



Politecnico
di Torino



An Overview of Concept-based XAI

Elena Baralis, Eliana Pastor, Gabriele Ciravegna,
Tania Cerquitelli, Eleonora Poeta

ECML PKDD 2025, Porto - September 15th, 2025



Future
Artificial
Intelligence
Research

ECML ▾
PKDD
2025

Meet the instructors and contributors



Eleonora Poeta
PhD student at Politecnico
di Torino, Italy.



Gabriele Ciravegna
Researcher at CENTAI
Institute, Italy.



Eliana Pastor
Assistant professor at
Politecnico di Torino, Italy.



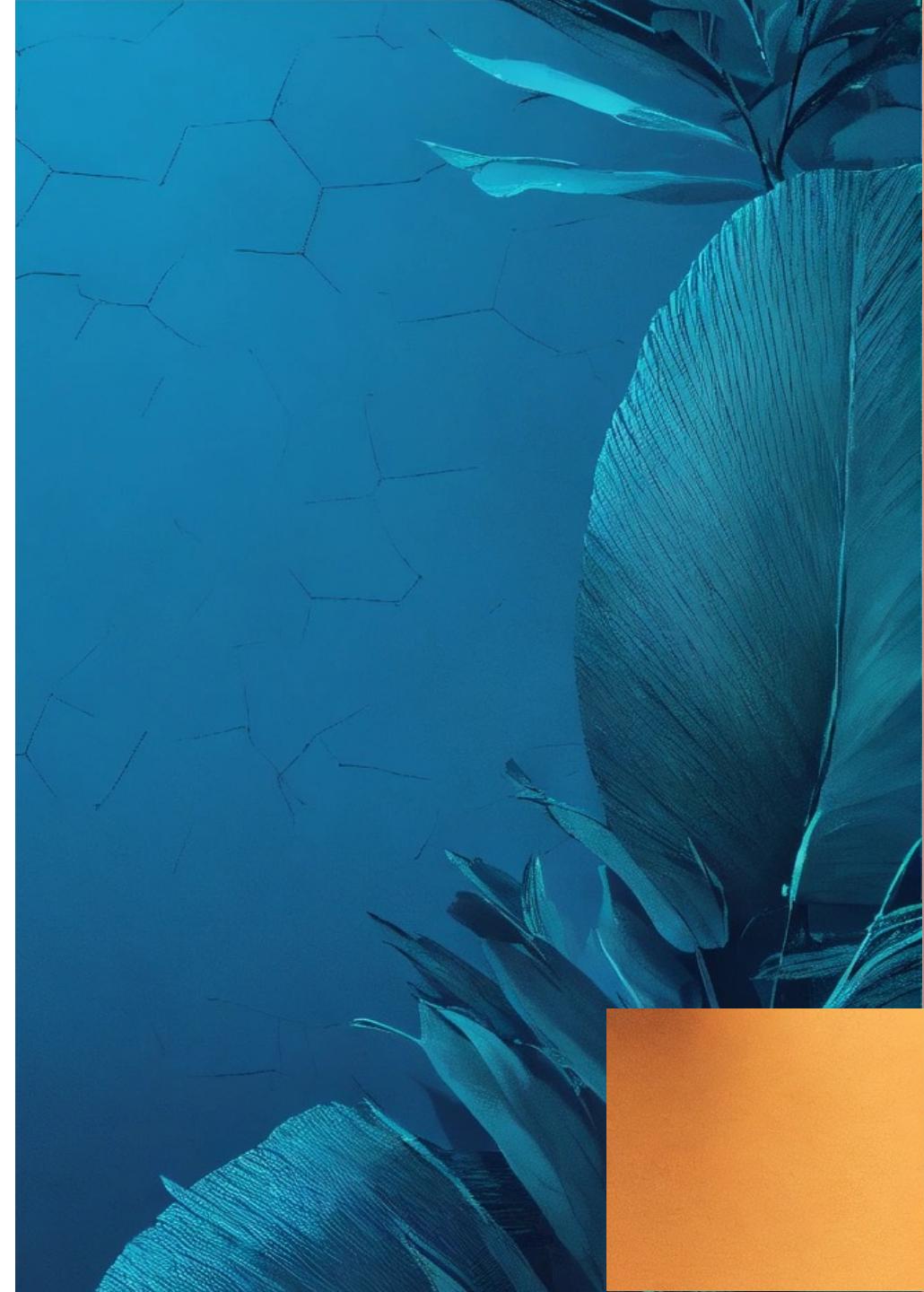
Tania Cerquitelli
Full professor at
Politecnico di Torino, Italy.



Elena Baralis
Full professor at Politecnico
di Torino, Italy.

