

Concept-based Explainable Artificial Intelligence: A Survey

ELEONORA POETA*, GABRIELE CIRAVEGNA*, ELIANA PASTOR, TANIA CERQUITELLI, and ELENA BARALIS, Politecnico di Torino, Italy

The field of explainable artificial intelligence emerged in response to the growing need for more transparent and reliable models. However, using raw features to provide explanations has been disputed in several works lately, advocating for more user-understandable explanations. To address this issue, a wide range of papers proposing Concept-based eXplainable Artificial Intelligence (C-XAI) methods have arisen in recent years. Nevertheless, a unified categorization and precise field definition are still missing. This paper fills the gap by offering a thorough review of C-XAI approaches. We define and identify different concepts and explanation types. We provide a taxonomy identifying nine categories and propose guidelines for selecting a suitable category based on the development context. Additionally, we report common evaluation strategies including metrics, human evaluations and dataset employed, aiming to assist the development of future methods. We believe this survey will serve researchers, practitioners, and domain experts in comprehending and advancing this innovative field.

CCS Concepts: • **Computing methodologies** → **Knowledge representation and reasoning**.

Additional Key Words and Phrases: Explainable AI, Concept-based XAI, Concept-based Explainability

1 INTRODUCTION

In recent years, the importance of Artificial Intelligence (AI) has surged due to its transformative impact on various aspects of society and industry. This success is predominantly attributed to the advancement of Deep Learning models [63]. However, the inherent complexity and opaque nature of Deep Neural Networks (DNN) hinders understanding the decision-making process underlying these models. This issue prevents the safe employment of these models in contexts that significantly affect users. Consequently, decision-making systems based on Deep learning have faced constraints and limitations from regulatory institutions [40, 71], which increasingly demand transparency in AI models [53].

To address the challenge of creating more trustworthy and transparent AI models, researchers have pursued eXplainable AI (XAI) methods [3]. Most XAI methods focus on a single prediction, highlighting which input features contributed the most to the prediction of a certain class. This is achieved through gradient-based analysis (Vanilla Gradient [102], CAM [124], Grad-CAM [100]), or local approximations with surrogate interpretable models (LIME, [91]) or game-theoretic approaches (SHAP [70]). Other methods, instead, focus on node-oriented explanations to enhance network transparency [31, 102, 119]. Finally, a few methods attempted to globally interpret a model behaviour via global feature importance [42] linear combination of non-interpretable ones [6] or global surrogate models [43, 91].

However, the reliability of standard XAI methods has been contested: [4] showed that gradient-based explanations may not change even after randomly re-parametrizing the network or the data labels, while [47] showed that methods based on surrogate models are unable to identify the most important features. Furthermore, [36, 59] demonstrated that feature importance methods can be misled by simple data modification, which does not affect the model’s prediction. Even if these challenges were addressed through the development of more robust methods, a fundamental concern remains: standard XAI techniques provide explanations at the feature level, which may lack meaningful interpretation, especially for non-expert users [86]. Understanding *where* the network is looking is not sufficient to explain *what* the

*These authors contributed equally to this work

Authors’ address: Eleonora Poeta*, eleonora.poeta@polito.it; Gabriele Ciravegna*, gabriele.ciravegna@polito.it; Eliana Pastor, eliana.pastor@polito.it; Tania Cerquitelli, tania.cerquitelli@polito.it; Elena Baralis, elena.baralis@polito.it, Politecnico di Torino, Torino, Italy.



Fig. 1. Word cloud generated from the titles, abstracts, and keywords of the reviewed C-XAI papers. The trend towards human understanding and interpretability is noticeable, in contrast to standard XAI approaches.

network is seeing in a given input [2, 93]. To really achieve this, there is a growing consensus that XAI techniques should provide explanations in terms of higher-level attributes [2, 33, 56], commonly referred to as *concepts*.

Concept-based explanations offer a compelling alternative by providing a more holistic view of the model’s inner workings. By explaining the model’s predictions in terms of human-understandable attributes or abstractions, concept-based explanations better resemble the way humans reason and explain [56, 58]. They also empower users to gain deeper insights into the underlying reasoning, helping them to detect model biases and improve the classification model [16, 51, 126]. Finally, they are more stable to perturbations than standard XAI techniques [8] and can create models inherently more robust to adversarial attacks [24].

This paper proposes the first systematic review of Concept-based eXplainable Artificial Intelligence (C-XAI) methods. While there exists an abundant literature of surveys in the realm of explainability [3, 9, 15, 27, 30, 44, 85], none of them has distinctly concentrated on the concept-based domain and only few [15, 85] identified C-XAI as a separate category. We analyzed papers from 2017 to July 2023 appearing in the proceedings of NeurIPS, ICLR, ICML, AAAI, IJCAI, ICCV and CVPR conferences or appeared in important journals of the field such as Nature Machine Intelligence, Transaction on Machine Learning Research, Artificial Intelligence, or also cited in the related work of these papers. To gain a visual insight into the examined papers, please refer to Figure 1. This illustration showcases the 100 most frequently occurring words found in the titles, abstracts, or keywords of the papers under review. The figure already indicates a noteworthy trend in the C-XAI field, where the emphasis is shifting towards interpretable models and explanations more in line with human understanding.

To summarize, this work contributes with the following:

- We properly define the key terms in the concept-based explanation field, including what is a concept and a concept-based explanation.
- We present a taxonomy that outlines 9 different C-XAI categories based on the utilization of concepts during training, the specific type of concept used, and the required explanation.
- We propose guidelines for choosing a suited approach, describing each category’s applicability, advantages, and disadvantages.
- We analyze the methods according to 13 dimensions of analysis.
- We provide for each method a brief outline of its main aspects and how it differs from the other methods in the same category, including papers of critics and comparisons.

- We report metrics and datasets available for assessing and developing novel C-XAI methods.
- We highlight first C-XAI applications and promising future directions.

We hope this paper may serve as a reference for AI researchers and practitioners seeking to enhance the reliability of explanation methods and models. The suggested guidelines are a valuable tool for domain experts, enabling them to select the most fitting explanation method tailored to the specific type of explanation needed or the concepts at hand, for instance. Policymakers may also find concept-based models to be sufficiently accountable and trustworthy for the deployment of neural models in safety-critical contexts, making this paper relevant to a wide spectrum of stakeholders.

The paper is structured as follows: in Section 2 we provide the definitions of the most important terms in the concept-based explanation field; in Section 3, we define a high-level categorization of C-XAI approaches and provide readers with guidelines for selecting the suitable category; in Section 4 we resume the most important post-hoc concept-based explainability methods, while in Section 5 we review explainable-by-design concept-based models following the previously introduced categorization; in Section 6, we report valuable tools developed in the literature to develop and assess novel methods, including metrics, human evaluations, and datasets; in Section 7, we report the first applications of concept-based explainability methods; in Section 8, we hypothesize some promising future directions of the concept-XAI field; finally, Section 9 concludes the paper.

2 DEFINITIONS

In Figure 2, we provide a common pipeline of a C-XAI method. The input is projected in the latent space of the network, where the method aims to extract or directly represent a set of concepts (normally of the same type). For some concepts, additional knowledge may be required to employ them properly. Finally, various types of explanations can be generated in addition to the model prediction. Due to the recent genesis of the C-XAI literature, it is challenging to establish consensus on the definitions of certain terms. Specifically, the definitions of “concept” and “explanation” are not consistently uniform, and similar studies may adopt varying interpretations.

In the following, we first propose four typologies of concepts (Section 2.1); secondly, we identified three different concept-based explanations with different types of explainability/interpretability (Section 2.2). Further, we recall the difference between a post-hoc method and an explainable-by-design model, particularly within the C-XAI context (Section 2.3). Finally, in Appendix A, we compiled a glossary of the most commonly used terms, while in Table 5, we report the acronyms and abbreviations used throughout the paper.

2.1 What is a Concept?

In the concept-based explainability literature, the term *concept* has been defined in very different ways. Citing [78], “A concept can be any abstraction, such as a colour, an object, or even an idea”. As we resume in the top part of Table 1, we propose a categorization of concepts into four typologies: *Symbolic Concepts*, *Unsupervised Concept Bases*, *Prototypes*, and *Textual Concepts*.

Symbolic Concepts are human-defined symbols [24]. They can be high-level attributes of the task under consideration or interpretable abstractions, such as the color or the shape of the predicted object [12]. Figure 2 shows that *beak* is a suitable symbolic concept for a bird identification task. Since these concepts are pre-defined by humans, we generally require auxiliary data equipped with concept annotation, particularly when we are dealing with non-symbolic features (e.g., image pixels vs tabular data). The representation in the network of symbolic concepts can be analyzed post-hoc or forced during training to create an explainable-by-design model.

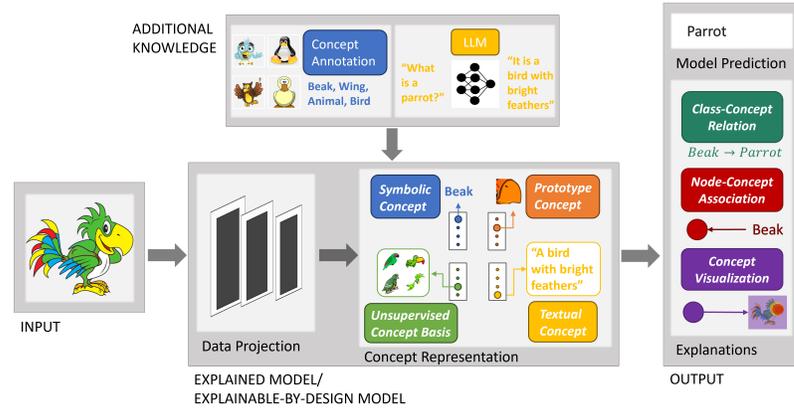


Fig. 2. Concepts and explanations provided by C-XAI methods and models. Some concepts require additional knowledge to be implemented.

Unsupervised Concept Bases are clusters of samples the network learns. Even if they are not built to resemble human-defined concepts, these unsupervised representations may still capture abstractions that are more understandable to humans than individual features or pixels [37]. As shown in Figure 2, a network may learn a cluster of green birds. To extract such a concept, we require the employment of a clustering algorithm either post-hoc in the latent space of the model or during training, e.g., by means of a reconstruction loss. As in the previous case, these concepts can be employed by both explainable-by-design models or post-hoc methods.

Prototypes are representative examples of peculiar traits of the training samples. They can be a training sample or only a part of a training sample. As shown in Figure 2, a part prototype of a training sample might be a *beak*. The set of prototypes is, in general, representative of the whole data set [66]. We regard prototypes as concepts following recent literature [75] since they are still higher-level terms than the input features. Also, they are different from unsupervised concept bases since they explicitly encode the peculiar traits in the network weights, and consequently, they cannot be employed in post-hoc methods but only in explainable-by-design models.

Textual Concepts are textual descriptions of the main classes. As reported in Figure 2, a textual concept might be ‘A bird with bright feathers’ for a sample of a parrot. Textual concepts are provided at training time by means of an external generative model [80], and they are employed inside a concept-based model in the form of a numerical embedding of the text. This type of concept is gaining prominence due to the recent development of Large-Language Models (LLMs). So far, it has been only employed in explainable-by-design models.

For a more formal definition of a Concept, following category theory principles, please refer to Appendix B.

2.2 What is a Concept-based Explanation?

Even when considering a fixed type of concept, the definition of a concept-based explanation can be elusive. Existing literature generally agrees that concept-based explanations should explain how DNNs make particular decisions using concepts [117] and adhere to specific criteria such as being meaningful, coherent, and relevant to the final class [37] and also explicit and faithful [8]. As we report in the bottom part of Table 1, we define three distinct categories of concept-based explanations: *Class-Concept Relation*, *Node-Concept Association*, and *Concept Visualization*. Each method may offer one or more explanations, serving specific purposes in different scenarios.

Concept Type	Definition
Symbolic Concept	Human-defined attribute or abstraction
Unsupervised Concept Basis	Cluster of similar samples
Prototype	(Part-of) a training sample
Textual Concept	Short textual description of a main class
Explanation Type	
Class-Concept Relation	Relation among a concept and an output class of a model
Node-Concept Association	Explicit association of a concept with a hidden node of the network
Concept-Visualization	Visualization of a learnt concept in terms of the input features

Table 1. Concept and explanation types.

Class-Concept Relation explanations consider the relationship between a specific concept and an output class within the model. This relationship can be expressed either as the importance of a concept for a particular class or through the use of a logic rule that involves multiple concepts and their connection to an output class. For example, in Figure 2, the image in input is classified as a Parrot because it is mainly related to the concept of *beak*. In general, this type of explanation can be applied to all types of concepts: for instance, with prototypes, we have class-concept explanations like $parrot := 0.8 prototype_1 + 0.2 prototype_2$. This explanation can be extracted by analyzing post-hoc the correlation of concepts and classes in the latent space of a network or with the employment of an interpretable model from the concepts to the tasks.

Node-Concept Association explicitly assigns a concept to an internal unit (or a filter) of a neural network. This explanation enhances the transparency of deep learning models, highlighting what internal units see in a given sample. It can be defined post-hoc by considering the hidden units maximally activating on input samples representing a concept. Otherwise, it can also be forced during training by requiring a unit to predict a concept. For instance, consider Figure 2: some nodes of the network can automatically learn to recognize the beak since it is a crucial part of the prediction of a bird. Otherwise, an expert can also define a set of attributes in advance and require their explicit representation in the network intermediate layers.

Finally, **Concept Visualization** explanations highlight the input features that best represent a specific concept. When symbolic concepts are used, this explanation closely resembles the saliency map of standard XAI methods. However, when non-symbolic concepts are employed (such as unsupervised concept basis, prototypes, or even textual concepts), the focus shifts towards understanding which unsupervised attributes or prototypes the network has learned. This form of explanation is often combined with one of the previous explanations, and it enables the understanding of what are the concepts associated with a specific class or node. For instance, consider again Figure 2, where we reported a model employing prototype concepts as an example of Concept Visualization. A visualization of the prototype in the input sample is crucial to understand which concept the network has learned.

2.3 C-XAI Post-hoc Methods vs Explainable-by-design Models

In the XAI literature, the difference between *post-hoc explanation methods* and *explainable-by-design models* is well-defined: the latter are inherently designed for transparency, such as decision trees or linear models; in contrast, the former are tools applied post-training to explain complex models, such as DNNs. On the contrary, in the field of

concept-based approaches, the difference becomes more subtle since C-XAI explainable-by-design models also work with DNNs but explicitly represent some concepts within the architecture. In the following, we try to clarify this difference.

Post-hoc Concept-based Explanation Methods operate on concepts rather than dissecting input features. They explain a prediction, an entire class, or an internal network node given a set of concepts. Typically, this involves projecting the samples representing the concepts in the model’s latent space and analyzing their relation to the prediction and the hidden node activations. Symbolic and unsupervised concept bases have been both employed to provide all types of explanations.

Explainable-by-design Concept-based Model, on the other hand, are neural models employing an explicit representation of a set of concepts as an intermediate layer. In this way, the predicted concepts influence the task predictions, and thus, they are closely tied to the provided explanations. Since they generally provide Node-Concept association by-design, they can be regarded as inherently transparent models, at least to some extent. In this case, the whole spectrum of concepts can be used, either explicitly annotated, extracted in an unsupervised manner, encoding prototypes, or even generated by a separate model.

3 C-XAI TAXONOMY

To assist users in understanding the C-XAI landscape, in this section, we provide a taxonomy of concept-based methods and models together with some guidelines for selecting an appropriate method for the development constraints and the desired outcomes. In Section 3.1, we describe the dimensions of analysis that we employed to categorize and characterize C-XAI methods, while in Section 3.2, we present a procedure to follow to select the desired C-XAI method for different requirements.

3.1 Dimensions of analysis

We characterize each concept-based explainability method by 13 dimensions of analysis. We categorize these dimensions into three groups: (i) concepts and explanations characteristics, (ii) the applicability of the method, and (iii) resources employed. The dimensions related to the first two groups will be analyzed in Sections 4, 5 and reported in Tables 2, 3, respectively analyzing methods following two different concept training (first dimension of analysis). The third group of dimensions will be analyzed vertically (per resource/evaluation) in Section 6, but we also report a horizontal analysis (per method) in Section 6.4, and in the Appendix in Tables 6 and 7.

- (1) **Concept and Explanations Characteristics.** These dimensions capture key aspects of concepts and explanations, such as how concepts are integrated into models, concept annotation strategies, concept types, the form of explanations, and their scope.
 - the *Concept Training*, an approach is first categorized according to whether it employs concepts only to provide explanations of an existing model (Post-hoc methods), or also while training the same model (Explainable-by-design models);
 - the *Concept Annotation*, according to whether the method employs an annotated set of concepts (Supervised), it does not (Unsup.), it employs only a few concepts (Hybrid), or it generates concept annotation (Generative);
 - the *Concept Type* (described in Section 2.1), which can be either symbolic (Symb.), unsupervised concept bases (Uns. Basis), prototype-based (Proto.), or textual;

- the *Explanation Type*, (defined in Section 2.2), either Class-concept relation (C-CR), Node-concept association (N-CA), Concept Visualization (C-Viz) or a combination of these;
- the *Explanation Scope*, according to whether the explanation is local (LO) (i.e., it explains an individual prediction), global (GL) (i.e., it provides insights into the overall model behavior), or both.

Regarding Concept-based Models, we also characterize the methods with the following:

- the *Concept Employment*, how an explainable by-design supervised model employs the concepts during training, either through joint training (Joint) or concept instillation (Instil);
 - the *Performance Loss*, whether there is a performance loss (✓) or not (✗) compared to a black-box model.
- (2) **Applicability of the approach.** These dimensions describe the approach’s applicability by detailing the types of data it can handle, the primary tasks it is designed for, and the specific neural architectures utilized.
- the *Data*, i.e., the data type the method can handle, generally images (IMG), but also text (TXT), graph (GRA), tabular data (TAB), videos (VID) and time series (TS);
 - the *Task* the method is designed for, primarily classification (CLF) but also regression (REG) and clustering (CLU);
 - the *Network* dimension describes the type of neural network employed, including Convolutional Neural Network (CNN), 3D-CNN, Auto-Encoder (AE), Graph Neural Network (GNN), Multi-Layer Perceptron (MLP) and Transformer (TRANSF).
- (3) **Resources and Evaluation conducted.** These dimensions focus on the resources provided (data releases and code/models availability) and the conducted evaluation.
- the *Data Release* indicates whether a new dataset annotated with concepts has been made available;
 - The *Code/Models Availability* indicates whether the authors have released the code and/or the model;
 - the *New metric* dimension indicates if the paper introduces a new evaluation metric;
 - the *Human evaluation* dimension indicates whether the paper includes a user study to evaluate the method.

