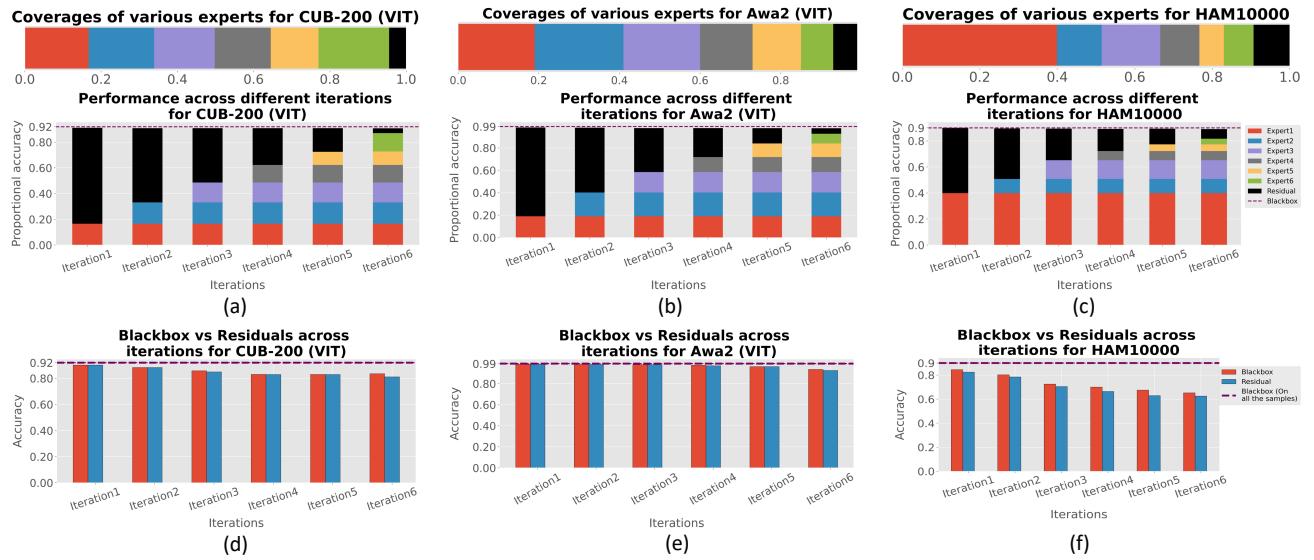


Figure 3. The performance of experts and residuals across iterations. (a-c) Coverage and proportional accuracy of the experts and residuals. (d-f) We route the samples covered by the residuals across iterations to the initial Blackbox f^0 and compare the accuracy of f^0 (red bar) with the residual (blue bar). Figures d-f show the progressive decline in performance of the residuals across iterations as they cover the samples in the increasing order of ‘hardness’. We observe the similar abysmal performance of the initial blackbox f^0 for these samples.

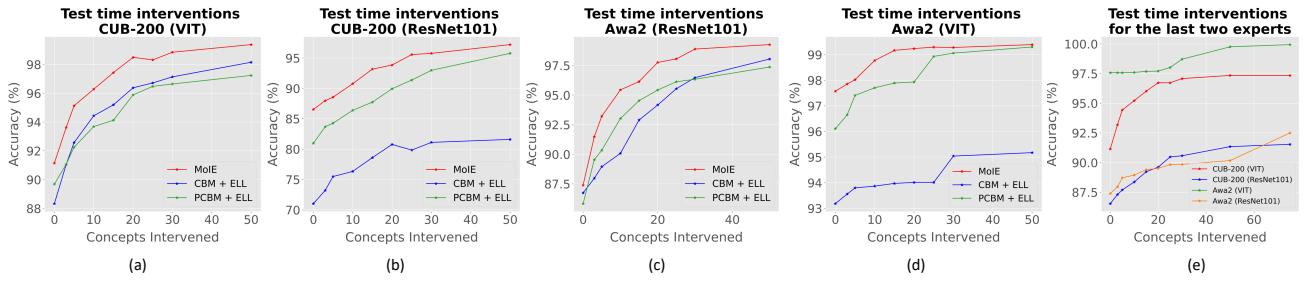


Figure 4. Across architectures test time interventions of concepts on all the samples (a-d), on the ‘hard’ samples (e), covered by only the last two experts of MoIE.

tifying the top concepts from g using the attention weights, as in the last section, we set these concepts’ values to zero, compute the model’s accuracy drop, and plot in Figure 5(b). When zeroing out the top 10 essential concepts for VIT based CUB-200 models, MoIE records a drop of 53% compared to 28% and 42% for the CBM + ELL and PCBM + ELL baselines, respectively, showing the faithfulness of the identified concepts to the prediction.

In both of the last experiments, MoIE outperforms the baselines as the baselines mark the same concepts as significant for all samples of each class. However, MoIE leverages various experts specializing in different subsets of samples of different classes. For results of Awa2, refer to Appendix A.10.2.

4.1.2. IDENTIFICATION OF HARDER SAMPLES BY SUCCESSIVE RESIDUALS

Figure 3 (a-c) display the proportional accuracy of the experts and the residuals of our method per iteration. The proportional accuracy of each model (experts and/or residuals) is defined as the accuracy of that model times its coverage. Recall that the model’s coverage is the empirical mean of the samples selected by the selector. Figure 3a shows that the experts and residual cumulatively achieve an accuracy ~ 0.92 for the CUB-200 dataset in iteration 1, with more contribution from the residual (black bar) than the expert1 (blue bar). Later iterations cumulatively increase and worsen the performance of the experts and corresponding residuals, respectively. The final iteration carves out the entire interpretable portion from the Blackbox f^0 via all the

Table 2. MoIE does not hurt the performance of the original Blackbox using a held-out test set. We provide the mean and standard errors of AUROC and accuracy for medical imaging (*e.g.*, HAM10000, ISIC, and Effusion) and vision (*e.g.*, CUB-200 and Awa2) datasets respectively over 5 random seeds. For MoIE, we also report the percentage of test set samples covered by all experts as “coverage”. Here, MoIE + Residual represents the experts with the final residual. Following the setting (Zarlunga et al., 2022), we only report the performance of the convolutional CEM, leaving the construction of VIT-based CEM as a future work. Recall that interpretable-by-design models can not be constructed for HAM10000 and ISIC as they have no concept annotation; we learn the concepts from Derm7pt dataset. For all the datasets, MoIE covers a significant portion of data (at least 90%) cumulatively. We boldfaced our results and the Blackbox.

MODEL	CUB-200 (RESNET101)	CUB-200 (VIT)	AWA2 (RESNET101)	DATASET	HAM10000	SIIM-ISIC	EFFUSION
BLACKBOX	0.88	0.92	0.89	0.99	0.96	0.85	0.91
INTERPRETABLE-BY-DESIGN							
CEM (Zarlunga et al., 2022)	0.77 ± 0.22	-	0.88 ± 0.50	-	NA	NA	0.76 ± 0.00
CBM (Sequential) (Koh et al., 2020)	0.65 ± 0.37	0.86 ± 0.24	0.88 ± 0.35	0.94 ± 0.28	NA	NA	0.79 ± 0.00
CBM + ELL (Koh et al., 2020; Barbiero et al., 2022)	0.71 ± 0.35	0.88 ± 0.24	0.86 ± 0.35	0.93 ± 0.25	NA	NA	0.79 ± 0.00
POSTHOC							
PCBM (Yüsekgonul et al., 2022)	0.76 ± 0.01	0.85 ± 0.20	0.82 ± 0.23	0.94 ± 0.17	0.93 ± 0.00	0.71 ± 0.01	0.81 ± 0.01
PCBM-h (Yüsekgonul et al., 2022)	0.85 ± 0.01	0.91 ± 0.18	0.87 ± 0.20	0.98 ± 0.17	0.95 ± 0.00	0.79 ± 0.05	0.87 ± 0.07
PCBM + ELL (Yüsekgonul et al., 2022; Barbiero et al., 2022)	0.80 ± 0.36	0.89 ± 0.26	0.85 ± 0.25	0.96 ± 0.18	0.94 ± 0.02	0.73 ± 0.01	0.81 ± 0.01
PCBM-h + ELL (Yüsekgonul et al., 2022; Barbiero et al., 2022)	0.85 ± 0.30	0.91 ± 0.28	0.88 ± 0.24	0.98 ± 0.20	0.95 ± 0.03	0.82 ± 0.05	0.87 ± 0.03
OURS							
MoIE (COVERAGE)	$0.86 \pm 0.01 (0.9)$	$0.91 \pm 0.00 (0.95)$	$0.87 \pm 0.02 (0.91)$	$0.97 \pm 0.00 (0.94)$	$0.95 \pm 0.00 (0.9)$	$0.84 \pm 0.00 (0.94)$	$0.87 \pm 0.00 (0.98)$
MoIE + RESIDUAL	0.84 ± 0.01	0.90 ± 0.01	0.86 ± 0.020	0.94 ± 0.004	0.92 ± 0.00	0.82 ± 0.01	0.86 ± 0.00

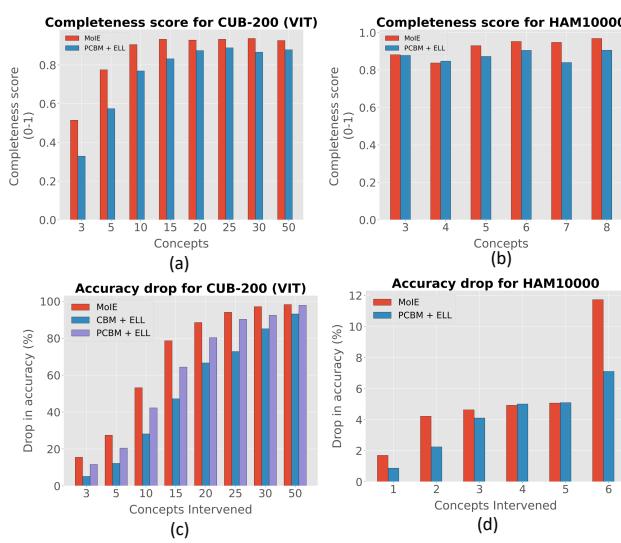


Figure 5. Quantitative validation of the extracted concepts and test time interventions. (a-b) Completeness scores of the models for varying number of top concepts. (c-d) Drop in accuracy compared to the original model after zeroing-out the top significant concepts iteratively. The highest drop for MoIE indicates that MoIE selects more instance specific concepts than generic ones by the baselines.

experts, resulting in their more significant contribution to the cumulative performance. The residual of the last iteration covers the “hardest” samples, achieving low accuracy. Tracing these samples back to the original Blackbox f^0 , it also classifies these samples poorly (Figure 3(d-f)). As shown in the coverage plot, this experiment reinforces Figure 1, where the flow through the experts gradually becomes thicker compared to the narrower flow of the residual with every iteration. Refer to Figure 9 in the Appendix A.10.1 for the results of the ResNet-based MoIEs.

4.1.3. QUANTITATIVE ANALYSIS OF MOIE WITH THE BLACKBOX AND BASELINE

Comparing with the interpretable-by-design baselines: Table 2 shows that MoIE achieves comparable performance to the Blackbox. Recall that “MoIE” refers to the mixture of all interpretable experts (g) only excluding any residuals. MoIE outperforms the interpretable-by-design baselines for all the datasets except Awa2. Since Awa2 is designed for zero-shot learning, its rich concept annotation makes it appropriate for interpretable-by-design models. In general, VIT-derived MoIEs perform better than their ResNet-based variants.

Comparing with the PCBMs: Table 2 shows that interpretable MoIE outperforms the interpretable posthoc baselines – PCBM and PCBM + ELL for all the datasets, especially by a significant margin for CUB-200 and ISIC. We also report “MoIE + Residual” as the mixture of interpretable experts plus the final residual to compare with the residualized PCBM, *i.e.*, PCBM-h. Table 2 shows that PCBM-h performs slightly better than MoIE + Residual. Note that PCBM-h learns the residual by fitting the complete dataset to fix the interpretable PCBM’s mistakes to replicate the performance of the Blackbox, resulting in better performance for PCBM-h than PCBM. However, we assume the Blackbox to be a combination of interpretable and uninterpretable components. So, we train the experts and the final residual to cover the interpretable and uninterpretable portions of the Blackbox respectively. Each iteration, our method learns the residuals to focus on the samples, which are not covered by the respective interpretable experts. Therefore, residuals are not designed to fix the mistakes made by the experts. In doing so, the final residual in MoIE + Residual covers the “hardest” examples, lowering

385 its overall performance compared to MoIE.
 386
 387

4.1.4. TEST TIME INTERVENTIONS

388 Figure 4(a-d) shows effect of test time interventions. Any
 389 concept-based models (Koh et al., 2020; Zarlenaga et al.,
 390 2022) allow test time interventions for datasets with concept
 391 annotation (e.g., CUB-200, Awa2). We identify the signifi-
 392 cant concepts via their attention scores in g , as during the
 393 computaion of completeness scores and set their values with
 394 the ground truths, considering the ground truth concepts as
 395 an oracle. As MoIE identifies a more diverse set of concepts
 396 by focusing on different subsets of classes, MoIE outper-
 397 forms the baselines in terms of accuracy for such test time
 398 interventions. Instead of manually deciding the samples
 399 to intervene, it is generally preferred to intervene on the
 400 “harder” samples, making the process efficient. As per Sec-
 401 tion 4.1.2, experts of different iterations cover samples with
 402 increasing order of “hardness”. To intervene efficiently, we
 403 perform identical test-time interventions with varying num-
 404 bers of concepts for the “harder” samples covered by the
 405 final two experts and plot the accuracy in Figure 4(e). For
 406 the VIT-derived MoIE of CUB-200, intervening only on 20
 407 concepts enhances the accuracy of MoIE from 91% to 96%
 408 ($\sim 6.1\% \uparrow$). We cannot perform the same for the baselines
 409 as they cannot directly estimate “harder” samples. Also, Fig-
 410 ure 4 shows relatively higher gain for ResNet-based models
 411 in general.
 412

4.1.5. APPLICATION IN THE REMOVAL OF SHORTCUTS

413 First, we create the Waterbirds dataset as in (Sagawa et al.,
 414 2019)by using forest and bamboo as the spurious land con-
 415 cepts of the Places dataset for landbirds of the CUB-200
 416 dataset. We do the same by using oceans and lakes as the
 417 spurious water concepts for waterbirds. We utilize ResNet50
 418 as the Blackbox f^0 to identify each bird as a Waterbird or
 419 a Landbird. The Blackbox quickly latches on the spurious
 420 backgrounds to classify the birds. As a result, the black
 421 box’s accuracy differs for land-based versus aquatic subsets
 422 of the bird species, as shown in Figure 6a. The Waterbird on
 423 the water is more accurate than on land ((96% vs. 67 % in
 424 the Figure 6a). The FOL from the biased Blackbox-derived MoIE captures the spurious concept *forest* a waterbird,
 425 misclassified as a landbird. Assuming the background
 426 concepts as metadata, we minimize the back-
 427 ground bias from the representation of the Blackbox using
 428 Metadata normalization (MDN) layers (Lu et al., 2021)
 429 between two successive layers of the convolutional backbone
 430 to fine-tune the biased Blackbox. Next, we train t , using
 431 the embedding Φ of the robust Blackbox, and compare the
 432 accuracy of the spurious concepts with the biased blackbox
 433 in Figure 6d. The validation accuracy of all the spurious
 434 concepts retrieved from the robust Blackbox falls well short
 435 of the predefined threshold 70% compared to the biased
 436
 437
 438
 439

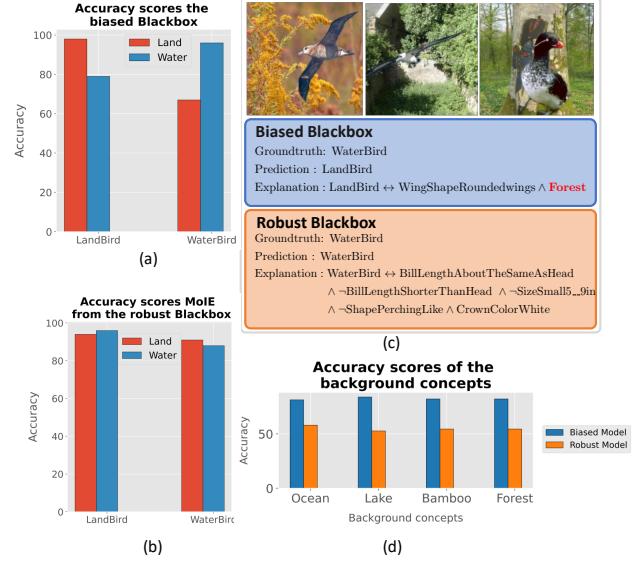


Figure 6. MoIE can fix shortcuts. (a) Performance of the biased Blackbox. (b) Performance of final MoIE extracted from the robust Blackbox after removing the shortcuts using MDN. (c) Examples of samples (top-row) and their explanations by the biased (middle-row) and robust Blackboxes (bottom-row). (d) Comparison of accuracies of the spurious concepts extracted from the biased vs the robust Blackbox .