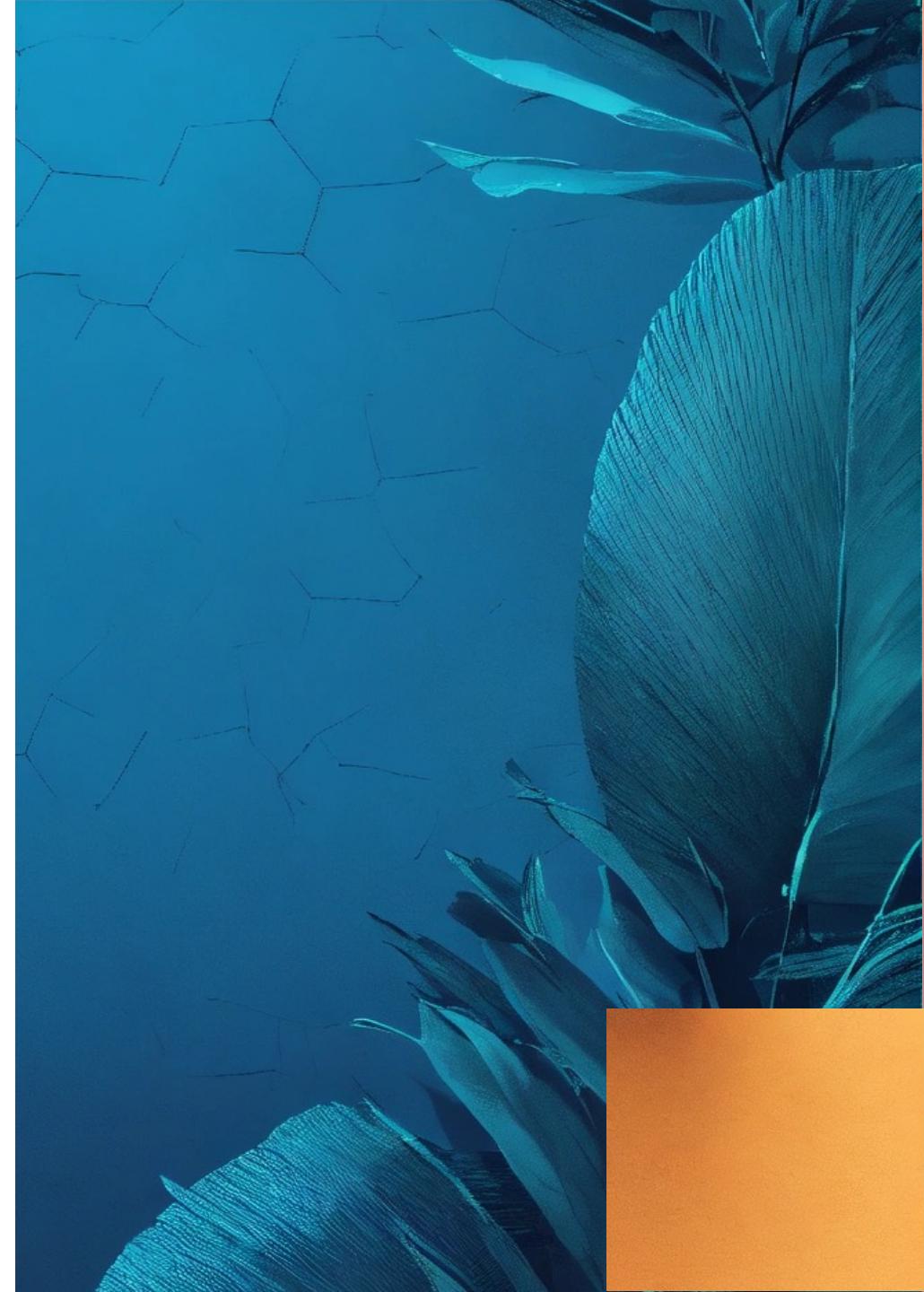
Tutorial structure

- Introduction to XAI techniques
- Taxonomy of C-XAI
- Post-hoc vs By-design C-XAI approaches
- Guidance for selecting C-XAI approaches
- Hands-on session



Tutorial structure

- **Introduction to XAI techniques**
- Taxonomy of C-XAI
- Post-hoc vs By-design C-XAI approaches
- Guidance for selecting C-XAI approaches
- Hands-on session



Why Explainability Matters in AI

- Deep learning systems are everywhere
- These systems influence critical decisions
- Demand is growing for AI that is *trustworthy* and *understandable*



Call for transparency

Increasing adoption of AI for medical diagnosis
Accurate results, even 'outperforming doctors'

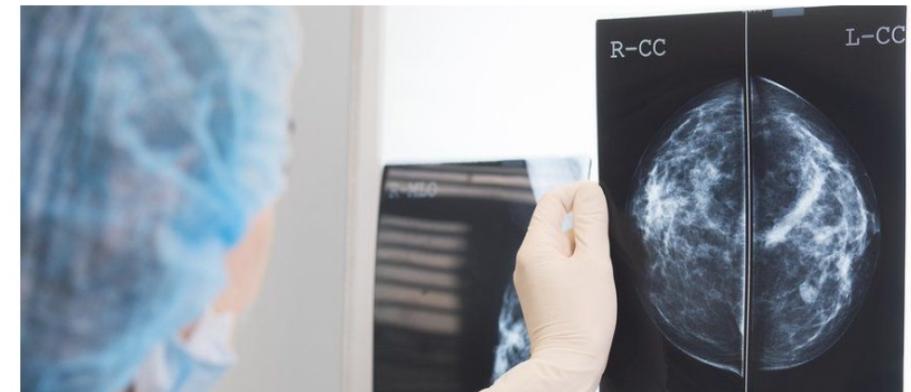
But...

AI 'outperforms' doctors diagnosing breast cancer



Fergus Walsh
Medical correspondent
[@BBCFergusWalsh](#)

2 January 2020



Article | [Published: 01 January 2020](#)

International evaluation of an AI system for breast cancer screening

[Scott Mayer McKinney](#) , [Marcin Sieniek](#), [Varun Godbole](#), [Jonathan Godwin](#)

+ Show authors

[Nature](#) **577**, 89–94 (2020) | [Cite this article](#)

98k Accesses | **1288** Citations | **3927** Altmetric | [Metrics](#)

Call for transparency

Response article

Researchers call for ***transparent and reproducible AI research***

Matters Arising | Published: 14 October 2020

Transparency and reproducibility in artificial intelligence

Benjamin Haibe-Kains , George Alexandru Adam, Ahmed Hosny, Farnoosh Khodakarami, Massive Analysis Quality Control (MAQC) Society Board of Directors, Levi Waldron, Bo Wang, Chris McIntosh, Anna Goldenberg, Anshul Kundaje, Casey S. Greene, Tamara Broderick, Michael M. Hoffman, Jeffrey T. Leek, Keegan Korthauer, Wolfgang Huber, Alvis Brazma, Joelle Pineau, Robert Tibshirani, Trevor Hastie, John P. A. Ioannidis, John Quackenbush & Hugo J. W. L. Aerts

Nature 586, E14–E16 (2020) | [Cite this article](#)

15k Accesses | 54 Citations | 520 Altmetric | [Metrics](#)

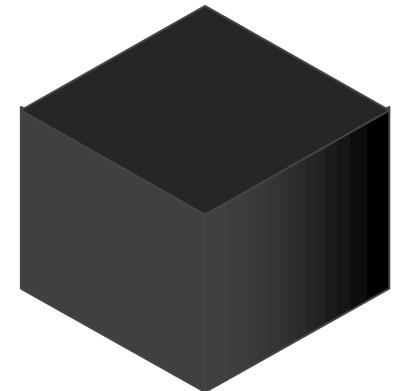
*[...] The lack of details of the methods and algorithm code undermines its scientific value. Here, we identify obstacles that hinder ***transparent and reproducible AI research*** and provide solutions to these obstacles with implications for the broader field.'*

Transparency & Explainability

Most AI models are **black boxes** and lack interpretability

Opaqueness of a model can be at **multiple levels**, involving

- Data
- Model/System/Algorithm
- Learned function and pattern → reasons behind its functioning
- Intention and business model of the AI product



These elements should be transparent - clear, disclosed - to the end users

A definition of Explainability/Interpretability

*The ability to **explain** or to present the **reasoning** behind the decisions **of AI systems** or its technical process **in understandable terms to humans***

Multiple terms – not monolithic concepts

- **Interpretability:** Focus more on transparency, models that are **inherently interpretable**
- **Explainability:** Focus more on models that are not comprehensible by design to humans
- + understandability, comprehensibility, intelligibility, mental fit

Often used interchangeably

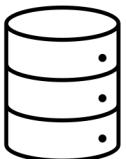
A vertical decorative strip on the left side of the slide featuring abstract, organic shapes in shades of blue, teal, and orange against a white background.

A taxonomy of Explainable AI

- **At which step of the ML Pipeline?**
 - Pre-modeling, modelling, post-modelling explainability
- **Is the explanation method general?**
 - Model dependent vs model agnostic
- **What do we explain?**
 - The global model, subgroups, a single prediction
- **How do we represent the explanations?**
 - Feature importance, rules, visualization-based explanation, explanation by example
- **How can we derive explanations?**
 - Explaining by removing, local models, gradient-based

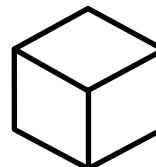
Explainability Stages

Explainability involves the entire AI development pipeline



Pre-modelling explainability

- Before building the model
 - Data exploration
 - Data selection
 - Feature engineering



Explainable modeling

- Build inherently interpretable models
 - Manage the accuracy and interpretability trade-off



Post-modelling explainability

- After model development
 - Explaining predictions and behavior of trained models

Generalizability of Explainability

Is the explanation method general?

