

# Now You See Me (CME): Concept-based Model Extraction

Dmitry Kazhdan<sup>a,c</sup>, Botty Dimanov<sup>a,c</sup>, Mateja Jamnik<sup>a</sup>, Pietro Liò<sup>a</sup> and Adrian Weller<sup>a,b</sup>

<sup>a</sup>The University of Cambridge, UK

<sup>b</sup>The Alan Turing Institute, London, UK

<sup>c</sup>Denotes equal contribution

## Abstract

Deep Neural Networks (DNNs) have achieved remarkable performance on a range of tasks. A key step to further empowering DNN-based approaches is improving their explainability. In this work we present CME: a concept-based model extraction framework, used for analysing DNN models via concept-based extracted models. Using two case studies (dSprites, and Caltech UCSD Birds), we demonstrate how CME can be used to (i) analyse the concept information learned by a DNN model (ii) analyse how a DNN uses this concept information when predicting output labels (iii) identify key concept information that can further improve DNN predictive performance (for one of the case studies, we showed how model accuracy can be improved by over 14%, using only 30% of the available concepts).

## Keywords

interpretability, concept extraction, concept-based explanations, model extraction, latent space analysis, xai

## 1. Introduction

The black-box nature of Deep Neural Networks (DNNs) hinders their widespread adoption, especially in industries under heavy regulation with high-cost of error [1]. As a result, there has recently been a dramatic increase in research on Explainable AI (XAI), focusing on improving explainability of DL systems [2, 3].

Currently, the most widely used XAI methods are feature importance methods (also referred to as saliency methods) [4]. For a given data point, these methods provide scores showing the importance of each feature (e.g., pixel, patch, or word vector) to the algorithm’s decision. Unfortunately, feature importance methods have been shown to be fragile to input perturbations [5, 6] or model parameter perturbations [7, 8]. Human experiments also demonstrate that feature importance explanations do not necessarily increase human understanding, trust, or ability to correct mistakes in a model [9, 10].

As a consequence, two other types of XAI approaches are receiving increasing attention: *model extraction* approaches, and *concept-based* explanation approaches. Model extraction methods (also referred to as model translation methods) approximate black-box models with simpler models to increase model explainability.

Concept-based explanation approaches provide model explanations in terms of human-understandable units, rather than individual features, pixels, or characters (e.g., the concepts of a *wheel* and a *door* are important for the detection of cars) [10, 11, 12].

In this paper we introduce CME<sup>1</sup>: a (C)oncept-based (M)odel (E)xtraction framework<sup>2</sup>. Figure 1 depicts how CME can be used to analyse DNN models via explainable concept-based extracted models, in order to explain and improve performance of DNNs, as well as to extract useful knowledge from them. Although this example focuses on a CNN model, CME is model-agnostic, and can be applied to *any* DNN architecture.

In particular, we make the following contributions:

- We present the novel CME framework, capable of analysing DNN models via concept-based extracted models
- We demonstrate, using two case-studies, how CME can analyse (both quantitatively and qualitatively) the concept information a DNN model has learned, and how this information is represented across the DNN layers
- We propose a novel metric for evaluating the quality of concept extraction methods
- We demonstrate, using two case-studies, how CME can analyse (both quantitatively and qualitatively) how a DNN uses concept information when predicting output labels
- We demonstrate how CME can identify key concept information that can further improve DNN predictive performance

*Title of the Proceedings: "Proceedings of the CIKM 2020 Workshops"*

*Editors of the Proceedings: Stefan Conrad, Ilaria Tiddi*

EMAIL: dk525@cam.ac.uk (D. Kazhdan); btd26@cam.ac.uk (B. Dimanov); mateja.jamnik@cl.cam.ac.uk (M. Jamnik);

pietro.lio@cl.cam.ac.uk (P. Liò); adrian.weller@eng.cam.ac.uk (A. Weller)

ORCID:



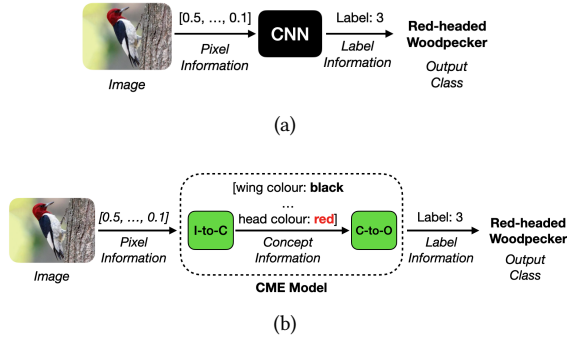
© 2020 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).



CEUR Workshop Proceedings (CEUR-WS.org)

<sup>1</sup>Pronounced "See Me."

<sup>2</sup>All relevant code is available at <https://github.com/dmitrykazhdan/CME>



**Figure 1:** CME extracted model example. (a) Given an input image, a CNN uses the image’s pixel information as input, and returns class information as output (in this case, class label 3, corresponding to the *Red-headed Woodpecker* class), performing data processing in a non-explainable, black-box fashion. (b) Given an input image, a CME extracted model uses an *Input-to-Concept* function (*I-to-C*) to compute *concept information* from the pixel data (e.g. bird wing color, or head color values). Next, the model uses a *Concept-to-Output* function (*C-to-O*) to compute the output class label from this concept information.

## 2. Related Work

### 2.1. Concept-based Explanations

Concept-based explanations have been used in a wide range of different ways, including: inspecting what a model has learned [12, 13], providing class-specific explanations [14, 10], and discovering causal relations between concepts [15]. Similarly to CME, these approaches typically seek to explain model behaviour in terms of high-level concepts, extracting this concept information from a model’s latent space.

Importantly, existing concept-based explanation approaches are typically capable of handling binary-valued concepts only, which implies that multi-valued concepts have to be binarised first. For instance, given a concept such as “shape”, with possible values ‘square’ and ‘circle’, these approaches have to convert “shape” into two binary concepts ‘is\_square’, and ‘is\_circle’. This makes such approaches (i) computationally expensive, since the binarised concept space usually has a high cardinality, (ii) error-prone, since mutual exclusivity of concept values is now not enforced (e.g., a single data point can now have both ‘is\_square’ and ‘is\_circle’ concepts being true). In contrast, our approach is capable of handling multi-valued concepts directly, without binarisation.

Furthermore, concept-based explanation approaches typically rely on the latent space of a single layer when

extracting concept information. DNNs have been shown to perform hierarchical feature extraction, with layers closer to the output utilising higher-level data representations, compared to layers closer to the input [16, 17]. This implies that choosing a single layer imposes an *unnecessary* trade-off between low- and high-level concepts. On the other hand, CME is capable of efficiently combining latent space information from multiple layers, thereby avoiding this constraint.

Finally, existing methods typically represent concept explanations as a list of concepts, with their relative importance with respect to the classification task. In contrast, our approach describes the functional relationship between concepts and outputs, thereby showing in more detail how the model utilises concept information when making predictions.

### 2.2. Concept Bottleneck Models

Recent work on concept-based explanations relies on models that use an intermediate concept-based representation when making predictions [18, 19]. Work in [18] refer to these types of models as *concept bottleneck models* (CBMs). A concept bottleneck model is a model which, given an input, first predicts an intermediate set of human-specified concepts, and then uses only this concept information to predict the output task label. Work in [18] proposes a method for turning any DNN into a concept bottleneck model given concept annotations at training time. This is achieved by resizing one of the layers to match the number of concepts provided, and re-training with an added intermediate loss that encourages the neurons in that layer to align component-wise to the provided concepts.