3.2 Guidelines for selecting a suitable method

We now provide some guidelines to facilitate the method selection, employing Figure 3 as a visual reference. In particular, let us consider a user interested in a C-XAI method. We propose a set of a maximum of three questions that he should answer to choose a category of methods. These questions regard the possibility of modifying the learning model, the availability of an annotated set of concepts, and, according to the subcategory, the desired outcome regarding explanation or concept types or the type of training set employed. We adopted a conversational approach to help the user identify a solution that is suitable for his development conditions. As we report in the figure, this categorization will also characterize the structure of the following sections.

Can I modify the model? This represents the first coarse-grained question a user should answer. According to the possibility of intervening in the model’s development, he can choose among two macro-categories in which the C-XAI literature is divided. He must look at *Post-hoc* concept-based explanation methods (Section 4) if he needs to consider an already trained model that he cannot (or does not want to) modify. On the other hand, if he can create a model from scratch, he can employ a *Explainable-by-design* concept-based models (Section 5), which allows an explicit representation of the concepts within the model architecture.

Do I have annotated concepts? Secondly, the user should check whether he can find a dataset annotated with concepts that are related to the task at hand. In case of a positive response, he should look at supervised methods;

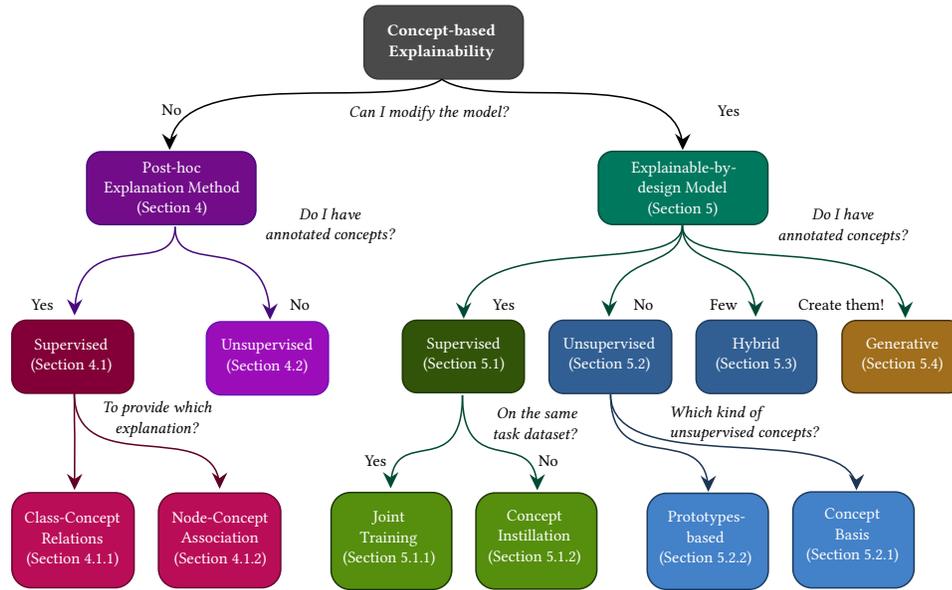


Fig. 3. Categorization of the C-XAI methods with guidelines for selecting a suitable approach.

otherwise, he should consider unsupervised ones that can extract concepts from the same data automatically. For explainable-by-design approaches, he can pick from two more categories. He can employ a hybrid solution that leverages a few supervised concepts together with unsupervised extracted ones since, recently, he can also use a generative method that creates the concepts by means of an external model.

For some sub-categories in Figure 3, we identify a third-level of fine-grained categories (and associated questions) according to whether the explanation, the concept employed, or the concept annotation varies among the methods within the subcategory. **To provide which Explanation?** - If he selected a post-hoc supervised approach, these differ in the type of explanations, either class-concept relations or node-concept associations. **On the same task Dataset?** - Explainable-by-design supervised approaches either require concept annotation on the same dataset of the task at hand or allow them on a separate one. **Which kind of unsupervised concepts?** - Finally, if he chose an unsupervised explainable-by-design approach, they differ in the representation of unsupervised concepts adopted. Instead, the methods in the remaining sub-categories share these main characteristics; hence, we have not split them further.

4 POST-HOC CONCEPT-BASED EXPLANATION METHODS

Concept-based explanation methods explain an existing model without modifying its internal architecture. For this reason, they are post-hoc methods, i.e., they analyze the model only after it has finished the training process.

As we can observe from Table 2, *Post-hoc* concept-based explanation methods primarily diverge on the employment of a dataset containing concept annotations (Supervised, Section ??) or not (Unsupervised methods, Section 4.2). While in the first case, they employ Symbolic Concepts, in the latter, they extract explanations based on unsupervised concept basis. Furthermore, supervised methods differ in focusing on either explaining class-concept relations (C-CR) or node-concept associations (N-CA). In both cases, they can also provide concept visualizations (C-Viz). Unsupervised methods, instead, generally focus on class-concept relations as well as on visualizing the extracted concepts. The explanation

Table 2. Post-hoc Concept-based Explainability methods. We characterize the approaches based on the dimensions describing concepts and their explainability and the applicability of the approach. A full description of each category and of the acronyms is provided in Section 3.1.

	Method	Concept Type	Expl. Type	Expl. Scope	Data Type	Task	Network Type	
Concept Annotation	Supervised	Symb.	T-CAV [56]	C-CR	GL	IMG	CLF	CNN
			CAR [25]	C-CR	GL	IMG + TS	CLF	CNN
			IBD [126]	C-CR, C-Viz	LO & GL	IMG	CLF	CNN
			CaCE [41]	C-CR	GL	IMG	CLF	CNN
			CPM [111]	C-CR	LO	TXT	CLF	TRANSF
		Symb.	Obj. Det. [123]	N-CA, C-Viz	GL	IMG	CLF	CNN
			ND [13]	N-CA, C-Viz	GL	IMG	CLF	CNN
			Net2Vec [34]	N-CA, C-Viz	GL	IMG	CLF	CNN
			GNN-CI	N-CA, C-CR, C-Viz	GL	GRAPH	CLF	GNN
			ACE [37]	C-CR, C-Viz	GL	IMG	CLF	CNN
	Unsupervised	C.Basis	Compl. Aware [117]	C-CR	GL	IMG, TXT	CLF	CNN
			ICE [122]	C-CR, C-Viz	LO & GL	IMG	CLF	CNN
			MCD [107]	C-CR, C-Viz	LO & GL	IMG	CLF	CNN, TRANSF
			CRAFT [33]	C-CR	GL	IMG	CLF	CNN
DMA & IMA [64]			-	-	IMG	CLF	CNN	
STCE [52]			C-CR, C-Viz	GL	VID	CLF	3D-CNN	

scope of most methods is global (GL), with few methods providing both local and global explanations (LO&GL) and one only local (LO). All the reviewed methods have been developed to explain classification tasks (CLF), and many have been designed for CNNs and image (IMG) data types. A few methods have been developed for working with texts (TXT), graphs (GRAPH), or videos (VID) and employing transformers (TRANS) (both for texts and images), GNNs, and 3D-CNN architectures. To characterize the key aspects of post-hoc methods, here we discuss the main positive aspects versus the negative ones.

Advantages. Post-hoc explanation methods are the only viable option in scenarios where a pre-trained model exists or a predefined model is necessary. Also, Concept-based explainability methods are the preferred choice when there can be no compromise on the predictive and generalization capabilities of the model. They allow for enhanced interpretability without affecting the model’s performance.

Disadvantages. The main drawback of these methods is that they do not guarantee that the model truly comprehends or employs the adopted concepts. This is because the model was not exposed to these concepts during its training process. Concept-based explainability methods can only be regarded as a better way of explaining the network behavior by employing terms that are more comprehensible to humans than raw input features. We will further discuss the issues and challenges post-hoc methods face at the end of this section (Section 4.3).

4.1 Supervised Concept-based Explainability Methods

Post-hoc supervised concept-based explanation methods analyze network behavior over samples annotated with *symbolic concepts*. The underlying idea is that a network automatically learns to recognize some concepts. These methods aim to assess which concepts have been learned, where, and how they influence the model. This analysis is conducted differently according to the type of explanations delivered. In particular, to provide Class-concepts relations, some methods analyze how the current prediction or an entire class weight correlates with the projection of a concept in

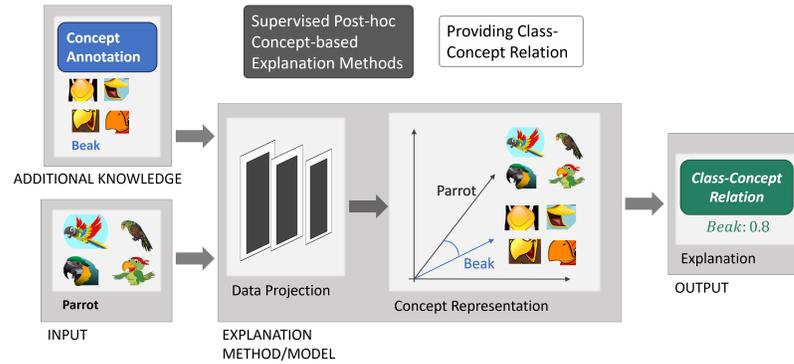


Fig. 4. Post-hoc Supervised Methods providing explanations in terms of class-concept relations. These approaches provide a set of samples annotated with concepts to the explained network to determine their influence on the output class.

the latent space of the model (Section 4.1.1). On the contrary, methods extracting explanations in terms of Node-concept associations study the activations of single hidden nodes to associate them to a given concept (Section 4.1.2).

4.1.1 Class-Concept Relation methods. Supervised post-hoc methods offering explanations examine the connection between output classes and a set of symbolic concepts. As illustrated in Figure 4, this is generally achieved by presenting samples representing concepts to the explained network and analyzing its behaviour. In the figure, we report the case of a model trained to recognize birds (e.g., a Parrot) and tested with samples representing a related concept (e.g., the Beak). The effect of the concepts on the final class is studied by either (i) analyzing the representation of concepts in the model’s latent space and its relation with the output classes (TCAV, TCAR, IBD) or (ii) exploring the causal effect of inserting or removing a concept on the prediction of a particular class (CaCE, CPM). The samples representing the symbolic concepts are drawn from an external dataset, which is then a requirement for this class of methods. As previously mentioned, the set of concepts must be related to the output classes to extract meaningful explanations.

Testing with Concept Activation Vector (**T-CAV [56]**) has been the first work in this category. It introduces the idea of employing a linear probe to assess the interpretability of a model. It consists of (i) freezing the model’s internal architecture and (ii) training a linear layer in the latent space of the model to predict the concepts. For each user-defined concept, it requires a positive set comprising images representing the concept and a complementary one with negative samples. The Concept Activation Vectors (CAVs) represent the orthogonal vectors to the weights of the linear probe, separating positive examples of a given concept from negative ones. Then, T-CAV computes the inner product between the directional derivatives of the class and the CAVs to quantify the models’ conceptual sensitivity, i.e., how much each concept positively influences the prediction of a given class. The T-CAV score is the fraction of the class’s inputs with a positive conceptual sensitivity. The higher the score, the higher the influence of the concept. As noted by the authors, CAVs require concepts to be linearly separable. However, concepts may be semantically linked and mapped to latent space portions that are close to each other. In this case, the authors propose to employ Relative CAVs, i.e., CAVs obtained when considering as negative samples the positive concepts of other concepts. The authors showed that T-CAV enables humans to identify the most important concepts for a network trained over standard and corrupted features, while several standard XAI methods do not.

Concept Activation Regions (**CAR [25]**) relaxes the linear separability assumption of concepts imposed by T-CAV. CAR only requires concept examples to be grouped in clusters in the model’s latent space such that they are identifiable by a non-linear probe (e.g., a kernel-based SVM). Each concept is thus represented by a *region* in the latent space and no longer by a single vector. To this aim, the authors define the notion of concept smoothness, i.e., the possibility to separate a *positive region* in which the concept is predominantly present and a *negative region* in which the concept is not present. A concept is important for a given example if its projection lies in the cluster of concept positive representations. Equivalently to the TCAV score, the authors propose the TCAR score, quantifying the proportion of instances belonging to a specific class whose representations fall within the positive concept region. In some settings, TCAR returns high scores expressing significant associations between classes and concepts, while TCAV returns low scores. The authors claim that this is due to employing a more powerful concept classifier and, consequently, that TCAR scores reflect the relations between classes and concepts more accurately.

Interpretable Basis Decomposition (**IBD [126]**) provides class-concept relations by decomposing the evidence supporting a prediction into the most semantically related components. Concretely, IBD considers the classes and the concepts as vectors in the latent space of the penultimate layer of the network. As for T-CAV, to represent the concepts, IBD trains a set of linear classifiers over the projections in the latent space. IBD then performs an optimization to assess the relevance of each concept for a class. For all samples in the concept dataset, they try to reconstruct a given class’s prediction using a linear regression. In this, the variables are the inner product between each concept vector and the projection of the sample in the latent space. Instead, the coefficients represent each concept’s associated importance for the given class. Interestingly, the *residual* of the optimization assesses the variability in the class predictions that cannot be represented by means of the learned concepts. Furthermore, each concept is accompanied by a heat map showing the features with the highest impact on the concept prediction. By means of human evaluation, the authors demonstrated that IBD explanations enhance standard XAI techniques (like CAM and Grad-CAM) regarding model assessment quality.

The Causal Concept Effect (**CaCE [41]**) investigates the causal effect made by the presence or absence of a symbolic concept on a network prediction. To measure how much a concept can change the classifier’s output, they propose operating over the original sample, removing the concept and generating a counterfactual sample. The CaCE measure is defined as the effect of the presence of the binary concept on the classifier’s output. To effectively compute it, we must control the generative process of the samples annotated with concepts. To achieve this, the authors propose two approaches. (i) Using a controlled environment, it is possible to intervene in the sample generation process, manually inserting or removing a concept and calculating the exact CaCE value for this intervention. (ii) Training a generative model to represent images with and without the concepts to approximate the generative process. In the second scenario, a Variational Auto-Encoder (VAE) architecture can be used in two different variants, either generating images with one concept or another, or generating the counterfactual of an image with/without a certain concept. The authors apply their approach to four datasets, comparing CaCE with a non-causal version of CaCE and TCAV [56], showing that CaCE better understands and identifies the concept as having a causal relation with the class prediction.

Similarly to CaCE, the Causal Proxy Model (**CPM [111]**) provides causal concept-based explanations, this time tailored for Natural Language Processing (NLP) models. CPM builds upon the CEBaB benchmark [1], a dataset studying the presence/absence of a concept on textual data and its causal effect on a model. Starting from a dataset of factual restaurant reviews, it provides counterfactual examples in which a concept (food, service, ambiance, etc.) has been modified. CPM is first initialized with the weights of the black-box model to explain. For both factual and counterfactual

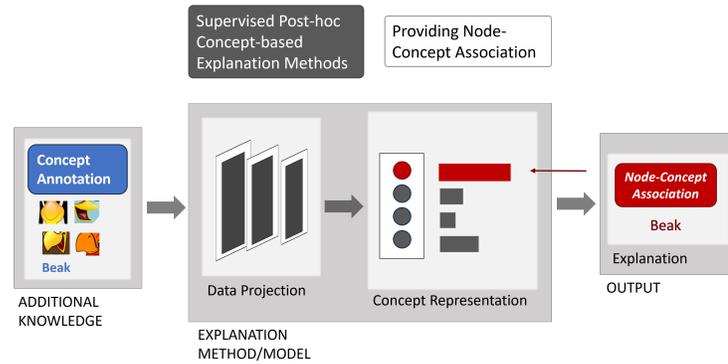


Fig. 5. Post-hoc supervised methods providing an explanation in terms of Node-Concept Associations. This class of approaches associates concepts with internal nodes or filters of the network to increase the transparency of the network decision-making process.

inputs, CPM is then trained to mimic the behaviour of the black-box model. Since counterfactual examples are not provided at test time, CPM is still prompted with the factual sample and information coming from the counterfactual one for counterfactual training. CPM explanations are compared with TCAV [56], the baseline of CaCE [41], and others. The comparison is performed by measuring the distance between the ground-truth (i.e., the actual prediction variation when prompted with a counterfactual sample) and the relative importance of the concept for the predicted class provided by each explainer. The results show that both versions of the CPM significantly outperform previous methods in assessing the causality of a concept.

4.1.2 Node-Concept Association methods. Supervised post-hoc methods may also aim to associate symbolic concepts with internal nodes or filters of the neural model. In this case, the idea is to understand where the model has automatically learned some concepts. This allows us to shed some light on the decision process of a given model. As shown in Figure 5, these methods supply Node-Concept Association explanations by checking which nodes activate the most when the network is presented with input image representing the concepts (e.g., the beak) (ND, Net2Vec) or input graphs (GNN-CI). However, this process can also be performed by human annotators analyzing a node’s common activation areas (Object Detectors).

Object Detectors [123] is a precursor to what has been termed concept-based explainability since 2017. It aims to shed some light on the learned representations in the hidden layers of a neural network. In particular, the study reveals that hidden neurons automatically learn to recognize some objects when training a network to classify scenes. More precisely, to understand and quantify the precise semantics learned from each unit, this work uses AMT workers. Workers have to uniquely identify the concept given the input activation of each unit, as well as define the unit’s precision by selecting all the activations that do not correspond to the concept. According to the AMT workers, the results show that a hidden unit automatically learns to recognize a concept with a precision higher than 70% on average across all hidden layers. Interestingly, they show that the types of concepts learned gradually shift from colours and textures in lower layers to objects and object parts in higher ones. It can be said that this is the first example of human-supervised annotation of concepts.

Similarly to Object Detectors, Network Dissection (**ND [13]**) aims to explain single network neurons by associating them with a unique concept. Rather than employing AMT workers, they propose a new dataset, BRODEN (BROad and

DENSEly labeled dataset) [7], which includes concept labels from diverse data sources. To associate each node to a specific concept, the authors propose to compute the Intersection over Unions (IoU) of the bounding boxes representing the concepts and the activation maps of the node over the BRODEN dataset. The node is then associated with the concept with the highest IoU. In output, ND reports the units with the highest average IoU and the samples with the highest overlap for each concept. For each concept, we can have more than a unit in the network that is associated with it. It is essential to notice that Network Dissection’s effectiveness relies on the quality of the underlying dataset. If a network unit corresponds to a concept understandable to humans but not present in the given dataset, it will receive a lower interpretability score IoU. As a result, they show that individual units may automatically learn to represent high-level semantics objects, concluding that the units of a deep representation may be more interpretable than expected.

The work Network to Vector (**Net2Vec** [34]) argues that semantic representations of concepts might be distributed across filters. Hence, filters must be studied in conjunction, not individually, as in ND [7]. Net2Vec aligns concepts to filters of a CNN via a two-step approach. For each concept, Net2Vec first collects filter activations of the network to explain when it is probed with input samples from a dataset annotated with concepts (e.g., BRODEN [7]). Then, Net2Vec performs an optimization process to either classify or segment the same concept starting from these activations. The resulting weights tell us how much each node is relevant for predicting the given concept. The results showed that combining filters allows to identify better concepts than using a single filter in both tasks, indicating that multiple filters better represent concepts. Moreover, filters are not concept-specific, but they can encode multiple ones. As a limitation, Net2Vec models the relationship between concepts and filters as linear, potentially failing to capture more complex and non-linear concept alignments.