**Model dependent
solutions**

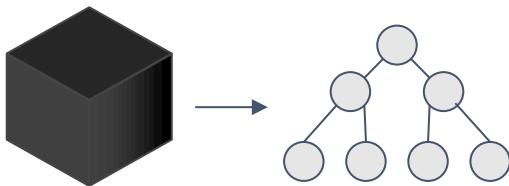
- Only applicable for specific models
- Rely on the model structure/properties

**Model agnostic
solutions**

- Applicable to any model
- Rely on the model as an oracle (model predictions, output probabilities)

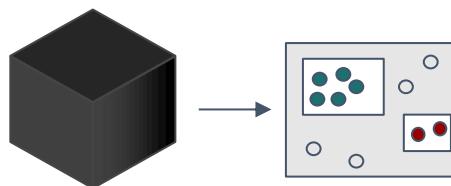
Scope of Explainability

What do we explain?



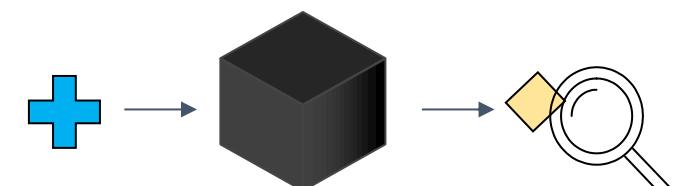
Global

How the model
works globally



Subgroup

How the model
behaves in data
subgroups



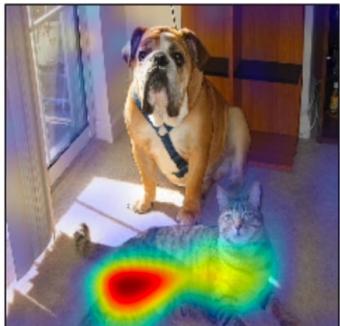
Individual/local

Explaining the
reasons behind
individual predictions

Multiple forms of explanation representation

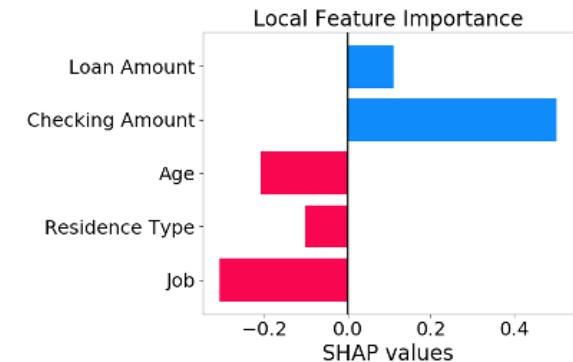
How do we represent the explanations?

Input attributions



(c) Grad-CAM ‘Cat’

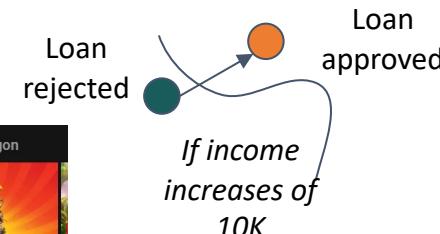
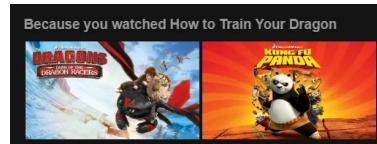
I am really happy



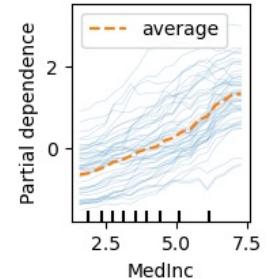
Rules

If Country is US, married, work hours > 45 → Income>50K

Explanations by example



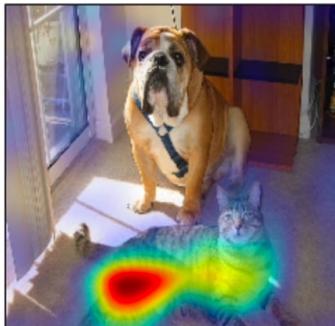
Visualizations



Multiple forms of explanation representation

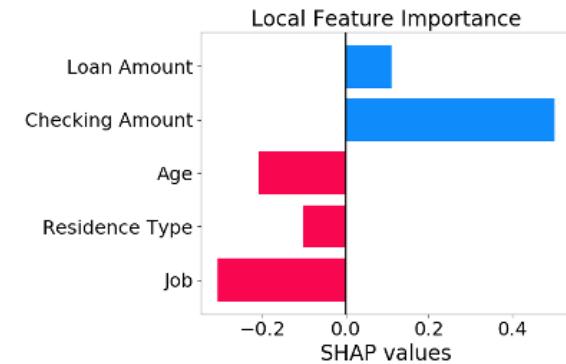
How do we represent the explanations?

Input attributions



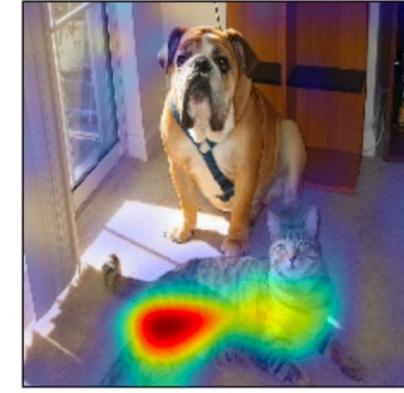
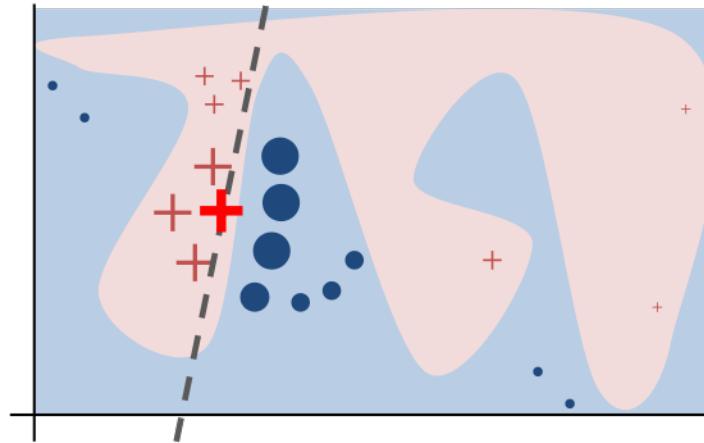
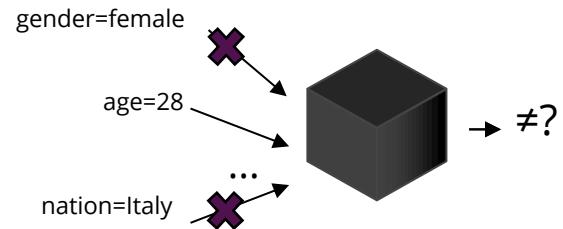
(c) Grad-CAM ‘Cat’

I am really happy



Identify which input features (e.g., pixels, words, attributes) contribute most to the model's output