Blackbox. Finally, we re-train the MoIE distilling from the new robust Blackbox. Figure 6b illustrates similar accuracies of MoIE for Waterbirds on water vs. Waterbirds on land (91% -88 %). The FOL from the robust Blackbox does not include any background concepts (6c, bottom row). Refer to 8 in Appendix A.9 for the flow diagram of this experiment.

5. Discussion & Conclusions

This paper proposes a novel method to iteratively extract a mixture of interpretable models from a flexible Blackbox. The comprehensive experiments on various datasets demonstrate that our method 1) captures more meaningful instance-specific concepts with high completeness score than baselines without losing the performance of the Blackbox, 2) does not require explicit concept annotation, 3) identifies the “harder” samples using the residuals, 4) achieves significant performance gain than the baselines during test time interventions, 5) eliminate shortcuts effectively. In the future, we aim to apply our method to other modalities, such as text or video. Also, as in the prior work, MoIE-captured concepts may not reflect a causal effect. The assessment of causal concept effects necessitates estimating inter-concept interactions, which will be the subject of future research.

References

- 440
441
442
443
444
- 445
446
447
448
449
- 450
451
452
453
454
- 455
456
457
458
459
- 460
461
462
463
464
- 465
466
467
468
469
- 470
471
- 472
473
474
475
476
- 477
478
479
- 480
481
482
483
484
- 485
486
487
488
- 489
490
491
492
493
494
- 495
496
497
498
499
500
- 501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
988
989
989
990
991
992
993
994
995
996
997
998
999
1000
- Abid, A., Yuksekgonul, M., and Zou, J. Meaningfully explaining model mistakes using conceptual counterfactuals. *arXiv preprint arXiv:2106.12723*, 2021.
- Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., and Kim, B. Sanity checks for saliency maps. *Advances in neural information processing systems*, 31, 2018.
- Alharbi, R., Vu, M. N., and Thai, M. T. Learning interpretation with explainable knowledge distillation. In *2021 IEEE International Conference on Big Data (Big Data)*, pp. 705–714. IEEE, 2021.
- Barbiero, P., Ciravegna, G., Giannini, F., Lió, P., Gori, M., and Melacci, S. Entropy-based logic explanations of neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6046–6054, 2022.
- Belle, V. Symbolic logic meets machine learning: A brief survey in infinite domains. In *International Conference on Scalable Uncertainty Management*, pp. 3–16. Springer, 2020.
- Besold, T. R., Garcez, A. d., Bader, S., Bowman, H., Domingos, P., Hitzler, P., Kühnberger, K.-U., Lamb, L. C., Lowd, D., Lima, P. M. V., et al. Neural-symbolic learning and reasoning: A survey and interpretation. *arXiv preprint arXiv:1711.03902*, 2017.
- Binder, A., Montavon, G., Lapuschkin, S., Müller, K.-R., and Samek, W. Layer-wise relevance propagation for neural networks with local renormalization layers. In *International Conference on Artificial Neural Networks*, pp. 63–71. Springer, 2016.
- Breiman, L., Friedman, J., Stone, C., and Olshen, R. Classification and regression trees (crc, boca raton, fl). 1984.
- Cheng, X., Rao, Z., Chen, Y., and Zhang, Q. Explaining knowledge distillation by quantifying the knowledge. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12925–12935, 2020.
- Ciravegna, G., Barbiero, P., Giannini, F., Gori, M., Lió, P., Maggini, M., and Melacci, S. Logic explained networks. *arXiv preprint arXiv:2108.05149*, 2021.
- Daneshjou, R., Vodrahalli, K., Liang, W., Novoa, R. A., Jenkins, M., Rotemberg, V., Ko, J., Swetter, S. M., Bailey, E. E., Gevaert, O., et al. Disparities in dermatology ai: Assessments using diverse clinical images. *arXiv preprint arXiv:2111.08006*, 2021.
- Garcez, A. d., Besold, T. R., De Raedt, L., Földiak, P., Hitzler, P., Icard, T., Kühnberger, K.-U., Lamb, L. C., Miikkulainen, R., and Silver, D. L. Neural-symbolic learning and reasoning: contributions and challenges. In *2015 AAAI Spring Symposium Series*, 2015.
- Geifman, Y. and El-Yaniv, R. Selectivenet: A deep neural network with an integrated reject option. In *International conference on machine learning*, pp. 2151–2159. PMLR, 2019.
- Hastie, T. and Tibshirani, R. Generalized additive models: some applications. *Journal of the American Statistical Association*, 82(398):371–386, 1987.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hinton, G., Vinyals, O., Dean, J., et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Jain, S., Agrawal, A., Saporta, A., Truong, S. Q., Duong, D. N., Bui, T., Chambon, P., Zhang, Y., Lungren, M. P., Ng, A. Y., et al. Radgraph: Extracting clinical entities and relations from radiology reports. *arXiv preprint arXiv:2106.14463*, 2021.
- Johnson, A., Lungren, M., Peng, Y., Lu, Z., Mark, R., Berkowitz, S., and Horng, S. Mimic-cxr-jpg-chest radiographs with structured labels.
- Kawahara, J., Daneshvar, S., Argenziano, G., and Hamarneh, G. Seven-point checklist and skin lesion classification using multitask multimodal neural nets. *IEEE journal of biomedical and health informatics*, 23(2):538–546, 2018.
- Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., Viegas, F., and Sayres, R. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). (2017). *arXiv preprint arXiv:1711.11279*, 2017.
- Koh, P. W., Nguyen, T., Tang, Y. S., Mussmann, S., Pierson, E., Kim, B., and Liang, P. Concept bottleneck models. In *International Conference on Machine Learning*, pp. 5338–5348. PMLR, 2020.
- Letham, B., Rudin, C., McCormick, T. H., and Madigan, D. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics*, 9(3):1350–1371, 2015.

- 495 Lu, M., Zhao, Q., Zhang, J., Pohl, K. M., Fei-Fei, L.,
 496 Niebles, J. C., and Adeli, E. Metadata normalization.
 497 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10917–10927,
 498 2021.
- 500 Lucieri, A., Bajwa, M. N., Braun, S. A., Malik, M. I., Dengel, A., and Ahmed, S. On interpretability of deep learning based skin lesion classifiers using concept activation vectors. In *2020 international joint conference on neural networks (IJCNN)*, pp. 1–10. IEEE, 2020.
- 506 Lundberg, S. M. and Lee, S.-I. A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*, pp. 4768–4777, 2017.
- 511 Mendelson, E. *Introduction to mathematical logic*. Chapman and Hall/CRC, 2009.
- 514 Montavon, G., Samek, W., and Müller, K.-R. Methods
 515 for interpreting and understanding deep neural networks.
 516 *Digital signal processing*, 73:1–15, 2018.
- 518 Ribeiro, M. T., Singh, S., and Guestrin, C. ”why should
 519 i trust you?” explaining the predictions of any classifier.
 520 In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp.
 521 1135–1144, 2016.
- 523 Rosenzweig, J., Sicking, J., Houben, S., Mock, M., and
 525 Akila, M. Patch shortcuts: Interpretable proxy models
 526 efficiently find black-box vulnerabilities. In *Proceedings
 527 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 56–65, 2021.
- 530 Rotemberg, V., Kurtansky, N., Betz-Stablein, B., Caffery,
 531 L., Chousakos, E., Codella, N., Combalia, M., Dusza, S.,
 532 Guitera, P., Gutman, D., et al. A patient-centric dataset
 533 of images and metadata for identifying melanomas using
 534 clinical context. *Scientific data*, 8(1):1–8, 2021.
- 535 Rudin, C. Stop explaining black box machine learning
 536 models for high stakes decisions and use interpretable
 537 models instead. *Nature Machine Intelligence*, 1(5):206–
 538 215, 2019.
- 540 Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P.
 541 Distributionally robust neural networks for group shifts:
 542 On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- 545 Samek, W., Binder, A., Montavon, G., Lapuschkin, S., and
 546 Müller, K.-R. Evaluating the visualization of what a deep
 547 neural network has learned. *IEEE transactions on neural
 548 networks and learning systems*, 28(11):2660–2673, 2016.
- 549 Sarkar, A., Vijaykeerthy, D., Sarkar, A., and Balasubramanian, V. N. Inducing semantic grouping of latent concepts for explanations: An ante-hoc approach. *arXiv preprint arXiv:2108.11761*, 2021.
- 550 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R.,
 551 Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- 556 Shrikumar, A., Greenside, P., Shcherbina, A., and Kundaje, A. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.
- 561 Simonyan, K., Vedaldi, A., and Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- 566 Singla, S., Pollack, B., Chen, J., and Batmanghelich, K. Explanation by progressive exaggeration. *arXiv preprint arXiv:1911.00483*, 2019.
- 571 Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- 576 Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. In *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- 581 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- 586 Tschandl, P., Rosendahl, C., and Kittler, H. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018.
- 591 Wadden, D., Wennberg, U., Luan, Y., and Hajishirzi, H. Entity, relation, and event extraction with contextualized span representations. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5784–5789, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1585. URL <https://aclanthology.org/D19-1585>.
- 596 Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The caltech-ucsd birds-200-2011 dataset. 2011.

- 550 Wan, C., Belo, R., and Zejnilovic, L. Explainability's gain is
551 optimality's loss? how explanations bias decision-making.
552 In *Proceedings of the 2022 AAAI/ACM Conference on AI,
553 Ethics, and Society*, pp. 778–787, 2022.
- 554 Wang, J., Yu, X., and Gao, Y. Feature fusion vision
555 transformer for fine-grained visual categorization. *arXiv
556 preprint arXiv:2107.02341*, 2021.
- 558 Xian, Y., Lampert, C. H., Schiele, B., and Akata, Z. Zero-
559 shot learning—a comprehensive evaluation of the good,
560 the bad and the ugly. *IEEE transactions on pattern anal-
561 ysis and machine intelligence*, 41(9):2251–2265, 2018.
- 562
- 563 Yeh, C.-K., Kim, B., Arik, S., Li, C.-L., Ravikumar, P., and
564 Pfister, T. On concept-based explanations in deep neural
565 networks. 2019.
- 566
- 567 Yu, K., Ghosh, S., Liu, Z., Deible, C., and Batmanghelich, K.
568 Anatomy-guided weakly-supervised abnormality local-
569 ization in chest x-rays. *arXiv preprint arXiv:2206.12704*,
570 2022.
- 571 Yuksekgonul, M., Wang, M., and Zou, J. Post-hoc concept
572 bottleneck models. *arXiv preprint arXiv:2205.15480*,
573 2022.
- 574
- 575 Zarlenga, M. E., Barbiero, P., Ciravegna, G., Marra, G., Gi-
576 annini, F., Diligenti, M., Shams, Z., Precioso, F., Melacci,
577 S., Weller, A., et al. Concept embedding models. *arXiv
578 preprint arXiv:2209.09056*, 2022.
- 579
- 580
- 581
- 582
- 583
- 584
- 585
- 586
- 587
- 588
- 589
- 590
- 591
- 592
- 593
- 594
- 595
- 596
- 597
- 598
- 599
- 600
- 601
- 602
- 603
- 604

605 **A. Appendix**

606 **A.1. Reproducibility**

608 The code with all the experiments is available in the anonymous GitHub repository [https://github.com/AI09-guy/](https://github.com/AI09-guy/ICML-Submission)
 609 [ICML-Submission](https://github.com/AI09-guy/ICML-Submission).

610

611 **A.2. Background of First-order logic (FOL) and Neuro-symbolic-AI**

613 FOL is a logical function that accepts predicates (concept presence/absent) as input and returns a True/False output being a
 614 logical expression of the predicates. The logical expression, which is a set of AND, OR, Negative, and parenthesis, can be
 615 written in the so-called Disjunctive Normal Form (DNF) (Mendelson, 2009). DNF is a FOL logical formula composed of a
 616 disjunction (OR) of conjunctions (AND), known as the “sum of products”.

617 Neuro-symbolic AI is an area of study that encompasses deep neural networks with symbolic approaches to computing
 618 and AI to complement the strengths and weaknesses of each, resulting in a robust AI capable of reasoning and cognitive
 619 modeling (Belle, 2020). Neuro-symbolic systems are hybrid models that leverage the robustness of connectionist methods
 620 and the soundness of symbolic reasoning to effectively integrate learning and reasoning (Garcez et al., 2015; Besold et al.,
 621 2017).

622

623

624

625 **A.3. Learning the concepts**

626 As discussed in Section 2, $f^0 : \mathcal{X} \rightarrow \mathcal{Y}$ is a pre-trained Blackbox. Also, $f^0(\cdot) = h^0 \circ \Phi(\cdot)$. Here, $\Phi : \mathcal{X} \rightarrow \mathbb{R}^l$ is the image
 627 embeddings, transforming the input images to an intermediate representation and $h^0 : \mathbb{R}^l \rightarrow \mathcal{Y}$ is the classifier, classifying
 628 the output \mathcal{Y} using the embeddings, Φ . Our approach is applicable for both datasets with and without human-interpretable
 629 concept annotations. For datasets with the concept annotation $\mathcal{C} \in \mathbb{R}^{N_c}$ (N_c being the number of concepts per image \mathcal{X}), we
 630 learn $t : \mathbb{R}^l \rightarrow \mathcal{C}$ to classify the concepts using the embeddings. Per this definition, t outputs a scalar value c representing
 631 a single concept for each input image. We adopt the concept learning strategy in PosthocCBM (PCBM) (Yuksekgonul
 632 et al., 2022) for datasets without concept annotation. Specifically, we leverage a set of image embeddings with the concept
 633 being present and absent. Next, we learn a linear SVM (t) to construct the concept activation matrix (Kim et al., 2017) as
 634 $\mathbf{Q} \in \mathbb{R}^{N_c \times l}$. Finally we estimate the concept value as $c = \frac{\langle \Phi(x), q^i \rangle}{\|q_i\|_2^2} \in \mathbb{R}$ utilizing each row q^i of \mathbf{Q} . Thus, the complete
 635 tuple of j^{th} sample is $\{x_j, y_j, c_j\}$, denoting the image, label, and learned concept vector, respectively.

636

637 **A.4. Optimization**

638 In this section, we will discuss the loss function used in distilling the knowledge from the blackbox to the symbolic model.
 639 We remove the superscript k for brevity. We adopted the optimization proposed in (Geifman & El-Yaniv, 2019). Specifically,
 640 we convert the constrained optimization problem in Equation (2) as

641
$$\begin{aligned} \mathcal{L}_s &= \mathcal{R}(\pi, g) + \lambda_s \Psi(\tau - \zeta(\pi)) \\ 642 \Psi(a) &= \max(0, a)^2, \end{aligned} \tag{5}$$

643 where τ is the target coverage and λ_s is a hyperparameter (Lagrange multiplier). We define $\mathcal{R}(\cdot)$ and $\mathcal{L}_{g,\pi}(\cdot)$ in Equation (1)
 644 and Equation (3) respectively. ℓ in Equation (3) is defined as follows:

645
$$\ell(f, g) = \ell_{distill}(f, g) + \lambda_{lens} \sum_{i=1}^r \mathcal{H}(\beta^i), \tag{6}$$

646 where λ_{lens} and $\mathcal{H}(\beta^i)$ are the hyperparameters and entropy regularize, introduced in (Barbiero et al., 2022) with r being
 647 the total number of class labels. Specifically, β^i is the categorical distribution of the weights corresponding to each concept.
 648 To select only a few relevant concepts for each target class, higher values of λ_{lens} will lead to a sparser configuration of β . ℓ
 649 is the knowledge distillation loss (Hinton et al., 2015), defined as

650

651

652

653

654

655

656

657

658

659

$$\ell(f, g) = (\alpha_{KD} * T_{KD} * T_{KD})KL(\text{LogSoftmax}(g(.)/T_{KD}), \text{Softmax}(f(.)/T_{KD})) + (1 - \alpha_{KD})CE(g(.), y), \quad (7)$$

where T_{KD} is the temperature, CE is the Cross-Entropy loss, and α_{KD} is relative weighting controlling the supervision from the blackbox f and the class label y .

As discussed in (Geifman & El-Yaniv, 2019), we also define an auxiliary interpretable model using the same prediction task assigned to q using the following loss function

$$\mathcal{L}_{aux} = \frac{1}{m} \sum_{j=1}^m \ell_{distill}(f(\mathbf{x}_j), g(\mathbf{c}_j)) + \lambda_{lens} \sum_{i=1}^r \mathcal{H}(\beta^i), \quad (8)$$

which is agnostic of any coverage. \mathcal{L}_{aux} is necessary for optimization as the symbolic model will focus on the target coverage τ before learning any relevant features, overfitting to the wrong subset of the training set. The final loss function to optimize by g in each iteration is as follows:

$$\mathcal{L} = \alpha \mathcal{L}_f + (1 - \alpha) \mathcal{L}_{aux}, \quad (9)$$

where α is the can be tuned as a hyperparameter. Following (Geifman & El-Yaniv, 2019), we also use $\alpha = 0.5$ in all of our experiments.