In Graph Neural Networks Concept-based Interpretability (**GNN-CI** [114]), the authors extend the idea of node-concept explanations to the challenging graph classification scenario. In particular, they define concepts either as properties of a node in the graph (e.g., number of connected edges > 7) or as the class of the node (being an Oxygen atom) or of the neighbors (next to an Oxygen atom). For each neuron, they check whether its activation resembles the presence of any concepts. The resemblance is computed as the IoU between concept presence and neuron activation. Rather than considering a single concept, they allow each neuron to be described as a short composition of concepts, e.g., number edges $> 7 \wedge \neg$ is Oxygen. Furthermore, they produce an explanation of the final class in terms of the concepts as a linear classifier trained only on explained nodes to mimic the original classifier. Finally, they also propose a concept visualization in terms of a graph concept activation map, a variant of Grad-CAM applied to a graph. The experiments show that extracted concepts are meaningful for graph classification. For instance, they reveal the presence of functional groups which are known to be correlated with the graph class. Also, they show that interpretable neuron activations are positively correlated to correct class predictions.

4.2 Unsupervised Concept-based Explainability methods

Relying on a given set of supervised concepts is not the only way to assess whether a network is learning higher-level semantics. Indeed, parts of the samples may represent practical explanations of a given prediction or of a learned class. As shown in Figure 6, input pre-processing is normally necessary for these methods. Afterward, unsupervised methods analyze sample representation in the latent space to identify clusters composed of parts of the samples (i.e., *unsupervised concept basis*). Also, they analyze how these clusters contribute to the final prediction, thereby providing explanations in terms of Class-concept relations. How the clustering is conducted differentiates the methods: they either employ

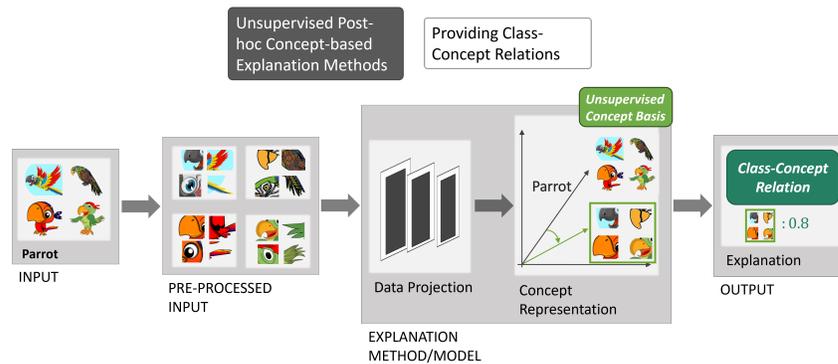


Fig. 6. Unsupervised post-hoc methods extract unsupervised concepts in terms of parts of the same input samples. In particular, these methods provide class-concept relations analyzing the representations of these parts with respect to the output class in the latent space of the network.

standard K-means (ACE), Non-Negative Matrix Factorization (ICE), recursion on different layers (CRAFT), subspace clustering (MCD), concept completeness maximization (Completeness-aware), or assuming concept independency, or, finally, working on voxels from videos. This category bears similarities to standard XAI methods. However, they differ in not selecting features relying solely on their saliency or internal similarity (as with superpixels) but rather on their capability to represent common patterns existing in other input samples.

The Automatic Concept-based Explanations (ACE [37]) method explains a class within a trained classifier, employing a probing strategy that does not require human supervision. The process involves first segmenting images belonging to a specific class at multiple resolutions, encompassing textures, object parts, and complete objects, thus capturing a diverse spectrum of concepts. Subsequently, the projection of image segments in the latent space of the network is clustered across different data. Outlier sample projections are filtered out since they possibly represent insignificant segments or underrepresented concepts. Finally, the T-CAV importance score [56] is calculated for each group (i.e., a concept) to estimate the average influence of a concept on class prediction. The final output comprises the most pertinent concepts based on their T-CAV score. As recognized by the authors, however, ACE might find meaningless or non-coherent concepts due to segmentation, clustering, or similarity score calculation errors. With a human-based evaluation, however, they tested that the samples assigned by ACE to the same concepts are very coherent and can be labeled by a human employing, in general, the same words.

While ACE allows the extraction of some concepts from a learned neural network, there is no guarantee that this set of concepts is sufficient to explain a prediction. In **Completeness-aware** [117], the authors propose a method for extracting complete concept-based explanations. They define a completeness score as the amount to which concept scores are sufficient for predicting model outcomes. To measure the concept completeness, they propose to assess the accuracy achieved by predicting the class label, given the concept scores. With this in mind, they propose a concept discovery algorithm is optimized to unveil maximally complete concepts. The algorithm also promotes proximity between each concept and its closest neighbouring patches to extract meaningful concepts. Furthermore, the authors propose ConceptSHAP, a metric measuring the importance of each concept for a given class. It modifies SHAP by replacing attribute importance with concept sufficiency in representing the class. Using a synthetic dataset with

complete concepts, the authors show that the proposed method excels in discovering the complete set of concepts, outperforming compared concept discovery methods, including ACE.

Invertible Concept-based Explanation (**ICE [122]**) builds upon ACE to improve its concept identification and importance assessment algorithm. ICE replaces K-means with Non-Negative Matrix Factorization (NMF) for extracting concept vectors. Using NMF over the feature maps uncovers a concept for each reduced dimension. This allows ICE to employ a classifier on top to approximate the model outputs and yield concept importance as the weights of the linear approximation instead of using T-CAV scores. The approximation error (or the completeness of the representation, as defined in [117]) is much lower when employing the richer NMF concept representation than when employing Principal Component Analysis (PCA) or K-Means, which can also be seen as a form of binary dimensionality reduction. Similarly to ACE, the authors perform a user study to assess the interpretability of the extracted concepts. They ask participants to assign a few-word description to a group of images belonging to the same extracted concept. The results show that ICE concepts are more interpretable than those derived from PCA or K-Means.

Concept Recursive Activation FacTORIZATION for Explainability (**CRAFT [33]**) extends the techniques introduced in ACE and ICE for unsupervised concept discovery. Also, it offers the possibility of visualizing the unsupervised concepts that contributed to the final prediction. In contrast to ACE and ICE approaches, in CRAFT, sub-regions are identified with a random crop technique. The activations from an intermediate layer of a neural network are calculated for each random crop. These are then factorized using the NMF technique to obtain concept representations similar to ICE. Selecting the right layer to analyze is critical because concepts can manifest across multiple layers. To address this, CRAFT introduces a recursive approach to identifying concepts and sub-concepts. Concept importance for the final prediction is determined using *Sobol* scores derived from sensitivity analysis. For concept visualization, CRAFT employs a saliency map starting from the concept representation and highlighting the most influential parts in the input image for the concept. A Human evaluation reports CRAFT's higher utility than standard XAI methods (e.g., Saliency, GradCAM) and ACE in two out of three scenarios. To assess the meaningfulness of extracted concepts, participants were required to detect intruder concepts and identify parent concept vs sub-concept identification to evaluate the effectiveness of the recursive approach.

Multi-dimensional Concept Discovery (**MCD [107]**) introduces multidimensionality in concept extraction by allowing concepts to span on a hyperplane across different convolutional channel directions. This is obtained by means of Sparse Subspace Clustering (SSC) for feature vector clustering and a subsequent PCA for cluster basis derivation. On top, a linear layer is employed (similarly to ICE) to mimic the final classification. This allows each sample to define a concept completeness score, quantifying the fraction of the model prediction that can be collectively recovered using the extracted concepts. The MCD framework offers two explanations: concept relevance maps, indicating a concept's importance for class predictions, and concept activation maps, for concept visualization. In the experiments, the authors show that MCD with an SSC clustering strategy overcomes other concept discovery methods (ICE, ACE) regarding concept faithfulness and conciseness. These are respectively defined as the reduction of the prediction when flipping the most important concepts and as the number of concepts required to reach a certain completeness score.

The objective of **DMA & IMA [64]** is to devise a concept discovery method characterized by *identifiability*, extending the idea from PCA and Independent Component Analysis (ICA) to conceptual explanations. The authors define a method as *identifiable* when it can provably recover all the concepts in a context where the components generating the data are known. They propose Disjoint Mechanism Analysis (DMA), which ensures identifiability in compositional data generation where distinct components control distinct pixel groups in an image. DMA guarantees the recoverability

even if ground truth components are correlated (unlike PCA or ICA). However, it requires components to affect disjoint parts of the input, which should be true as long as the explained encoder is faithful to the generative process generating input data. The second approach, Independent Mechanism Analysis (IMA), represents a less constrained version requiring concept vectors to be only orthogonal but not disjoint. The results demonstrate that IMA and DMA yield results comparable to PCA or ICA on datasets without correlated generative components and better on the contrary. Also, discovered concepts represent ground truth concepts more often than ACE and Completeness-aware. However, it’s worth noticing that DMA and IMA do not provide any of the explanations defined in Section 2.2: they only assess which concepts can be extracted from a neural representation.

Spatial-temporal Concept-based Explanation (**STCE [52]**) provides concept-based explanations for 3D ConvNets, a domain less explored due to the computational complexity of handling video data. STCE directly extends ACE’s 2D image strategy to the video domain without neglecting the spatio-temporal dimension. The method involves segmenting videos into supervoxels for each class, extracting features with a 3D ConvNet, and clustering similar supervoxels to define unsupervised spatial-temporal concepts. Concept importance is determined using T-CAV scores, considering, as negative samples, random videos from non-related datasets. Similarly to previous works, the quality of the STCE explanation is assessed in terms of faithfulness by measuring the reduction of accuracy of the explained model when removing concepts associated with high T-CAV scores. For qualitative analysis, the authors present video frames illustrating detected concepts, particularly by highlighting video frames that feature the most and least important concepts for a given classification

4.3 Discussion

In the following, we report the critics that have been moved to the concept-based methods reviewed in this section and how, in particular, they affect their performance.

Do standard networks learn concepts?. Several works in literature have disputed it [22, 61] whether standard models actually learn symbolic concepts when trained to accomplish a different task. In particular, [61] tested the concept accuracy of a linear probe. We remember that a linear probe is a linear layer trained to predict the concepts in the latent space of an already trained model, as proposed in [56]. They found that linear probes achieve a lower concept accuracy with respect to a model trained to predict the concepts (91% vs 97% on the CUB dataset [] and 0.68 vs. 0.53 error on the OAI dataset). This result confirms that post-hoc concept-based explainability methods cannot ensure that the model actually learns symbolic concepts, as introduced at the beginning of this section. On the other side, they help to understand what the network is learning: the still high accuracy of linear probes implies that the model learned many concepts.

Different probe datasets provide different explanations. In [90], the authors tested the importance of the dataset choice in extracting the explanations. They have shown that different probe datasets lead to different explanations even when explaining the same model using the same interpretability method. This issue arises for methods explaining both node-concept associations and class-concept relations. As an example of the first type, they show that ND label neurons differently according to the dataset employed (ADE20k or Pascal): while in some cases, the concepts are still similar (e.g., *plant vs potted-plant, computer vs tv*), in other cases are profoundly different (e.g., *chair vs horse, tent vs bus*). As an example of the second type, they show that the T-CAV concept activation vectors extracted from different datasets are projected in different directions. To numerically assess this, they compute the cosine similarity between the vectors representing the same concept but extracted from ADE20k and Pascal, and they show that while for same

concepts the similarity score is quite high (*ceiling* 0.27), for others it is very low (*bag* 0.01, *rock* -0.02), with a cosine similarity of common concept vector lower than 0.1 on average. As we also highlighted in Section 3, they suggest that the probe dataset should be chosen carefully and specifically recommend probe datasets that closely resemble the data distribution of the dataset on which the model was originally trained.

Post-hoc concept-based explanation methods are vulnerable to adversarial attack. In [18], Concept-based post-hoc explainability methods, much like their standard XAI counterparts [5, 36, 59, 103], are susceptible to adversarial perturbations. This study focuses primarily on supervised methods that explain class-concept relationships using linear probes, such as TCAV. The underlying idea is that a sample can be crafted to change the interpretation while maintaining the same classification. In particular, they show that an attack can induce a positive interpretation with respect to a concept, such as making a *Corgis* important for the classification of a *Honeycomb*. At the same time, they also provide examples of negative interpretation, such as making *Stripes* not important for classifying a *Zebra*. Notably, the authors also demonstrate that these attacks maintain some level of effectiveness even in black-box scenarios, i.e., when the adversarial samples are created for a different model to the one they are tested on, which is a more realistic scenario.

5 EXPLAINABLE-BY-DESIGN CONCEPT-BASED MODELS

Explainable-by-design concept-based models introduces a deviation from standard neural network training practices. They explicitly represent concepts within the same neural network architecture. We define them as explainable-by-design since they inherently provide Node-concept associations (N-CA) (except for some cases). They typically employ a dedicated hidden layer to predict concept scores to this aim. These scores represent numerical values quantifying the relevance or presence of specific concepts in a given input sample and condition the model’s output. The same concept scores can also be studied in relation to the output classes to define Class-Concept Relations (C-CR) and can be visualized through Concept-Visualization (C-Viz).

How the intermediate representation is defined varies according to the concept annotations and, consequently, the concept types. As shown in Table 3, Concept-based models either employ a dataset with annotated concepts (Supervised, Section 5.1), without concept annotations (Unsupervised, Section 5.2), with annotations for few concepts only (Hybrid, Section 5.3) or they generate the supervision by means of an external model (Generative, Section 5.4). Furthermore, if they employ annotated symbolic concepts (Symb.), they differ according to whether they employ them during training (Joint) or afterward to instill them in a trained network (Instil). If they automatically extract concepts, they differ for the type of concepts extracted, either clusters (Uns. Basis) or prototypes (Proto.). If they adopt a hybrid approach, they generally employ both symbolic and unsupervised basis concepts. Finally, if they generate concept annotations, these are normally Textual concepts. Concept-based models have been mostly developed to solve classification tasks (CLF), with one also performing regression (REG) and another clustering (CLU). Similar to the post-hoc method, in this case, the network’s backbone is mostly a CNN working on images (IMG), possibly reconstructing the data with an auto-encoder (AE). Some works employ a fully-connected network working on tabular data (TAB), GNN on graphs, and transformer (TRANS) working on both images and texts (TXT). To summarize the key aspects of concept-based models, we now discuss their main advantages and disadvantages compared to post-hoc methods.

Advantages. Concept-based Models offer the advantage of ensuring that the network learns a set of concepts explicitly. Furthermore, a domain expert can modify the predicted values for specific concepts and observe how the model’s output changes in response. This process is referred to as *concept intervention*, and it enables the generation of counterfactual explanations.

Table 3. Explainable by-Design Concept-based Models. We characterize the approaches based on the following dimensions. Full category description in Section 3.

	Method	Concept Employ.	Concept type	Scope	Expl Type	Data type	Loss	Task	Network Type
Supervised	CBM [61]	Joint	Symb.	N-CA, C-CR	GL	IMG	✓	CLF, REG	CNN
	LEN [12, 24]			N-CA, C-CR	LO & GL	TAB, IMG, TXT	✓	CLF, CLU	FCN, CNN
	CEM[32]			N-CA	GL	TAB, IMG	✗	CLF	FCN, CNN
	ProbCBM [57]			N-CA, C-CR	GL	IMG	✗/≈	CLF	CNN
	DCR[11]			N-CA, C-CR	LO	TAB, IMG, GRA	✗	CLF	FCN, CNN, GNN
	CW [22]	Instill	Symb.	N-CA	GL	IMG	✓	CLF	CNN
	CME [54]			N-CA, C-CR	GL	IMG	✓	CLF	CNN
	PCBM [118]			N-CA, C-CR	GL	IMG	✓/≈	CLF	CNN
	CT [92]			N-CA, C-CR	LO & GL	IMG	✗	CLF	TRANSF
	Interp. CNN [121]			N-CA, C-Viz	GL	IMG	✓/≈	CLF	CNN
Unsupervised	SENN [8]	-	Uns. Basis	N-CA, C-CR	LO	IMG	✗/≈	CLF	CNN + AE
	BotCL [109]			N-CA, C-CR	LO	IMG	✗/≈	CLF	CNN + AE
	SelfExplain [89]			C-CR	LO & GL	TXT	✗	CLF	TRANSF
	PrototypeDL [66]	-	Proto.	N-CA, C-CR, C-Viz	GL	IMG	✗	CLF	AE
	ProtoPNet [21]			N-CA, C-CR, C-Viz	GL	IMG	✓	CLF	CNN
	ProtoPool [94]			N-CA, C-CR, C-Viz	GL	IMG	✗	CLF	CNN
	Def. ProtoPNet [29]			N-CA, C-CR, C-Viz	GL	IMG	✗	CLF	CNN
	HPNet [45]			N-CA, C-CR, C-Viz	GL	IMG	✗	CLF	CNN
	ProtoPShare [95]			N-CA, C-CR, C-Viz	GL	IMG	✗	CLF	CNN
	ProtoPDebug [16]			C-Viz	GL	IMG	-	CLF	CNN
Hybrid	CBM-AUC [98]	-	Uns. Basis, Symb.	N-CA, C-CR	GL	IMG + VID	✗/≈	CLF	CNN
	Ante-hoc [97]			N-CA	LO & GL	IMG	✓	CLF	CNN + AE
	GlanceNets [75]			N-CA, C-CR	GL	IMG	✗	CLF	CNN + VAE
Gen.	LaBO [116]	-	Textual	C-CR	GL	IMG	✗/≈	CLF	CNN+LLM
	Label-free CBM [80]			C-CR	LO & GL	IMG	✓	CLF	CNN+LLM

Disadvantages. On the other hand, these methods can only be employed when there is flexibility to train a model from scratch, possibly tailoring it to the specific task at hand. Furthermore, in simpler solutions, the predictive accuracy of concept-based models may be lower than that of standard black-box models. Further issues and challenges of concept-based models will be analyzed at the end of the section (Section 5.5).

5.1 Supervised Concept-based Models

Supervised Concept-based Models employ a dataset annotated with symbolic concepts (like attributes of the output classes) to supervise an intermediate layer explicitly representing the concepts. When these annotations are readily available within the same training dataset, a *joint concept training* strategy can be employed (Section 5.1.1). Otherwise, they can also be embedded in the model through separate training on an external dataset dedicated solely to concept learning to perform a *concept instillation* (Section 5.1.2).

5.1.1 Joint concept training. As shown in Figure 7, supervised concept-based models employ concept annotations to supervise an intermediate representation and increase the transparency of the model. Unlike methods performing concept instillation, the training is performed *jointly* for the target task (e.g., image classification) and for concept representation. This requires employing a training dataset that includes annotations about both the main classes (e.g., Parrot) and the related attributes (e.g., Feather, Beak, Foot, Muzzle). Concepts can be represented by means of a single neuron (CBM, LEN) or through an embedding (CEM, ProbCBM, DCR), which avoids performance loss with respect to a black-box model. As shown in the picture, these methods can also analyze the relations between concepts and classes through concept importance (CBM, ProbCBM) or logic relations (LEN, DCR).