Crucially, CBM approaches provide ways for *generating DNN models*, which are explicitly encouraged to rely on specified concept information. In contrast, our approach is used for *analysing DNN models* (and is much cheaper computationally).

Furthermore, CBM approaches require concept annotations to be available at training time for *all* of the training data, which is often expensive to produce. In contrast, CME can be used with *partially-labelled* datasets in a semi-supervised fashion, as will be described in Section 3.

Finally, CBM approaches require the concepts themselves to be known beforehand. On the other hand, CME can efficiently utilise knowledge contained in pre-trained DNNs, in order to *learn* about which concepts are/aren’t required for a given task. Further details on CME/CBM comparison can be found in Appendix A.

### 2.3. Model Extraction

Model extraction techniques use rules [20, 21, 22], decision trees [23, 24], or other more readily explainable models [25] to approximate complex models, in order to study their behaviour. Provided the approximation quality (referred to as *fidelity*) is high enough, an extracted model can preserve many statistical properties of the original model, while remaining open to interpretation.

However, extracted models generated by existing methods represent their decision-making using the same input representation as the original model, which is typically difficult for the user to understand directly. Instead, our extracted models represent decision-making via human-understandable concepts, making them easier to interpret.

## 3. Methodology

In this section we present our CME approach, describing how it can be used to analyse DNN models using concept-based extracted models.

### 3.1. Formulation

We consider a pre-trained DNN classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , ( $\mathcal{X} \subset \mathbb{R}^n$ ,  $\mathcal{Y} \subset \mathbb{R}^0$ ), where  $f(\mathbf{x}) = y$  is mapping an input  $\mathbf{x} \in \mathcal{X}$  to an output class  $y \in \mathcal{Y}$ . For every DNN layer  $l$ , we denote the function  $f^l : \mathcal{X} \rightarrow \mathcal{H}^l$ , ( $\mathcal{H}^l \subset \mathbb{R}^m$ ) as a mapping from the input space  $\mathcal{X}$  to the hidden representation space  $\mathcal{H}^l$ , where  $m$  denotes the number of hidden units, and can be different for each layer.

Similarly to [18, 19], we assume the existence of a *concept representation*  $\mathcal{C} \subset \mathbb{R}^k$ , defining  $k$  distinct concepts associated with the input data.  $\mathcal{C}$  is defined such that every basis vector in  $\mathcal{C}$  spans the space of possible values for one particular concept. We further assume the existence of a function  $p^* : \mathcal{X} \rightarrow \mathcal{C}$ , where  $p^*(\mathbf{x}) = \mathbf{c}$  is mapping an input  $\mathbf{x}$  to its concept representation  $\mathbf{c}$ . Thus,  $p^*$  defines the concepts and their values (referred to as the *ground truth concepts*) for every input point.

### 3.2. CME

In this work, we define a DNN  $f$  as being *concept-decomposable*, if it can be well-approximated by a composition of functions  $p$  and  $q$ , such that  $f(\mathbf{x}) = q(p(\mathbf{x}))$ . In this definition, the function  $p : \mathcal{X} \rightarrow \mathcal{C}$  is an *input-to-concept* function, mapping data-points from their input representation  $\mathbf{x} \in \mathcal{X}$  to their concept representation  $\mathbf{c} \in \mathcal{C}$ . The function  $q : \mathcal{C} \rightarrow \mathcal{Y}$  is a *concept-to-*

*output* function, mapping data-points in their concept representation  $\mathcal{C}$  to output space  $\mathcal{Y}$ . Thus, when processing an input  $\mathbf{x}$ , a DNN  $f$  can be seen as converting this input into an interpretable concept representation using  $p$ , and using  $q$  to predict the output from this representation. The significance of this decomposition is further discussed in Appendix A.

CME explores whether a given DNN  $f$  is concept-decomposable, by attempting to approximate  $f$  with an extracted model  $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$ . In this case,  $\hat{f}$  is defined as  $\hat{f}(\mathbf{x}) = \hat{q}(\hat{p}(\mathbf{x}))$ , using input-to-concept  $\hat{p}$  and output-to-concept  $\hat{q}$  extracted by CME from the original DNN. We describe our approach to extracting  $\hat{p}$  and  $\hat{q}$  in the remainder of this section.

### 3.3. Input-to-Concept ( $\hat{p}$ )

When extracting  $\hat{p}$  from a pre-trained DNN, we assume we have access to the DNN training data and labels  $\{(\mathbf{x}^{(0)}, y^{(0)}), \dots, (\mathbf{x}^{(d)}, y^{(d)})\}$ . Furthermore, we assume partial access to  $p^*$ , such that a small set of  $i$  training points  $\{\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(i-1)}\}$  have concept labels  $\{\mathbf{c}^{(0)}, \dots, \mathbf{c}^{(i-1)}\}$  associated with them, while the remaining  $u$  points  $\{\mathbf{x}^{(i)}, \dots, \mathbf{x}^{(i+u)}\}$  do not (in this case  $u = d - i$ ). We refer to these subsets respectively as the *concept labelled dataset* and *concept unlabelled dataset*. Using these datasets, we generate  $\hat{p}$  by aggregating concept label predictions across multiple layers of the given DNN model, as described below.

Given a DNN layer  $l$  with  $m$  hidden units, we compute the layer's representation of the input data  $\mathbf{h} = f^l(\mathbf{x})$ , obtaining  $(\mathbf{h}^{(0)}, \dots, \mathbf{h}^{(i+u)})$ . Using this data and the concept labels, we construct a semi-supervised dataset, consisting of labelled data  $\{(\mathbf{h}^{(0)}, \mathbf{c}^{(0)}), \dots, (\mathbf{h}^{(i-1)}, \mathbf{c}^{(i-1)})\}$ , and unlabelled data  $\{\mathbf{h}^{(i)}, \dots, \mathbf{h}^{(i+u)}\}$ .

Next, we rely on Semi-Supervised Multi-Task Learning (SSMTL) [26], in order to extract a function  $g^l : \mathcal{H}^l \rightarrow \mathcal{C}$ , which predicts concept labels from layer  $l$ 's hidden space. In this work, we treat each concept as a separate, independent task. Hence,  $g^l(\mathbf{h})$  is decomposed into  $k$  separate tasks (one per concept), and is defined as  $g^l(\mathbf{h}) = (g_1^l(\mathbf{h}), \dots, g_k^l(\mathbf{h}))$  where each  $g_i^l(\mathbf{h})$  ( $i \in \{1..k\}$ ) predicts the value of concept  $i$  from  $\mathbf{h}$ .

Repeating this process for all model layers  $L$ , we obtain a set of functions  $G = \{g_i^l \mid l \in \{1..L\} \wedge i \in \{1..k\}\}$ . For every concept  $i$ , we define the "best" layer  $l^i$  for predicting that concept as shown in (1):

$$l^i = \arg \min_{l \in L} \ell(g_i^l, i) \quad (1)$$

Here,  $\ell$  is a loss function (in this case the error rate), computing the predictive loss of function  $g_i^l$  with re-

spect to a concept  $i$ . Finally, we define  $\hat{p}$  as shown in (2):

$$\hat{p}(\mathbf{x}) = (g_1^{l_1} \circ f^{l_1}(\mathbf{x}), \dots, g_k^{l_k} \circ f^{l_k}(\mathbf{x})) \quad (2)$$

Thus, given an input  $\mathbf{x}$ , the value computed by  $\hat{p}(\mathbf{x})$  for every concept  $i \in \{1..k\}$  is equal to the value computed by  $g_i^{l_i}$  from that input's representation in layer  $l_i$ . Overall,  $\hat{p}$  encapsulates concept information contained in a given DNN model, and can be used to analyse how this information is represented, as well as to predict concept values for new inputs.