A.5. Algorithm

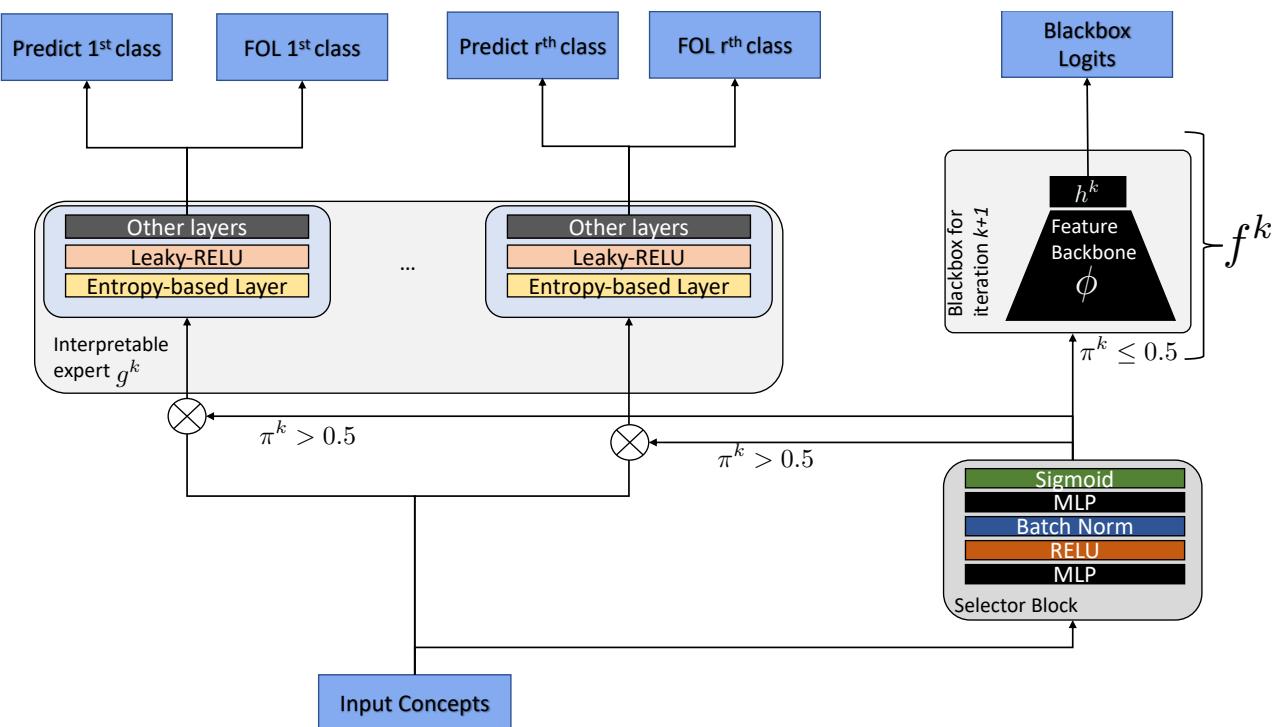


Figure 7. Architecture of MoIE. In an iteration k during inference, the selector routes the samples to go through the interpretable expert g^k if the probability $\pi^k \geq 0.5$. If $\pi^k < 0.5$, the selector routes the samples, through f^k , the Blackbox for iteration $k+1$. Note $f^k = h^k(\Phi(\cdot))$ is an approximation of the residual $r^k = f^{k-1} - g^k$.

715 **Algorithm 1** *Route, interpret and Repeat* algorithm to generate FOL explanations locally.

716 1: **Input:** Complete tuple: $\{x_j, y_j, c_j\}_{j=1}^n$; initial blackbox $f^0 = h^0(\Phi(\cdot))$; K as the total iterations; Coverages τ_1, \dots, τ_k .
 717 2: **Output:** Sparse mixture of experts and their selectors $\{g^k, \pi^k\}_{k=1}^K$ and the final residual $f^K = h^K(\Phi(\cdot))$
 718 3: Fix Φ .
 719 4: **for** $k = 1 \dots K$ **do**
 720 5: Fix $\pi^1 \dots \pi^{k-1}$.
 721 6: Swap x_i and x_{i+1} .
 722 7: Minimize \mathcal{L}^k using equation 9 to learn π^k and g^k .
 723 8: Calculate $r^k = f^{k-1}(\cdot) - g^k(\cdot)$.
 724 9: Minimize equation 4 to learn $f^k(\cdot)$, the new blackbox for the next iteration $k + 1$.
 725 10: **end for**
 726 11: **for** $k = 1 \dots K$ **do**
 727 12: **for** sample j in test-set **do**
 728 13: **repeat**
 729 14: Initialize sub_select_concept = True
 730 15: Initialize the percentile_threshold = 99.
 731 16: Retrieve the predicted class label of sample j from the expert k as: $\hat{y}_j = g^k(c_j)$
 732 17: Create a mask vector m_j . $m_j[i] = 1$ if $\tilde{\alpha}[\hat{y}_j][i] \geq \text{percentile}(\tilde{\alpha}[\hat{y}_j], \text{percentile_threshold})$ and 0 otherwise. Specifically, the i^{th} entry in m_j is one if the i^{th} value of the attention score $\tilde{\alpha}[\hat{y}_j]$ is greater than $(\text{percentile_attention})^{th}$ percentile.
 733 18: Subselect the concept vector as \tilde{c}_j as: $\tilde{c}_j = c_j \odot m_j$
 734 19: **if** $g^k(\tilde{c}_j) \neq \hat{y}_j$ **then**
 735 20: percentile_threshold = percentile_threshold - 1
 736 21: sub_select_concept = false
 737 22: **end if**
 738 23: **until** sub_select_concept is True
 739 24: Using the subselected concept vector \tilde{c}_j , construct the FOL expression of the j^{th} sample as suggested by (Barbiero et al., 2022).
 740 25: **end for**
 741 26: **end for**

742
 743
 744
 745
 746
 747 Algorithm 1 explains the overall training procedure of our method. Figure 7 displays the architecture of our model in
 748 iteration k .

A.6. Dataset

749
 750 **CUB-200** The Caltech-UCSD Birds-200-2011 ((Wah et al., 2011)) is a fine-grained classification dataset comprising 11788
 751 images and 312 noisy visual concepts. The aim is to classify the correct bird species from 200 possible classes. We adopted
 752 the strategy discussed in (Barbiero et al., 2022) to extract 108 denoised visual concepts. Also, we utilize training/validation
 753 splits shared in (Barbiero et al., 2022). Finally, we use the state-of-the-art classification models Resnet-101 ((He et al.,
 754 2016)) and Vision-Transformer (VIT) ((Wang et al., 2021)) as the blackboxes f^0 .
 755
 756

757
 758 **Animals with attributes2 (AwA2)** AwA2 dataset (Xian et al., 2018) consists of 37322 images of total 50 animals
 759 classes with 85 numeric attribute. We use the state-of-the-art classification models Resnet-101 ((He et al., 2016)) and
 760 Vision-Transformer (VIT) ((Wang et al., 2021)) as the blackboxes f^0 .
 761

762
 763 **HAM10000** HAM10000 ((Tschandl et al., 2018)) is a classification dataset aiming to classify a skin lesion benign or
 764 malignant. Following (Daneshjou et al., 2021), we use Inception (Szegedy et al., 2015) model, trained on this dataset as the
 765 blackbox f^0 . We follow the strategy in (Lucieri et al., 2020) to extract the 8 concepts from the Derm7pt ((Kawahara et al.,
 766 2018)) dataset.
 767

768
 769 **SIIM-ISIC** To test a real-world transfer learning use case, we evaluate the model trained on HAM10000 on a subset
 770 of the SIIM-ISIC(Rotemberg et al., 2021)) Melanoma Classification dataset. We use the same concepts described in the
 771

770 HAM10000 dataset.
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781

MIMIC-CXR We use 220,763 frontal images from the MIMIC-CXR dataset (Johnson et al.) aiming to classify effusion. We obtain the anatomical and observation concepts from the RadGraph annotations in RadGraph’s inference dataset ((Jain et al., 2021)), automatically generated by DYGIE++ ((Wadden et al., 2019)). We use the test-train-validation splits from (Yu et al., 2022) and Densenet121 (Huang et al., 2017) as the blackbox f^0 .
 782
 783
 784
 785

786 A.7. Architectural details of symbolic experts and hyperparameters

787 Table 3 demonstrates different settings to train the Blackbox of CUB-200, Awa2 and MIMIC-CXR respectively. For the
 788 VIT-based backbone, we used the same hyperparameter setting used in the state-of-the-art Vit-B_16 variant in (Wang et al.,
 789 2021). To train t , we flatten the feature maps from the last convolutional block of Φ using “Adaptive average pooling” for
 790 CUB-200 and Awa2 datasets. For MIMIC-CXR and HAM10000, we flatten out the feature maps from the last convolutional
 791 block. For VIT-based backbones, we take the first block of representation from the encoder of VIT. For HAM10000, we
 792 use the same Blackbox in (Yuksekgonul et al., 2022). Table 4, Table 5, Table 6, Table 7 enumerate all the different settings
 793 to train the interpretable experts for CUB-200, Awa2, HAM, and MIMIC-CXR respectively. All the residuals in different
 794 iterations follow the same settings as their blackbox counterparts.
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806

807 *Table 3.* Hyperparameter setting of different convolution-based Blackboxes used by CUB-200, Awa2 and MIMIC-CXR

Setting	CUB-200	Awa2	MIMIC-CXR
Backbone	ResNet-101	ResNet-101	DenseNet-121
Pretrained on ImageNet	True	True	True
Image size	448	224	512
Learning rate	0.001	0.001	0.01
Optimization	SGD	Adam	SGD
Weight-decay	0.00001	0	0.0001
Epcohs	95	90	50
Layers used as ϕ	till 4 th ResNet Block	till 4 th ResNet Block	till 4 th DenseNet Block
Flattening type for the input to t	Adaptive average pooling	Adaptive average pooling	Flatten

825
 826 *Table 4.* Hyperparameter setting of interpretable experts (g) trained on ResNet-101 (top) and VIT (bottom) blackboxes for the CUB-200
 827 dataset

Settings based on dataset	Expert1	Expert2	Expert3	Expert4	Expert5	Expert6
CUB-200 (ResNet-101)						
+ Batch size	16	16	16	16	16	16
+ Coverage (τ)	0.2	0.2	0.2	0.2	0.2	0.2
+ Learning rate	0.01	0.01	0.01	0.01	0.01	0.01
+ λ_{lens}	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001
+ α_{KD}	0.9	0.9	0.9	0.9	0.9	0.9
+ T_{KD}	10	10	10	10	10	10
+hidden neurons	10	10	10	10	10	10
+ λ_s	32	32	32	32	32	32
+ Temperature						
+ E-Lens (T_{lens})	0.7	0.7	0.7	0.7	0.7	0.7
CUB-200 (VIT)						
+ Batch size	16	16	16	16	16	16
+ Coverage (τ)	0.2	0.2	0.2	0.2	0.2	0.2
+ Learning rate	0.01	0.01	0.01	0.01	0.01	0.01
+ λ_{lens}	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001
+ α_{KD}	0.99	0.99	0.99	0.99	0.99	0.99
+ T_{KD}	10	10	10	10	10	10
+hidden neurons	10	10	10	10	10	10
+ λ_s	32	32	32	32	32	32
+E-Lens (T_{lens})	6.0	6.0	6.0	6.0	6.0	6.0

880
881 *Table 5.* Hyperparameter setting of interpretable experts (g) trained on ResNet-101 (top) and VIT (bottom) blackboxes for the Awa2
882 dataset

Settings based on dataset	Expert1	Expert2	Expert3	Expert4	Expert5	Expert6
Awa2 (ResNet-101)						
+ Batch size	30	30	30	30	-	-
+ Coverage (τ)	0.4	0.35	0.35	0.25	-	-
+ Learning rate	0.001	0.001	0.001	0.001	-	-
+ λ_{lens}	0.0001	0.0001	0.0001	0.0001	-	-
+ α_{KD}	0.9	0.9	0.9	0.9	-	-
+ T_{KD}	10	10	10	10	-	-
+hidden neurons	10	10	10	10	-	-
+ λ_s	32	32	32	32	-	-
+ Temperature						
+ E-Lens (T_{lens})	0.7	0.7	0.7	0.7	-	-
Awa2 (VIT)						
+ Batch size	30	30	30	30	30	30
+ Coverage (τ)	0.2	0.2	0.2	0.2	0.2	0.2
+ Learning rate	0.01	0.01	0.01	0.01	0.01	0.01
+ λ_{lens}	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001
+ α_{KD}	0.99	0.99	0.99	0.99	0.99	0.99
+ T_{KD}	10	10	10	10	10	10
+hidden neurons	10	10	10	10	10	10
+ λ_s	32	32	32	32	32	32
+ Temperature						
+ E-Lens (T_{lens})	6.0	6.0	6.0	6.0	6.0	6.0

908
909 *Table 6.* Hyperparameter setting of interpretable experts (g) for the dataset HAM10000
910
911
912
913
914
915
916
917
918
919

Settings based on dataset	Expert1	Expert2	Expert3	Expert4	Expert5	Expert6
HAM10000 (Inception-V3)						
+ Batch size	32	32	32	32	32	32
+ Coverage (τ)	0.4	0.2	0.2	0.2	0.1	0.1
+ Learning rate	0.01	0.01	0.01	0.01	0.01	0.01
+ λ_{lens}	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001
+ α_{KD}	0.9	0.9	0.9	0.9	0.9	0.9
+ T_{KD}	10	10	10	10	10	10
+hidden neurons	10	10	10	10	10	10
+ λ_s	64	64	64	64	64	64
+ Temperature						
+ E-Lens (T_{lens})	0.7	0.7	0.7	0.7	0.7	0.7

Table 7. Hyperparameter setting of interpretable experts (g) for the dataset MIMIC-CXR

Settings based on dataset	Expert1	Expert2	Expert3
Effusion-MIMIC-CXR (DenseNet-121)			
+ Batch size	1028	1028	1028
+ Coverage (τ)	0.6	0.2	0.15
+ Learning rate	0.01	0.01	0.01
+ λ_{lens}	0.0001	0.0001	0.0001
+ α_{KD}	0.99	0.99	0.99
+ T_{KD}	20	20	20
+hidden neurons	20, 20	20, 20	20, 20
+ λ_s	96	128	256
+E-Lens (T_{lens})	7.6	7.6	7.6

A.8. Estimation of completeness score

Let $f^0(x) = h^0(\Phi(x))$ is the initial Blackbox as per Section 2. The Concept completeness paper (Yeh et al., 2019) assumes $\Phi(\mathbf{x}) \in \mathbb{R}^l$ (s.t., $l = T.d$) to be a concatenation of $[\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_T)]$ s.t., $\phi(\mathbf{x}) \in \mathbb{R}^d$. Recall we utilize t to learn the concepts \mathcal{C} with N_c being the total number of concepts per image. So the parameters of t , represented by $\omega_1, \omega_2, \dots, \omega_{N_c}$ s.t., $\omega_i \in \mathbb{R}^d$ represent linear direction in the embedding space $\phi(\cdot) \in \mathbb{R}^d$. Next, we compute the concept product $v_c(\mathbf{x}_t)(<\phi(\mathbf{x}_t), \omega_j>)^{N_c}_{j=1}$, denoting the similarity between the image embedding and linear direction of j^{th} concept. Finally, we normalize $v_c(\cdot)$ to obtain the concept score as $v_v(\mathbf{x}) = (\frac{v_c(\mathbf{x}_t)}{\|v_c(\mathbf{x}_t)\|_2})^{T}_{t=1} \in \mathbb{R}^{T.N_c}$.

Next for a Blackbox $f^0(x) = h^0(\Phi(\mathbf{x}))$, set of concepts c_1, c_2, \dots, c_{N_c} and their linear direction $\omega_1, \omega_2, \dots, \omega_{N_c}$ in the embedding space and , we compute the completeness score as:

$$\eta_{f^0} = \frac{\sup_{\Gamma} \mathbb{P}_{\mathbf{x}, y \sim V}[y = \arg \max_{y'} h_{y'}^0(\Gamma(v_c(\mathbf{x}))) - a_r]}{\mathbb{P}_{\mathbf{x}, y \sim V}[y = \arg \max_{y'} f_{y'}^0(\mathbf{x})] - a_r}, \quad (10)$$

where V is the validation set and $\Gamma : \mathbb{R}^{T.m} \rightarrow \mathbb{R}^l$, projection from the concept score to the embedding space Φ . For CUB-200 and Awa2 we estimate $\mathbb{P}_{\mathbf{x}, y \sim V}[y = \arg \max_{y'} h_{y'}^0(\Gamma(v_c(\mathbf{x})))]$ as the best accuracy using the given concepts and a_r is the random accuracy. For HAM10000, we estimate the same as the best AUROC. Completeness score indicates the consistency between the prediction based just on concepts and the given Blackbox f^0 . If the identified concepts are sufficiently rich, label prediction will be similar to the Blackbox, resulting in higher completeness scores for the concept set. In all our experiments, Γ is a two-layer feedforward neural network with 1000 neurons.