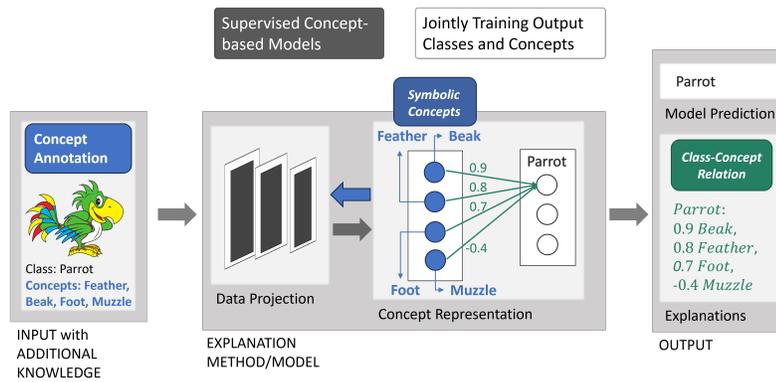


Fig. 7. Supervised Concept-based Models jointly training output classes and concept representation on a dataset where both class and concept annotation is provided.

The first work proposing a joint task-concept training strategy is attributed to Concept-Bottleneck Model (CBM [61]), where the authors modify the typical end-to-end neural architecture mapping from an input to a target by introducing an intermediate *concept bottleneck* layer. This layer has neurons dedicated to learning and predicting a set of human-specified concepts from input features. A task function then predicts the final target, relying only on concept prediction (concept scores). Specifically, a double loss is optimized, penalizing wrong predictions in both the output and concept bottleneck layers. CBM design enhances the model’s transparency and enables interaction with the model itself. A domain expert may intervene in the predicted concepts, allowing potential adjustments to the model’s predictions. This process is referred to as *concept intervention*. Concept interventions also facilitate the extraction of counterfactual explanations by manipulating the concept predictions and observing the model’s response on the counterfactual concept representation. CBMs, however, are affected by two main limitations. i) CBMs provide class-concept relations only when employing a single layer to model the task function. In this case, the importance of a concept corresponds to the weight connecting the concept to the final class, as shown in Figure 7. ii) The generalization capability of CBM is lower than standard end-to-end models. This is mainly due to the bottleneck layer where each concept is represented by a single neuron [74]. This limit becomes more severe in contexts with few concepts available and when employing a shallow task function.

In a series of works, the authors address the first issue by means of a framework named Logic Explained Networks (LENs [12, 24]). LENs improve the explainability of CBMs. These models, indeed, are trained to provide explanations in terms of First-Order-Logic (FOL) formulas associating task functions with the concept having the highest influence on it. To achieve this goal, they impose architectural constraints enforcing sparsity in the weights of the task function. As an example, on a bird classification task, a LEN may provide explanations of the type $billShapeHooked \wedge mediumSize \wedge \neg throatColorWhite \Leftrightarrow blackFootedAlbatross$. While most of the other works focus on vision tasks, LENs can be applied to tabular data and textual data [51] without a concept extractor, i.e., without requiring a module explicitly mapping the non-interpretable input data to a concept space. LEN can consider the same input data as concepts when they are associated with human-understandable symbols. Furthermore, LEN can solve clustering tasks and provide logic explanations about cluster assignment in terms of frequent concept relations [24]. Finally, logic explanations can be

used to defend a LEN from an adversarial attack by checking the consistency of the test-time predictions with the explanations extracted at training time.

Concept Embedding Models (**CEM [32]**) improve CBM performance by using an entire set of neurons to represent a concept (concept embedding). CEM obtains the same classification accuracy of black-box neural networks while retaining the possibility to employ and interact with a concept-based model. Unlike CBM, CEM performance does not decrease when the number of employed concepts decreases, showing a very high concept efficiency. A naive solution to the performance loss issue of CBM is adding a separate set of unsupervised neurons [74], encoding information about the task function. However, as shown in CEM, this hybrid solution loses the possibility of interacting with the network. According to CEM authors, in the hybrid configuration, the task functions only rely on unsupervised neurons, nullifying the effects of concept interventions on the task function. In CEM, the authors also introduce the Concept Alignment Score to measure the quality of a concept representation in describing a set of concepts. More precisely, it measures the entropy reduction in the concept space when a concept representation is provided. They show that CEM concept alignment to the real concept distribution is close to the one provided by CBM and much higher than the hybrid solution.

When concepts used to make the final predictions are not visible in the input sample, a deterministic prediction of the concepts may damage both the final task prediction and the explanation based on the concepts. To solve this problem, the Probabilistic Concept Bottleneck Model (**ProbCBM [57]**) employs probabilistic concepts to estimate uncertainty in both concept and class predictions. Similarly to CEM, ProbCBM also employs concept embedding. However, given an input, ProbCBM outputs mean, and variance vectors for each predicted concept. Concepts are then represented as multivariate normal distributions. The probability of each concept is calculated via Monte-Carlo estimation by sampling from the distribution and computing the mean distance from the true and false concept vectors. Class predictions are similarly computed, generating a class embedding from the concept embeddings. The uncertainty quantification in ProbCBM is done by calculating the determinant of the covariance matrix, which represents the volume of the probabilistic distribution. ProbCBM’s performance is compared with CBM [61] and CEM [32]. ProbCBM outperforms CBM in classification accuracy but falls short of CEM. Regarding concept prediction accuracy, it is on par with CBM and better than CEM, demonstrating that integrating uncertainty through probabilistic embedding doesn’t compromise performance.

Even though CEM and Prob-CBM enhance the performance of standard CBMs, they worsen the underlying explainability. When employing embedded concepts, even with the use of a shallow network, the interpretability of the task function is compromised since the dimensions within a conceptual embedding lack a direct association with a symbolic meaning. In response to this challenge, Deep Concept Reasoner (**DCR [11]**) is introduced as an interpretable concept embedding model. Starting from concept embeddings, DCR learn differentiable modules which output for each class and each sample a different fuzzy rule. These rules are then executed, taking into account the concept activations. As a result, the model outputs an interpretable prediction regarding the concepts, even though only locally. DCR achieves task accuracies on par with CEM and higher than CBM. The authors also show that in contexts where ground truths rule are available, DCR can recover them. Compared with LIME and LEN, DCR facilitates the generation of counterfactual examples, and extracted explanations are less sensitive to input perturbations. Finally, they showcase that DCR can also work on graphs employing GNN as a backbone. In this case, they extract concepts post-hoc from a trained GNN by means of a Graph-Concept Explainer [73], a variant of ACE for graph structure.

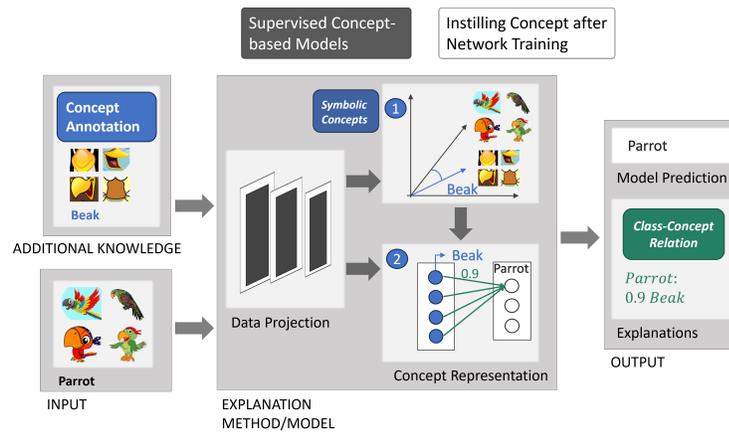


Fig. 8. Supervised Concept-based Models instilling concept in the network after a first end-to-end network training. Concept representation is either extracted from the network itself with post-hoc C-XAI techniques or from an external source (e.g., a knowledge graph).

5.1.2 Concept Instillation. This class of approaches considers the case where concept information is not available in the training set, but it is available from an external data source, such as a separate supervised dataset (CW, CT), a semi-supervised one (CME), or a knowledge graph (PCBM). As shown in Figure 8, the process is composed of two steps: 1) a black-box network is standardly trained end-to-end; 2) through *Concept Instillation*, a given layer of the network is modified to represent concepts. At test time, this process still allows the employment of a concept-based model. Similarly to jointly trained models, concepts can be represented by means of single nodes (CW, CME, PCBM) or through concept embeddings (CT). They also provide class-concept relations if they replace the classifier head with an interpretable task predictor (CME, PCBM, CT).

Concept Whitening (CW [22]) modifies a hidden layer of a neural network after it has been trained in such a way that it also predicts a given set of concepts. Usually, CW can be used to replace the Batch Normalization (BN) layer in a CNN because it combines batch whitening (decorrelating and normalizing each dimension) with an extra step involving a rotation matrix to align the concepts with the axes. The underlying idea is that axes of the latent space should align with concepts, allowing each point within the space to correspond to and be interpreted through known concepts. Further training on the standard classes is required to avoid loss in the final task performance. An interesting aspect of CW is its capability to determine how a network captures a concept across different layers, as it can be applied after any layer. The authors show that a network produces more abstract and fundamental representations in the lower layers. The authors compute the concept purity as the 1-vs-all AUC score on the test concept predictions to evaluate interpretability. A comparative analysis against concept-based post-hoc methods (TCAV, IBD) indicate that the concepts acquired with CW exhibit higher purity. This result is attributed to the inherent orthogonality of CW conceptual representations and the impact of the entire model fine-tuning process.

Similarly to CW, the Concept-based Model Extraction (CME [54]) also proposes to derive an interpretable model from a black-box. CME trains a concept extractor starting from the hidden representation of a given DNN. To reduce the concept annotation effort, CME employs only a small annotated dataset and trains the concept extractors in a semi-supervised fashion. Interestingly, the authors repeat the extraction process of each concept for different hidden

layers, retaining the one providing the higher concept accuracy. The predicted concept scores are employed as input for an interpretable task predictor, which is trained to mimic the task prediction of the original model. For the task function, the authors employ Decision Trees and Logistic Regression. Since they directly work on concepts, they both allow concept intervention at test time. CME’s performance is compared with Net2Vec [34] and CBM [61] frameworks. CME significantly outperforms Net2Vec in task accuracy while remaining comparable to CBM. The authors also introduce a new metric, the MisPrediction Overlap (MPO), to compute the portion of samples with at least m concepts inaccurately predicted. Based on MPO, CME performs similarly to Net2Vec but falls behind CBM.

Post-hoc Concept Bottleneck Models (**PCBM [118]**) implement two modalities to extract concepts: (i) from other datasets, similarly to CW [22], and (ii) from a knowledge graph and then encoded in the network by means of a textual encoder. In the first scenario, Concept Activation Vectors (CAVs) are employed to learn the concept representations in the latent space of a model pretrained on the final task. In the latter scenario, the authors employ an open-access knowledge graph (e.g., ConceptNet [105]), to gather concepts in relation with the queried classes. They then obtain a vectorial representation of the concepts by means of a text encoder. In both scenarios, the concepts are projected in the latent space of a pretrained model. The projections then train an interpretable task predictor, providing class-concept importance. In addition to concept intervention, PCBM also allow for *global model edits*, global modifications, which enables the model to adapt its behaviour to novel scenarios, including shifts in data distributions. Similarly to [74], the authors propose to employ an unsupervised set of neurons to improve network scaling. This time, however, they propose to employ a residual fitting, such that the hybrid prediction still has to rely on the interpretable one. With respect to the original model, PCBM performances are slightly worse, but still comparable.

ConceptTransformer (**CT [92]**) substitutes the task classifier of a DNN with an attention module enriched with concept embeddings. Instead of being extracted from a knowledge base as in PCBM, concept embeddings in CT come from a transformer model trained on a dataset annotated with user-defined symbolic concepts. The output of the attention module represents concept scores. A linear layer is employed on top to predict the final classes. CT provides Class-Concept Relation in terms of the attention scores between the output class and the concepts embeddings. However, CT does not provide Node-Concept Association since it does not explicitly represent concepts by means of a set of nodes but only with concept embeddings. The authors claim that their model is faithful by-design since it employs a linear layer from concepts to classes, similar to CBM. To obtain more plausible explanations, they also propose supervising attention scores so that concept-output relationships align with domain knowledge. CT classification performance matches end-to-end black box models. Leveraging the concepts with CT seems to even boost the performance in contexts where the set of concepts is well-defined, and the relationships with the final classes are known (as in the case of CUB-200-2011 dataset [110]).

5.2 Unsupervised Concept-based Models

Unsupervised concept models (Section 5.2), modify the internal representation of a network to autonomously extract concepts without explicitly associating them with pre-defined symbols. In some cases, these models jointly train an intermediate layer with an unsupervised learning technique to extract an *unsupervised concept basis*, i.e., a clusterized representation identifying groups of similar samples (Section 5.2.1). Otherwise, they explicitly require the network to represent in the layer weights *prototypes-concepts*, (parts of) frequently seen input samples (Section 5.2.2).

5.2.1 Unsupervised Concept Basis. The idea behind these approaches is to learn unsupervised disentangled representations in the latent space of a model. These representations must group samples with some characteristics, such as

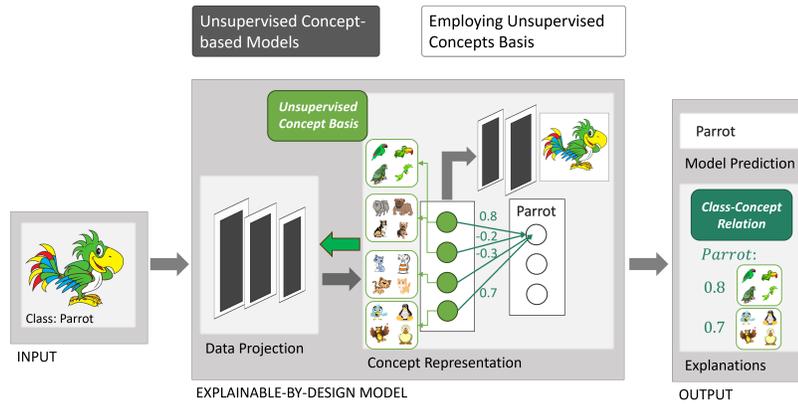


Fig. 9. Unsupervised concept-based model extracting unsupervised concept basis. How the concepts are extracted varies between the models. Here, we report the case in which a decoder branch is employed over the concept to reconstruct the input and extract meaningful concepts (SENN, BotCL).

a generative factor. This class of approaches takes inspiration and shares ideas with the disentangled representation learning field [14]. To extract unsupervised concepts, these models are jointly trained to reconstruct the input by means of a decoder branch (SENN, BotCL, and reported in Figure 9) to maximize the mutual information (Interpretable CNN) or via self-supervised loss over textual data (SelfExplain). When these models employ an interpretable task predictor over the concepts, they also provide class-concept relations (SENN, BotCL).

In **Interpretable CNN** [121], convolutional filters act as unsupervised concepts and are required to represent different object parts. The assignment occurs during the learning process without employing any concept annotation. The authors devise an unsupervised loss function to achieve this, maximizing the mutual information between the image and the filter activations. Also, the loss promotes low entropy in both the inter-category activations and the spatial distributions of neural activations. The idea is that filters should create a clustered representation of the concepts. Each filter should encode a distinct concept associated with a single object category, possibly not repeated in different image regions. In the evaluation, the authors measure the IoU score with ground-truth object bounding boxes, similar to ND. They also measure the part location stability, assessing the consistency in how a filter represents the same object part across various objects. The experimental findings show that Interpretable CNNs exhibit superior interpretability and location stability compared to standard CNNs. Regarding classification accuracy, Interpretable CNNs are on par with simple CNNs.

Self-explanatory neural network (**SENN** [8]), similarly to Interpretable CNN, also employs an unsupervised approach to create an explainable-by-design network, this time based on an auto-encoder architecture. SENN first employs a concept encoder to derive a clustered representation of its features from the input. Second, it employs an input-dependent parametrizer that generates class-concept relevance scores. The output is given by the linear combination of these basis concepts and their relevance scores. The degree of interpretability of SENN is intrinsically tied to the interpretability of the underlying concepts. SENN promotes concept diversity to this aim: inputs should be representable using a few non-overlapping concepts. This is achieved by making the concept encoder sparse. Also, to visualize the identified concepts, SENN provides a small set of training examples that maximally activate it. The authors show that SENN explanations compared to standard XAI techniques (i.e., LIME [91], SHAP[70]) result in more: i) explicit, SENN

explanations are more immediate and understandable; ii) faithful, SENN importance scores are assigned to features indicative of ground-truth importance; iii) stable, since similar examples yield similar explanations for SENN. SENN classification accuracy is slightly lower or on par with standard CNN according to the experimented scenarios.

Bottleneck Concept Learner (**BotCL [109]**) is based upon and improves the idea of SENN for extracting an unsupervised concept basis. In BotCL, the features extracted from the convolutional encoder are fed into a slot attention-based mechanism [65] predicting for concept scores. A set of vectors learnt by backpropagation represent the concept embeddings and are used as keys of the attention mechanism. A linear layer working over concept scores is employed as a task classifier, and it captures class-concept relationships. BotCL employs the reconstruction loss of SENN or a contrastive loss to learn meaningful concepts. The latter enforces similarity among concept activations belonging to the same class and dissimilarity otherwise. Also, BotCL employs two regularizations to facilitate training and ensure consistency in the learned concepts. The first is a consistency loss, enforcing similarity between the concept features that activate the concept in different images. The second is a mutual distinctiveness loss, requiring different concepts to cover different visual characteristics. BotCL task performance, on average, is higher than SENN, ProtoPNet, and Completeness-aware, particularly when employing the contrastive loss. BotCL extracted concepts, instead, are compared against ACE and Completeness-aware on a synthetic dataset annotated with concepts. They result in being good in terms of concept purity and efficiency and reconstruction error, particularly when employing many concepts.

SelfExplain [89] is a self-explaining model, providing local explanations within text classifiers as Class-Concept Relation. It extracts concepts in an unsupervised way directly from the input text. In particular, the concepts are here represented as the non-terminal leaves of the semantic tree parsing the text. SelfExplain incorporates two different layers into a neural text classifier to explain the predictions. Specifically, the Locally Interpretable Layer (LIL) compute the relevance score for all input concepts extracted from the input sample. The Globally Interpretable Layer (GIL), instead, interpret sample predictions by means of the concepts extracted from all the training data. The GIL layer brings the advantage of studying how concepts from the training set influenced the classifier’s decision. The concept scores of both layers are employed to re-predict the final class, but they are not employed for the standard prediction. As a baseline, RoBERTa and XLNet transformer encoders are used. GIL and LIL layers are incorporated into the transformers to make a comparison. The results reveal that incorporating GIL and LIL layers consistently leads to better classification across all evaluated datasets, probably due to the multitask loss employed. To gauge the efficacy of the explanations, the authors conducted a human evaluation, where explanations provided by SelfExplain are consistently acknowledged as more understandable, trustworthy, and adequately justifying model prediction than standard XAI techniques.