### 3.4. Concept-to-Label ( $\hat{q}$ )

We setup extraction of  $\hat{q}$  as a classification problem, in which we train  $\hat{q}$  to predict output labels  $y$  from concept labels  $c$  predicted by  $\hat{p}$ . We use  $\hat{p}$  to generate concept labels for all training data points, obtaining a set of concept labels  $\{c^{(0)}, \dots, c^{(i+u)}\}$ . Next, we produce a labelled dataset, consisting of concept labels and corresponding DNN output labels  $\{(c^{(0)}, y^{(0)}), \dots, (c^{(i+u)}, y^{(i+u)})\}$ , and use it to train  $\hat{q}$  in a supervised manner. We experimented with using Decision Trees (DTs), and Logistic Regression (LR) models for representing  $\hat{q}$ , as will be discussed in Section 5. Overall,  $\hat{q}$  can be used to analyse how a DNN uses concept information when making predictions.

## 4. Experimental Setup

We evaluated CME using two datasets: dSprites [27], and Caltech-UCSD birds [28]. All relevant code is publicly available at<sup>3</sup>.

### 4.1. dSprites

dSprites is a well-established dataset used for evaluating unsupervised latent factor disentanglement approaches. dSprites consists of 2D 64×64 pixel black-and-white shape images, procedurally generated from all possible combinations of 6 ground truth independent concepts (color, shape, scale, rotation, x and y position). Further details can be found in Appendix B, and the official dSprites repository.<sup>4</sup>

#### 4.1.1. Classification Tasks

We define 2 classification tasks, used to evaluate our framework:

- **Task 1:** This task consists of determining the shape concept value from an input image. For every image sample, we define its task label as the shape concept label of that sample.
- **Task 2:** This task consists of discriminating between all possible *shape* and *scale* concept value combinations. We assign a distinct identifier to each possible combination of the shape and scale concept labels. For every image sample, we define its task label as the identifier corresponding to this sample's shape and scale concept values.

Overall, Task 1 explores a scenario in which a DNN has to learn to recognise a specific concept from an input image. Task 2 explores a relatively more complex scenario, in which a DNN has to learn to recognise combinations of concepts from an input image.

#### 4.1.2. Model

We trained a Convolutional Neural Network (CNN) model [29] for each task. Both models had the same architecture, consisting of 3 convolutional layers, 2 dense layers with ReLUs, 50% dropout [30] and a softmax output layer. The models were trained using categorical cross-entropy loss, and achieved  $100.0 \pm 0.0\%$  classification accuracies on their respective held-out test sets. We refer to these models as the *Task 1 model* and the *Task 2 model* in the rest of this work.

#### 4.1.3. Ground-truth Concept Information

Importantly, the task and dataset definitions described in this section imply that we know precisely which concepts the models had to learn, in order to achieve  $100.0 \pm 0.0\%$  task performances (*shape* for Task 1, and *shape* and *scale* for Task 2). We refer to this as the *ground truth* concept information learned by these models.

### 4.2. Caltech-UCSD Birds (CUB)

For our second dataset, we used Caltech-UCSD Birds 200 2011 (CUB). This dataset consists of 11,788 images of 200 bird species with every image annotated using 312 binary concept labels (e.g. beak and wing colour, shape, and pattern). We relied on concept pre-processing steps defined in [18] (used for de-noising concept annotations, and filtering out outlier concepts), which produces a refined set of  $k = 112$  binary concept labels for every image sample.

<sup>3</sup><https://github.com/dmitrykazhdan/CME>

<sup>4</sup><https://github.com/deepmind/dsprites-dataset/>



#### 4.2.1. Classification Task.

We relied on the standard CUB classification task, which consists of predicting the bird species from an input image.

#### 4.2.2. Model

We used the Inception-v3 architecture [31], pretrained on ImageNet [32] (except for the fully-connected layers) and fine-tuned end-to-end on the CUB dataset, following the preprocessing practices described in [33]. The model achieved  $82.7 \pm 0.4\%$  classification accuracy on a held-out test set. We refer to this model as the *CUB model* in the rest of this work.

#### 4.2.3. Ground-truth Concept Information

Unlike dSprites, the CUB dataset does not explicitly define how the available concepts relate to the output task. Thus, we *do not* have access to the ground truth concept information learned by the CUB model.

### 4.3. Benchmarks

We compare performance of our CME approach to two other benchmarks, described in the remainder of this section.

#### 4.3.1. Net2Vec

We rely on work in [34] for defining benchmark  $\hat{p}$  functions for the three tasks. Work in [34] attempts to predict presence/absence of concepts from spatially-averaged hidden layer activations of convolutional layers of a CNN model. Given a binary concept  $c$ , this approach trains a logistic regressor, predicting the presence/absence of this concept in an input image from the latent representation of a given CNN layer. In case of multi-valued concepts, the concept space has to be binarised, as discussed in Section 2.2. In this case, the binarised concept value with the highest likelihood is returned.

Unlike CME, [34] does not provide a way of selecting the convolutional layer to use for concept extraction. We consider the best-case scenario by selecting, for all tasks, the convolutional layers yielding the best concept extraction performance. For all tasks, these layers were convolutional layers closest to the output (the 3rd conv. layer in case of dSprites tasks, and the final inception block output layer in case of the CUB task).

#### 4.3.2. CBM

As discussed in Section 4.2.3, we do not have access to ground truth concept information learned by the CUB model. Instead, we rely on the pre-trained *sequential bottleneck model* defined in [18] (referred to as CBM in the rest of this work). CBM is a bottleneck model, obtained by resizing one of the layers of the CUB model to match the number of concepts provided (we refer to this as the *bottleneck layer*), and training the model in two steps. First, the sub-model consisting of the layers between the input layer and the bottleneck layer (inclusive) is trained to predict concept values from input data. Next, the submodel consisting of the layers between the layer following the bottleneck layer and the output layer is trained to predict task labels from the concept values predicted by the first submodel. Hence, this bottleneck model is *guaranteed* to solely rely on concept information that is learnable from the data, when making task label predictions. Thus, this benchmark serves as an *upper bound* for the concept information learnable from the dataset, and for the task performance achievable using this information. Importantly, CBM does not attempt to approximate/analyse the CUB model, but instead attempts to solve the same classification task using concept information only.

We use the first CBM submodel as a  $\hat{p}$  benchmark, representing the upper bound of concept information learnable from the data. We use the second submodel as a  $\hat{q}$  benchmark, representing the upper bound of task performance achievable from predicted concept information only. Finally, we use the entire model as an  $\hat{f}$  benchmark. We make use of the saved trained model from [18], available in their official repository<sup>5</sup>.

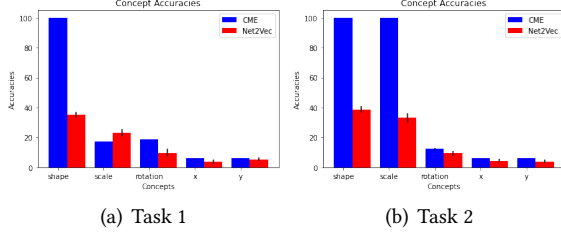
## 5. Results

We present the results obtained by evaluating our approach using the two case studies described above.