Multiple methodologies to derive input attribution explanations

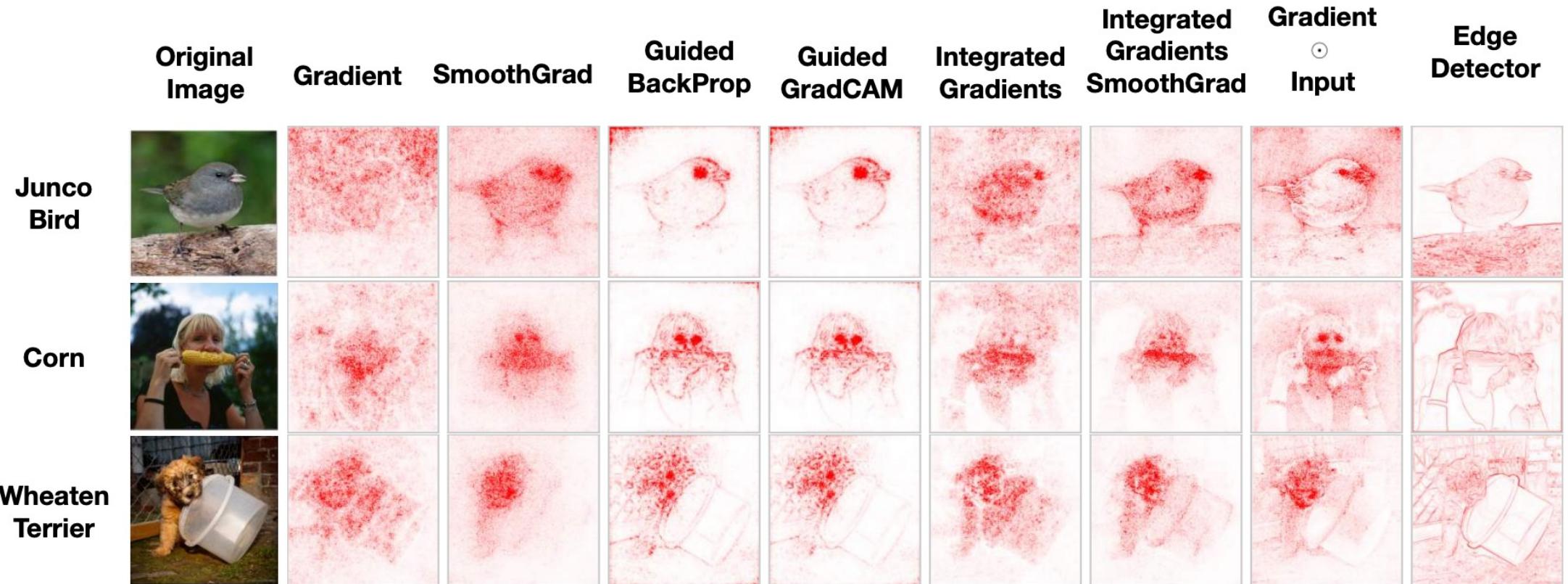


Explaining by removing
SHAP, ..

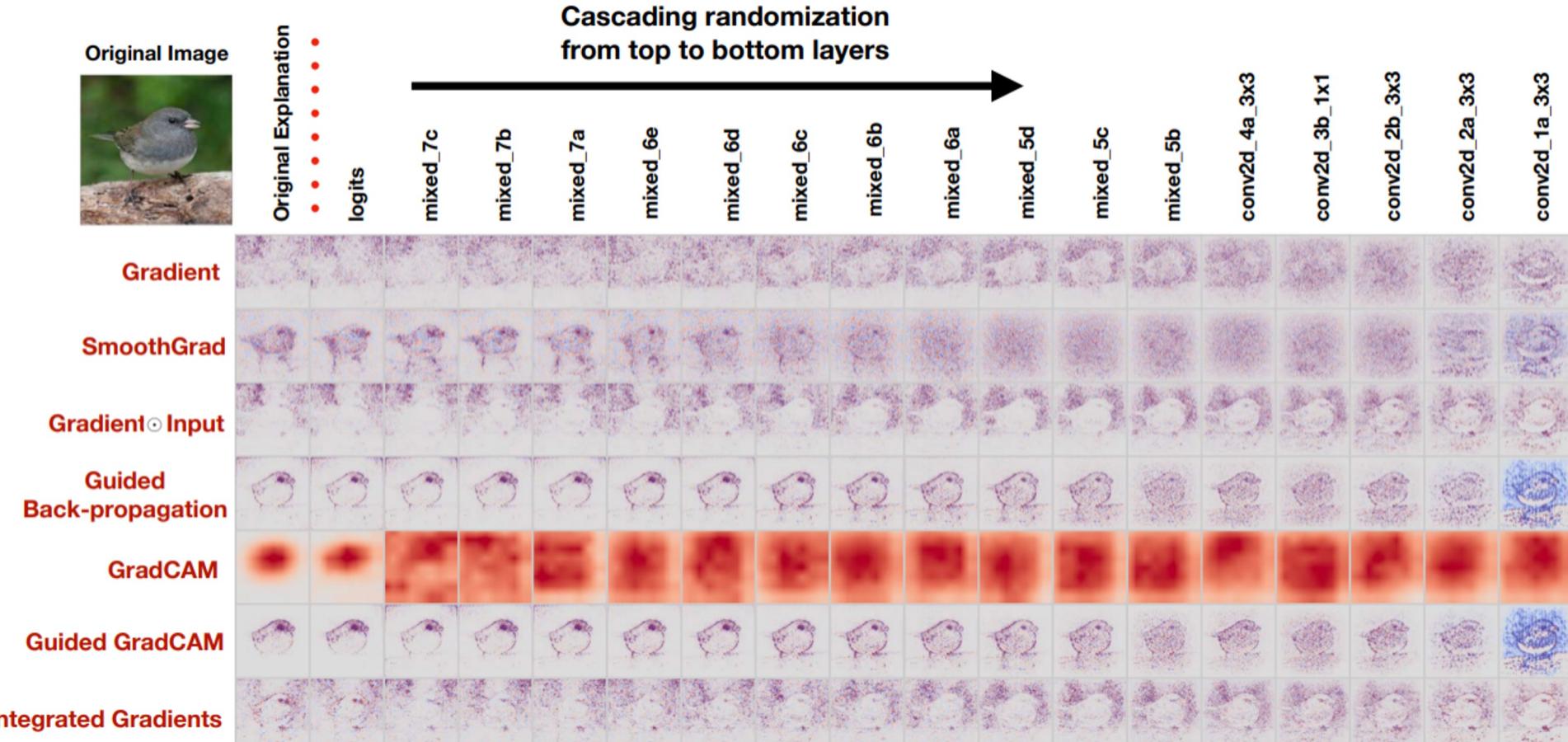
Local models
LIME, ..

Gradient-based
IntegratedGradient,
GradCAM, ..

Standard Explainable AI does not always work well- Edge detectors?



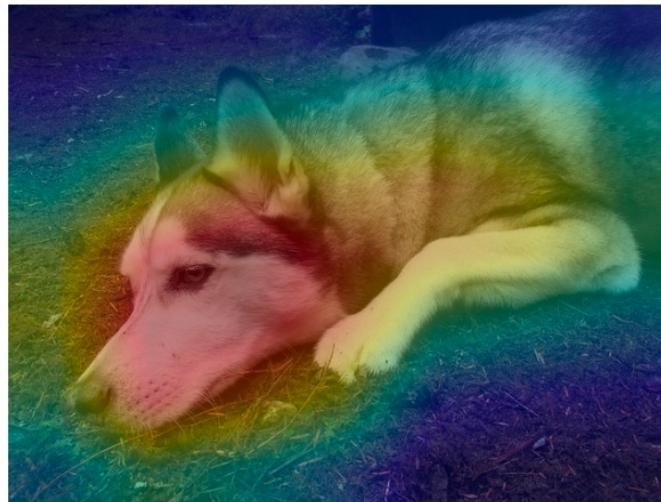
Standard Explainable AI does not always work well- Faithful to model behavior?