5.2.2 Prototype-based Concepts. Within the domain of unsupervised concept-based models, another possibility is to require the network to explicitly encode *prototype-based concepts*, which represent specific (parts of) train examples. As shown in Figure 10, the prototypes are compared to input samples, and their similarities are employed as sufficient statistics to provide the final prediction. How the prototypes are extracted is different among the methods, either by means of an autoencoder (Prototype DL), or per class by means of the same convolutional filters (ProtoPNets), sharing them among the classes (ProtoPShare, ProtoPool) with a hierarchy of prototypes (HPNet), with deformable filters (Deformable ProtoPNets) or through an interactive process with human experts (ProtoPDebug). In most cases, a linear layer is employed on top of the prototypes to classify the final class and provide an interpretable class-concept relation. All the methods provide different systems to visualize the concepts since prototype concepts particularly require a visualization.

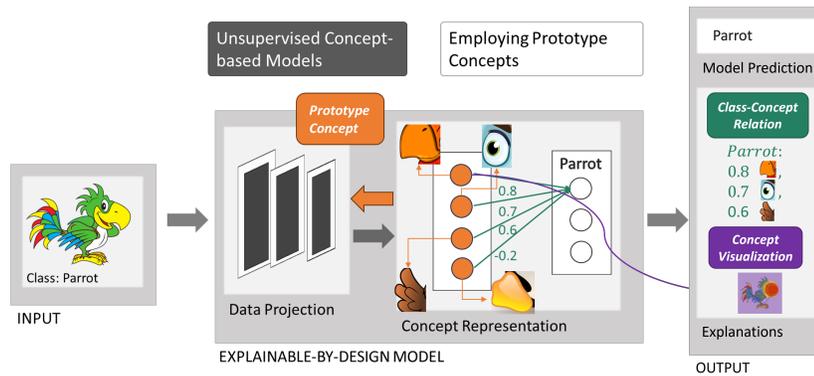


Fig. 10. Unsupervised concept-based models employing prototype concepts. The network explicitly encodes prototype concepts representing (part of) training examples. The explanations are typically provided as class-concept relation and/or concept visualization.

Prototype DL [66] introduces the use of prototypes in concept-based explainability. The proposed model explains predictions based on similarity to prototypical observations within the dataset. Its architecture comprises an autoencoder and a prototype-based classification network. The autoencoder aims to reduce input dimensionality and learn predictive features. The autoencoder's encoder facilitates comparisons within the latent space, while the decoder enables the visualization of learned prototypes. Given the encoded input, the prototype classification network generates a probability distribution across K classes. This network includes a prototype layer where each (prototype) unit stores a weight (prototype) vector mirroring an encoded training input. For model prediction interpretation, we can visualize the learned prototype by feeding prototype vectors into the decoder. Then, from the learned weights of the prototype layer, we can identify which prototypes are more representative of a class. The experimental results show that the proposed interpretable architecture does not compromise predictive ability compared to traditional and non-interpretable architecture.

Prototype Parts Network (ProtoPNet [21]) identifies different parts of the input image that look like some learned prototypical parts. Similarly to *Li et al. [66]*, ProtoPNet incorporates a prototype layer but employs convolutional layers for the prototype extraction process. Thus, ProtoPNet extracts prototypes based on segmented image parts rather than the whole images. Specifically, the convolutional layers extract a set of features for each image patch. ProtoPNet then computes the similarity between each prototype unit and all the image patches. This process yields a similarity score activation map (similar to a saliency map), where each value indicates the presence of a prototypical component for each part of the image. These similarity maps generated by individual prototype units are then condensed into a single similarity score for each prototype. The final prediction is based on a weighted combination of similarity scores according to the importance of a prototype for a given class. The resulting model provides explanations in three forms: **Class-Concept Relation** according to the weight connecting each prototype to a class, **Node-Concept Association** since each prototype represents a unit of the network, and **Concept Visualization** through the similarity activation map. ProtoPNet achieves a slightly lower accuracy (up to 3.5%) than non-interpretable baselines.

The **Hierarchical Prototype network (HPnet [45])** builds upon ProtoPNet but supports a hierarchical organization of prototypes. The model supports hierarchical classification, i.e., a classification based on a taxonomy at multiple granularities. The approach learns a set of prototypes for each taxonomy class, from coarser ones (e.g., animal) to finer

ones (e.g., cat). HPnet architecture, after the convolution layers as in [95], includes a prototype layer for each parent node of the taxonomy. Each layer includes and classifies the prototypes of the finer classes of the corresponding parent one. The design of the prototype layer is close to that of ProtoPNet and provides for the computation of the similarity scores of the learned prototype and patches of the input image. The approach can provide Concept Visualization via a heat map showing localized portions of the input image that strongly activate the prototype at multiple levels of the hierarchy. Hence, it enables the understanding of the prediction at multiple granularities. Thanks to its ability to handle hierarchical classification, HPNet addresses the novel class detection problem by classifying data from previously unseen classes at the appropriate level in the taxonomy. The experiments show that HPnet experiences a slight drop in fine-grained accuracy with respect to a black-box model.

Prototypical Part Shared network (**ProtoPShare** [95]) is an extension of ProtoPNet that introduces the idea of sharing prototypes between classes to reduce the number of prototypes. A reduced number of prototypes, which are, as in ProtoPnet, representations of prototypical parts, improves the interpretability of predictions. ProtoPshare addresses two issues of ProtoPNet. First, ProtoPNet derives many prototypes since each is assigned to only one class. Second, two prototype representations of two distinct classes with similar semantics can be far apart in latent space, making the final predictions unstable. ProtoPShare addresses these issues by introducing prototype sharing between the classes via prototype pruning. ProtoPShare has two phases: i) initial training and ii) prototype pruning. The initial training phase follows the one of ProtoPNet. The second phase consists of the *data-dependent merge-pruning*. Given the set of prototypes after initial training, this step merges a fixed percentage of most similar prototypes. Specifically, the work introduces a data-dependent similarity function to identify prototypes with different semantics, even if they are apart in the latent space. The pruning process proceeds by computing the similarity of pairs of prototypes. If a pair is among a given percent of most similar pairs, one prototype of the pair is removed. Even with fewer prototypes, ProtoPShare outperforms ProtoPNet in terms of accuracy.

ProtoPool [94], inspired by ProtoPNet [21], tackles some of its limitations and those of subsequent works. One limitation it addresses is the assumption of separate prototypes for each class in ProtoPNet. Like ProtoPShare [95], ProtoPool reduces the number of prototypes by sharing them across different data classes. However, ProtoPShare’s merge-pruning step leads to increased training time. ProtoPool avoids the pruning by adopting a soft assignment of prototype represented by a distribution and automatic, fully differentiable prototype assignment, thus simplifying the training process. Furthermore, ProtoPool introduces a focal similarity function that generates saliency maps more focused on salient features compared to using ProtoPnet similarity. Additionally, ProtoPool introduces Concept Visualization via prototype projection. This involves replacing the abstract prototype learned by the model with the representation of the nearest training patch. Compared to ProtoPNet and ProtoPShare, ProtoPool achieves the highest accuracy on one dataset using fewer prototypes and comparable but lower accuracy results on another one. A user study rated ProtoPool’s visualization as more interpretable than other compared methods.

Also **Deformable ProtoPNet** [29] builds upon ProtoPNet [21]. The work addresses the limitation of ProtoPNet and its subsequent works of using spatially rigid prototypes, which cannot explicitly consider geometric transformations or pose variations of objects. Deformable ProtoPNet uses spatially flexible deformable prototypes. Each prototype consists of adaptive prototypical parts that dynamically adjust their spatial positions based on the input image. The work builds upon previous work modeling object deformations and geometric transformations to deform prototypes, specifically on Deformable Convolutional Network [26]. The experimental results show that Deformable ProtoPNet achieves state-of-the-art accuracy.

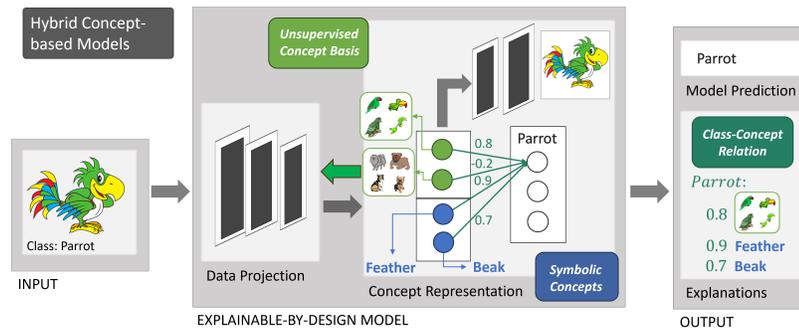


Fig. 11. Hybrid concept-based model employing both supervised and unsupervised concepts. This approach normally involves employing an unsupervised learning strategy (such as reconstructing the input image) to learn an unsupervised concept basis and a few supervisions over some neurons representing symbolic concepts.

ProtoPNets [21] may provide confounder factors as part-prototype explanations (e.g., textual metadata in image X-rays or background parts), limiting the model’s predictive performance and generalizability. **ProtoPDebug** [16] addresses this issue by proposing a human-in-the-loop concept-level debugger for ProtoPNets architectures. Human experts identify which part-prototypes are confounders by examining ProtoPNets’ explanations. This information is provided as a feedback to the model and it is fine-tuned to align with the user supervision. ProtoPDebug assesses all learned part-prototypes during each iteration, retrieves maximally activated training samples, and prompts user judgment on confounding. Confirmed confounders are added to the *forbidden concepts* set. Conversely, high-quality part-prototypes are incorporated into the *valid concepts* set. For the fine-tuning step, ProtoPDebug proposes two losses for the fine-tuning step. The *forgetting loss* penalizes forbidden concepts while the *remembering loss* encourages remembering valid ones. ProtoPDebug outperforms a predecessor debugger and ProtoPNet, enhancing test classification accuracy.

5.3 Hybrid Concept-based Models

Hybrid concept-based models propose the joint employment of both supervised and unsupervised concepts, as shown in Figure 11. The strategy normally involves the employment of an unsupervised learning strategy, either reconstructing the input image with an AE (Ante-hoc, and reported in Figure11) or a VAE (GlanceNet) or maximizing the mutual information (CBM-AUC), together with some supervision over some concepts. The goal of this integration is generally to allow the employment of concept-based models also in those scenarios where few supervised concepts are available for the task at hand. This approach has been adopted to improve the performance of standard CBM (CBM-AUC, Ante-hoc Explainable Model) and its capability to represent concepts (GlanceNet).

Concept Bottleneck Model with Additional Unsupervised Concepts (**CBM-AUC** [98]) is the first approach designed to merge CBM [61] with SENN [8], employing supervised and unsupervised concepts at the same time. CBM-AUC addresses the issue of CBM’s restricted concept efficiency. To address this, the authors suggest enabling the model to leverage the representation capability associated with unsupervised concepts. The first contribution is the introduction of M-SENN, a variant of SENN reducing the computational complexity and improving the scalability of SENN. First, it replaces the decoder of the AE architecture employed by SENN with a discriminator. This module now optimizes the mutual information between the input and the unsupervised concepts. Second, M-SENN parametrizer now shares the

intermediate network with the encoder, reducing the computational complexity. For integrating CBM into SENN, they add supervision on some concept basis, which then comprises unsupervised and supervised neurons. The results show that the modifications made to join SENN with CBM are effective. CBM-AUC outperforms the task accuracy of both SENN and CBM. Also, they analyze both supervised and unsupervised concepts by visualizing their activation through saliency maps. They show that both types of concepts look at semantically meaningful areas of the input images.

Similarly to CBM-AUC, also Learning **Ante-hoc Explainable Models via Concepts** [97] integrated supervised and unsupervised concepts starting from SENN and CBM. Differently from CBM-AUC, however, it still employs SENN concept decoder to learn unsupervised concepts. However, the most important difference between both models is that the classifier is added on top of the base encoder and not the concepts. This entails that the final classification is not based on the concepts. To minimize this issue, they add a classifier to the concept encoder and minimize predictions' divergence via a fidelity loss. As presented, the model is designed to operate in unsupervised concept learning mode. A loss function is employed to align the learned concepts with the available annotations to allow the model to learn supervised concepts. Otherwise, self-supervision can also be integrated as an auxiliary task, such as predicting the rotation of the input images starting from the encoded concepts. The proposed framework demonstrates competitive performance across both annotated and not annotated datasets, improving the task accuracy of CBM and SENN, respectively.

Unlike previous work, **GlanceNets** [75] aims to improve the concept representations of SENN and CBM rather than the performance. It improves the interpretability of SENN by employing a VAE instead of a standard AE to learn disentangled concepts. Due to the variational approach, learned concepts can better represent sample generative factors and respond to their alterations. When supervised concepts are available, they force a mapping with the learned concepts to preserve known semantics. They also improve the representations of CBM by identifying a solution to the concept leakage issue (see Section 5.5 for a description). They claim it is a deficiency in the out-of-distribution generalization of concept-based models. Therefore, GlanceNets employs an open-set recognition technique at inference time to detect instances that do not belong to the training distribution. They show this solution substantially improves the concept alignment with respect to CBM. Task performances are on par with CBM and SENN.

5.4 Generative Concept-based Models

Rather than employing a set of symbolic supervised concepts or extracting an unsupervised concept basis, generative concept-based models have recently shown that they can directly create concept representations. As shown in Figure 12, this is generally performed by employing a Large Language Model (LLM) producing *Textual Concepts*, short textual descriptions of the final class. The embeddings representing the textual concepts are aligned to the latent input representation to output concept scores. These scores are then employed to provide for the final classification, possibly with an interpretable classifier providing class-concept relations. At test time, these models predict both the final class and the most suitable descriptions among those that have been learned. In the following, we outline two methods (LaBO [116] and Label-free CBM [80]) of this recent line of work. The two methods are similar in their main components since they both employ an LLM in combination with a vision-language one and provide class-concept relation explanations.

Language-guided Bottlenecks (**LaBO** [116]) is the first model proposing this idea. It exploits the Generative Pre-trained Transformer 3 (GPT-3) language model [19] to generate factual sentences to derive concepts. This approach circumvents the requirement for costly manual annotation of concepts. In the LaBO framework, the process initiates by prompting the GPT-3 model with queries designed to generate a collection of *candidate concepts* for each class. An

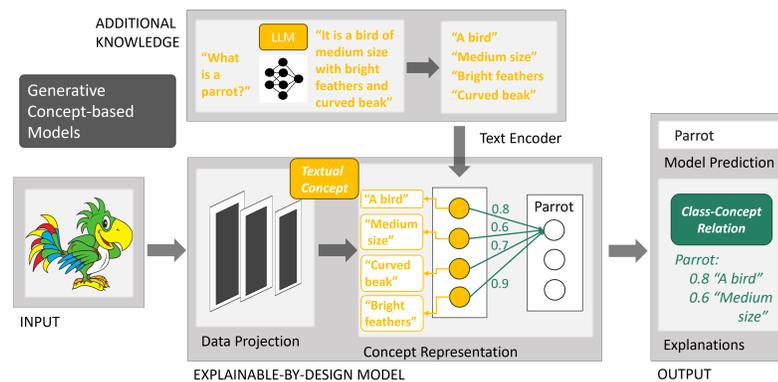


Fig. 12. Generative concept-based models employ an external LLM to generate textual concepts, i.e., short descriptions of the final classes without requiring a concept annotation. The embeddings representing the textual concepts are aligned with the input latent representation. The alignment produces the concept scores, which are used to provide the final classification.

example of a prompt is: "Describe what the <class_name> looks like." Subsequently, the approach applies a submodular optimization technique to select a subset of the generated sentences to maximize the discriminability and diversity of the concepts. The generated textual concepts are then aligned to images using a pretrained alignment model to form a bottleneck layer. Specifically, LaBo uses CLIP [88] to derive embeddings of concepts in textual form and images. Subsequently, a linear layer is applied to the similarity scores of concepts and images to learn a weight matrix whose values represent the importance of the concepts in the classification task. The effectiveness of LaBo is intrinsically linked to the representation of the classes and concepts within the adopted LLM. LaBo retains accuracy in limited training data scenarios, but as data volume grows, its performance slightly decreases while remaining competitive compared to black-box models. Compared to PCBM [118] as an interpretable baseline, LaBo exhibits a substantial performance improvement.

Label-free Concept Bottleneck Model (**LabelFree-CBM [80]**) is a framework close to LaBo [116]. As LaBo, LF-CBM can convert a neural network into an interpretable CBM without requiring annotating concepts but by generating them using a language model. Its main steps also resemble LaBo ones. It generates the concepts using GPT-3 [19] and filters them to enhance the quality and reduce the concept set's size. Then, it computes embeddings from the network backbone of input data and the concept matrix. The concept matrix is derived using a visual-text alignment model (CLIP-Dissect [81]). Each element of the matrix is the product between the CLIP encoding of the images of the training set and the concepts. The framework then involves learning the concept weights to create the Concept Bottleneck Layer. Lastly, the final predictor is learned by training a fully connected sparse layer. The authors assessed the interpretability of the model with a user study. Regarding performance, LF-CBM demonstrates only a small loss in task accuracy compared to the original neural network backbones while it outperforms the interpretable PCBM [118].

5.5 Discussion

In the following, we analyze critical aspects of explainable-by-design models. We focus on three perspectives: (i) performance loss, (ii) information leakage in concept representation, and (iii) pitfalls of concept intervention.

Performance loss. Concept-based models explicitly represent a set of concepts within the model architecture and rely on these concepts for the final task predictions. On the other side, the higher transparency of these architectures comes at the cost of a performance loss compared to traditional, black-box architectures. We analyze this aspect in Table 3. Still, we note that several approaches (that we report with \times) match the performance of complex black-box prediction models, with the advantage of enabling decision explanations through concepts.

Information leakage in supervised concept representation. Recent works revealed that supervised concept-based models encode additional information in the concept space other than the concept information itself [46, 74, 76]. Concept representations may include representations of the task labels when trained jointly with the downstream task. This phenomenon is denoted as *information leakage in concept representation*, and compromises the model’s interpretability and intervenability (i.e., the possibility to perform concept intervention) [46, 74, 76]. From the interpretability perspective, we cannot directly interpret the association between the concept representations and the model decision since the outcomes may depend on other encoded information. From the intervenability one, interventions may remove the extra information encoded and can decrease the final model performance. Mahinpei et al. [74] argue this issue is shared among concept-based models that use soft or fuzzy representations to represent concepts and experimentally demonstrate it for CBM [61] and CW [22]). Mahinpei et al. [74] also show that information leakage is not avoided even when applying mitigation strategies, including (i) training concept representations and the concept-to-task sequentially (rather than jointly), (ii) adding unsupervised concept dimensions to account for additional task-relevant information, or (iii) via concept whitening [22] which explicitly decorrelate concept representations during training [74]. Hence, they conclude that soft representations encode concept information and data distribution. Recent works specifically address leakage issues in concept-based models, such as [46] for CBMs and GlanceNets [75].