We obtain the concept labelled dataset by returning the ground-truth concept values for a random set of samples in the model training data. For dSprites, we found that a concept labelled dataset of a 100 samples or more worked well in practice for both tasks. Thus, we fix the size of the concept labelled dataset to 100 in all of the dSprites experiments. For CUB, we found that a concept labelled dataset containing 15 or more samples per class worked well in practice. Thus, we fix the size of the concept labelled dataset to 15 samples per class in all of the CUB experiments. In the future, we intend to explore the variation of model extraction performance

---

<sup>5</sup><https://github.com/yewsiang/ConceptBottleneck>



**Figure 2:** Predictive accuracy of CME and Net2Vec  $\hat{p}$  functions for all concepts

with the size of the concept labelled dataset in more detail.

## 5.1. Concept Prediction Performance

First, we evaluate the quality of  $\hat{p}$  functions produced by CME, Net2Vec, and CBM. For both dSprites tasks, we relied on the *Label Spreading* semi-supervised model [35], provided in scikit-learn [36], when learning the  $g_i^l$  functions for CME. For CUB, we used logistic regression functions instead, as they gave better performance.

### 5.1.1. dSprites

Figure 2 shows predictive performance of the  $\hat{p}$  functions on all concepts for the two dSprites tasks (averaged over 5 runs). As discussed in Section 4.1.1, we have access to the ground truth concept information learned by these models (*shape* concept information for Task 1, and *shape* and *scale* concept information for Task 2). For both tasks,  $\hat{p}$  functions extracted by CME successfully achieved high predictive accuracy on concepts relevant to the tasks, whilst achieving a low performance on concepts irrelevant to the tasks. Thus, CME was able to successfully extract the concept information contained in the task models. For both tasks,  $\hat{p}$  functions extracted by Net2Vec achieved a much lower performance on the relevant concepts.

### 5.1.2. CUB

As discussed in Section 4.2.3, the CUB dataset does not explicitly define how the concepts relate to the output task labels. Thus, we do not know how relevant/important different concepts are, with respect to task label prediction. In this section, we make the conservative assumption that all concepts are relevant, when evaluating  $\hat{p}$  functions, and explore relative concept importance in more detail in Section 5.3.

Firstly, we relied on the ‘average-per-concept’ metrics introduced in [18] when evaluating the  $\hat{p}$  function

performances, by computing their  $F1$  predictive scores for each concept, and then averaging over all concepts. We obtained  $F1$  scores of  $92 \pm 0.5\%$ ,  $86.3 \pm 2.0\%$ , and  $85.9 \pm 2.3\%$  for CBM, CME, and Net2Vec  $\hat{p}$  functions, respectively (averaged over 5 runs).

Importantly, we argue that in case of a large number of concepts, it is crucial to measure how concept mispredictions are distributed across the test samples. For instance, consider a dSprites Task 2  $\hat{p}$  function that achieves 90% predictive accuracy on both *shape* and *scale* concepts. The average predictive accuracy on relevant concepts achieved by this  $\hat{p}$  will therefore be 90%. However, if the two concepts are mis-predicted for strictly different samples (i.e. none of the samples have both *shape* and *scale* predicted incorrectly at the same time), this means that 20% of the test samples will have one relevant concept predicted incorrectly. Given that both concepts need to be predicted correctly when using them for task label prediction, this implies that consequent task label prediction will not be able to achieve over 80% task label accuracy. This effect becomes even more pronounced in case of a larger number of relevant concepts.

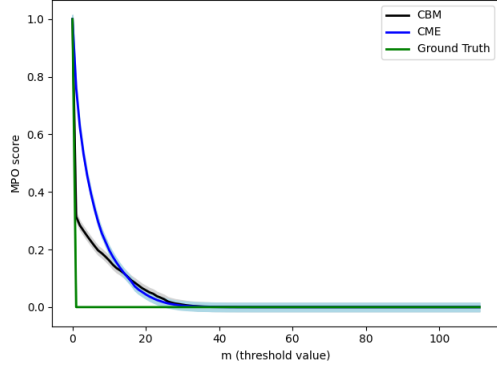
Consequently, we defined a novel cumulative misprediction error metric, which we refer to as the ‘mis-prediction-overlap’ (MPO) metric. Given a test set  $T = \{(\mathbf{x}^{(0)}, \mathbf{c}^{(0)}), \dots, (\mathbf{x}^{(n)}, \mathbf{c}^{(n)})\}$  consisting of  $n + 1$  input samples  $\mathbf{x}$  with corresponding concept labels  $\mathbf{c}$ , and a prediction set  $P = \{(\hat{\mathbf{c}}^{(0)}), \dots, \hat{\mathbf{c}}^{(n)}\}$ ,  $MPO$  computes the fraction of samples in the test set, that have *at least*  $m$  relevant concepts predicted incorrectly, as shown in Equation 3 (where  $\mathbb{I}(\cdot)$  denotes the indicator function):

$$MPO(T, P, m) = \frac{1}{n} \sum_{i=0}^n \mathbb{I}(\text{err}(\mathbf{c}_i, \hat{\mathbf{c}}_i) \geq m) \quad (3)$$

Here,  $\text{err}$  can be used to specify which concepts to measure the mis-prediction error on (i.e. in case some of the provided concepts are irrelevant). Under our assumption of all concepts being relevant, we defined  $\text{err}$  as shown in Equation 4:

$$\text{err}(\mathbf{c}_i, \hat{\mathbf{c}}_i) = \sum_{j=0}^k \mathbb{I}(c_{i,j} \neq \hat{c}_{i,j}) \quad (4)$$

Using a held-out test set, we plot the  $MPO$  metric values for  $m \in [0, \dots, 112]$ , as shown in Figure 3 (averaged over 5 runs). Importantly,  $\hat{p}$  function performances can be evaluated by observing their  $MPO$  scores for different values of  $m$ . A larger  $MPO$  score implies a bigger proportion of samples had at least  $m$  relevant concept predicted incorrectly.



**Figure 3:** Performances of  $\hat{p}$  functions, evaluated using the *MPO* metric. The green line plots the case for perfect prediction, when the predicted concepts are equivalent to the ground truth concepts (i.e. the  $p^*$  performance), in which case  $MPO = 1$  for  $m = 0$ , and  $MPO = 0$  otherwise. Net2Vec obtained values within 1% deviation from the corresponding CME values for all  $m$ , and is therefore omitted here for simplicity

Overall, CME performed almost identically to Net2Vec, and worse than CBM according to the *MPO* metric. Similar performance to Net2Vec is likely caused by (i) concepts being binary (requiring no binarisation) (ii) the Inception-v3 model having a relatively large number of convolutional layers, implying that the final convolutional layer likely learned higher-level features, relevant to concept prediction.

Importantly, *MPO* showed that both CBM and CME  $\hat{p}$  functions had a significant proportion of test samples with incorrectly-predicted relevant concepts (e.g. CME had an *MPO* score of 0.25 at  $m = 4$ , implying that 25% of all test samples have at least 4 concepts predicted incorrectly). In practice, these mispredictions can have a significant impact on consequent task label predictive performance, as will be further explored in the next section.