Standard Explainable AI may be difficult to interpret- Which class are we explaining?



Evidence for Siberian Husky



Evidence for Transverse Flute



The where vs the what

Showing *where* a network is looking
does not tell us what the network is seeing
in a given input

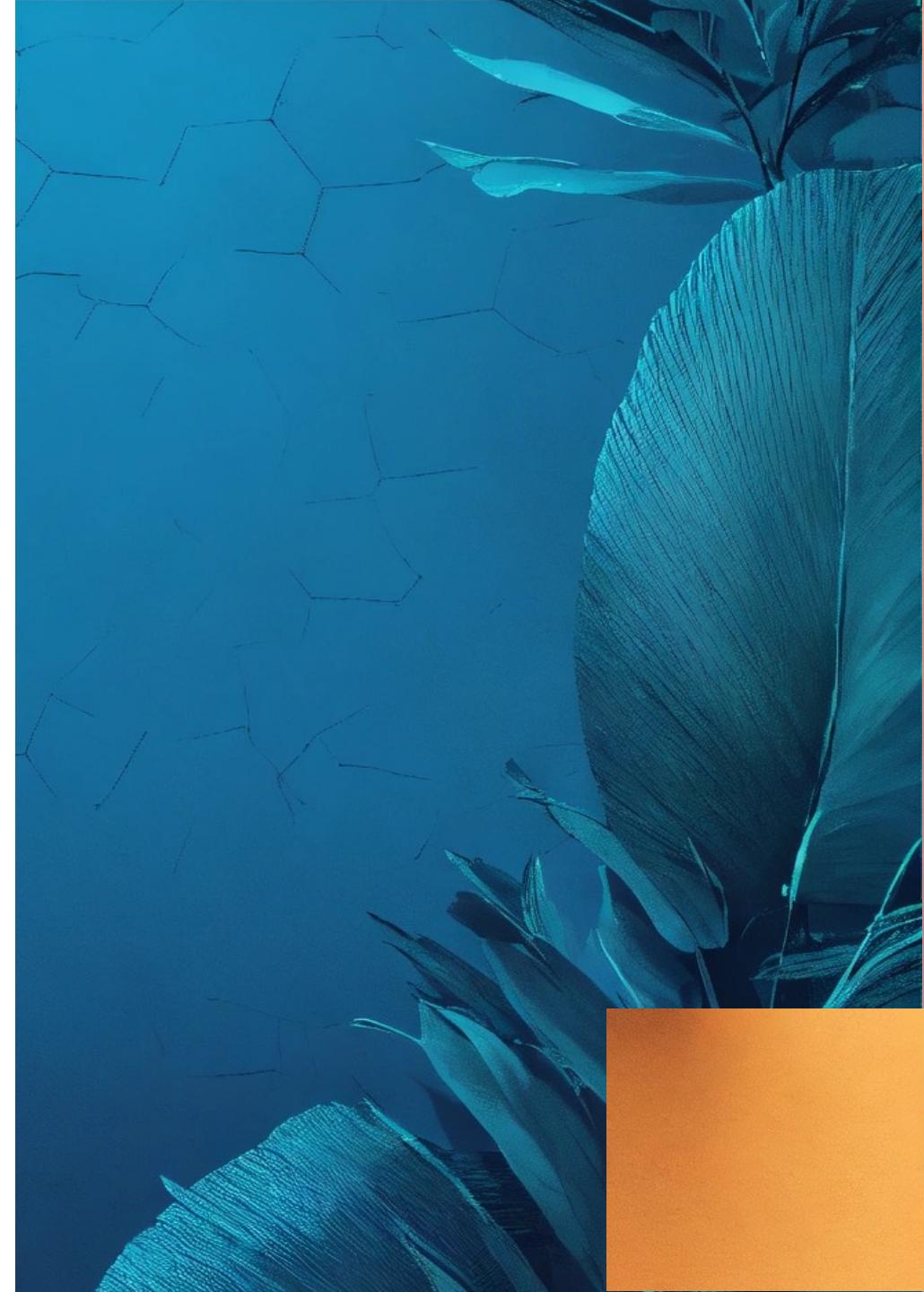
Beyond Input Attribution

Concept-Based XAI (C-XAI)

Explain model behavior via *high-level, human-interpretable* concepts

Tutorial structure

- Introduction to XAI techniques
- **Taxonomy of C-XAI**
- Post-hoc vs By-design C-XAI approaches
- Evaluation, Resources and applications
- Hands-on session





Taxonomy of C-XAI

- **Definition of concept**
- **Types of concepts**
 - Symbolic Concepts, Unsupervised Concept Basis, Prototypes, Textual Concepts
- **Types of concept-based explanations**
 - Class-Concept Relations , Node-Concept Association , Concept-Visualization
- **Stages of Explainability**
 - Post-hoc, By design
- **Annotation availability**
 - Supervised/annotated concepts, Unsupervised, Hybrid, Generated
- **Overview of C-XAI approaches**

What is a concept?

A concept can be any abstraction, such as a colour, an object, or even an idea



Taxonomy of C-XAI

- Definition of concept
- **Types of concepts**
 - Symbolic Concepts, Unsupervised Concept Basis, Prototypes, Textual Concepts
- Types of concept-based explanations
 - Class-Concept Relations , Node-Concept Association , Concept-Visualization
- Stages of Explainability
 - Post-hoc, By design
- Annotation availability
 - Supervised/annotated concepts, Unsupervised, Hybrid, Generated
- Overview of C-XAI approaches

Types of concepts

- **Symbolic Concepts**

Human-defined attributes

“BEAK”

“GREEN”

“WINGS”



- **Unsupervised Concept Basis**

Cluster of similar samples

- **Prototypes**

(Part-of) a training sample



- **Textual Concepts**

Textual representation

“A bird with bright feathers”

Symbolic Concepts

“BEAK”

“GREEN”

“WINGS”

- **Human-defined attributes or abstractions**
 - Of the final classes
 - e.g., bird --> the beak of the bird, the color of the bird
- **Require auxiliary data & annotations**
 - **Instance-level annotation**
 - Annotate the presence of a concept for each instance
 - More expensive
 - **Class-level annotation**
 - All samples belonging to a class are annotated as having a certain attribute
 - Less expensive but less precise (e.g., attribute could not be visible)

Unsupervised Concept Basis



- **Cluster of similar samples**
 - Extracted from the network representation (the latent space)
 - e.g., a cluster of green birds.
- Can capture abstractions more understandable to humans than individual features or pixels
- Clustering algorithms employed to extract unsupervised concepts

Prototypes



- **Explanation by example**
- **Representative examples** of peculiar traits of the training samples
 - **Entire samples**
 - **Parts of a training sample** (e.g., a hooked beak)
- Different from unsupervised concept basis
 - Represent a single example instead of a group of examples

Textual Concepts

“A bird
with bright
feathers”