To plot the completeness score in Figure 5a-c, we select the topN concepts iteratively representing the $N < N_c$ concepts most significant to the prediction of the interpretable model g . Recall we follow Entropy based linear neural network (Barbiero et al., 2022) as g . So each concept has an associated attention score, α in g (Barbiero et al., 2022), denoting the importance of the concept for the prediction. We select the topN concepts based on the N concepts with highest attention weights. We get the linear direction of these topN concepts from the parameters of the learned t and project it to the embedding space Φ using Γ . If Γ reconstructs the discriminative features from the concepts successfully, the concepts achieves high completeness scores, showing faithfulness with the Blackbox. Recall Figure 5a-c demonstrate that MoIE outperforms the baselines in terms of the completeness scores. This suggests that MoIE identifies rich instance-specific concepts than the baselines, being consistent with the Blackbox.

A.9. Flow diagram to eliminate shortcut

Figure 8 shows the flow diagram to eliminate shortcut.

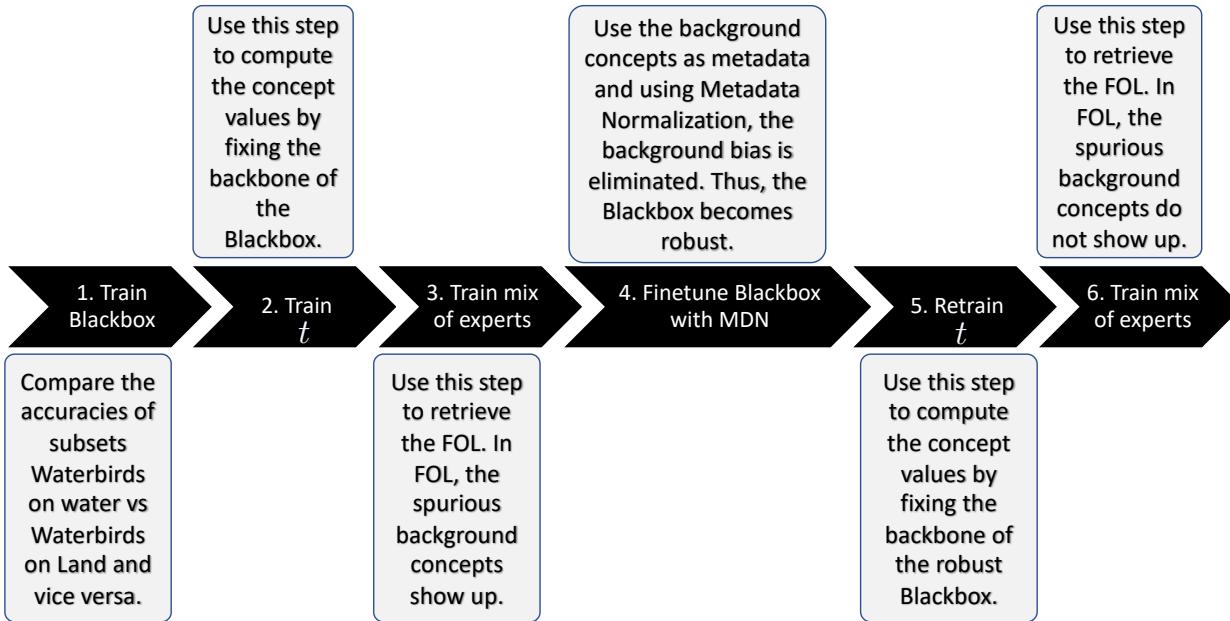


Figure 8. The flow diagram to eliminate the shortcut from vision datasets using FOL by MoIE.

A.10. More Results

A.10.1. PERFORMANCE OF EXPERTS AND RESIDUAL FOR RESNET-DERIVED EXPERTS OF AWA2 AND CUB-200 DATASETS

Figure 9 shows the coverage (top row), performances (bottom row) of each expert and residual across iterations of - (a) ResNet101-derived Awa2 and (b) ResNet101-derived CUB-200 respectively.

A.10.2. CONCEPT VALIDATION OF AWA2

Figure 10 shows the completeness scores and the drop in accuracy by zeroing out the concepts for Awa2.

A.11. Computational performance

Figure 11 shows the computational performance compared to the blackbox. Though in MoIE, we sequentially learn the experts and the residuals, they take less computational resources than the blackbox. The experts are simple one or two layer neural networks. Also we only update the classification layer (h) for the residuals, so it takes such a less time. The Flops in the Y axis is computed as Flop of (forward propagation + backward propagation) \times (minibatch size) \times (no of training epochs).

A.11.1. VALIDATING CONCEPTS FROM PURELY THE FOL RULES

To validate the concepts using the extracted FOL explanations, we intervene on the concepts in the derived FOL by setting the values of those concepts to zero for each sample. For example, the FOL explanation of instances in expert4 for the class “Bay Breasted Warbler” include concepts *s.t.*, *back_pattern_stripped* and *leg_color_grey* in Figure 2. For intervention, we set the values of these concepts to zero while values of other concepts remain unchanged. Next we pass the complete intervened concept vector as input to the associated expert, compute the performance and summarize the results in Table 8. We discover that MoIE is highly susceptible to such interventions, and its performance drops significantly. For example, the performance of MoIE deteriorates from 0.91 to 0.42 % (a 53.8 % drop) for CUB-200 VIT-derived MoIE. We compare the results with 1) the CBM + ELL baseline corresponding to CUB-200 (both Resnet101 and VIT), Awa2 (both Resnet101 and VIT) and Effusion of MIMIC-CXR; 2) PCBM + ELL baseline for the skin datasets. For CUB-200 VIT-based baseline model, the performance of the baseline drops by 26.5 % drop, from 0.90 to 0.66. We observe a similar trend for other datasets as

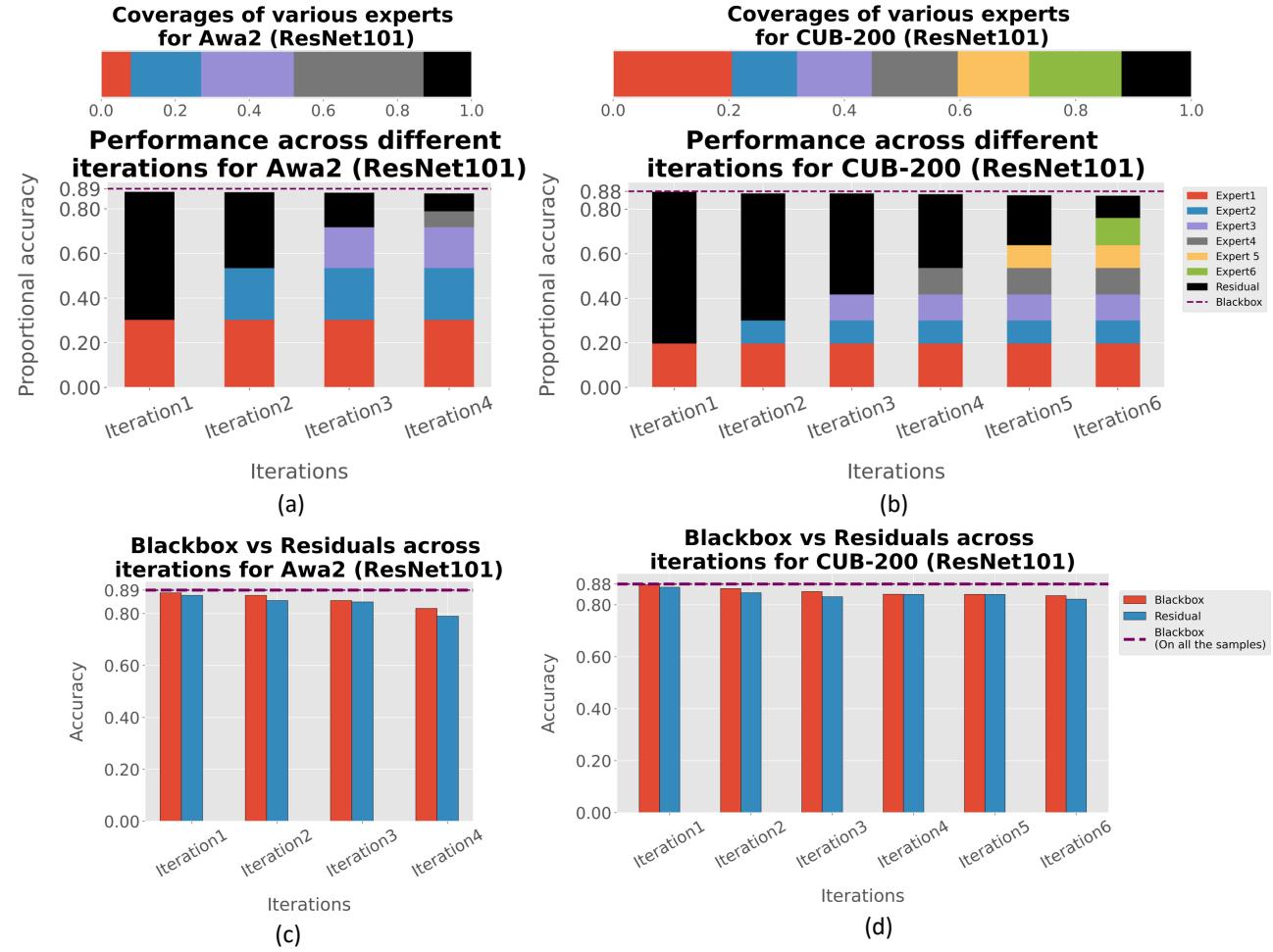


Figure 9. The performances of experts and residuals across iterations for ResNet derived MoIE for CUB-200 and Awa2. (a-c) Coverage and proportional accuracy of the experts and residuals. (e-g) We route the samples covered by the residuals across iterations to the initial Blackbox f^0 and compare the accuracy of f^0 (red bar) with the residual (blue bar).

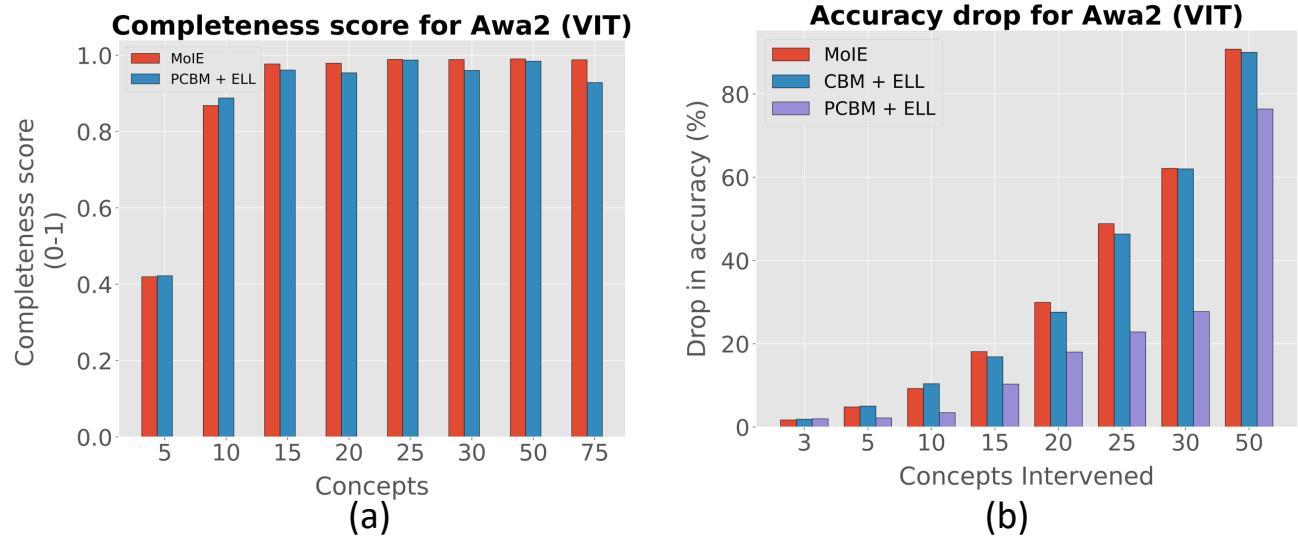


Figure 10. (a): Completeness scores for different significant concepts of Awa2. (a): Drop in accuracy by zeroing out the concepts for Awa2.

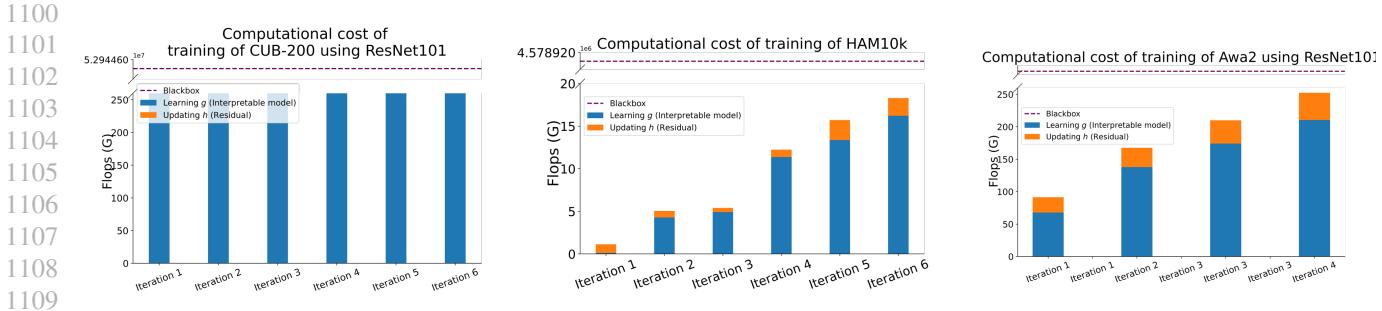


Figure 11. Flops vs iteration for MoIE and the Blackbox. The dotted line in the figure represents the flops taken by the blackbox.

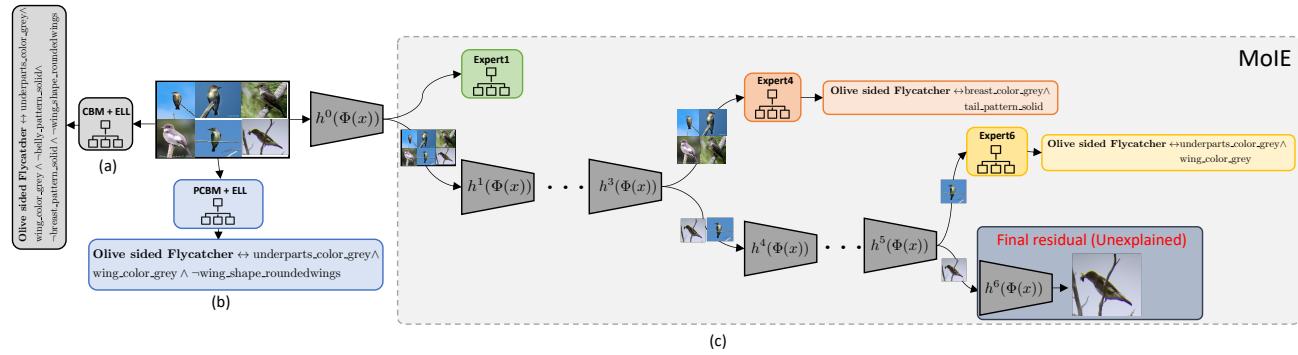


Figure 12. Construction logical explanations of the samples of a category, “Olive sided Flycatcher” in the CUB-200 dataset for (a) VIT-based sequential CBM + ELL as an *interpretable by design* baseline, (b) VIT-based PCBM + ELL as a posthoc based baseline, (c) various experts in MoIE at *inference*. This is an example where the final residual covers the unexplained sample, which is “harder” to explain (indicated in red). Also, MoIE can capture more instance-specific concepts than generic ones by the baselines.

well. As MoIE selects more concrete, instance specific concepts than the baseline by covering subsets of data using various experts, its performance degrades severely compared to the baseline.

Table 8. Validating the concepts purely by extracted FOL rules from MoIE and the baselines. Performance using original explanation refers to the performance of the model employing the concepts from the FOL explanations per sample. We intervene with these concepts by setting their values to zero and reporting it as “Performance using intervened explanation” in the table. In addition, we also show the drop in performance for the two scenarios. The larger drop in performance illustrates the model to be more sensitive to such intervention of the derived concepts.