## 5.2. Task Performance

In this section, we evaluate the fidelity and performance of the extracted  $\hat{f}$  models. For all CME and Net2Vec  $\hat{p}$  functions evaluated in the previous section, we trained output-to-concept functions  $\hat{q}$ , predicting class labels from the  $\hat{p}$  concept predictions. Next, for every  $\hat{p}$ , we defined its corresponding  $\hat{f}$  as discussed in Section 3, via a composition of  $\hat{p}$  and its associated  $\hat{q}$ . For every  $\hat{f}$ , we evaluated its fidelity and its task performance, using a held-out sample test set. Table 1 shows the fidelity of extracted models, and Table 2 shows the task

performance for these models (averaged over 5 runs). The original Task 1, Task 2, and CUB models achieved task performances of  $100 \pm 0\%$ ,  $100 \pm 0\%$ , and  $82.7 \pm 0.4\%$ , respectively, as described in Section 4.

**Table 1**

Fidelity of extracted  $\hat{f}$  models

	CME	CBM	Net2Vec
<b>Task 1</b>	$100.0 \pm 0.0\%$	–	$24.5 \pm 3.6\%$
<b>Task 2</b>	$99.3 \pm 0.5\%$	–	$38.3 \pm 4.0\%$
<b>CUB</b>	$74.42 \pm 3.1\%$	$77.5 \pm 0.2\%$	$73.8 \pm 2.8\%$

**Table 2**

Task performance of extracted  $\hat{f}$  models

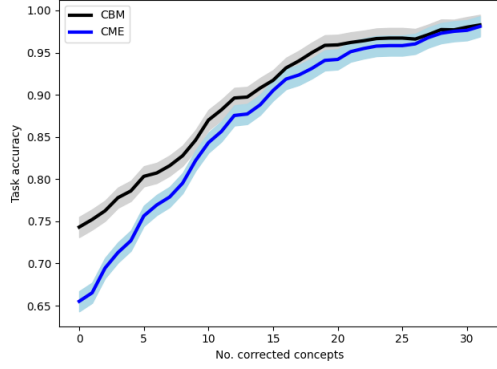
	CME	CBM	Net2Vec
<b>Task 1</b>	$100.0 \pm 0\%$	–	$24.5 \pm 3.6\%$
<b>Task 2</b>	$99.3 \pm 0.5\%$	–	$38.3 \pm 4.0\%$
<b>CUB</b>	$70.8 \pm 1.8\%$	$75.7 \pm 0.6\%$	$69.8 \pm 1.5\%$

For both dSprites tasks, CME  $\hat{f}$  models achieved high (99%+) fidelity and task performance scores, indicating that CME successfully approximated the original dSprites models. Furthermore, these scores were considerably higher than those produced by the Net2Vec  $\hat{f}$  models.

For the CUB task, both CME and Net2Vec  $\hat{f}$  models achieved relatively lower fidelity and task performance scores (in this case, performance of CME was very similar to that of Net2Vec). Crucially, the CBM model *also* achieved relatively low fidelity and accuracy scores (as anticipated from our *MPO* metric analysis). This implies that concept information learnable from the data is insufficient for achieving high task accuracy. Hence the relatively high CUB model accuracy has to be caused by the CUB model relying on other non-concept information. Thus, the low fidelity of CME and Net2Vec is a consequence of the CUB model being *non-concept-decomposable*, implying that its behaviour cannot be explained by the desired concepts. The next section discusses possible approaches to fixing this issue.

## 5.3. Intervening

In the previous section, we demonstrated how CME can be used to identify whether a model relies on desired concepts during decision-making. In this section, we demonstrate how CME can be used to suggest *model*



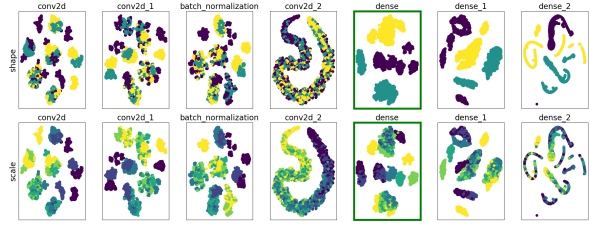
**Figure 4:** The task accuracy of  $\hat{q}$  functions, trained on concepts predicted by  $\hat{p}$  functions, with top # *No. corrected concepts* set to their ground truth values. Performance of Net2Vec was very similar to that of CME, and is thus omitted here for simplicity.

*improvements*, aligning model behaviour with the desired concepts.

We trained a logistic regression  $\hat{q}$  model predicting task labels from ground-truth concept labels for the CUB task, obtaining an accuracy score of  $96.4 \pm 0.5\%$  on a held-out test set (averaged over 5 runs). Using this model’s coefficient magnitudes as a measure of concept importance, we discovered that the 32 most important concepts identified this way were sufficient for achieving over 96% task accuracy using logistic regression.

Using this reduced concept set, we inspected how our CUB  $\hat{q}$  function performances would change, if their corresponding  $\hat{p}$  functions extracted these concepts perfectly. This was achieved by taking the  $\hat{p}$  concept predictions of these concepts on the test and training sets, setting the values of the top  $i$  most important concepts to their ground truth values, training logistic regression  $\hat{q}$  functions on these modified training sets, and measuring their accuracies on the modified test sets (this approach is referred to as *concept intervention* in the rest of this work). The results are shown in Figure 4, with  $i$  ranging from 0 to 32.

These results demonstrate that concept information from only 32 concepts is sufficient for achieving over 96% task performance. Thus, predictive performance of the CUB model can be significantly improved (up to 14%) by ensuring that the model is able to learn and use this concept information. Crucially, these results show that CME concept intervention also significantly improves CBM model performance, indicating that the necessary concept information is *not learnable from the data*. Hence, undesired CUB model behaviour is likely



**Figure 5:** t-SNE plots for the relevant Task 2 concepts. Each row corresponds to a different concept, and each column corresponds to a different layer of the Task 2 model. Each plot is colored with respect to the concept’s values. For every concept row, the subplot with a green border indicates the layer CME selected for predicting the value of that concept.

arising due to data properties (e.g. the data not being representative with respect to key concepts), not model properties (e.g. architecture, or training regime).

Overall, we demonstrated how CME can be used to identify the key concept information that can be used to improve performance of DNN models, and ensure that they are closer aligned with the desired concept-based behaviour. Furthermore, we demonstrated how CME can be used to identify whether undesired model behaviour is caused by model properties, or data properties.

## 5.4. Explainability

By studying CME-extracted  $\hat{p}$  and  $\hat{q}$  functions separately, we can gain additional insights into what concept information the original model learned and how this concept information is used to make predictions. We give examples of how these sub-models can be inspected in the remainder of this section.

### 5.4.1. Input-to-Concept ( $\hat{p}$ )

CME extraction of  $\hat{p}$  functions from a DNN model is highly complementary to existing approaches on latent space analysis. For example, Figure 5 shows a t-SNE [37] 2D projected plot of every layer’s hidden space of the dSprites Task 2 model, highlighting different concept values of the two relevant concepts, as well as the layers used by CME to predict them. Figure 5 demonstrates several important ways in which CME concept extraction can be combined with existing latent space analysis approaches, which will be discussed in the remainder of this section. Further examples are given in Appendix C.

**Manifold Types** Using ground-truth concept information and hidden space visualisation, it is possible



to inspect the nature of latent space manifolds, with respect to specific concepts. Firstly, this inspection allows to build an intuition of how concept information is represented in a particular latent space. Secondly, it is possible to use this information when selecting the types of  $\hat{p}$  functions to use during concept extraction. For instance, some manifolds consist of “blobs” encoding distinct concept values (e.g. row shape, columns dense, dense\_1), suggesting that the latent space is clustered with respect to a concept’s values.