- **Textual descriptions of main classes**
 - From an individual description, distinctive parts are extracted
- **Provided at training time by an external generative model**
 - It requires a Large-Language Models with knowledge of the given task



Taxonomy of C-XAI

- Definition of concept
- Types of concepts
 - Symbolic Concepts, Unsupervised Concept Basis, Prototypes, Textual Concepts
- **Types of concept-based explanations**
 - Class-Concept Relations, Node-Concept Association, Concept-Visualization
- Stages of Explainability
 - Post-hoc, By design
- Annotation availability
 - Supervised/annotated concepts, Unsupervised, Hybrid, Generated
- Overview of C-XAI approaches

Types of concept-based Explanations

- **Class-Concept Relations**

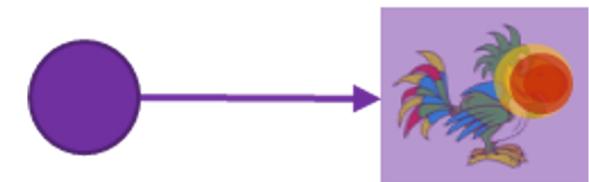
Relation among a concept and an output class of a model

Beak → Parrot



- **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

Class-Concept Relations

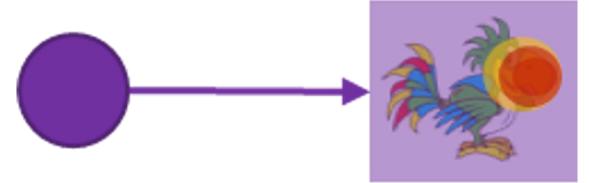
- **Relationship between a concept and an output class of the model**
 - Concept importance
 - Logic rule involving multiple concepts and their connection to an output class
- Can be applied to all type of concepts
 - e.g., **symbolic**, $\text{parrot} = 0.3 \text{ beak} + 0.2 \text{ green} + 0.3 \text{ wings} + 0.2 \text{ tropical}$
 - e.g., with **prototypes**, $\text{parrot} := 0.8 \text{ prototype}_1 + 0.2 \text{ prototype}_2$

Node-Concept Association



- Assign a concept to an internal unit (or a filter) of a network
- It 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
- It can also be forced during training
 - by requiring a unit to predict a concept

Concept Visualization



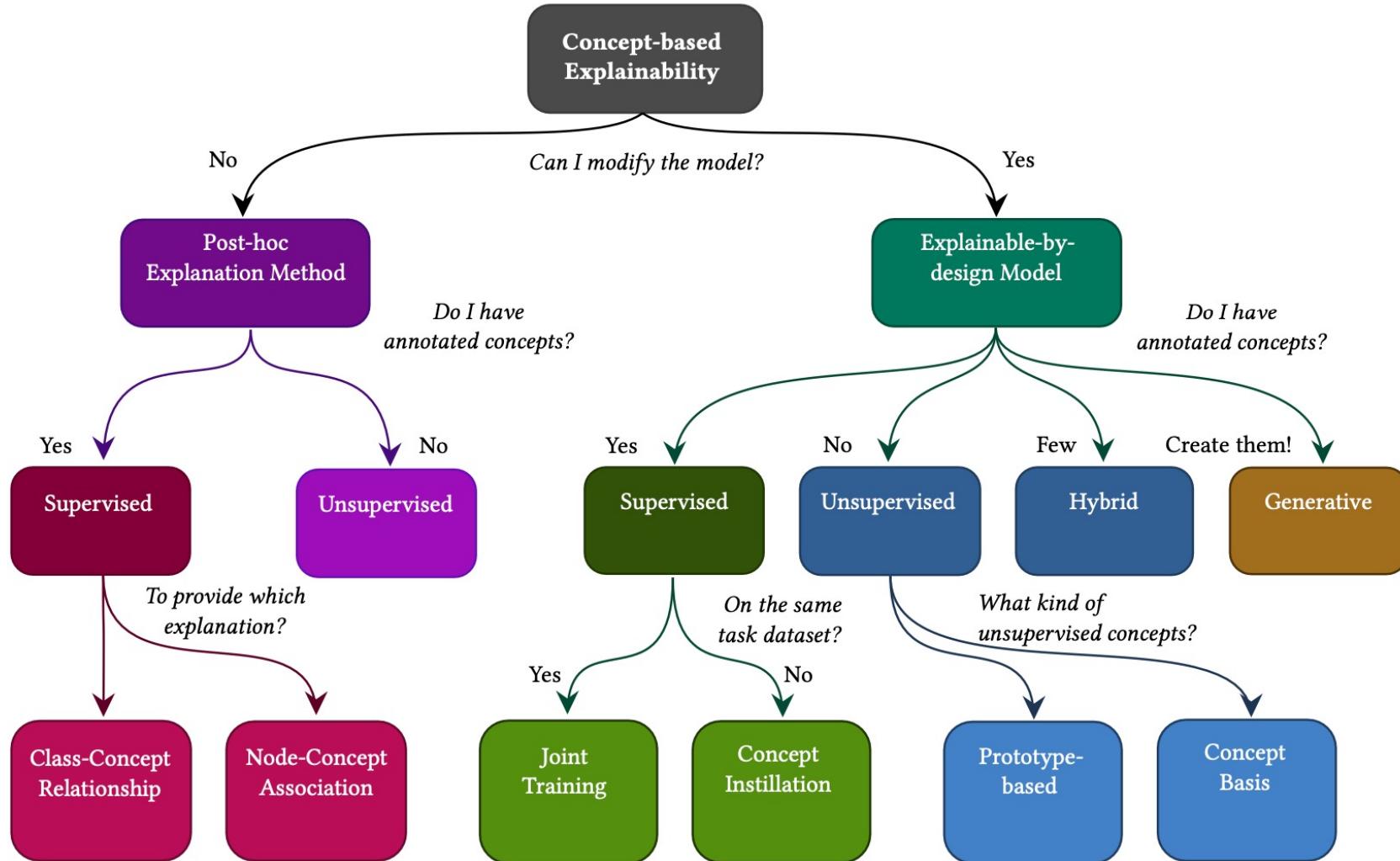
- **Highlight the input features that best represent a specific concept**
 - Similar to saliency map but for concepts
- **Crucial when non-symbolic concepts are employed**
 - Need to understand which **unsupervised concept basis or prototypes** the network has learned
- Often combined with one of the previous explanations
 - Enable **understanding the concepts associated** with a specific class or node