Model	Performance using original explanation → Performance using intervened explanation (drop %)						
	CUB-200 (ResNet101)	CUB-200 (VIT)	Awa2 (ResNet101)	Awa2 (VIT)	HAM 10000	ISIC	Effusion
MoIE	0.88 → 0.65 (26.1)	0.91 → 0.42 (53.8)	0.87 → 0.54 (37.9)	0.97 → 0.90 (7.2)	0.95 → 0.92 (3.1)	0.82 → 0.79 (3.6)	0.87 → 0.82 (5.7)
Baseline	0.71 → 0.56 (21.1)	0.90 → 0.66 (26.5)	0.86 → 0.56 (34.8)	0.94 → 0.92 (2.1)	0.94 → 0.92 (2.1)	0.83 → 0.80 (3.6)	0.73 → 0.72 (1.3)

A.11.2. DIVERSITY OF EXPLANATIONS FOR CUB

Figure 12 shows the construction of FOL explanations of a category, “Olive sided Flycatcher” in the CUB-200 dataset for the VIT-based baselines and MoIE. In this example, the final expert6 covers the relatively “harder” sample. Figure 13, Figure 14, Figure 15, Figure 16 shows more such FOL explanations. All these examples demonstrate MoIE’s high capability to identify more meaningful instance-specific concepts in FOL explanations. In contrast, the baselines identify the generic concepts for all samples in a class.

1155			
1156			
1157			
1158			
1159			
1160			
1161	Baseline (CBM + ELL)		
1162		• • •	
1163			
1164			
1165	Baseline (PCBM + ELL)		
1166		• • •	
1167			
1168			
1169			
1170			
1171			

Figure 13. Construction logical explanations of the samples of a category, “Harris Sparrow” in the CUB-200 dataset for (a) VIT-based sequential CBM + ELL as an *interpretable by design* baseline, (b) VIT-based PCBMB + ELL as a posthoc based baseline, (c) various experts in MoIE at inference.

A.11.3. DIVERSITY OF EXPLANATIONS FOR AWA2

Figure 17 and 18 demonstrates the flexibility of FOL explanations by VIT-derived MoIE compared to the different baselines for Awa2 dataset qualitatively.

A.11.4. VIT-BASED EXPERTS COMPOSE OF LESS CONCEPTS THAN THE RESNET-BASED COUNTERPARTS

Figure 19 shows the summary statistics for multiclass classification vision datasets. For both datasets, we observe that the VIT-based MoIE uses fewer concepts for explanation than their ResNet-based counterparts. For example, for the CUB-200 dataset, expert6 of VIT-backbone requires 25 concepts compared to 105 by expert6 of ResNet-101-backbone (Figure 19a). The 105 concepts by expert6 is the highest number of concepts utilized by any expert for CUB-200. Similarly, for Awa2, the highest number concept used by an expert is 8 for the VIT-based backbone compared to 80 for the ResNet-101-based backbone (Figure 19b). As mentioned before, the average number of concepts for class $j = \frac{\sum_{\text{all concepts for the samples belong to class } j}}{\# \text{ samples of class } j}$. We can see that for ResNet-101, on average 80 concepts are required to explain a sample correctly for the class “Rhinoceros_Auklet” (expert3 in Figure 24 a). However, for VIT, only 6 concepts are needed to explain a sample correctly “Rhinoceros_Auklet” (expert3 in Figure 24 a). From both of these figures, we can see that different experts require a different number of concepts to explain the same class. For example, Figure 20 (b) and Figure 22 (b) reveal that experts 2 and 6 require 25 and 58 concepts on average to explain “Artic_Tern” correctly respectively.

Figure 26, Figure 27, Figure 28 display the average number of concepts required to predict an animal species correctly in the Awa2 dataset for VIT as backbones. Similarly Figure 29 and Figure 30 display the average number of concepts required to predict an animal species correctly in the Awa2 dataset for ResNet101 as backbones. We can see that for ResNet101, on average, 80 concepts are required to explain a sample correctly for the class “Weasel” (Expert 1 in Figure 29 a). However, for VIT, only three concepts are needed to explain a sample correctly for “Weasel” (Expert 6 in Figure 28 f). Also from both of these figures, we can see that different experts require different number concepts to explain same class. For example Figure 28 (e) and (f) reveal that experts 5 and 6 require 4 and 30 concepts on average to explain “Wolf” correctly.

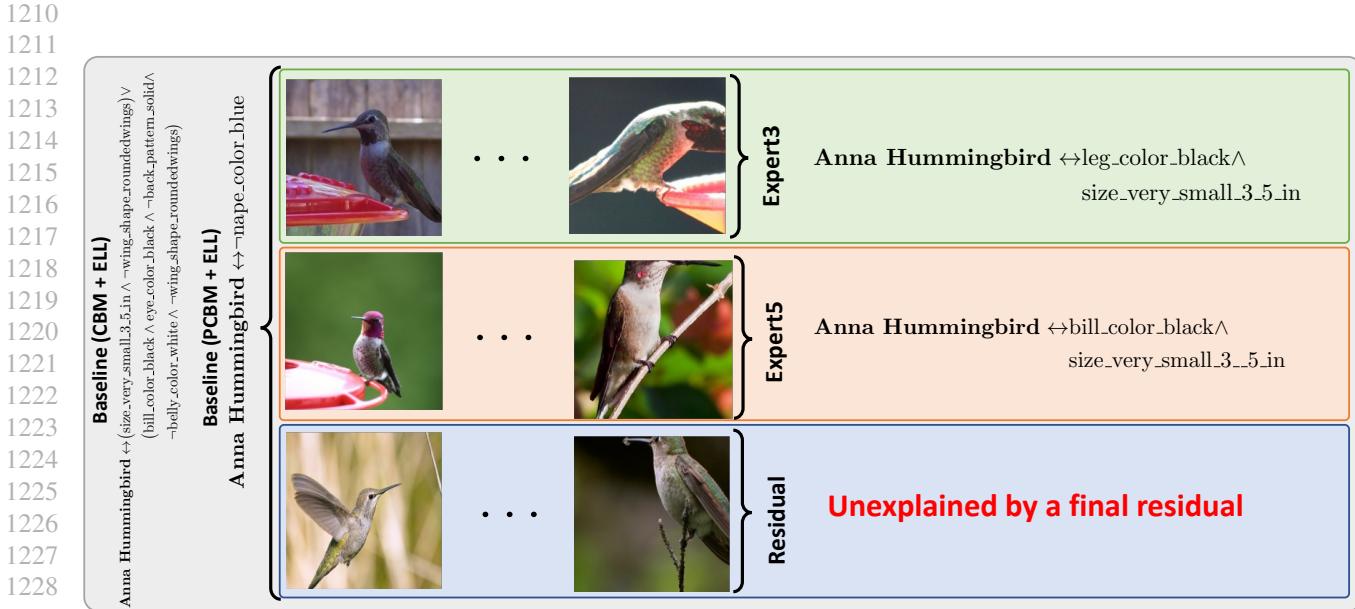


Figure 14. Construction logical explanations of the samples of a category, “Anna Hummingbird” in the CUB-200 dataset for (a) VIT-based sequential CBM + ELL as an *interpretable by design* baseline, (b) VIT-based PCBM + ELL as a posthoc based baseline, (c) various experts in MoIE at inference.

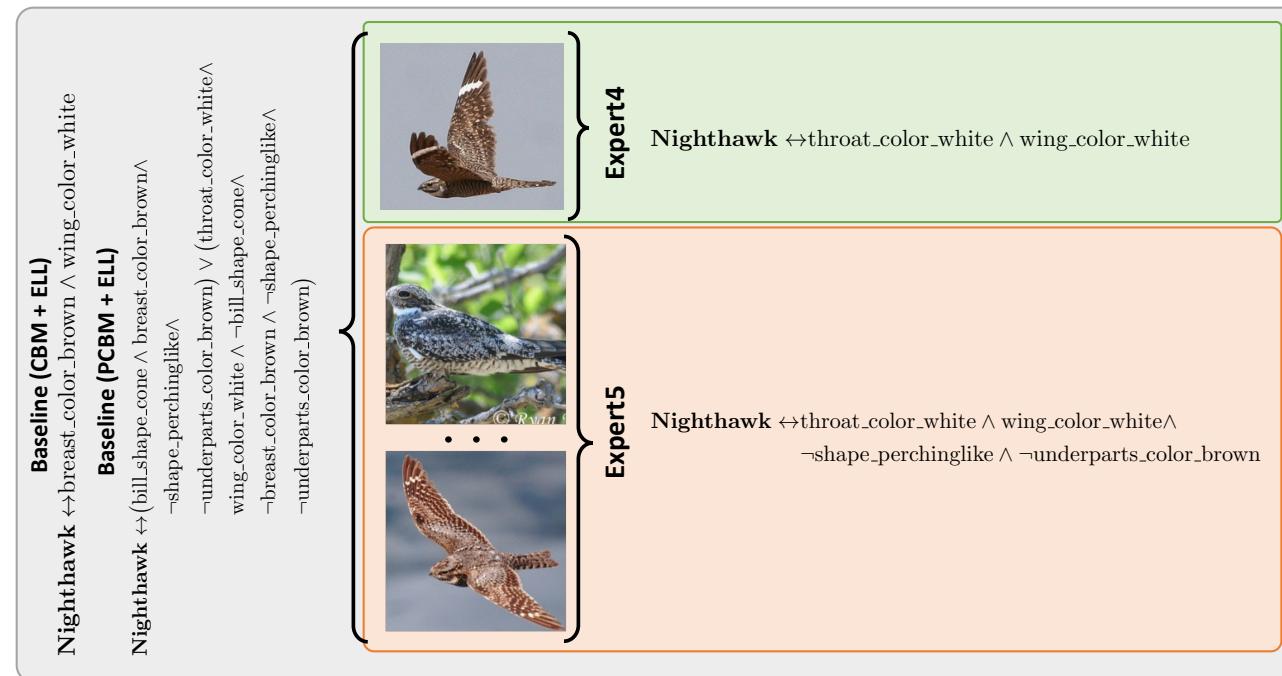


Figure 15. Construction logical explanations of the samples of a category, “Nighthawk” in the CUB-200 dataset for (a) VIT-based sequential CBM + ELL as an *interpretable by design* baseline, (b) VIT-based PCBM + ELL as a posthoc based baseline, (c) various experts in MoIE at inference.

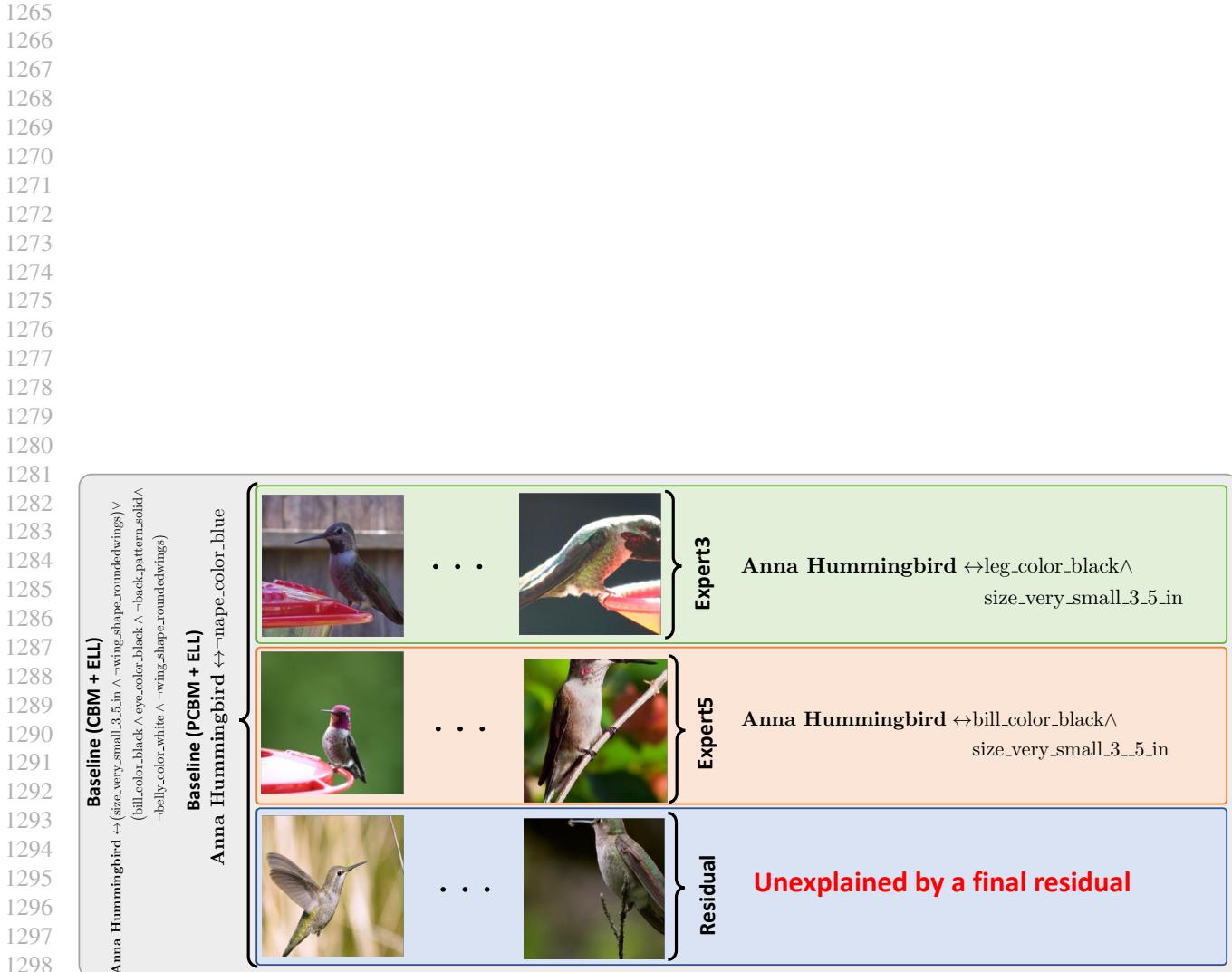


Figure 16. Construction logical explanations of the [samples of](#) a category, “Painted Bunting” in the CUB-200 dataset for (a) VIT-based sequential CBM + ELL as an *interpretable by design* baseline, (b) VIT-based PCBM + ELL as a posthoc based baseline, (c) various experts in MoIE [at inference](#).

1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330

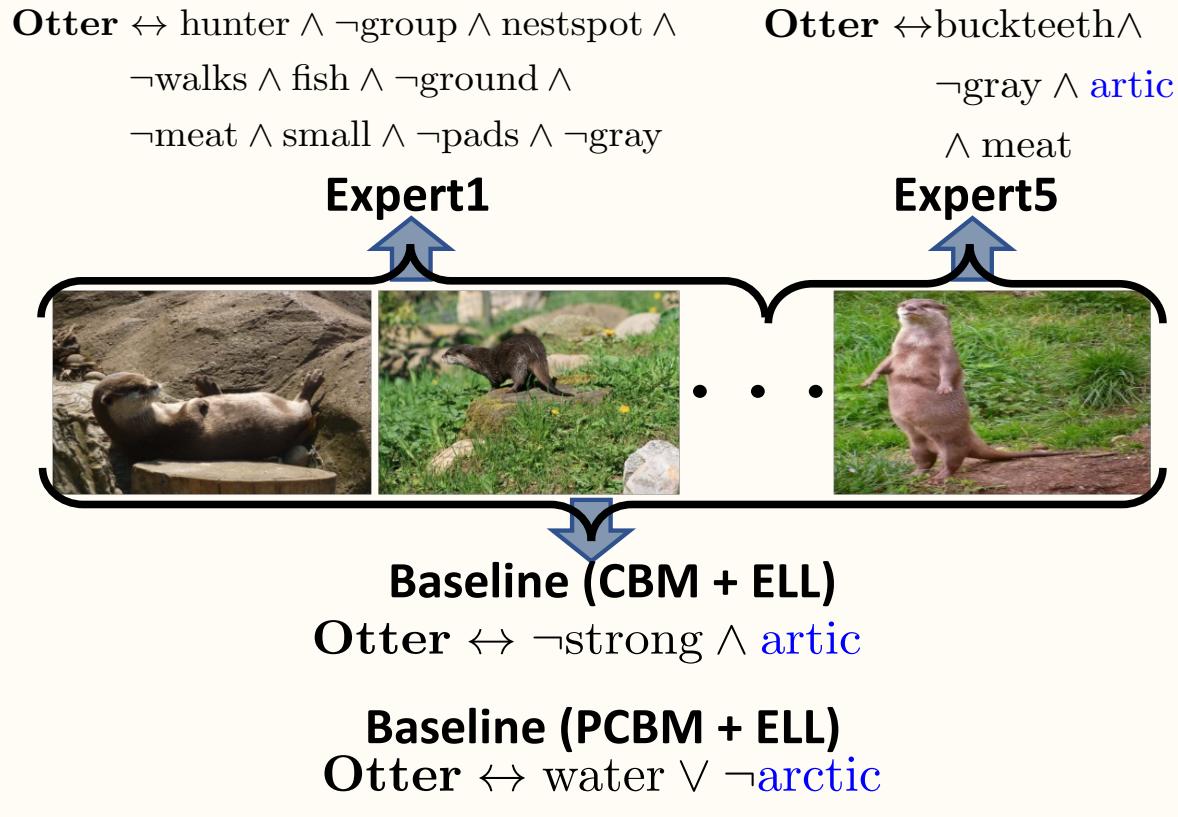


Figure 17. Flexibility of FOL explanations by VIT-derived MoIE MoIE and the CBM + ELL and PCBM + ELL baselines for Awa2 dataset to classify “Otter” at inference. Both the baseline’s FOL constitutes identical concepts to distinguish all the samples. However, expert1 classifies “Otter” with *hunter*, *group etc.* as the identifying concept for the instances covered by it. Similarly expert5 classifies “Otter” using *buckteeth*, *gray etc.*. Note that, *meat* and *gray* are shared between the two experts. We highlight the shared concepts (*artic*) between the experts and the baselines as blue.