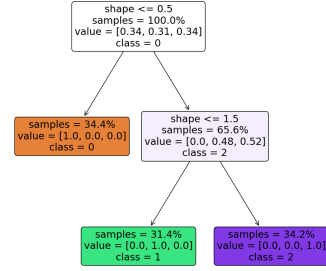
**Variation Across Layers** Using ground-truth concept information and hidden space visualisation, it is also possible to inspect how concept information representation varies across layers of a DNN model. Firstly, this inspection allows to build an intuition of how concept-related information is transformed by the DNN. Secondly, it is possible to use this information to identify the ‘best’ layers to extract concept information from. For instance, both rows shape and scale illustrate that the manifolds of higher layers become more unimodal (separating concept values) with respect to the relevant concepts. Importantly, this analysis, together with the definition of  $\hat{p}$  allows using *different* layers for extracting different concepts.

Overall, we argue that CME concept extraction can be well-integrated with existing latent space analysis approaches, in order to study which concept information is learned by a DNN, and how this information is represented across DNN layers. This type of inspection can have numerous applications, including: (i) inspecting which concepts a model has learned, and verifying whether it has learned the desired concepts (useful for *model explanations* and *model verification*), (ii) inspecting how concept information is represented across different layers (useful for fine-grained *model analysis*), (iii) extracting concept predictions from a DNN (useful for *knowledge extraction*). Further examples and analysis of extracted  $\hat{p}$  functions can be found in Appendix C.

#### 5.4.2. Concept-to-Output ( $\hat{q}$ )

$\hat{q}$  functions encapsulate how a DNN uses concept information when making predictions. Hence, these functions can be inspected directly, in order to analyse model behaviour *represented in terms of concepts*. An example is given in Figure 6, in which we plot the decision tree  $\hat{q}$  function extracted by CME from the Task 1 model. Further examples are given in Appendix D.

Overall, inspection of  $\hat{q}$  functions can be used for (i) verifying that a DNN uses concept information correctly during decision-making, and that it’s high-level



**Figure 6:** Visualisation of a decision tree  $\hat{q}$  extracted from the Task 1 model. The model has correctly learned to differentiate between classes based on the shape concept values.

behaviour is consistent with user expectations (*model verification*), (ii) identifying specific concepts or concept interactions (if any) causing incorrect behaviour (*model debugging*), (iii) extracting new knowledge about how concept information can be used for solving a particular task (*knowledge extraction*). Further examples and analysis of extracted  $\hat{q}$  functions can be found in Appendix D.

## 6. Conclusions

We present CME: a concept-based model extraction framework, used for analysing DNN models via concept-based extracted models. Using two case-studies, we demonstrate how CME can be used to (i) analyse concept information learned by DNN models (ii) analyse how DNNs use concept information when making predictions (iii) identifying key concept information that can further improve DNN predictive performance. CME is a model-agnostic, general-purpose framework, which can be combined with a wide variety of different DNN models and corresponding tasks.

In this work, we assume a fixed set of concept labels available to CME before model extraction begins (i.e. the concept-labelled dataset). In the future, we intend to explore active-learning based approaches to obtaining maximally-informative concept labels in an interactive fashion. Consequently, these approaches will improve extracted model fidelity by retrieving the most informative concept labels, and reduce manual concept labelling effort.

Given the rapidly-increasing interest in concept-based explanations of DNN models, we believe our approach can play an important role in providing granular concept-based analyses of DNN models.

## Acknowledgements

AW acknowledges support from the David MacKay Newton research fellowship at Darwin College, The Alan Turing Institute under EPSRC grant EP/N510129/1 & TU/B/000074, and the Leverhulme Trust via the Leverhulme Centre for the Future of Intelligence (CFI). BD acknowledges support from EPSRC Award #1778323. DK acknowledges support from EPSRC ICASE scholarship and GSK. DK and BD acknowledge the experience at Tenyks as fundamental to developing this research idea.

## References

- [1] B. Goodman, S. Flaxman, European union regulations on algorithmic decision-making and a “right to explanation”, *AI magazine* 38 (2017) 50–57.
- [2] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Benetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al., Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai, *Information Fusion* 58 (2020).
- [3] A. Adadi, M. Berrada, Peeking inside the black-box: A survey on explainable artificial intelligence (xai), *IEEE Access* 6 (2018).
- [4] U. Bhatt, A. Xiang, S. Sharma, A. Weller, A. Taly, Y. Jia, J. Ghosh, R. Puri, J. M. Moura, P. Eckersley, Explainable machine learning in deployment, in: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020, pp. 648–657.
- [5] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, B. Kim, The (un) reliability of saliency methods, in: *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, Springer, 2019, pp. 267–280.
- [6] D. A. Melis, T. Jaakkola, Towards robust interpretability with self-explaining neural networks, in: *Advances in Neural Information Processing Systems*, 2018, pp. 7775–7784.
- [7] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, B. Kim, Sanity checks for saliency maps, in: *Advances in Neural Information Processing Systems*, 2018, pp. 9505–9515.
- [8] B. Dimanov, U. Bhatt, M. Jamnik, A. Weller, You shouldn’t trust me: Learning models which conceal unfairness from multiple explanation methods, in: *European Conference on Artificial Intelligence*, 2020.
- [9] F. Poursabzi-Sangdeh, D. G. Goldstein, J. M. Hofman, J. W. Vaughan, H. Wallach, Manipulating and measuring model interpretability, *arXiv preprint arXiv:1802.07810* (2018).
- [10] B. Kim, M. Wattenberg, J. Gilmer, C. J. Cai, J. Wexler, F. B. Viégas, R. Sayres, Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV), in: J. G. Dy, A. Krause (Eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm*, Sweden, July 10–15, 2018, volume 80 of *Proceedings of Machine Learning Research*, PMLR, 2018, pp. 2673–2682. URL: <http://proceedings.mlr.press/v80/kim18d.html>.
- [11] B. Zhou, Y. Sun, D. Bau, A. Torralba, Interpretable basis decomposition for visual explanation, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 119–134.
- [12] A. Ghorbani, J. Wexler, J. Y. Zou, B. Kim, Towards automatic concept-based explanations, in: *Advances in Neural Information Processing Systems*, 2019.
- [13] C.-K. Yeh, B. Kim, S. O. Arik, C.-L. Li, P. Ravikumar, T. Pfister, On concept-based explanations in deep neural networks, *arXiv preprint arXiv:1910.07969* (2019).
- [14] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viégas, R. Sayres, Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav), *arXiv preprint arXiv:1711.11279* (2017).
- [15] Y. Goyal, U. Shalit, B. Kim, Explaining classifiers with causal concept effect (cace), *arXiv preprint arXiv:1907.07165* (2019).
- [16] G. E. Hinton, Learning multiple layers of representation, *Trends in cognitive sciences* 11 (2007) 428–434.
- [17] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Object detectors emerge in deep scene cnns, *arXiv preprint arXiv:1412.6856* (2014).
- [18] P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, P. Liang, Concept bottleneck models, in: *Proceedings of Machine Learning and Systems 2020, International Conference on Machine Learning*, 2020, pp. 11313–11323.
- [19] F. D.-V. Isaac Lage, Human-in-the-loop learning of interpretable and intuitive representations, in: *ICML Workshop on Human Interpretability*, 2020. URL: [http://whi2020.online/static/pdfs/paper\\_31.pdf](http://whi2020.online/static/pdfs/paper_31.pdf).
- [20] R. Andrews, J. Diederich, A. B. Tickle, Survey and critique of techniques for extracting rules from trained artificial neural networks, *Knowledge-*