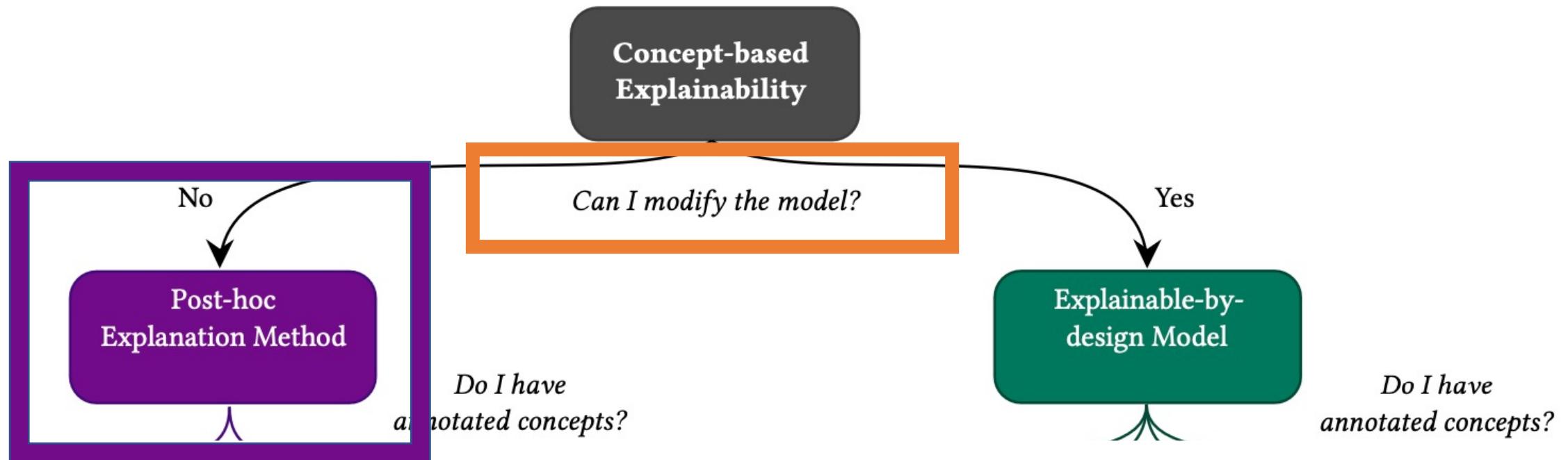
Taxonomy of C-XAI

- **Definition of concept**
- **Types of concepts**
 - Symbolic Concepts, Unsupervised Concept Basis, Prototypes, Textual Concepts
- **Types of concept-based explanations**
 - Class-Concept Relations , Node-Concept Association , Concept-Visualization
- **Stages of Explainability**
 - Post-hoc, By design
- **Annotation availability**
 - Supervised/annotated concepts, Unsupervised, Hybrid, Generated
- **Overview of C-XAI approaches**

C-XAI Taxonomy



C-XAI Taxonomy



Post-hoc Concept-based Explanation methods

Explain via concepts an **already trained model** that we cannot (or do not want to) modify

Underlying idea is that **a network automatically learns to recognize some concepts**

Post-hoc Concept-based Explanation methods

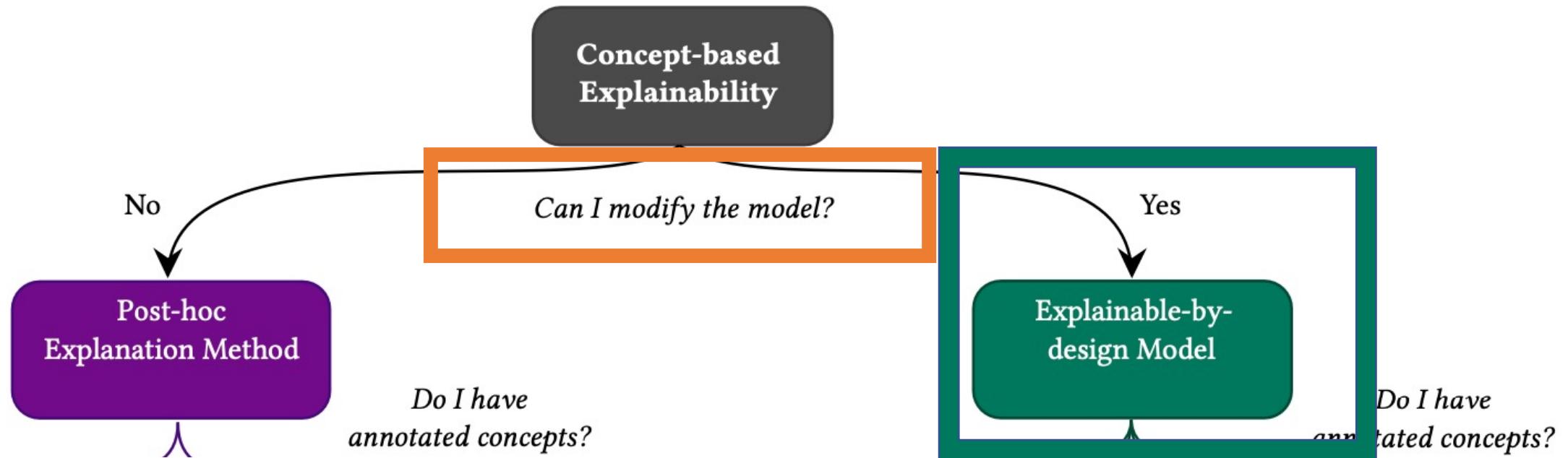
Pros.

- They **do not compromise the learning capacity** of a model
- They can provide more interpretable explanations than standard post-hoc methods

Cons.

- **Cannot ensure the network really knows the concepts**
 - It has not been trained for that
 - Not interpretable by design

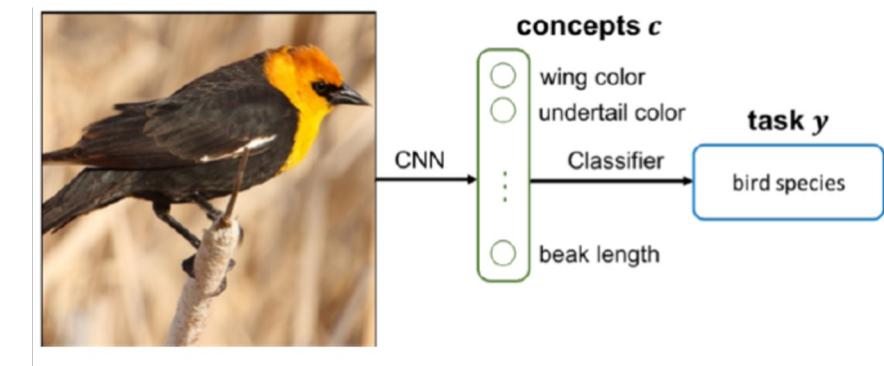
C-XAI Taxonomy



Explainable-by-design Concept-based Models

Directly **train concept-aware models**

- Neural models with an **explicit concept representation**, typically as an intermediate layer
- Task predictions are influenced by the predicted concepts



Explainable-by-design Concept-based Models

Pros.

- They can be regarded as **inherently transparent models** as they provide **node-concept association** by-design



Cons.