1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374

1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427

Horse \leftrightarrow smelly

Expert4



Horse \leftrightarrow \neg longneck \wedge fields

Expert5



Baseline (CBM + ELL)

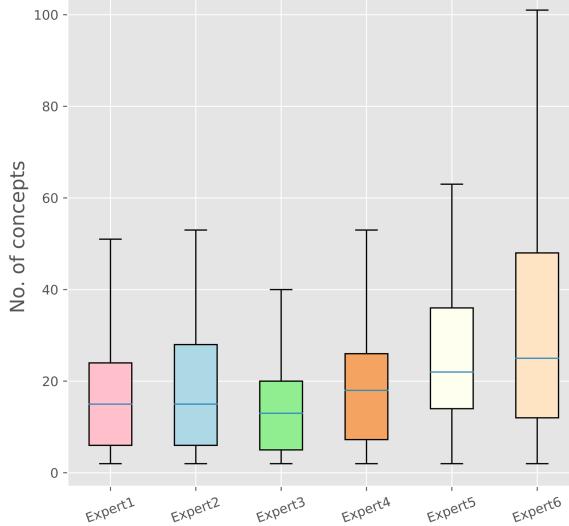
Horse \leftrightarrow (buckteeth \wedge longneck) \vee (longneck \wedge smelly) \vee (longleg \wedge smelly \wedge \neg buckteeth)

Baseline (PCBM + ELL)

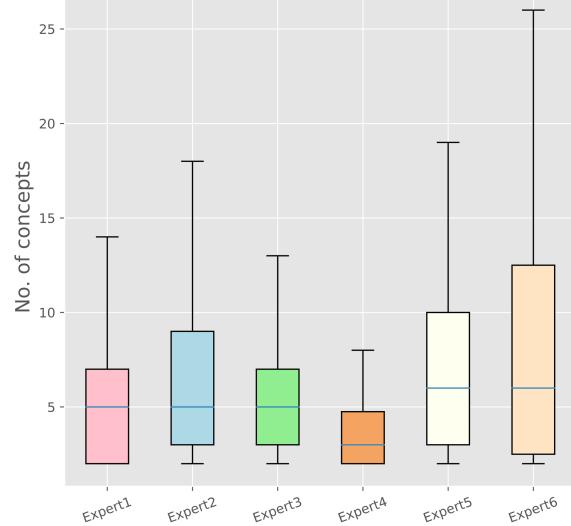
Horse \leftrightarrow (buckteeth \wedge longneck) \vee (\neg buckteeth \wedge \neg longneck) \vee (buckteeth \wedge bulbous \wedge gray \wedge longleg \wedge longneck \wedge \neg forager \wedge \neg solitary \wedge \neg spots) \vee (buckteeth \wedge gray \wedge longleg \wedge longneck \wedge \neg bulbous \wedge \neg forager \wedge \neg solitary \wedge \neg spots) \vee (active \wedge buckteeth \wedge chewteeth \wedge hooves \wedge horns \wedge lean \wedge longleg \wedge longneck \wedge muscle \wedge oldworld \wedge patches \wedge smelly \wedge tail \wedge timid \wedge toughskin \wedge \neg bulbous \wedge \neg bush \wedge \neg forager \wedge \neg forest \wedge \neg gray \wedge \neg hairless \wedge \neg inactive \wedge \neg meatteeth \wedge \neg mountains \wedge \neg nestspot \wedge \neg paws \wedge \neg small \wedge \neg solitary \wedge \neg spots) \vee (active \wedge big \wedge black \wedge bulbous \wedge chewteeth \wedge furry \wedge grazer \wedge ground \wedge hooves \wedge horns \wedge inactive \wedge longleg \wedge longneck \wedge muscle \wedge oldworld \wedge patches \wedge quadrupedal \wedge slow \wedge smelly \wedge strong \wedge tail \wedge timid \wedge toughskin \wedge walks \wedge white \wedge \neg agility \wedge \neg arctic \wedge \neg buckteeth \wedge \neg bush \wedge \neg claws \wedge \neg coastal \wedge \neg fast \wedge \neg fierce \wedge \neg fish \wedge \neg flippers \wedge \neg forager \wedge \neg forest \wedge \neg gray \wedge \neg hairless \wedge \neg hibernate \wedge \neg hunter \wedge \neg jungle \wedge \neg lean \wedge \neg meat \wedge \neg meatteeth \wedge \neg mountains \wedge \neg nestspot \wedge \neg nocturnal \wedge \neg ocean \wedge \neg pads \wedge \neg paws \wedge \neg small \wedge \neg smart \wedge \neg solitary \wedge \neg spots \wedge \neg stripes \wedge \neg swims \wedge \neg tunnels \wedge \neg water \wedge \neg weak) \vee (active \wedge big \wedge black \wedge bulbous \wedge chewteeth \wedge furry \wedge grazer \wedge ground \wedge hooves \wedge horns \wedge inactive \wedge longleg \wedge longneck \wedge muscle \wedge oldworld \wedge patches \wedge quadrupedal \wedge slow \wedge smelly \wedge strong \wedge tail \wedge timid \wedge toughskin \wedge walks \wedge white \wedge \neg agility \wedge \neg arctic \wedge \neg buckteeth \wedge \neg bush \wedge \neg claws \wedge \neg coastal \wedge \neg fast \wedge \neg fierce \wedge \neg fish \wedge \neg flippers \wedge \neg forager \wedge \neg forest \wedge \neg gray \wedge \neg hairless \wedge \neg hibernate \wedge \neg hunter \wedge \neg jungle \wedge \neg lean \wedge \neg meat \wedge \neg meatteeth \wedge \neg mountains \wedge \neg nestspot \wedge \neg nocturnal \wedge \neg ocean \wedge \neg pads \wedge \neg paws \wedge \neg small \wedge \neg smart \wedge \neg solitary \wedge \neg spots \wedge \neg stripes \wedge \neg swims \wedge \neg tunnels \wedge \neg water \wedge \neg weak) \vee (active \wedge big \wedge black \wedge bulbous \wedge chewteeth \wedge furry \wedge grazer \wedge ground \wedge hooves \wedge horns \wedge inactive \wedge longleg \wedge longneck \wedge muscle \wedge oldworld \wedge patches \wedge quadrupedal \wedge slow \wedge smelly \wedge strong \wedge tail \wedge timid \wedge toughskin \wedge walks \wedge white \wedge \neg agility \wedge \neg arctic \wedge \neg buckteeth \wedge \neg bush \wedge \neg claws \wedge \neg coastal \wedge \neg fast \wedge \neg fierce \wedge \neg fish \wedge \neg flippers \wedge \neg forager \wedge \neg forest \wedge \neg gray \wedge \neg hairless \wedge \neg hibernate \wedge \neg hunter \wedge \neg jungle \wedge \neg lean \wedge \neg meat \wedge \neg meatteeth \wedge \neg mountains \wedge \neg nestspot \wedge \neg nocturnal \wedge \neg ocean \wedge \neg pads \wedge \neg paws \wedge \neg small \wedge \neg smart \wedge \neg solitary \wedge \neg spots \wedge \neg stripes \wedge \neg swims \wedge \neg tunnels \wedge \neg water \wedge \neg weak))

Figure 18. Flexibility of FOL explanations by VIT-derived MoIE MoIE and the CBM + ELL and PCBM + ELL baselines for Awa2 dataset to classify “Horse” at inference. Both the baseline’s FOL constitutes identical concepts to distinguish all the samples. However, expert4 classifies “Horse” with *smelly* as the identifying concept for the instances covered by it. Similarly expert5 classifies the same “Horse” using *longneck* and *fields*. We highlight the shared concepts between the experts and the baselines as blue.

1430
1431
1432
1433
1434
1435 **Summary statistics of the no. of concepts used**
1436 **by an expert for explanation for CUB (Resnet-101)**

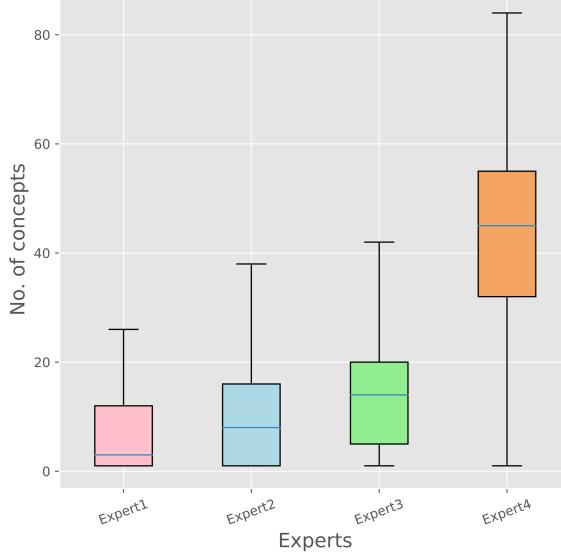


1437
1438
1439
1440
1441
1442 **Summary statistics of the no. of concepts used**
1443 **by an expert for explanation for CUB (ViT)**

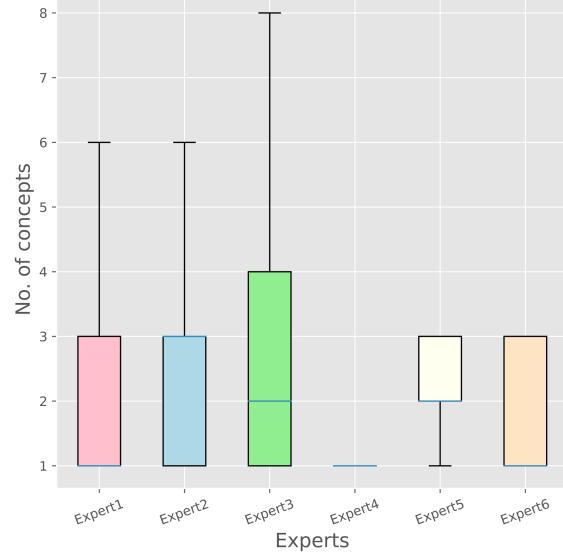


(a)

1455 **Summary statistics of the no. of concepts used**
1456 **by an expert for explanation for AWA2 (ResNet-101)**



1457 **Summary statistics of the no. of concepts used**
1458 **by an expert for explanation for Awa2 (ViT)**



(b)

1477 **Figure 19.** Summary statistics of the number of concepts utilized by various experts of datasets (a) CUB -200 (top row) and (b) Awa2 (bottom row). In general, we can see that experts carving out the explanations from ViT often uses less number of concepts.

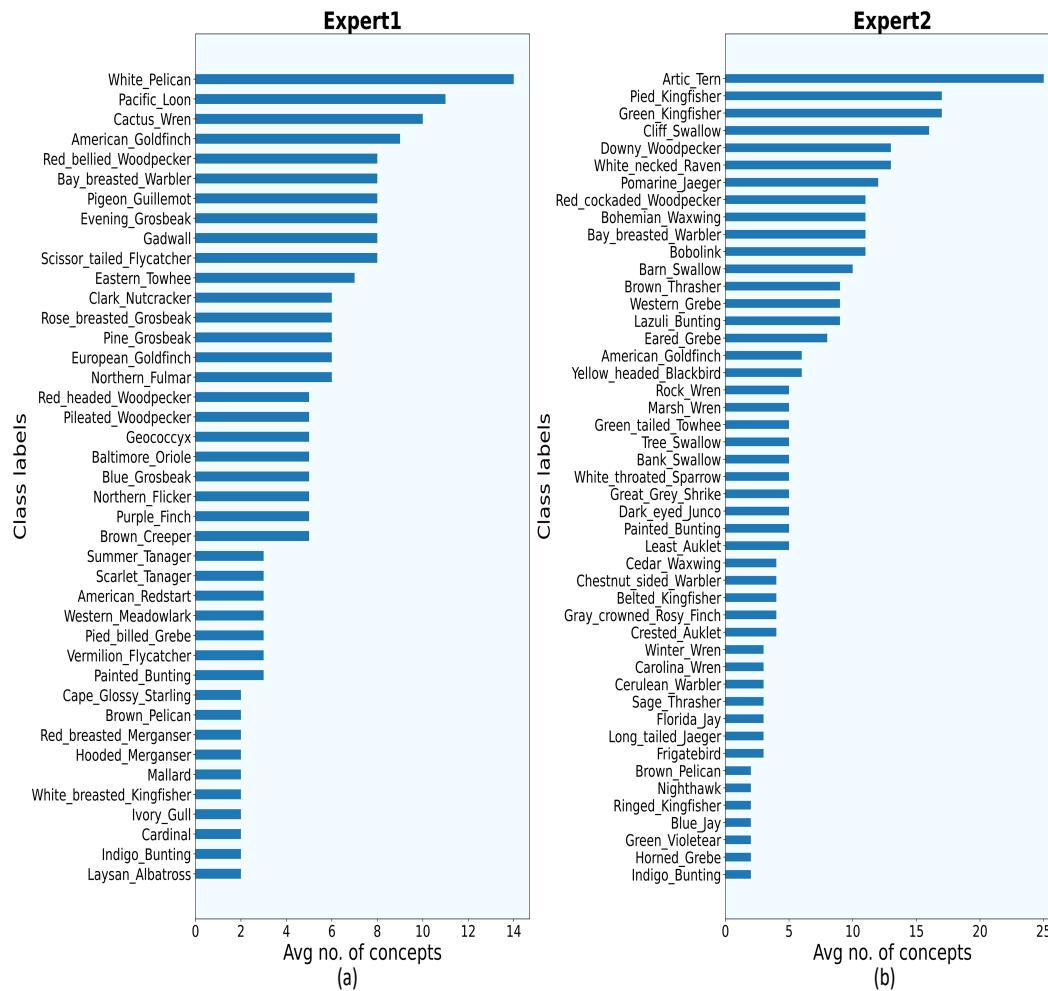


Figure 20. Class labels (Bird species) vs avg concepts using VIT as backbone for CUB-200 by (a) Expert1 (b) Expert2. Each bar in this plot indicates the average number concepts required to explain each sample of that bird species correctly. For example according to (a) expert1 requires 14 concepts to explain an instance of “White Pelican”.

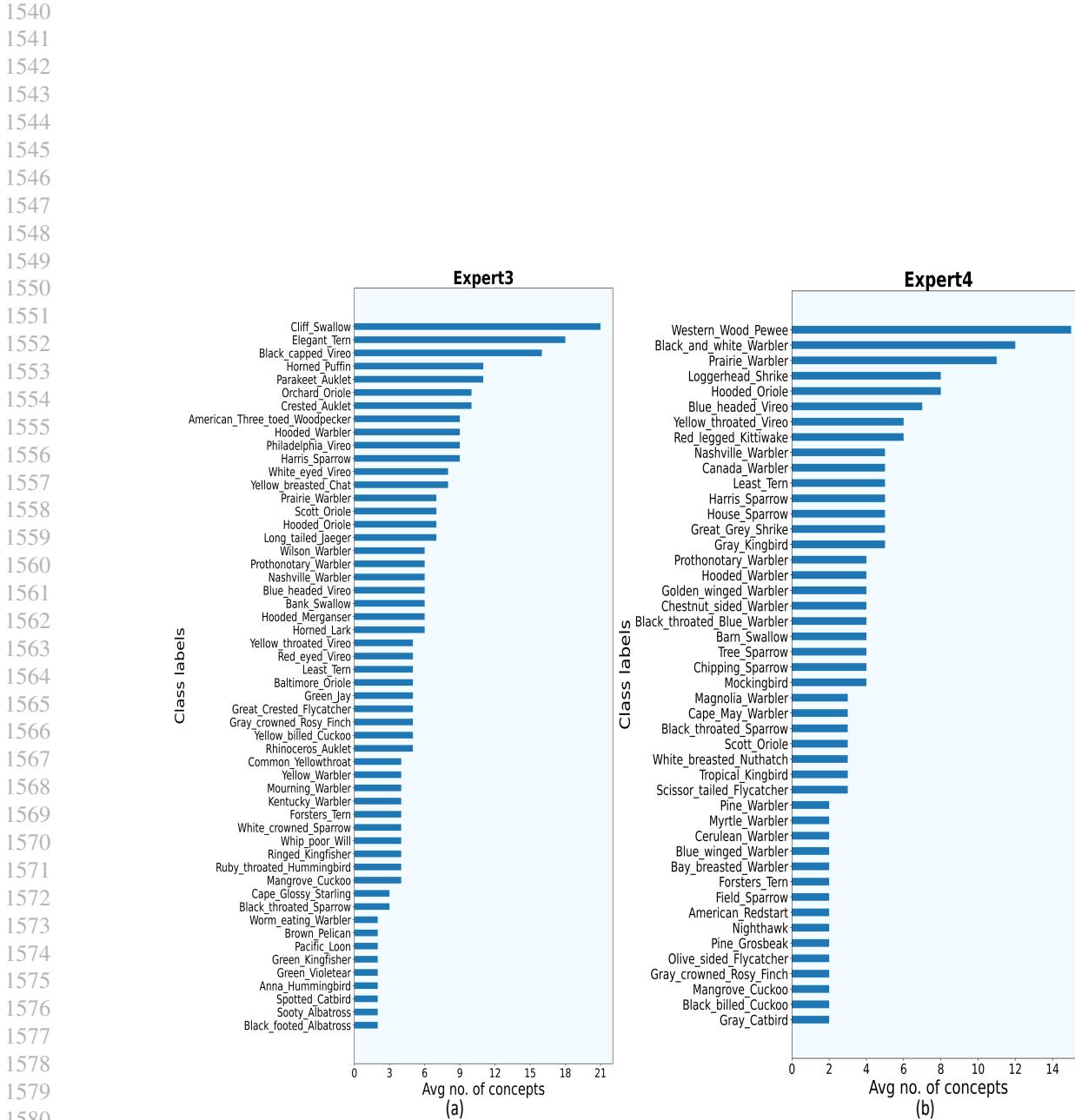


Figure 21. Class labels (Bird species) vs avg concepts using ViT as backbone for CUB-200 by (a) Expert3 (b) Expert4. Each bar in this plot indicates the average number concepts required to explain each sample of that bird species correctly. For example according to (a) expert3 requires 21 concepts to explain an instance of “Cliff Swallow”.