Pitfalls of concept intervention. Concept-based models enable human experts to intervene in the concept representation and modify them at test time. Concept intervention improves final task performance when intervening on mispredicted concepts and generating counterfactual explanations. However, concept intervention techniques might be affected by major pitfalls. Shin et al. identify two major issues [101] on CBMs [61]. First, intervention could increase the task error, compromising the reliability of intervention techniques. This decrease in performance could happen when users nullify void concepts by setting unsure concepts to zero. The authors justify this phenomenon, noting that the absence of a concept in a specific image does not entail the concept is not associated with its class. We argue this is also linked to the soft representation of concepts and leakage issues. The second issue arises from the practice of using majority voting to associate concept values to classes to improve task performance (e.g., in [32, 46, 61]). Majority voting affects a fair representation of minority samples and makes task outcomes biased toward the majority. The authors argue that intervention can exacerbate this unfair treatment since it would increase minority performance only when changing the concept value to the majority one.

6 EVALUATION

This section overviews the quantitative and qualitative metrics, the datasets, and the resources mainly used in the literature to assess the performance of C-XAI methods. In Section 6.1, we delve into the quantitative metrics employed to evaluate the methodologies. This includes measures of the effect of concepts on task predictions and performance, and the quality of concepts. Conversely, Section 6.2 elaborates on the qualitative evaluations which leverage human assessments. Section 6.3 provides an overview of the datasets utilized to assess the effectiveness of C-XAI methods. Finally, in Section 6.4, we provide details about the resources accessible for each method.

6.1 Quantitative evaluation

We categorize the quantitative metrics used and proposed in the reviewed papers into two categories according to their final purpose and what they measure. The metrics of the first category assess how well concepts contribute to class predictions and task performance. The second category focuses on the intrinsic quality of concepts. For a more detailed explanation of each single metric, please refer to AppendixC.

Concept effect on class prediction and task performance. These metrics assess how concepts contribute to (i) the final class prediction or (ii) the task performance.

(i) *Concept effect on class prediction.* This set of metrics assesses the effect of the concepts on the final class prediction. Concept-based methods often naturally provide the concept relevance (importance) by means of the concept weight connecting the concept to the final classes. Hence, these metrics have been mostly employed for post-hoc explanation methods providing Class-Concept Relation to study the impact of the identified concepts on task prediction. Among these metrics we find the T-CAV score [37, 56], CaCE [41, 111], ConceptSHAP [117], Smallest Sufficient Concepts (SSC) [37], Smallest Destroying Concepts (SDC) [37, 107], STCE Concept Importance [52], T-CAR score [25], Sobol score [33], Concept Weight [116, 118], Concept Relevance [8, 11, 92].

(ii) *Concept effect on task performance.* These metrics assess how well concepts contribute to the predictive capacity of the model. They are employed in post-hoc methods, to assess the quality of the extracted concept for predicting the final class; at the same time, they are also studied in concept-based models, to assess how well concepts contribute to the predictive capacity within the given model. Among these metrics we find the Completeness Score [117], Fidelity [33, 97, 122], Concept Efficiency [32], Conciseness [107].

Quality of concepts. This second category focuses on the quality of concepts, studying (i) their inherent properties, (ii) how they relate to internal network representations, and (iii) their impact on prediction errors. In the following, we review these metrics with respect to these three targets of focus.

(i) *Concepts properties.* This group of metrics centres on the intrinsic characteristics of concepts themselves. They evaluate the retained information, distinctiveness, completeness, and purity of concepts, providing insights into the information content and coverage of concepts. Among these metrics we find the Mutual Information [32], Concept Distinctiveness [109], Concept Purity [109].

(ii) *Relation of Concepts with internal network representation.* This group of metrics quantifies the relationship between a node (or filter) of the model and the concept it represents. These metrics are adopted for methods that provide explanations as Node-Concept Association. They assess the spatial overlap (IoU) and location stability of concepts within activation maps. Among these metrics we find Intersection over Union [13, 34, 114] (and adopted [121] with the name of *Part Interpretability*), Location stability [121].

(iii) *Concept prediction error.* This group of metrics measures the alignment between learned concepts and concept ground truth. Hence, they are applicable when concept labels are available. These metrics evaluate and quantify mispredictions and the purity of concepts. Among these metrics we find the Concept Error [61] (and used in [97] with the name Explanation Error), AUC score [22], Misprediction-overlap metric (MPO) [54].

6.2 Qualitative evaluation

As common in XAI assessment [15], several C-XAI works perform human evaluations to assess their method through human feedback. We report which methods perform human evaluations in the “Human evaluation” column of Tables 6 (post-hoc methods) and 7 (by-design approaches). Half of the post-hoc approaches conduct qualitative assessments,

while only a quarter of explainable-by-design ones perform them. In the following, we overview the assessment based on the nature of the aspect under study: understandability, coherence, and utility of concept explanations.

Understandability of explanations. When evaluating explanations, the predominant focus involves assessing how understandable and reasonable the explanations are to humans. Various measures and evaluation settings have been introduced to capture specific nuances within this context. One relevant property of explanation is their *plausibility*, denoting the extent to which explanations align with human reasoning. One direct setting for evaluating the explanation plausibility consists of asking users to determine if the explanation is justifiable and understandable based on their expectation [13, 89]. To evaluate *human understanding* of concepts, another evaluation approach requires users to label the identified concepts by selecting some phrases from a predefined vocabulary that best describes the concept [109]. This method provides a structured measure for summarizing users' comprehension of identified concepts. Other approaches as [122] measure the interpretability of a concept by asking users to provide free-text labels for concepts and identify, among a set, the candidate explanations matching the concept and measure their accuracy as the percentage of correctly identified concept explanations. Close to this setting, [123] focus on quantifying the semantics learned by each network unit and evaluating the accuracy of human-assigned text annotations to concepts. The *trustability* of explanations is another critical aspect, where users provide feedback on their level of trust in the explanation [89]. The *reasonability* of explanations evaluates the effectiveness of conveying meaningful information through explanations. A way to evaluate it is by investigating what information saliency maps can communicate to humans [56]. Another setting consists of providing different explanations for the same prediction and asking human raters to indicate which is more reasonable, as adopted by IBD [126]. *Factuality* assesses the accuracy of concepts by having annotators judge their alignment with ground truth images [116].

Coherence of concepts. Another aspect to evaluate is the coherence of concepts as perceived by humans. This involves scrutinizing and evaluating the characteristics of the discovered concepts. An adopted evaluation method for measuring concept coherence is the intruder detection experiment. In this setup, for example, users are asked to identify an image among a set that is conceptually different from the rest [33, 37]. Another setting involves examining the nearest neighbor images associated with each discovered concept vector and selecting the most prevalent and cohesive concept [117]. The work in [116] defines the *groundability* to measure the consistency of the vision-language model adopted in their framework to ground concepts to images in a manner that aligns with human interpretations.

Utility of explanations . The *utility* of an explanation refers to its practical usefulness and effectiveness in providing insights or guidance to users, typically in the context of the specific task. It is often measured by its ability to assist users in predicting, performing a task, or making informed decisions in a given context. For example, in [51], users evaluate the utility of explanations in improving classifier generalization through feature engineering. They identify problematic features using global explanations and evaluate the ease of use of the explanations for this task. In addition, users select the most general classifier based solely on the global explanations of various methods. Other user study settings, such as in [33], measure how well the explanations assist users in inferring the insights guiding classifications. The utility is evaluated by measuring users' accuracy in predicting model decisions on novel images given the explanations.

6.3 Datasets

We now analyze the most interesting datasets employed in the literature. Among the datasets experimented by the method reviewed in this survey, we report in Table 4 only those where the authors either explicitly describe the concept

Table 4. Datasets employed by the reviewed methods in terms of type of data, size, task label and concepts. We report only the dataset in which the concepts are explicitly defined. All datasets have been defined for classification tasks. We report with *, the datasets suitable for regression too.

	Dataset	Used in	Size	Label	Concepts		
Type of Data	IMG	MNIST [62]	[25] [8] [109] [66]	60000	Digits	Geometrical attributes	
		MNIST even/odd [12]	[92]	60000	Even/odd	Digit parity	
		Fashion MNIST [112]	[66]	60000	Clothes	Details of clothes	
		colored-MNIST [55]	[41]	60000		Digits	Digit's colours
		Caltech-UCSD	[25] [122] [64] [61]	11788		Bird species	Visual attributes of the birds
		Birds-200 (CUB) [108]	[57] [54] [121] [109]				
			[21] [95] [94] [29]				
			[98] [97] [116] [80]				
			[32]				
			[13] [126] [34]	63305		Concepts	1197 visual concepts.
	[41] [75] [32] [11]	202599		Person's visual attributes	Non-overlapping person visual attributes		
	[61]	36369		Osteoarthritis levels	Knee conditions relative to the osteoarthritis		
	[123] [22] [80]	10million		Scenes	Objects landmarks		
	[117] [57] [97]	37322		Animals	85 numeric attribute		
	[118] [116]	10000		Malignant/Benign	Characteristics of malignant skin lesions		
	[118]	33126		Malignant/Benign	Characteristics of malignant melanoma		
	[75]	>1million		Objects	Object attributes and robot's sensory measurements		
Type of Data	TXT	CEBaB [1]	[111]	>15000	Restaurant reviews' sentiment	Aspect-level sentiment	
		IMBD [48]	[117]	50000	Movie reviews' sentiment (pos/neg)	Extracted	
		SST-2 [104]	[89]	>70000	Movie reviews' sentiment (pos/neg)	Extracted	
		SST-5 [104]	[89]	>11000	Movie reviews' sentiment (5 classes)	Extracted	
		TREC-6/50 [67]	[89]	>5000	Question types (6 classes/50 classes)	Extracted	
		SUBJ [83]	[89]	9000	Subjective/objective.	Extracted	
		Type of Data	VID	Kinetics-770 [20]	[52]	65000	Human actions (700 classes)
KTH Action [99]	[52]			>2000	Human actions (6 classes)	Extracted	
BDD-OIA [113]	[98]			>22000	Human actions (4 classes)	21 concepts (e.g. Traffic light is green, etc.)	
Type of Data	TAB	COMPAS [87]	[8]	>60000	Recidivism (yes/no)	Table attributes	
		V-Dem [84]	[12, 24]	202	State Democracy level	Table attributes	
		MIMIC-II [96]	[12, 24]	>40000	Patient survival (dead/alive)	Table attributes	
Type of Data	TS	MIT-BIH Electrocardiogram (ECG) [79]	[25]	47 (187 timesteps)	Heartbeat (normal/abnormal)	Heart issues	
		MUTAG [28]	[114]	188	Node labels (7 classes)	Extracted	
Type of Data	GRAPH	Reddit-binary [60]	[114]	2000	Type of graphs (question/answer-based community or a discussion-based community)	Extracted	
		PROTEINS [17]	[114]	>1000	Type of protein (enzymes/not enzymes)	Extracted	
		IMBD-Binary [115]	[114]	1000	Movies	Extracted	

annotations or the procedure followed to extract them (Extracted). We classified them based on the data type: images (IMG), text (TXT), videos (VID), tabular (TAB), time series (TS), and graph data (Graph).

On top of Table 4, we report the image datasets. We observe that the CUB dataset is the predominant choice in the C-XAI literature. This dataset is very suitable for such methodologies due to its extensive annotation of concepts, representing the visual attributes inherent to different bird species. We also note that the BRODEN dataset [13] was introduced explicitly for concept-based XAI. This dataset encompasses over 63,000 images and 1,197 concepts distributed across six categories (i.e., textures, colors, materials, parts, objects, and scenes). Finally, OAI is the only dataset that is suitable for regression tasks and not only for classification (the level of osteoarthritis may be both considered as a class or as a continuous level). In the second part, we find the text datasets. These datasets predominantly revolve

around binary or multi-class sentiment classifications. In most datasets, the concepts are not already present but are extracted from the input texts using various techniques. For example, in [89], concepts are represented by non-terminal phrases of the parsing tree. Within the third part of the table reporting videos dataset, BDD-OIA is the only dataset with pre-established concepts, while in the other two cases, concepts are extracted through specific processes. In tabular data sets, data attributes are used concepts. Finally, in graph datasets, concepts are always extracted by the method.

6.4 Resources

We now report information and insights regarding the resources associated with the reviewed methods. Specifically, when we refer to *resources*, we are addressing whether the authors of the respective papers have made the code, models utilized in their research publicly accessible or have released a new dataset. Among all methods discussed, nearly all but four have made their code publicly available on GitHub, two post-hoc methods (CaCE, Object Detection) and two concept-based models (SENN, CBM-AUC). Overall, there is a common trend of releasing code on public platforms (e.g., GitHub), which benefits the community by providing a starting point for testing and development. Concerning data release, four methods released novel datasets, equally distributed between post-hoc methods (ND, DMA & IMA) and concept-based models (CME, Glance-Nets). The release of new datasets is infrequent due to the resource-intensive nature of creating concept annotated datasets. It may require $C \times N$ extra annotations where C represents the number of concepts and N the number of samples. In other cases, a class-level annotation of the concept is preferred to decrease this effort. It consists in employing the same concept annotation for all the samples belonging to the same class, therefore reducing the number of annotations to $N \times Y$, where Y represents the number of final classes. For additional information on the works with shared code, models, or datasets released, please consult Appendix D.

7 APPLICATIONS

The recent emergence of C-XAI methods has led to some interesting applications. We report first novel algorithms based on concepts in the field within the AI area (outside explainability), followed by real-world applications.

Support to AI algorithm’s development. The first application of concept-based methods within AI came out in the image generation field. In VAEL [77], the authors proposed to condition the generation of images by means of a neuro-symbolic program working over concepts. For example, given two pairs of digits from MNIST, they train the concepts to represent the digits and the logic program to predict the addition of the digits. When decoding, they require instead that the generated images correspond to samples solving the given query (e.g., digits whose sum is 13). Interestingly, they show that their method also generalizes to queries unseen in training (e.g., they generate the samples whose difference or factorization is provided).

Other concept-based algorithms appeared in the out-of-distribution sample detection field. In particular, the author of [72] proposed a concept-based approach to characterize outliers. They propose first mapping the sample into an interpretable feature space comprising human-defined concepts like textures, colors, and geometry. In this space, they analyze the rarity of the features to detect outliers and explain the rarest concepts identified in the sample. The author of [23], instead, employs an unsupervised approach to extract concepts that can sufficiently characterize the decisions of an out-of-distribution detector. Furthermore, they require extracted concepts to separate in-distribution and out-of-distribution samples. Interestingly, they show that extracting concepts to explain the classifier is insufficient to explain the out-of-distribution detector’s prediction.

Real-world applications. Outside the AI world, the medical sector is where we find most applications of C-XAI already in the reviewed papers (CAR, CBM, PCBM, ProtoPDebug, LaBO). Another example is [69], where the authors employ a post-hoc concept-based method in the Computer-Aided Diagnosis (CAD) systems in skin cancer diagnosis. The authors want to elucidate the basic principles that guide the model’s predictions, discerning whether the model operates on similar concepts to those used by dermatologists to describe and diagnose the disease. Specifically, using TCAV [56] method, they show a strong correlation between DNN’s learned representation of concepts and those routinely used by dermatologists. Hence, listing the concepts that positively influenced the final outcome may be sufficient for trusting the predictions of a CAD system.

In the field of education, we identify RIPPLE [10] as a first application. This work introduces a pipeline to make predictions from raw time series data about students’ success that is also interpretable. Specifically, they use a combination of a graph-based neural network approach for classifying raw time series of student interactions and an adaptation of concept activation vectors TCAV [56] for interpreting NN internal state in terms of human-interpretable concepts. In this case, concepts are represented by six learning dimensions from educational scenarios defined a-priori.

8 PROMISING DIRECTIONS

The analysis of the C-XAI field literature allows highlighting several promising directions existing for both researchers and practitioners. Addressing these frontiers is important not only for enhancing the interpretability of artificial intelligence systems but also for fostering their responsible and effective deployment across diverse domains. Herein, we outline key directions that we believe merit focused explorations, encompassing both methodological advancements and improvements from an evaluation and application point of view.

From a methodological point of view, we believe that generative approaches could also be explored for post-hoc concept-based methods. On one side, Generative models could be harnessed to create concepts akin to their role in concept-based models. On the other, LLMs could also be employed to enhance supervised post-hoc methods in assessing model causality when prompted with specific concepts. Still considering the post-hoc methods, the analysis of Table 2 reveals a gap in exploring node-concept association explanations for unsupervised methods. Even though it may be difficult to associate nodes to extract concepts, it represents a valuable avenue for exploration. Concerning concept-based models, extending concept representations to accommodate multi-modal data, such as text, images, and audio, could greatly enhance the versatility of C-XAI across different domains. Also, the novel generative approach seems open to several extensions. Moreover, the new generative approach seems to be open to different types of fusion following the same approach: a concept generator that creates representations (not necessarily textual) that can be used in an external (not necessarily visual) model.

On the evaluation front, the field currently lacks benchmark datasets and standardized metrics, hindering systematic comparisons. Developing these resources would facilitate fair assessments of different methods.

From an application point of view, it would be important to develop scalable concept-based methods and models that can handle and be integrated into large-scale models and applied to large-scale datasets. Finally, properly integrating concept-based models is crucial to creating accountable AI-powered decision-making systems, particularly in critical domains like healthcare, finance, and autonomous systems.

9 CONCLUSION

In conclusion, our comprehensive review of Concept-based explainable Artificial Intelligence highlights the growing interest in concept-based explanations in addressing the demand for transparent and interpretable AI models. The

ongoing debate regarding the limitations of raw feature-based explanations has led to a substantial body of work on concept-based explanations. By defining the various types of concepts and concept-based explanations and introducing a systematic categorization framework, we hope to have facilitated a more structured understanding of the C-XAI landscape. Our analysis of C-XAI papers, which we reviewed and categorized with the proposed taxonomy into nine categories, may provide valuable insights into the strengths and challenges of different categories. Furthermore, by highlighting key resources within the C-XAI literature, such as datasets and metrics, we aspire to catalyse future research and facilitate comparisons between different methodologies.