- based systems 8 (1995) 373–389.
- [21] J. R. Zilke, E. L. Mencia, F. Janssen, Deepred-rule extraction from deep neural networks, in: International Conference on Discovery Science, Springer, 2016, pp. 457–473.
  - [22] D. Chen, S. P. Fraiberger, R. Moakler, F. Provost, Enhancing transparency and control when drawing data-driven inferences about individuals, *Big data* 5 (2017) 197–212.
  - [23] R. Krishnan, G. Sivakumar, P. Bhattacharya, Extracting decision trees from trained neural networks, *Pattern recognition* 32 (1999).
  - [24] M. Sato, H. Tsukimoto, Rule extraction from neural networks via decision tree induction, in: IJCNN’01. International Joint Conference on Neural Networks. Proceedings (Cat. No. 01CH37222), volume 3, IEEE, 2001, pp. 1870–1875.
  - [25] D. Kazhdan, Z. Shams, P. Liò, Marleme: A multi-agent reinforcement learning model extraction library, *arXiv preprint arXiv:2004.07928* (2020).
  - [26] Q. Liu, X. Liao, L. Carin, Semi-supervised multi-task learning, in: Advances in Neural Information Processing Systems, 2008.
  - [27] L. Matthey, I. Higgins, D. Hassabis, A. Lerchner, dsprites: Disentanglement testing sprites dataset, <https://github.com/deepmind/dsprites-dataset/>, 2017.
  - [28] C. Wah, S. Branson, P. Welinder, P. Perona, S. Belongie, The caltech-ucsd birds-200-2011 dataset (2011).
  - [29] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, L. D. Jackel, Handwritten digit recognition with a back-propagation network, in: Advances in neural information processing systems, 1990, pp. 396–404.
  - [30] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: A simple way to prevent neural networks from overfitting, *Journal of Machine Learning Research* 15 (2014) 1929–1958. URL: <http://jmlr.org/papers/v15/srivastava14a.html>.
  - [31] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the inception architecture for computer vision, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2818–2826.
  - [32] A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet classification with deep convolutional neural networks, in: Advances in neural information processing systems, 2012, pp. 1097–1105.
  - [33] Y. Cui, Y. Song, C. Sun, A. Howard, S. Belongie, Large scale fine-grained categorization and domain-specific transfer learning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 4109–4118.
  - [34] R. Fong, A. Vedaldi, Net2vec: Quantifying and explaining how concepts are encoded by filters in deep neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 8730–8738.
  - [35] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, B. Schölkopf, Learning with local and global consistency, in: Advances in Neural Information Processing Systems 16, 2004.
  - [36] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* 12 (2011).
  - [37] L. v. d. Maaten, G. Hinton, Visualizing data using t-sne, *Journal of Machine Learning Research* 9 (2008) 2579–2605.

## A. Concept Decomposition

The results and findings presented in existing work on concept-based explanations suggests that users often think of tasks in terms of concepts and concept interactions (see Section 2.1 for further details). For instance, consider the task of determining the species of a bird from an image. A user will typically perform this task by first identifying relevant concepts (e.g. wing color, head color, and beak length) present in a given image, and then using the values of these concepts to infer the bird species, in a bottom-up fashion.

On the other hand, Machine Learning (ML) models usually rely on high-dimensional data representations, and infer task labels directly from these high-dimensional inputs (e.g. a CNN produces a class label from raw input pixels of an image).

Consequently, *Concept Decomposition* (CD) approaches attempt to explain the behaviour of such ML models by decomposing their processing into two distinct steps: concept extraction, and label prediction. In concept extraction, concept information is extracted from the high-dimensional input data. In label prediction, concept information is used to produce the output label. Hence, CD approaches attempt to explain ML model behaviour in terms of human-understandable concepts and their interactions in a bottom-up fashion, paralleling human-like reasoning more closely.

Importantly, whilst this work focuses on CNN models and tasks, the notion of CD can in principle be

applied to *any* ML model and task.

### A.1. CBMs

CBMs can be seen as a special case of models performing CD, in which CD behaviour is *enforced by design*. Hence, these models explicitly consist of two submodels, with the first submodel extracting concept information, and the second submodel using this concept information for producing task labels. Importantly, non-CBM models can still demonstrate CD behaviour. For instance, the dSprites Task 2 model was shown to have CD behaviour, with relevant concept information extracted in the dense layers, and used for classification decisions.

### A.2. CBMs & CME

The utility of CBMs is that they produce models explicitly encouraged to use CD. Consequently, these models are much more likely to rely on the desired concepts during decision-making, and be more aligned with a user’s mental model of the corresponding task.

However, a given DNN model can already exhibit CD behaviour, and use the desired concept information (e.g. as was the case with both dSprites task models). In this case, costly modifications and model re-training are unnecessary. As discussed in Section 3, CME can extract concept information from pre-trained DNNs by training  $L * k$  concept predictors (where  $L$  denotes the number of DNN layers used in concept extraction, and  $k$  denotes the number of concepts). As demonstrated in Section 5, these concept predictors can consist of simpler models (e.g. LR), trained on only a fraction of the DNN training data. Thus, the computational cost of training these concept predictors is significantly smaller, compared to training a bottleneck model on all the training data, as done in the case of CBMs.

More importantly, CBM models require knowledge of existing concepts and available concept annotations. In practice, these annotations are often expensive to produce, especially for large datasets and/or a large number of concepts. Furthermore, information about which concepts are relevant and/or sufficient for solving a given task is often not fully available either. Instead, CME is capable of using existing DNN models to extract this information automatically in a semi-supervised fashion, making concept discovery (identifying the relevant concepts), and concept annotation both faster and cheaper.

Overall, CME permits efficient interaction with pre-trained DNN models, which can be used to leverage

concept-related knowledge stored in these models. Consequently, we believe that CME will be invaluable in situations where concept-related information is expensive/difficult to obtain, or is only partially-known. In these cases, a user may interact with existing DNN models via CME, in order to refine existing concept-related knowledge.

It should be noted that a CBM can trivially be approximated using CME, by defining  $\hat{p}$  as the output of a CBM’s concept bottleneck layer, and defining  $\hat{q}$  as the CBM’s submodel producing task labels from the bottleneck layer output.

### A.3. Further Discussion

As discussed in Section 3, CME explores whether a DNN is concept-decomposable, by attempting to approximate it with an extracted model that is concept-decomposable by design (i.e. explicitly consists of two separate stages). Intuitively, if a given DNN learns and relies on concept information of the specified concepts during label prediction, this concept information will be contained in the DNN latent space. Hence, the DNN decision process could be separated into two steps: concept information extraction, and consequent task label prediction.

Importantly, existing CD-based approaches (such as those discussed in Section 2.2) require the set of concepts and their values to be (i) *sufficient* to solve the corresponding classification task (i.e. the class labels can be predicted from concept information with high accuracy) (ii) *learnable* from the data (i.e. the DNN model will be able to learn concept information from the given dataset), in order to achieve high task performance.