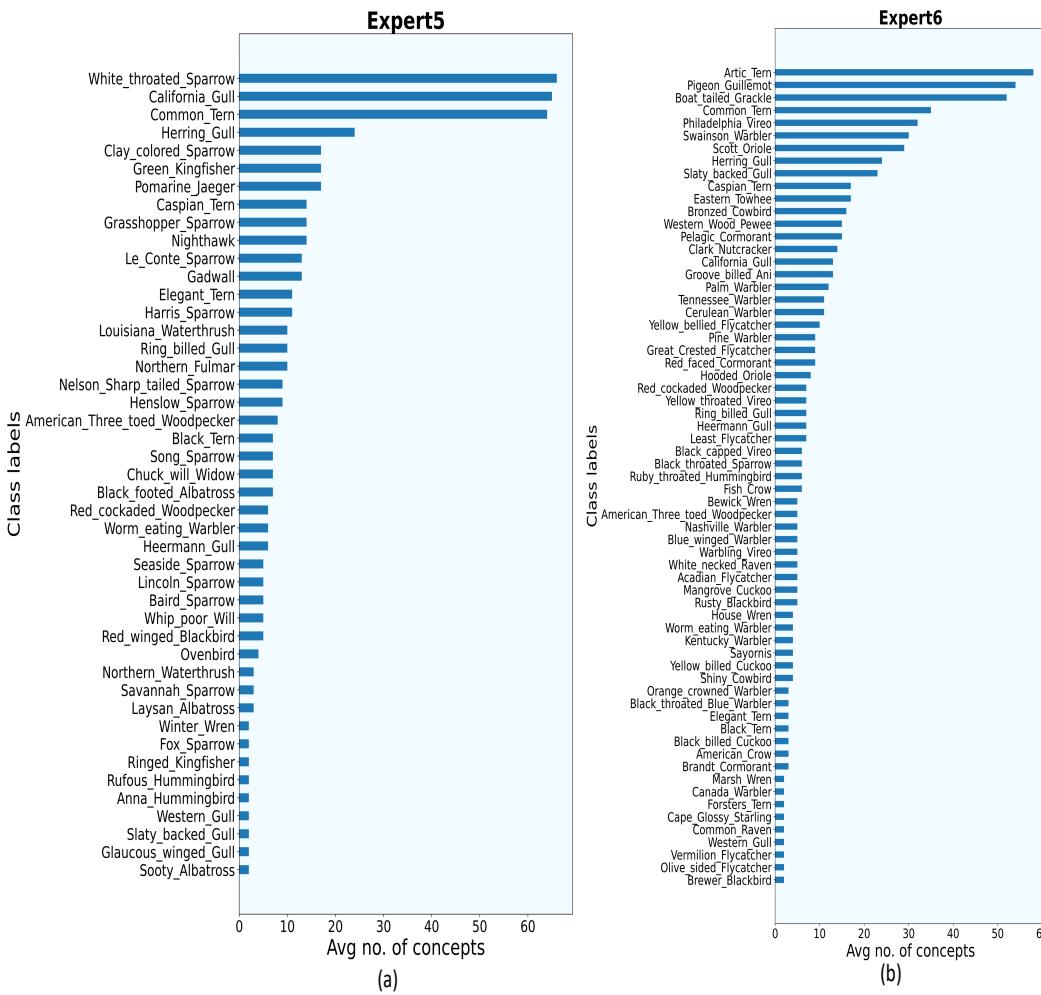


Figure 22. Class labels (Bird species) vs avg concepts using ViT as backbone for CUB-200 by (a) Expert5 (b) Expert6. Each bar in this plot indicates the average number concepts required to explain each sample of that bird species correctly. For example according to (a) expert5 requires approximately 65 concepts to explain an instance of “White throated sparrow”.

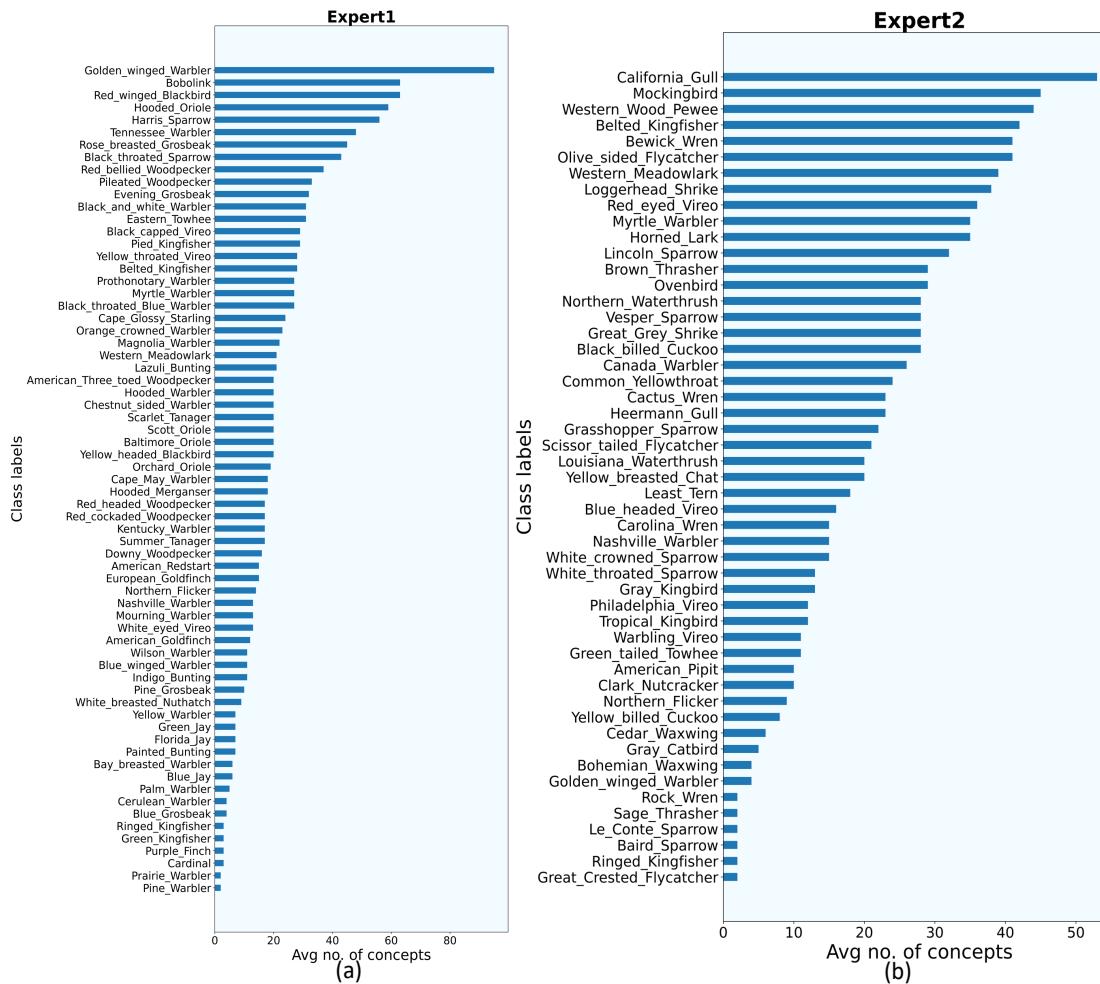


Figure 23. Class labels (Bird species) vs avg concepts using ResNet-101 as backbone for CUB-200 by (a) Expert1 (b) Expert2. Each bar in this plot indicates the average number concepts required to explain each sample of that bird species correctly. For example according to (a) expert1 requires approximately 85 concepts to explain an instance of “Golden winged warbler”.

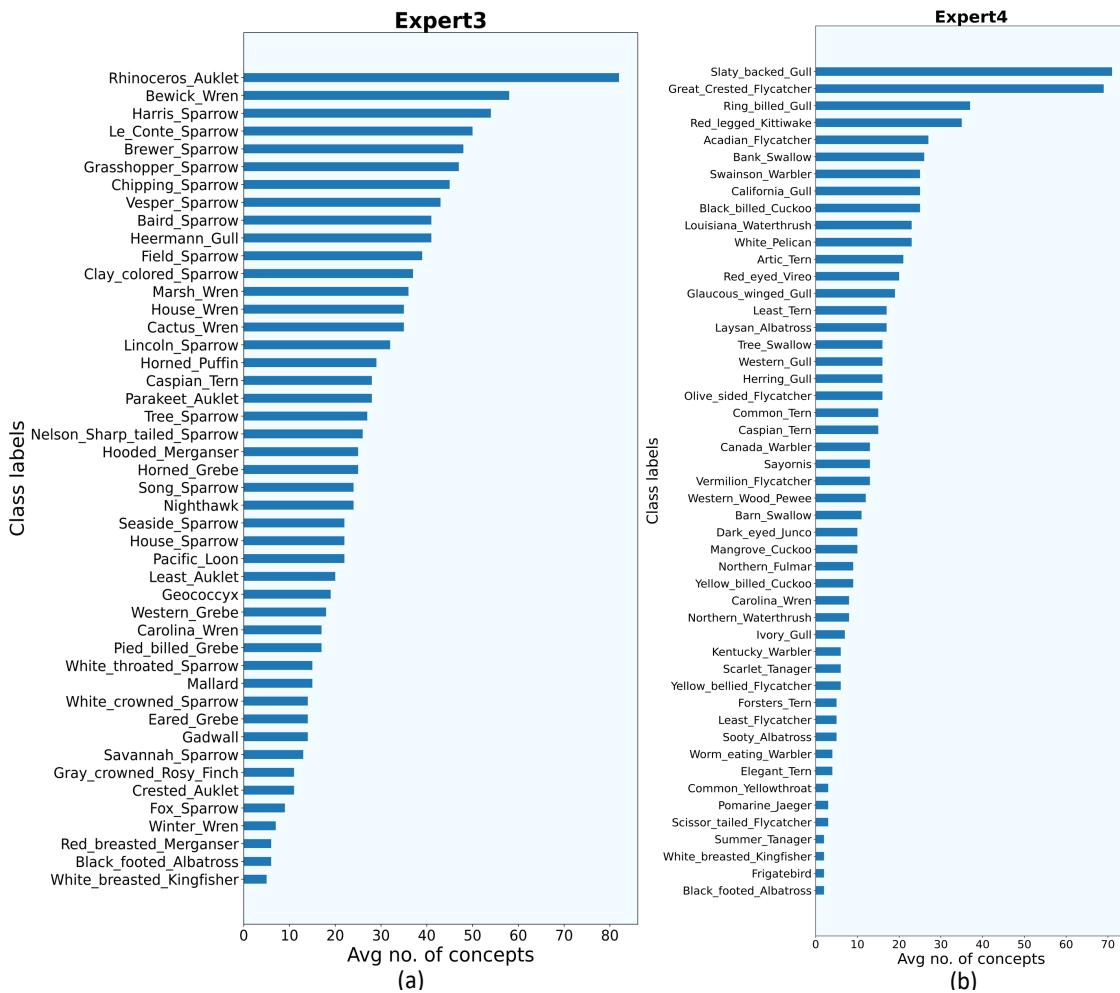


Figure 24. Class labels (Bird species) vs avg concepts using ResNet-101 as backbone for CUB-200 by (a) Expert3 (b) Expert4. Each bar in this plot indicates the average number concepts required to explain each sample of that bird species correctly. For example according to (a) expert3 requires approximately 82 concepts to explain an instance of “Rhinoceros auklet”.

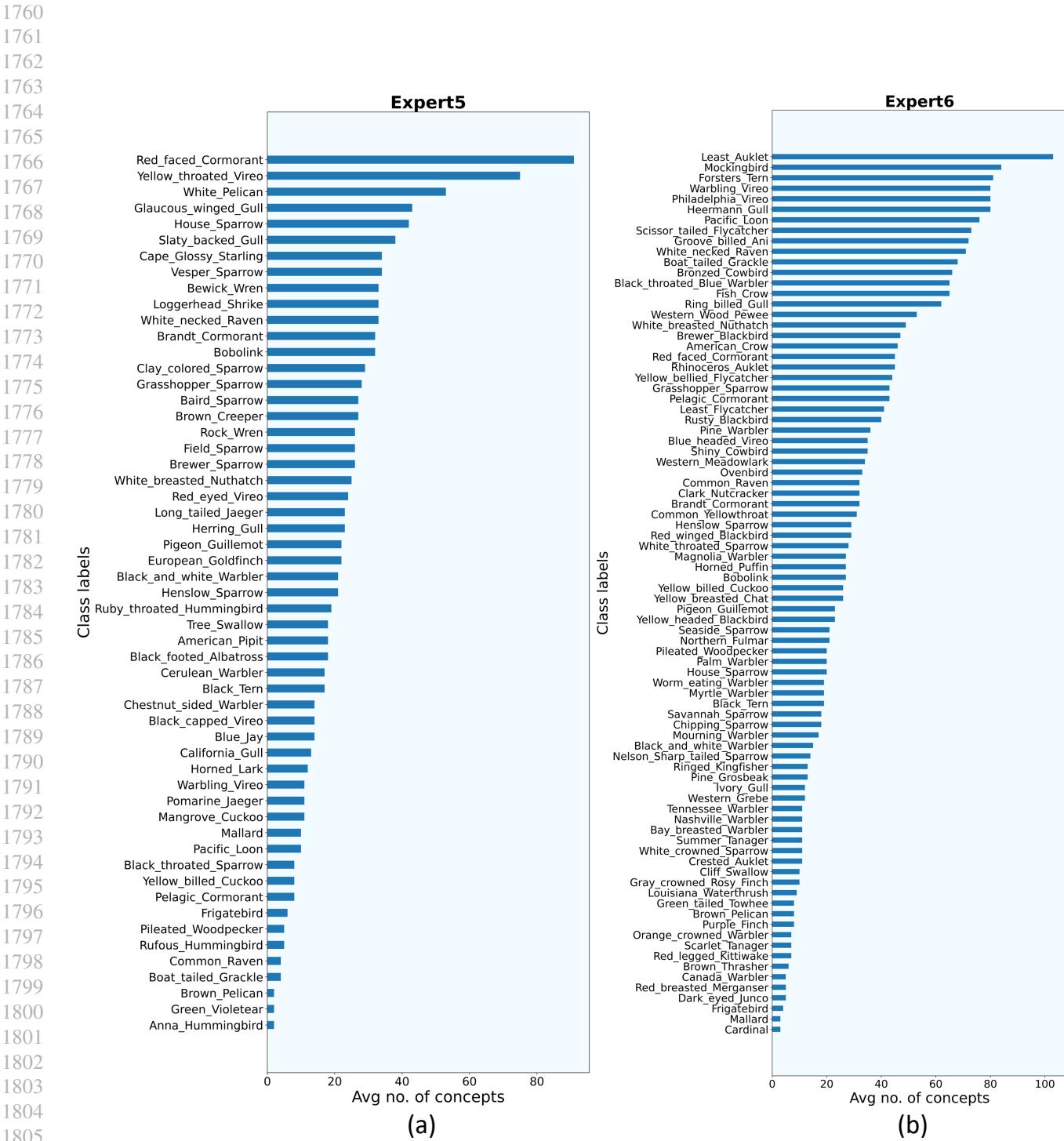


Figure 25. Class labels (Bird species) vs avg concepts using ResNet-101 as backbone for CUB-200 by (a) Expert5 (b) Expert6. Each bar in this plot indicates the average number concepts required to explain each sample of that bird species correctly. For example according to (a) expert5 requires approximately 85 concepts to explain an instance of "Red faced cormorant".