REFERENCES

- [1] ABRAHAM, E. D., D'OOSTERLINCK, K., FEDER, A., GAT, Y., GEIGER, A., POTTS, C., REICHART, R., AND WU, Z. Cebab: Estimating the causal effects of real-world concepts on nlp model behavior. *Advances in Neural Information Processing Systems* 35 (2022), 17582–17596.
- [2] ACHTIBAT, R., DREYER, M., EISENBRAUN, I., BOSSE, S., WIEGAND, T., SAMEK, W., AND LAPUSCHKIN, S. From attribution maps to human-understandable explanations through concept relevance propagation. *Nature Machine Intelligence* 5, 9 (2023), 1006–1019.
- [3] ADADI, A., AND BERRADA, M. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access* 6 (2018), 52138–52160.
- [4] 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).
- [5] ADEBAYO, J., MUELLY, M., ABELSON, H., AND KIM, B. Post hoc explanations may be ineffective for detecting unknown spurious correlation. In *International conference on learning representations* (2021).
- [6] AGARWAL, R., MELNICK, L., FROSST, N., ZHANG, X., LENGERICH, B., CARUANA, R., AND HINTON, G. E. Neural additive models: Interpretable machine learning with neural nets. *Advances in neural information processing systems* 34 (2021), 4699–4711.
- [7] AGRAWAL, P., GIRSHICK, R., AND MALIK, J. Analyzing the performance of multilayer neural networks for object recognition. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VII 13* (2014), Springer, pp. 329–344.
- [8] ALVAREZ MELIS, D., AND JAAKKOLA, T. Towards robust interpretability with self-explaining neural networks. *Advances in neural information processing systems* 31 (2018).
- [9] ARRIETA, A. B., DÍAZ-RODRÍGUEZ, N., DEL SER, J., BENNETOT, A., TABIK, S., BARBADO, A., GARCÍA, S., GIL-LÓPEZ, S., MOLINA, D., BENJAMINS, R., ET AL. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information fusion* 58 (2020), 82–115.
- [10] ASADI, M., SWAMY, V., FREJ, J., VIGNOUD, J., MARRAS, M., AND KÄSER, T. Ripple: concept-based interpretation for raw time series models in education. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2023), vol. 37, pp. 15903–15911.
- [11] BARBIERO, P., CIRAVEGNA, G., GIANNINI, F., ESPINOSA ZARLENGA, M., MAGISTER, L. C., TONDA, A., LIO, P., PRECIOSO, F., JAMNIK, M., AND MARRA, G. Interpretable neural-symbolic concept reasoning. In *Proceedings of the 40th International Conference on Machine Learning* (23–29 Jul 2023), vol. 202 of *Proceedings of Machine Learning Research*, PMLR, pp. 1801–1825.
- [12] 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* (2022), vol. 36, pp. 6046–6054.
- [13] BAU, D., ZHOU, B., KHOSLA, A., OLIVA, A., AND TORRALBA, A. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 6541–6549.
- [14] BENGIO, Y., COURVILLE, A., AND VINCENT, P. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence* 35, 8 (2013), 1798–1828.
- [15] BODRIA, F., GIANNOTTI, F., GUIDOTTI, R., NARETTO, F., PEDRESCHI, D., AND RINZIVILLO, S. Benchmarking and survey of explanation methods for black box models. *Data Mining and Knowledge Discovery* (2023), 1–60.
- [16] BONTEPELLI, A., TESO, S., TENTORI, K., GIUNCHIGLIA, F., AND PASSERINI, A. Concept-level debugging of part-prototype networks. In *The Eleventh International Conference on Learning Representations* (2023).
- [17] BORGWARDT, K. M., ONG, C. S., SCHÖNAUER, S., VISHWANATHAN, S., SMOLA, A. J., AND KRIEDEL, H.-P. Protein function prediction via graph kernels. *Bioinformatics* 21, suppl_1 (2005), i47–i56.
- [18] BROWN, D., AND KVINGE, H. Making corgis important for honeycomb classification: Adversarial attacks on concept-based explainability tools. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 620–627.
- [19] BROWN, T., MANN, B., RYDER, N., SUBBLAH, M., KAPLAN, J. D., DHARIWAL, P., NEELAKANTAN, A., SHYAM, P., SASTRY, G., ASKELL, A., ET AL. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [20] CARREIRA, J., NOLAND, E., HILLIER, C., AND ZISSERMAN, A. A short note on the kinetics-700 human action dataset. *arXiv preprint arXiv:1907.06987* (2019).
- [21] CHEN, C., LI, O., TAO, D., BARNETT, A., RUDIN, C., AND SU, J. K. This looks like that: deep learning for interpretable image recognition. *Advances in neural information processing systems* 32 (2019).
- [22] CHEN, Z., BEL, Y., AND RUDIN, C. Concept whitening for interpretable image recognition. *Nature Machine Intelligence* 2, 12 (2020), 772–782.

- [23] CHOI, J., RAGHURAM, J., FENG, R., CHEN, J., JHA, S., AND PRAKASH, A. Concept-based explanations for out-of-distribution detectors. In *International Conference on Machine Learning* (2023), PMLR, pp. 5817–5837.
- [24] CIRAVEGNA, G., BARBIERO, P., GIANNINI, F., GORI, M., LIÓ, P., MAGGINI, M., AND MELACCI, S. Logic explained networks. *Artificial Intelligence* 314 (2023), 103822.
- [25] CRABBÉ, J., AND VAN DER SCHAAR, M. Concept activation regions: A generalized framework for concept-based explanations. *Advances in Neural Information Processing Systems* 35 (2022), 2590–2607.
- [26] DAI, J., QI, H., XIONG, Y., LI, Y., ZHANG, G., HU, H., AND WEI, Y. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision* (2017), pp. 764–773.
- [27] DAS, A., AND RAD, P. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371* (2020).
- [28] DEBNATH, A. K., LOPEZ DE COMPADRE, R. L., DEBNATH, G., SHUSTERMAN, A. J., AND HANSCH, C. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity. *Journal of medicinal chemistry* 34, 2 (1991), 786–797.
- [29] DONNELLY, J., BARNETT, A. J., AND CHEN, C. Deformable protopnet: An interpretable image classifier using deformable prototypes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 10265–10275.
- [30] DWIVEDI, R., DAVE, D., NAIK, H., SINGHAL, S., OMER, R., PATEL, P., QIAN, B., WEN, Z., SHAH, T., MORGAN, G., ET AL. Explainable ai (xai): Core ideas, techniques, and solutions. *ACM Computing Surveys* 55, 9 (2023), 1–33.
- [31] ERHAN, D., BENGIO, Y., COURVILLE, A., AND VINCENT, P. Visualizing higher-layer features of a deep network. *University of Montreal* 1341, 3 (2009), 1.
- [32] ESPINOSA ZARLENGA, M., BARBIERO, P., CIRAVEGNA, G., MARRA, G., GIANNINI, F., DILIGENTI, M., SHAMS, Z., PRECIOSO, F., MELACCI, S., WELLER, A., LIÓ, P., AND JAMNIK, M. Concept embedding models: Beyond the accuracy-explainability trade-off. In *Advances in Neural Information Processing Systems* (2022), S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35, Curran Associates, Inc., pp. 21400–21413.
- [33] FEL, T., PICARD, A., BETHUNE, L., BOISSIN, T., VIGOUROUX, D., COLIN, J., CADÈNE, R., AND SERRE, T. Craft: Concept recursive activation factorization for explainability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 2711–2721.
- [34] FONG, R., AND VEDALDI, A. 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] GANTER, B., AND WILLE, R. *Formal concept analysis: mathematical foundations*. Springer Science & Business Media, 2012.
- [36] GHORBANI, A., ABID, A., AND ZOU, J. Interpretation of neural networks is fragile. In *Proceedings of the AAAI conference on artificial intelligence* (2019), vol. 33, pp. 3681–3688.
- [37] GHORBANI, A., WEXLER, J., ZOU, J. Y., AND KIM, B. Towards automatic concept-based explanations. *Advances in neural information processing systems* 32 (2019).
- [38] GOGUEN, J. What is a concept? In *International Conference on Conceptual Structures* (2005), Springer, pp. 52–77.
- [39] GONDAL, M. W., WUTHRICH, M., MILADINOVIC, D., LOCATELLO, F., BREIDT, M., VOLCHKOV, V., AKPO, J., BACHEM, O., SCHÖLKOPF, B., AND BAUER, S. On the transfer of inductive bias from simulation to the real world: a new disentanglement dataset. *Advances in Neural Information Processing Systems* 32 (2019).
- [40] GOODMAN, B., AND FLAXMAN, S. European union regulations on algorithmic decision-making and a “right to explanation”. *AI magazine* 38, 3 (2017), 50–57.
- [41] GOYAL, Y., FEDER, A., SHALIT, U., AND KIM, B. Explaining classifiers with causal concept effect (cace). *arXiv preprint arXiv:1907.07165* (2019).
- [42] GREENWELL, B. M., BOEHMKE, B. C., AND MCCARTHY, A. J. A simple and effective model-based variable importance measure. *arXiv preprint arXiv:1805.04755* (2018).
- [43] GUIDOTTI, R., MONREALE, A., RUGGIERI, S., PEDRESCHI, D., TURINI, F., AND GIANNOTTI, F. Local rule-based explanations of black box decision systems. *arXiv preprint arXiv:1805.10820* (2018).
- [44] GUIDOTTI, R., MONREALE, A., RUGGIERI, S., TURINI, F., GIANNOTTI, F., AND PEDRESCHI, D. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)* 51, 5 (2018), 1–42.
- [45] HASE, P., CHEN, C., LI, O., AND RUDIN, C. Interpretable image recognition with hierarchical prototypes. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* (2019), vol. 7, pp. 32–40.
- [46] HAVASI, M., PARBHOO, S., AND DOSHI-VELEZ, F. Addressing leakage in concept bottleneck models. *Advances in Neural Information Processing Systems* 35 (2022), 23386–23397.
- [47] HOOKER, S., ERHAN, D., KINDERMANS, P.-J., AND KIM, B. A benchmark for interpretability methods in deep neural networks. *Advances in neural information processing systems* 32 (2019).
- [48] IMDB. IMDb Non-Commercial Datasets. <https://developer.imdb.com/non-commercial-datasets/>. Accessed: December 21, 2023.
- [49] INTERNATIONAL SKIN IMAGING COLLABORATION (ISIC). ISIC Challenge 2020. <https://challenge2020.isic-archive.com/>. Accessed: December 21, 2023.
- [50] IST AUSTRIA. AWA2: An Attribute Database for Animal Behavior Analysis. <https://cvml.ist.ac.at/AWA2/>. Accessed: December 21, 2023.
- [51] JAIN, R., CIRAVEGNA, G., BARBIERO, P., GIANNINI, F., BUFFELLI, D., AND LIO, P. Extending logic explained networks to text classification. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing* (2022), Association for Computational Linguistics, pp. 8838–8857.
- [52] JI, Y., WANG, Y., AND KATO, J. Spatial-temporal concept based explanation of 3d convnets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 15444–15453.

- [53] KAMINSKI, M. E. The right to explanation, explained. *Berkeley Technology Law Journal* 34, 1 (2019), 189–218.
- [54] KAZHDAN, D., DIMANOV, B., JAMNIK, M., LIÒ, P., AND WELLER, A. Now you see me (cme): concept-based model extraction. *arXiv preprint arXiv:2010.13233* (2020).
- [55] KIM, B., KIM, H., KIM, K., KIM, S., AND KIM, J. Learning not to learn: Training deep neural networks with biased data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 9012–9020.
- [56] KIM, B., WATTENBERG, M., GILMER, J., CAI, C., WEXLER, J., VIEGAS, F., ET AL. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning* (2018), PMLR, pp. 2668–2677.
- [57] KIM, E., JUNG, D., PARK, S., KIM, S., AND YOON, S. Probabilistic concept bottleneck models. *arXiv preprint arXiv:2306.01574* (2023).
- [58] KIM, S. S., WATKINS, E. A., RUSSAKOVSKY, O., FONG, R., AND MONROY-HERNÁNDEZ, A. "help me help the ai": Understanding how explainability can support human-ai interaction. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (2023), pp. 1–17.
- [59] KINDERMANS, P.-J., HOOKER, S., ADEBAYO, J., ALBER, M., SCHÜTT, K. T., DÄHNE, S., ERHAN, D., AND KIM, B. The (un) reliability of saliency methods. *Explainable AI: Interpreting, explaining and visualizing deep learning* (2019), 267–280.
- [60] KIPF, T. N., AND WELLING, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [61] 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* (2020), PMLR, pp. 5338–5348.
- [62] LECUN, Y. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/> (1998).
- [63] LECUN, Y., BENGIO, Y., AND HINTON, G. Deep learning. *nature* 521, 7553 (2015), 436–444.
- [64] LEEMANN, T., KIRCHHOF, M., RONG, Y., KASNECI, E., AND KASNECI, G. When are post-hoc conceptual explanations identifiable? In *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence* (31 Jul–04 Aug 2023), R. J. Evans and I. Shpitser, Eds., vol. 216 of *Proceedings of Machine Learning Research*, PMLR, pp. 1207–1218.
- [65] LI, L., WANG, B., VERMA, M., NAKASHIMA, Y., KAWASAKI, R., AND NAGAHARA, H. Scouter: Slot attention-based classifier for explainable image recognition. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 1046–1055.
- [66] LI, O., LIU, H., CHEN, C., AND RUDIN, C. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2018), vol. 32.
- [67] LI, X., AND ROTH, D. Learning question classifiers. In *COLING 2002: The 19th International Conference on Computational Linguistics* (2002).
- [68] LIU, Z., LUO, P., WANG, X., AND TANG, X. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)* (December 2015).
- [69] 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)* (2020), IEEE, pp. 1–10.
- [70] LUNDBERG, S. M., AND LEE, S.-I. A unified approach to interpreting model predictions. *Advances in neural information processing systems* 30 (2017).
- [71] MACCARTHY, M. An examination of the algorithmic accountability act of 2019. *Algorithms* (2020).
- [72] MADEIRA, P., CARREIRO, A., GAUDIO, A., ROSADO, L., SOARES, F., AND SMAILAGIC, A. Zebra: Explaining rare cases through outlying interpretable concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 3781–3787.
- [73] MAGISTER, L. C., KAZHDAN, D., SINGH, V., AND LIÒ, P. Gcexplainer: Human-in-the-loop concept-based explanations for graph neural networks. *arXiv preprint arXiv:2107.11889* (2021).
- [74] MAHINPEI, A., CLARK, J., LAGE, I., DOSHI-VELEZ, F., AND PAN, W. Promises and pitfalls of black-box concept learning models. *arXiv preprint arXiv:2106.13314* (2021).
- [75] MARCONATO, E., PASSERINI, A., AND TESO, S. Glancenets: Interpretable, leak-proof concept-based models. *Advances in Neural Information Processing Systems* 35 (2022), 21212–21227.
- [76] MARGELOIU, A., ASHMAN, M., BHATT, U., CHEN, Y., JAMNIK, M., AND WELLER, A. Do concept bottleneck models learn as intended? *CoRR abs/2105.04289* (2021).
- [77] MISINO, E., MARRA, G., AND SANSONE, E. Vael: Bridging variational autoencoders and probabilistic logic programming. *Advances in Neural Information Processing Systems* 35 (2022), 4667–4679.
- [78] MOLNAR, C. *Interpretable machine learning*. Independently published, 2020.
- [79] MOODY, G. B., AND MARK, R. G. The impact of the mit-bih arrhythmia database. *IEEE engineering in medicine and biology magazine* 20, 3 (2001), 45–50.
- [80] OIKARINEN, T., DAS, S., NGUYEN, L. M., AND WENG, T.-W. Label-free concept bottleneck models. *arXiv preprint arXiv:2304.06129* (2023).
- [81] OIKARINEN, T., AND WENG, T.-W. Clip-dissect: Automatic description of neuron representations in deep vision networks. *arXiv preprint arXiv:2204.10965* (2022).
- [82] OSTEOARTHRITIS INITIATIVE. Osteoarthritis Initiative (OAI) Data. <https://nda.nih.gov/oai/>. Accessed: December 21, 2023.
- [83] PANG, B., AND LEE, L. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *arXiv preprint cs/0506075* (2005).
- [84] PEMSTEIN, D., MARQUARDT, K. L., TZELGOV, E., WANG, Y.-T., KRUSELL, J., AND MIRI, F. The v-dem measurement model: latent variable analysis for cross-national and cross-temporal expert-coded data. *V-Dem working paper 21* (2018).
- [85] POCHÉ, A., HERVIER, L., AND BAKKAY, M.-C. Natural example-based explainability: a survey.
- [86] POURSABZI-SANGDEH, F., GOLDSTEIN, D. G., HOFMAN, J. M., WORTMAN VAUGHAN, J. W., AND WALLACH, H. Manipulating and measuring model

- interpretability. In *Proceedings of the 2021 CHI conference on human factors in computing systems* (2021), pp. 1–52.
- [87] PROPUBLICA. COMPAS Analysis Repository. <https://github.com/propublica/compas-analysis>, Year Accessed. Accessed: December 21, 2023.
- [88] RADFORD, A., KIM, J. W., HALLACY, C., RAMESH, A., GOH, G., AGARWAL, S., SASTRY, G., ASKELL, A., MISHKIN, P., CLARK, J., ET AL. Learning transferable visual models from natural language supervision. In *International conference on machine learning* (2021), PMLR, pp. 8748–8763.
- [89] RAJAGOPAL, D., BALACHANDRAN, V., HOVY, E. H., AND TSVETKOV, Y. Selfexplain: A self-explaining architecture for neural text classifiers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* (2021), pp. 836–850.
- [90] RAMASWAMY, V. V., KIM, S. S., FONG, R., AND RUSSAKOVSKY, O. Overlooked factors in concept-based explanations: Dataset choice, concept learnability, and human capability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 10932–10941.
- [91] RIBEIRO, M. T., SINGH, S., AND GUESTRIN, C. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (2016), pp. 1135–1144.
- [92] RIGOTTI, M., MIKSOVIC, C., GIURGIU, I., GSCHWIND, T., AND SCOTTON, P. Attention-based interpretability with concept transformers. In *International Conference on Learning Representations* (2022).
- [93] RUDIN, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence* 1, 5 (2019), 206–215.
- [94] RYMARCZYK, D., STRUSKI, Ł., GÓRSZCZAK, M., LEWANDOWSKA, K., TABOR, J., AND ZIELIŃSKI, B. Interpretable image classification with differentiable prototypes assignment. In *European Conference on Computer Vision* (2022), Springer, pp. 351–368.
- [95] RYMARCZYK, D., STRUSKI, Ł., TABOR, J., AND ZIELIŃSKI, B. Protoshare: Prototypical parts sharing for similarity discovery in interpretable image classification. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining* (2021), pp. 1420–1430.
- [96] SAEED, M., VILLARROEL, M., REISNER, A. T., CLIFFORD, G., LEHMAN, L.-W., MOODY, G., HELDT, T., KYAW, T. H., MOODY, B., AND MARK, R. G. Multiparameter intelligent monitoring in intensive care ii (mimic-ii): a public-access intensive care unit database. *Critical care medicine* 39, 5 (2011), 952.
- [97] SARKAR, A., VIJAYKEERTHY, D., SARKAR, A., AND BALASUBRAMANIAN, V. N. A framework for learning ante-hoc explainable models via concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 10286–10295.
- [98] SAWADA, Y., AND NAKAMURA, K. Concept bottleneck model with additional unsupervised concepts. *IEEE Access* 10 (2022), 41758–41765.
- [99] SCHULTZ, C., LAPTEV, I., AND CAPUTO, B. Recognizing human actions: a local svm approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.* (2004), vol. 3, IEEE, pp. 32–36.
- [100] SELVARAJU, R. R., COGSWELL, M., DAS, A., VEDANTAM, R., 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* (2017), pp. 618–626.
- [101] SHIN, S., JO, Y., AHN, S., AND LEE, N. A closer look at the intervention procedure of concept bottleneck models. *arXiv preprint arXiv:2302.14260* (2023).
- [102] SIMONYAN, K., VEDALDI, A., AND ZISSERMAN, A. Deep inside convolutional networks: visualising image classification models and saliency maps. In *Proceedings of the International Conference on Learning Representations (ICLR)* (2014), ICLR.
- [103] SLACK, D., HILGARD, S., JIA, E., SINGH, S., AND LAKKARAJU, H. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* (2020), pp. 180–186.
- [104] SOCHER, R., PERELYGIN, A., WU, J., CHUANG, J., MANNING, C. D., NG, A. Y., AND POTTS, C. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing* (2013), pp. 1631–1642.
- [105] SPEER, R., CHIN, J., AND HAVASI, C. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI conference on artificial intelligence* (2017), vol. 31.
- [106] 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 (2018), 1–9.
- [107] VIELHABEN, J., BLUECHER, S., AND STRODTHOFF, N. Multi-dimensional concept discovery (mcd): A unifying framework with completeness guarantees. *Transactions on Machine Learning Research* (2023).
- [108] WAH, C., BRANSON, S., WELINDER, P., PERONA, P., AND BELONGIE, S. The caltech-ucsd birds-200-2011 dataset.
- [109] WANG, B., LI, L., NAKASHIMA, Y., AND NAGAHARA, H. Learning bottleneck concepts in image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 10962–10971.
- [110] WELINDER, P., BRANSON, S., MITA, T., WAH, C., SCHROFF, F., BELONGIE, S., AND PERONA, P. Caltech-ucsd birds 200.
- [111] WU, Z., D'OOSTERLINCK, K., GEIGER, A., ZUR, A., AND POTTS, C. Causal proxy models for concept-based model explanations. In *International Conference on Machine Learning* (2023), PMLR, pp. 37313–37334.
- [112] XIAO, H., RASUL, K., AND VOLLGRAF, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747* (2017).
- [113] XU, Y., YANG, X., GONG, L., LIN, H.-C., WU, T.-Y., LI, Y., AND VASCONCELOS, N. Explainable object-induced action decision for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), pp. 9523–9532.
- [114] XUANYUAN, H., BARBIERO, P., GEORGIEV, D., MAGISTER, L. C., AND LIÒ, P. Global concept-based interpretability for graph neural networks via neuron analysis. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2023), vol. 37, pp. 10675–10683.
- [115] YANARDAG, P., AND VISHWANATHAN, S. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining* (2015), pp. 1365–1374.