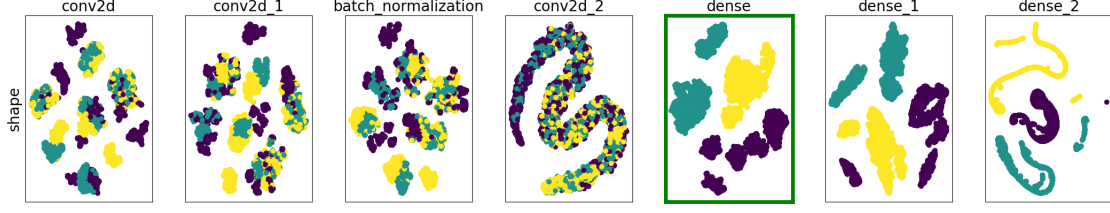
However, these works do not discuss how to handle cases where these assumptions do not hold (e.g. as was the case with the CUB task). Thus, exploring ways of efficiently *discovering* relevant concepts sufficient for solving a given task, as well as ways of ensuring whether this concept information is *learnable* from the data are both important research directions for future work.

## B. dSprites Dataset

### B.1. Description

dSprites is a dataset of 2D shapes, procedurally generated from 6 ground truth independent concepts (color, shape, scale, rotation, x and y position). Table 3 lists the concepts, and corresponding values. dSprites consists





**Figure 7:** t-SNE plots for the relevant Task 1 concept. Each column corresponds to a different layer of the Task 1 model. Each plot is colored with respect to the concept’s values. The subplot with a green border indicates the layer  $\hat{p}$  uses for predicting the value of that concept

of  $64 \times 64$  pixel black-and-white images, generated from all possible combinations of these concepts, for a total of  $1 \times 3 \times 6 \times 40 \times 32 \times 32 = 737280$  total images.

**Table 3**  
dSprites concepts and values

Name	Values
Color	white
Shape	square, ellipse, heart
Scale	6 values linearly spaced in $[0.5, 1]$
Rotation	40 values in $[0, 2\pi]$
Position X	32 values in $[0, 1]$
Position Y	32 values in $[0, 1]$

## B.2. Pre-processing

We select 16 of the 32 values for *Position X* and *Position Y* (keeping every other value only), and select 8 of the 40 values for *Rotation* (retaining every 5th value). This step makes the dataset size more manageable (reducing it from 737280 to  $3 \times 6 \times 8 \times 16 \times 16 = 36864$  samples), whilst preserving its characteristics and properties, such as concept value ranges and diversity.

## C. Input-to-Concept Functions

Figure 7 shows a t-SNE 2D projected plot of every layer’s hidden space of the dSprites Task 1 model, highlighting different concept values of the relevant shape concept, and which layers were used by CME to predict it.

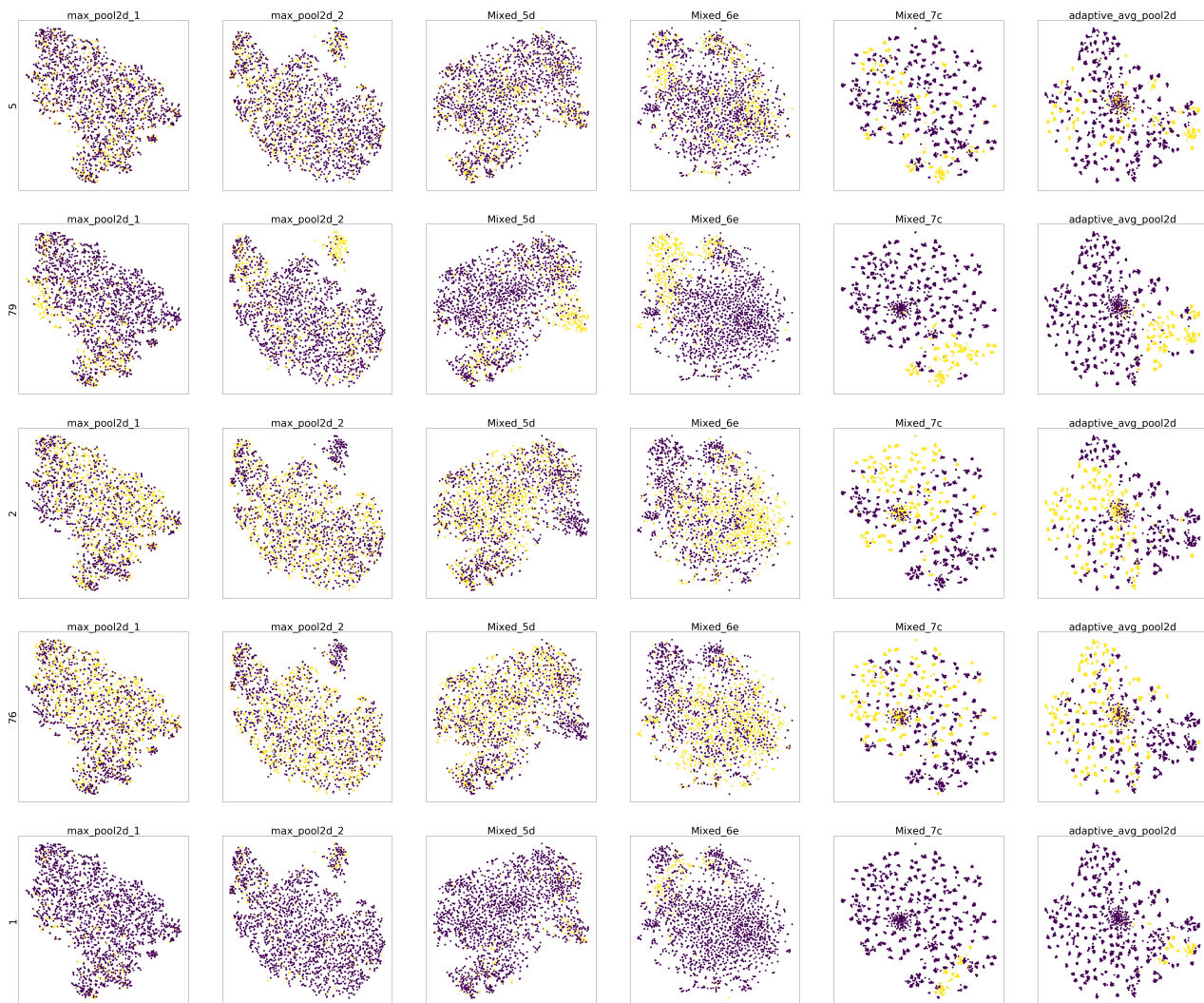
The CUB model has a considerably larger number of layers, and a considerably larger number of task concepts. Hence, for the sake of space, we demonstrate an example here using only 6 different model layers of the CUB model, and showing only the top 5 important concepts identified in Section 5.3. In this Figure, the concepts are named using their indices, and the layers are named following the naming convention used

in [18]. Further details regarding layer naming and/or concept naming can be found in <sup>6</sup>. For all concepts, concept values become significantly better-separated after the *Mixed\_7c* layer. However, the figure shows that concept values are still quite mixed together for some of the points, even for later layers. This low separability indicates that concept values will still be mis-predicted for some of the points, and that concept extraction for the CUB task will likely perform suboptimally.

## D. Concept-to-Output Functions

Figure 9 shows the decision tree extracted for dSprites Task 2. Overall, this model has correctly learned to differentiate between classes based on the shape and scale concepts (note: there are  $3 \times 6$  shape and scale concept values, for a total of 18 output classes).

<sup>6</sup><https://github.com/yewsiang/ConceptBottleneck/tree/master/CUB>



**Figure 8:** t-SNE plots for the top 5 CUB concepts. Each column corresponds to a different layer of the CUB model. Each plot is colored with respect to the concept's values.