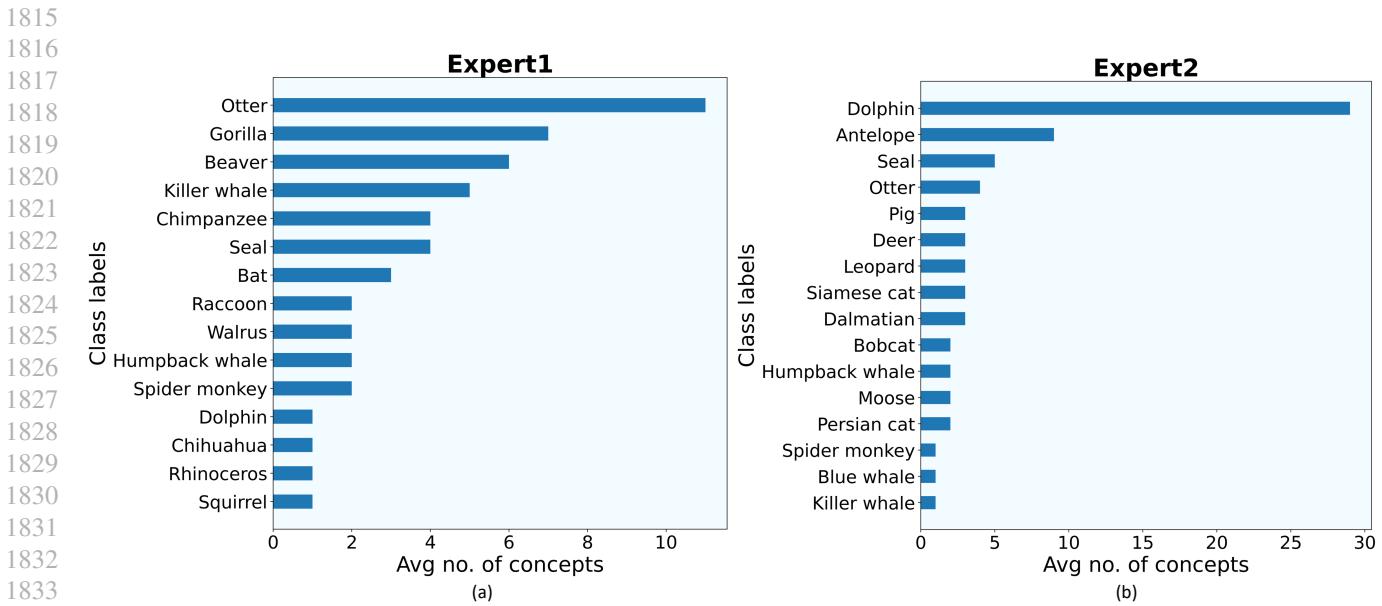


Figure 26. Class labels (Animal species) vs avg concepts using ViT as backbone for Awa2. Each bar in this plot indicates the average number concepts required to explain each sample of that animal species correctly. For example according to (c) expert1 requires approximately 12 concepts to explain an instance of “Otter”.

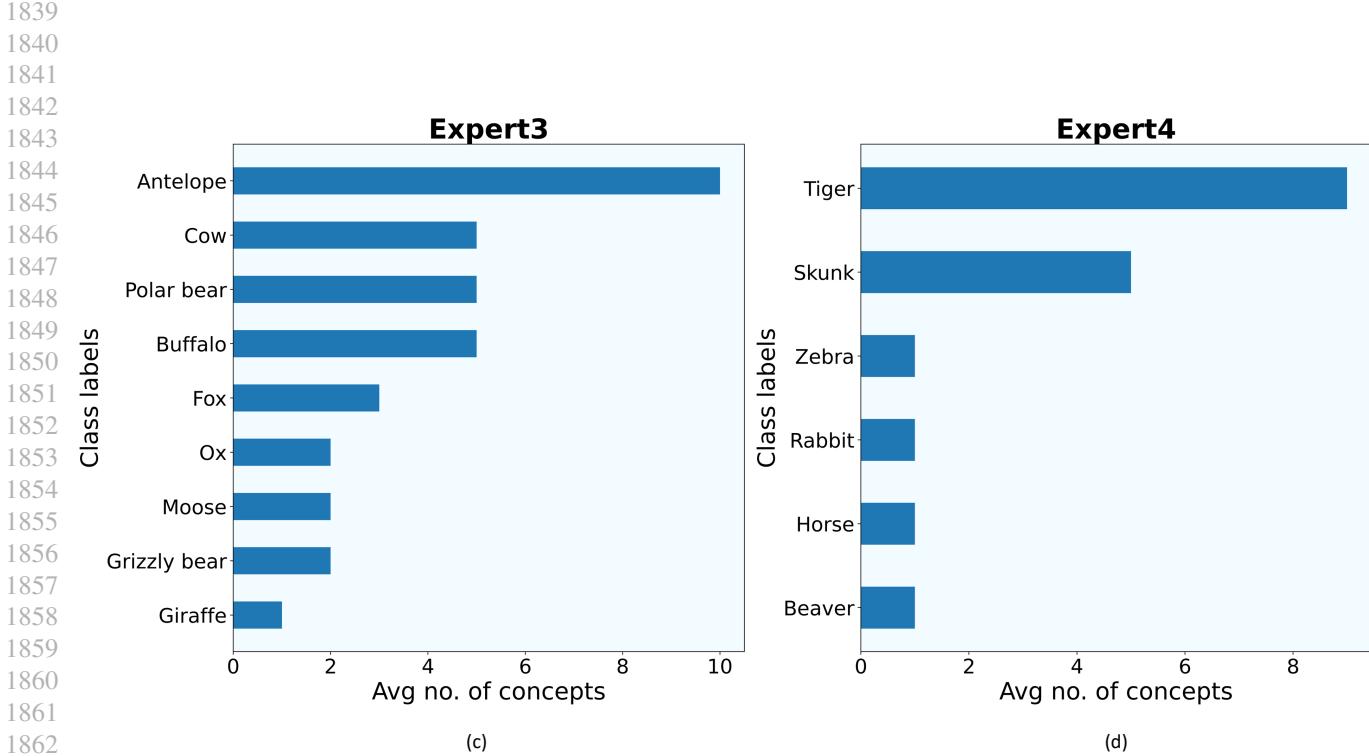


Figure 27. Class labels (Animal species) vs avg concepts using ViT as backbone for Awa2. Each bar in this plot indicates the average number concepts required to explain each sample of that animal species correctly. For example according to (c) expert3 requires approximately 10 concepts to explain an instance of “Antelope”.

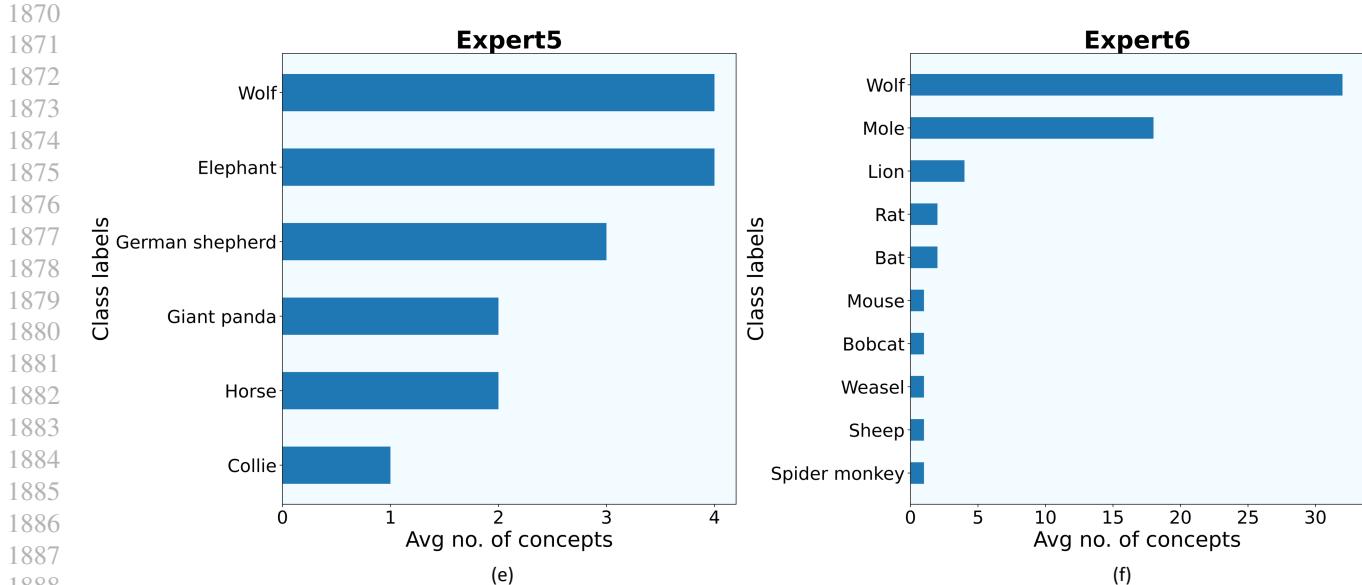


Figure 28. Class labels (Animal species) vs avg concepts using ViT as backbone for Awa2. Each bar in this plot indicates the average number concepts required to explain each sample of that animal species correctly. For example according to (e) expert5 requires approximately 4 concepts to explain an instance of “Antelope”.

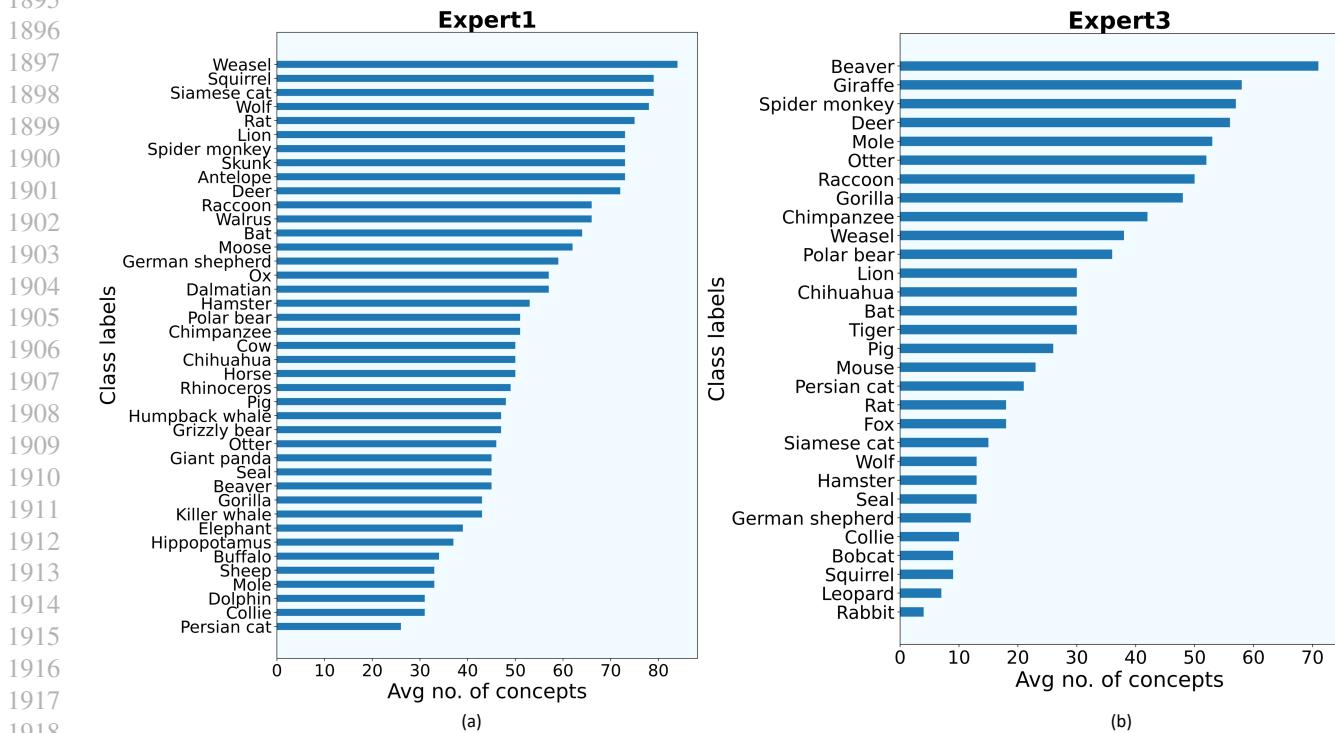


Figure 29. Class labels (Animal species) vs avg concepts using ResNet-101 as backbone for Awa2. Each bar in this plot indicates the average number concepts required to explain each sample of that animal species correctly. For example according to (a) expert1 requires approximately 80 concepts to explain an instance of “Weasel”.

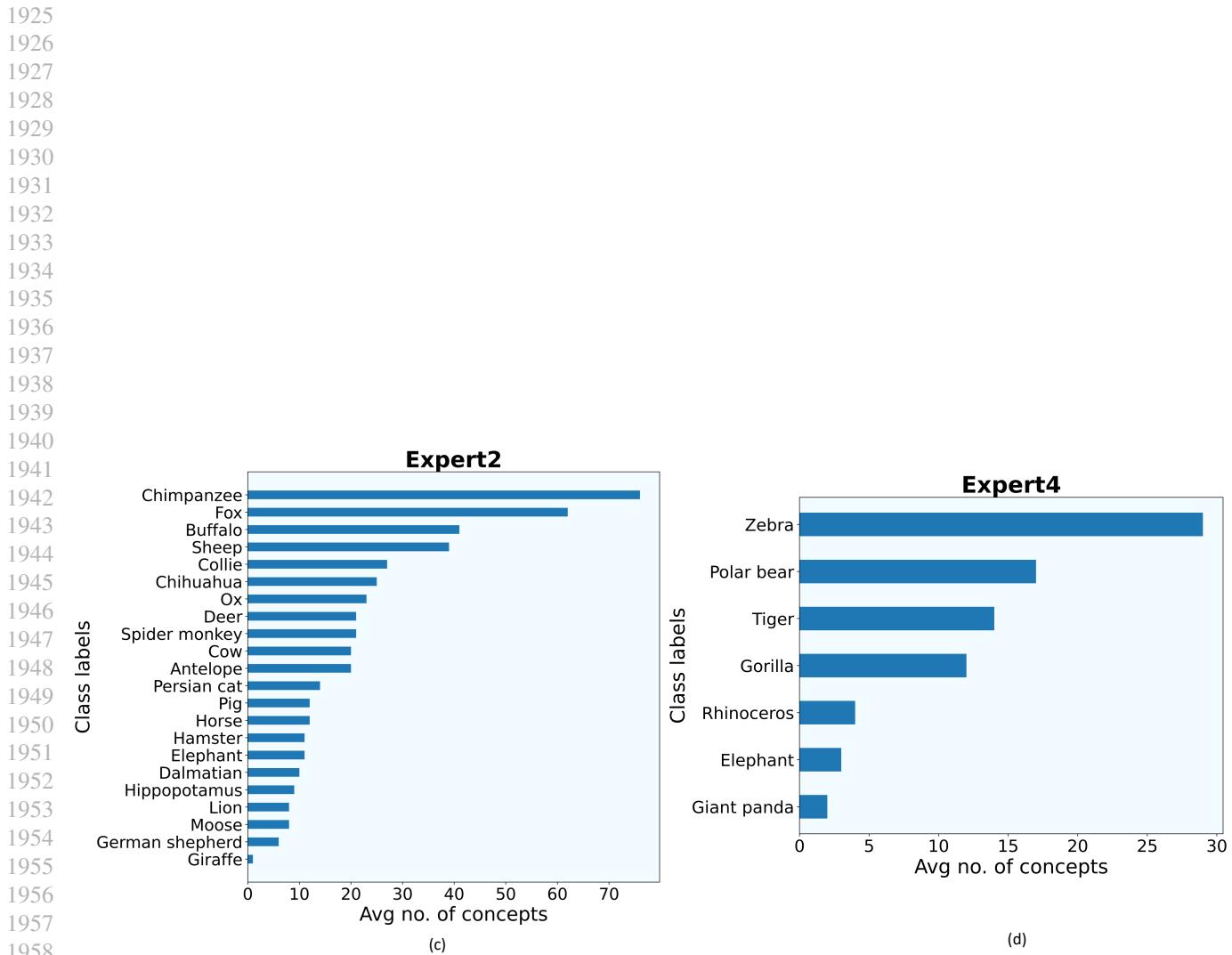


Figure 30. Class labels (Animal species) vs avg concepts using ResNet-101 as backbone for Awa2. Each bar in this plot indicates the average number concepts required to explain each sample of that animal species correctly. For example according to (b) expert2 requires approximately 72 concepts to explain an instance of “Chimpanzee”.