- [116] YANG, Y., PANAGOPOULOU, A., ZHOU, S., JIN, D., CALLISON-BURCH, C., AND YATSKAR, M. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)*, pp. 19187–19197.
- [117] YEH, C.-K., KIM, B., ARIK, S., LI, C.-L., PFISTER, T., AND RAVIKUMAR, P. On completeness-aware concept-based explanations in deep neural networks. *Advances in neural information processing systems* 33 (2020), 20554–20565.
- [118] YUKSEKONUL, M., WANG, M., AND ZOU, J. Post-hoc concept bottleneck models. In *ICLR 2022 Workshop on PAIR’2Struct: Privacy, Accountability, Interpretability, Robustness, Reasoning on Structured Data (2022)*.
- [119] ZEILER, M. D., AND FERGUS, R. Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I 13 (2014)*, Springer, pp. 818–833.
- [120] ZHANG, G.-Q., AND SHEN, G. Approximable concepts, chu spaces, and information systems. *Theory and Applications of Categories* 17 (2006), 79–102.
- [121] ZHANG, Q., WU, Y. N., AND ZHU, S.-C. Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition (2018)*, pp. 8827–8836.
- [122] ZHANG, R., MADUMAL, P., MILLER, T., EHINGER, K. A., AND RUBINSTEIN, B. I. Invertible concept-based explanations for cnn models with non-negative concept activation vectors. In *Proceedings of the AAAI Conference on Artificial Intelligence (2021)*, vol. 35, pp. 11682–11690.
- [123] ZHOU, B., KHOSLA, A., LAPEDRIZA, A., OLIVA, A., AND TORRALBA, A. Object detectors emerge in deep scene cnns. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings (2015)*, Y. Bengio and Y. LeCun, Eds.
- [124] ZHOU, B., KHOSLA, A., LAPEDRIZA, A., OLIVA, A., AND TORRALBA, A. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition (2016)*, pp. 2921–2929.
- [125] ZHOU, B., LAPEDRIZA, A., KHOSLA, A., OLIVA, A., AND TORRALBA, A. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence* 40, 6 (2017), 1452–1464.
- [126] ZHOU, B., SUN, Y., BAU, D., AND TORRALBA, A. Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV) (2018)*, pp. 119–134.

A GLOSSARY OF TERMS

In the following, we provide a glossary of the most important term in the field of Concept-based XAI. We hope this glossary will serve readers as a unique reference source, allowing future literature to use the different terms consistently. For the sake of readability, we will always reintroduce these terms upon their first mention. In Table 5, instead, we provide a list of the acronyms and abbreviations.

- Explainability The capacity of a method to provide intelligible explanations for model predictions or decisions. It should unveil the inner workings and factors that contribute to the model’s outcomes in a manner that is understandable to humans.
- Transparency The characteristics of a model whose inner processes, mechanisms, and logic are clear and comprehensible to users. A transparent model should reveal how it reaches specific predictions, enhancing user trust and aiding in identifying potential biases or errors.
- Interpretability It is the degree to which the outcomes and decision-making of a machine learning model can be grasped and reasoned about by humans. An interpretable model reveals the relationships between input data and predictions, facilitating a deeper understanding of its functioning.
- Local Expl. It focuses on interpreting the behavior of a model for a specific individual prediction. It provides insights into why the model made a particular decision for a given input instance.
- Node/Filter Expl. It involves detailing the behavior of individual nodes or neurons within a neural network. It helps decipher the specific features or relationships that targeted nodes consider during predictions.
- Global Expl. It provides insights into a model’s behavior across the entire dataset. It identifies patterns, trends, and overall tendencies to understand better how the model functions holistically.
- Post-hoc Expl. Explanations are generated after the model has made a prediction. They do not interfere with the model standard working.

- Model Agnostic** Technique that is not dependent on the specific algorithm or model architecture used. Model-agnostic approaches can be applied to various machine learning models by only analyzing their input and outputs.
- Counterfactual Expl.** Explanations showing how altered input data could change model predictions. They create modified scenarios to show why a model made a specific prediction and how input changes would result in different outcomes.
- Feature Importance** It measures the contribution of each input feature to a model's predictions. It helps identify which features have the most significant influence on the model's outcomes (e.g., pixel images).
- Saliency Map** A saliency map is a visual representation highlighting the important regions or features of an input that significantly affect a model or a node's behavior.
- Concept** High-level, human-understandable abstraction that captures key features or patterns within data (e.g., object parts, colors, textures, etc.). Concepts are used to simplify complex information, enabling easier interpretation and communication of AI model behavior.
- Concept-based Expl.** Type of explanations aiming to elucidate model predictions by decomposing them into human-understandable concepts. Rather than providing explanations regarding feature importance, these methods provide explanations in terms of high-level entity attributes (i.e., concepts). These explanations aim to bridge the gap between complex model decisions and intuitive human reasoning.
- Concept-based Model** Type of machine learning model designed to make predictions while employing high-level concepts or abstractions rather than working directly with raw data. These models enhance interpretability by relying on human-understandable concepts.
- Concept Bottleneck** Design principle of certain concept-based models. It involves incorporating an intermediate layer within the model's architecture to represent a given set of concepts. This layer constrains the flow of information by forcing the model to represent its predictions using these predefined concepts, possibly decreasing the model classification performance.
- Concept Score** Numerical value that quantifies the relevance or presence of a specific concept within the context of a machine learning model's decision or prediction.
- Concept Embedding** Representation of concepts into a numerical form. Concept embeddings map complex human-understandable concepts into numerical representations, providing a model of more information than what can be described by a single score.
- Logic Expl.** An explanation outlining the step-by-step logical reasoning process leading to a particular decision or prediction. It presents a sequence of logical conditions (i.e., a rule) used by the model to reach its conclusion, e.g., $paws(x) \wedge tail(x) \wedge muzzle(x) \rightarrow dog(x)$.
- Probing Method** Technique to assess how the latent representations captured by a given model can discriminate a set of concepts. This method involves training an auxiliary concept-specific model on top of the neural network's latent representation. Probing methods enable researchers to uncover the underlying concepts learned by deep models.
- Concept Intervention** Technique to assess and modify a machine learning prediction based on modification of the predicted concepts. They enable domain experts to test hypothetical scenarios to improve the model's prediction and provide concept-based counterfactual explanations.

Acronym	Full name
ACE [37]	Automatic Concept-based Explanations
AE	Auto Encoder
AI	Artificial Intelligence
CAM [124]	Class Activation Map
CAV [56]	Concept Activation Vector
CAR [25]	Concept Activation Region
CBM [61]	Concept Bottleneck Models
CEM [32]	Concept Embedding Models
CNN	Convolution Neural Network
CT [92]	Concept Transformer
CW [22]	Concept Whitening
C-XAI	Concept-based XAI
DCR [11]	Deep Concept Reasoning
DeconvNet [119]	Deconvolutional Network
DL	Deep Learning
DNN	Deep Neural Network
GDPR	General Data Protection Regulation
Grad-CAM [100]	Gradient-weighted Class Activation Mapping
GNN	Graph Neural Network
IBD [126]	Interpretable Basis Decomposition
ICE [122]	Invertible Concept-based Explanation
LaBO [116]	Language in a bottle
LEN [24]	Logic Explained Networks
LIME [91]	Local Interpretable Model-agnostic Explanation
LLM	Large Language Model
ML	Machine Learning
MLP	Multi-Layer Perceptron
ND [13]	Network Dissection
Net2Vec [34]	Network to Vector
PCA	Principal Component Analysis
PCBM [118]	Post-hoc Concept Bottleneck Models
ProtoPNets [21]	Prototype Parts Networks
SENN [8]	Self-Explaining Neural Network
SHAP [70]	SHAPley additive explanation
TCAV [56]	Testing with CAV
VAE	Variational Auto Encoder
XAI	EXplainable AI

Table 5. Table of acronyms and abbreviations.

B FORMAL CONCEPT ANALYSIS

We now formally define the notion of concepts, using as reference the category theory literature [35, 38, 120]. This area has been examined the topic from various perspectives, providing a robust base for our definition. A concept is a set of attributes that agrees on the intention of its extension [120], where the extension of a concept consists of all objects belonging to the concept and the intention of a concept is the set of attributes that are common to all these objects.

More precisely, consider a context defined as the triple (O, A, I) , where we indicate with O a set of objects, with A a set of attributes, and with I the satisfaction relation (also called “incidence”) which is a subset of $A \times O$. We also

define a subset of objects $\hat{O} \subset O$ and a subset of attributes $\hat{A} \subset A$, a derivation operator $\hat{A}' = \{o \in O \mid (o, a) \in I, \forall a \in \hat{A}\}$ and dually $\hat{O}' = \{a \in A \mid (o, a) \in I, \forall o \in \hat{O}\}$. We define (\hat{A}, \hat{O}) as a Formal Concept if and only if $\hat{B}' = \hat{A}$ and $\hat{A}' = \hat{B}$, or, also, if $\hat{A}'' = \hat{A}$ and $\hat{O}'' = \hat{O}$. In other words, a concept is a set of attributes that are shared among a set of objects and for which no other attribute exists that is shared among all the objects, even though other attributes can characterize some objects. According to [38], this is the broadest possible concept definition. Indeed, it allows us to consider all the previously identified typologies of concepts according to what we consider as the set of attributes A .

C METRICS DESCRIPTION

We now deep the analysis provided in Section 6.1 with a short description of each metric within each of the categories outlined.

C.1 Concept effect on class prediction and task performance

(i) Concept effect on class prediction.

- T-CAV score: it is the fraction of the class's inputs with a positive conceptual sensitivity [37, 56], where the conceptual sensitivity is how much each concept influences the prediction of a given class. So, it measures how many samples have concepts whose effects lie positively on the final prediction.
- CaCE: it measures the causal effect of the presence of concepts on the final class prediction [41, 111].
- ConceptSHAP: it assesses the effect of a concept on the final completeness score [117]. The completeness is defined as the amount to which concept scores are sufficient for predicting model outcomes.
- Smallest Sufficient Concepts (SSC): it computes the class accuracy when employing only a subset of the concepts, looking for the smallest set of concepts [37].
- Smallest Destroying Concepts (SDC): it looks for the smallest set of concepts deleting, which causes prediction accuracy to fall the most [37, 107].
- STCE Concept Importance: it computes the importance rank of each concept for a given class [52].
- T-CAR score: it quantifies the proportion of instances belonging to a specific class whose representations fall within the positive concept region [25].
- Sobol score: it measures the contribution of a concept and its interactions of any order with any other concepts to the model output variance [33].
- Concept Weight: in concept-based models employing a linear layer from concept to task, the same weight represents the importance of a concept for a given class [116, 118].
- Concept Relevance: the predicted importance of a concept for a final class in a concept-based model, determined differently for each sample [8, 11, 92].

(ii) Concept effect on task performance.

- Completeness score: it measures the amount to which concept scores are sufficient for predicting model outcomes [117]. This is based on the assumption that the concept scores of *complete* concepts are sufficient statistics of the model prediction.
- Fidelity: it counts the number of times in which the original and approximate model predictions are equal over the total number of images [33, 97, 122].
- Faithfulness: this measures the predictive capacity of the generated concepts [97]. It represents the capability of the overall concept vector to predict the ground truth task label.

- **Concept Efficiency:** in concept-based models, the concept capacity to predict the task depends on the number of concepts employed [32]. Richer concept representations mitigate this problem.
- **Conciseness:** similar to the concept efficiency but for post-hoc methods, it represents the number of concepts required to reach a certain grade of completeness and so a desired level of final accuracy [107].

C.2 Quality of Concepts

i) Concepts properties.

- **Mutual Information:** it quantifies the amount of information retained in a concept representation with respect to the input data [32].
- **Distinctiveness:** it quantifies the distinction between different concepts by considering the extent of their coverage [109].
- **Completeness:** Different from the previous definition of completeness given in [117], it measures how well a concept covers a specific associated class in the dataset [109].
- **Purity:** it quantifies the ability to discover concepts covering only a single class [109].

(ii) Relation of Concepts with internal network representation.

- **Intersection over Union (IoU):** it measures the degree of overlapping of the bounding boxes representing the concepts and the activation maps of the single node [13, 34, 114] (and adopted [121] with the name of *Part Interpretability*).
- **Location stability:** it assesses the consistency in how a filter represents the same object part across various objects [121].

(iii) Concept prediction error.

- **Concept Error:** it measures how close the concepts learned are to the concept ground truth [61] (and used in [97] with the name Explanation Error). It involves computing the L2 distance between the concepts learned and the ground truth concepts to measure the alignment.
- **AUC score:** this score measures whether the samples belonging to a concept are ranked higher than others [22]. That is, the AUC score indicates the purity of the concepts.
- **Misprediction-overlap metric (MPO):** it computes the sample fraction in the test set with at least m relevant concepts predicted incorrectly [54]. A larger MPO score implies a bigger proportion of incorrectly predicted concepts.

D METHODS RESOURCES

In this section, we extend Section 6, employing an orthogonal direction of analysis. Rather than describing the singular resources, here we report tables 6, 7 that contain, for each method, the type of resources employed in terms of metrics, dataset, and human evaluation, following the categorization provided in Section 3.1.

Received December 2023

Table 6. Post-hoc Concept-based Explainability methods. We characterize the approaches based on the following dimensions. The *Code/Models availability*, *Data release*: (✓), no (✗). *Human evaluation* (Human eval), if a user study is conducted for evaluation purposes: (✓), no (✗). Full category description in Section 3.

	Method	Code/Mod. Avail.	Data Release	New metric	Human eval
Supervised	T-CAV [56]	✓	✗	✓	✓
	CAR [25]	✓	✗	✓	✗
	IBD [126]	✓	✗	✗	✓
	CaCE [41]	✗	✗	✓	✗
	CPM [111]	✓	✗	✗	✗
	Obj. Det [123]	✗	✗	✗	✓
	ND [13]	✓	✓	✓	✓
	Net2Vec [34]	✓	✗	✗	✗
	GNN-CI [114]	✓	✗	✗	✗
Unsupervised	ACE [37]	✓	✗	✗	✓
	Compl. Aware [117]	✓	✗	✓	✓
	ICE [122]	✓	✗	✗	✓
	CRAFT [33]	✓	✗	✓	✓
	MCD [107]	✓	✗	✗	✗
	DMA & IMA [64]	✓	✓	✗	✗
	STCE [52]	✓	✗	✗	✗

Table 7. Explainable By-Design Concept-based Models. We characterize the approaches based on the following dimensions. The *Code/Models availability*, *Data release*: (✓), no (✗). *Human evaluation* (Human eval), if a user study is conducted for evaluation purposes: (✓), no (✗): (✓), no (✗). *Human evaluation* (Human eval), if a user study is conducted for evaluation purposes: (✓), no (✗). Full category description in Section 3.

	Method	Code/Mod Avail.	Data Release	New Metric	Human eval
Supervised	CBM [61]	✓	✗	✗	✗
	LEN [12, 24, 51]	✓	✗	✗	✓
	CEM[32]	✓	✗	✓	✗
	ProbCBM [57]	✓	✗	✗	✗
	DCR[11]	✓	✗	✗	✗
	CW [22]	✓	✗	✓	✗
	CME [54]	✓	✓	✓	✗
	PCBM [118]	✓	✗	✗	✗
	CT [92]	✓	✗	✗	✗
Unsupervised	Interp. CNN [121]	✓	✗	✗	✗
	SENN [8]	✗	✗	✗	✗
	SelfExplain [89]	✓	✗	✗	✓
	BotCL [109]	✓	✗	✗	✓
	PrototypeDL [66]	✓	✗	✗	✗
	ProtoPNet [21]	✓	✗	✗	✗
	ProtoPool [94]	✓	✗	✗	✓
	DeformableProtoPNet [29]	✓	✗	✗	✗
	HPnet [45]	✓	✗	✗	✗
	ProtoPShare [95]	✓	✗	✗	✓
Hybrid	CBM-AUC [98]	✗	✗	✗	✗
	Ante-hoc expl. [97]	✓	✗	✗	✗
	GlanceNets [75]	✓	✓	✗	✗
Generative	LaBO [116]	✓	✗	✗	✓
	Label-free CBM [80]	✓	✗	✗	✗