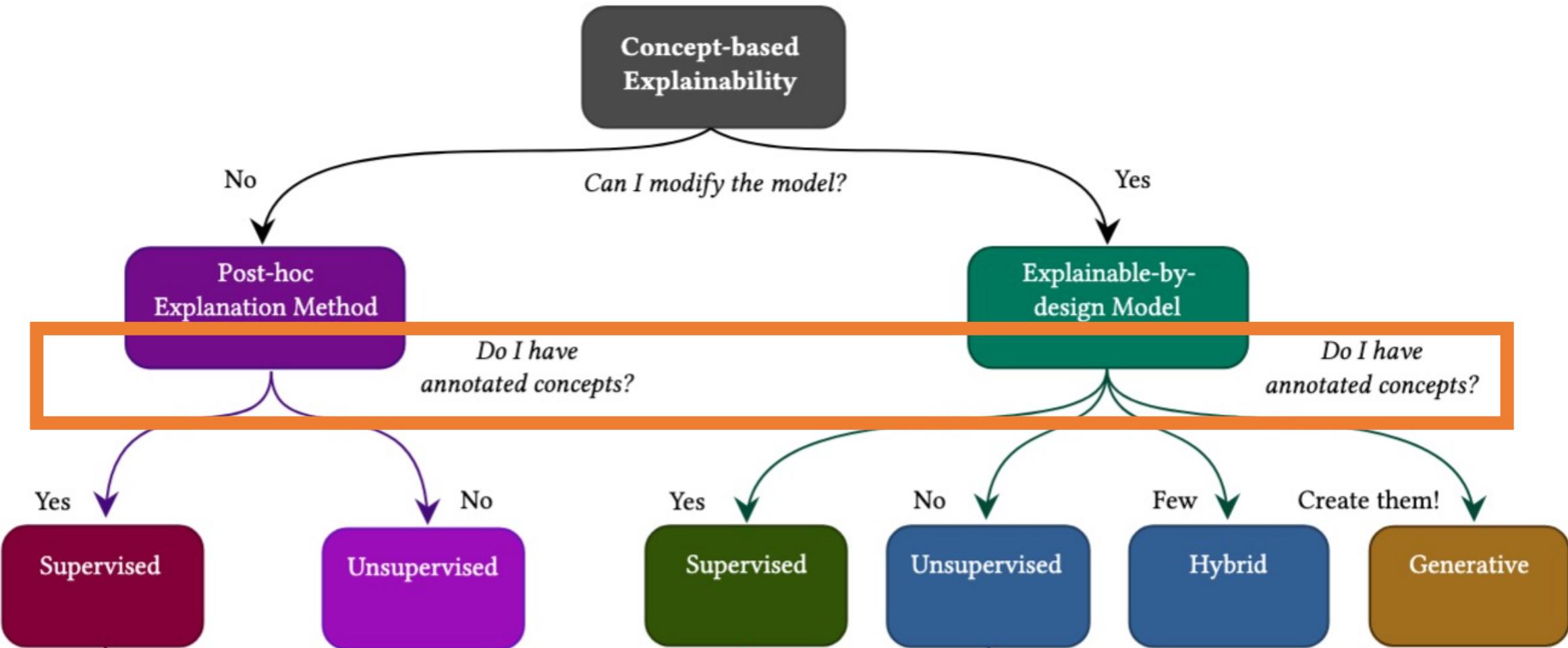
- They need ad-hoc training
- Predicting concepts **might reduce network task performance**



Taxonomy of C-XAI

- **Definition of concept**
- **Types of concepts**
 - Symbolic Concepts, Unsupervised Concept Basis, Prototypes, Textual Concepts
- **Types of concept-based explanations**
 - Class-Concept Relations , Node-Concept Association , Concept-Visualization
- **Stages of Explainability**
 - Post-hoc, By design
- **Annotation availability**
 - Supervised/annotated concepts, Unsupervised, Hybrid, Generated
- **Overview of C-XAI approaches**

C-XAI Taxonomy



*Do I have
annotations?*

Concept annotation

- **Supervised methods**

Use dataset **annotated with concepts** related to the task

Yes

- **Unsupervised approaches**

Extract concepts from the same data **automatically**

No

- **Hybrid (by-design) approaches**

Leverage **few supervised** concepts **and unsupervised** (extracted) ones;

Few

- **Generative methods**

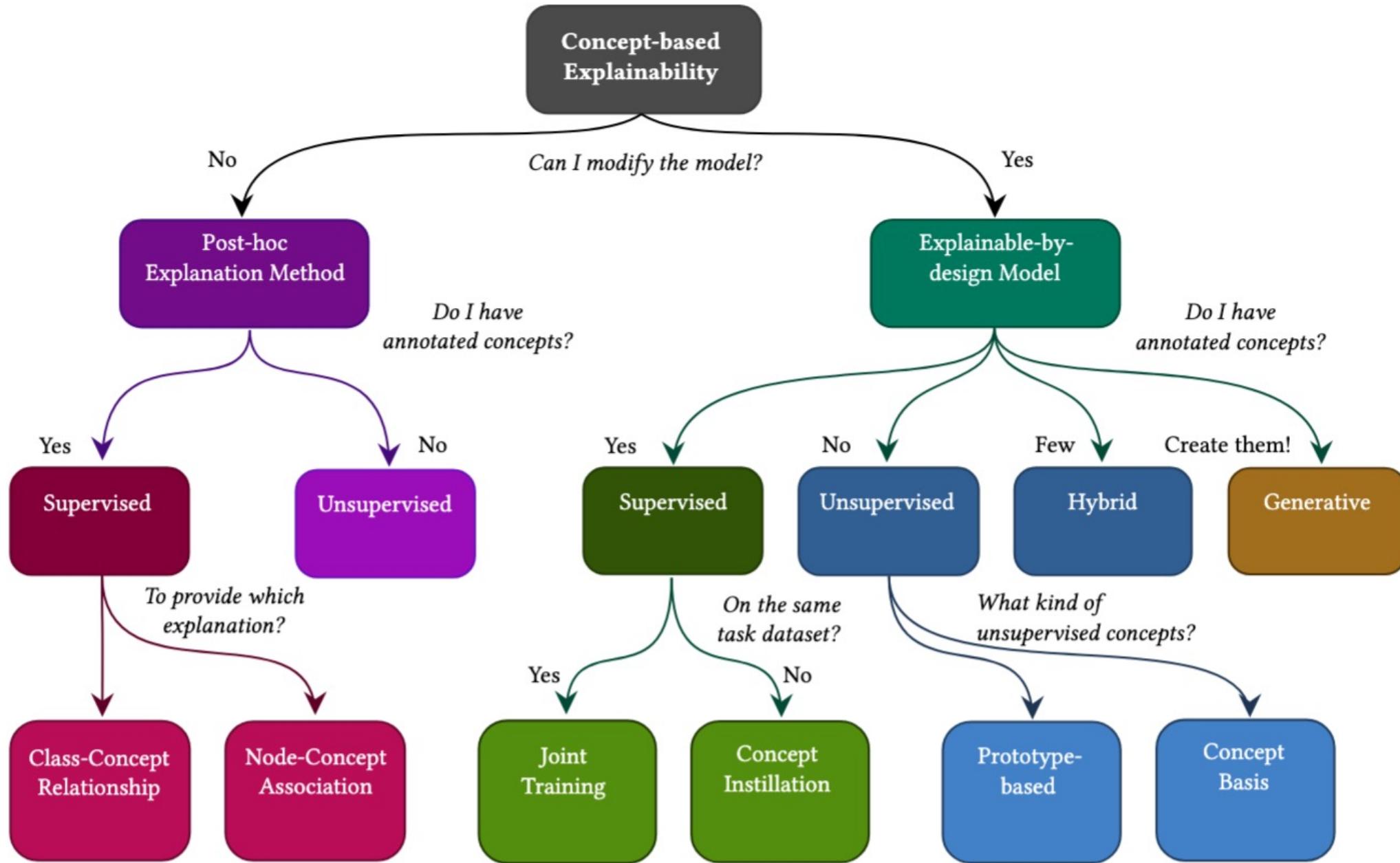
Create concepts by means of an external generative model

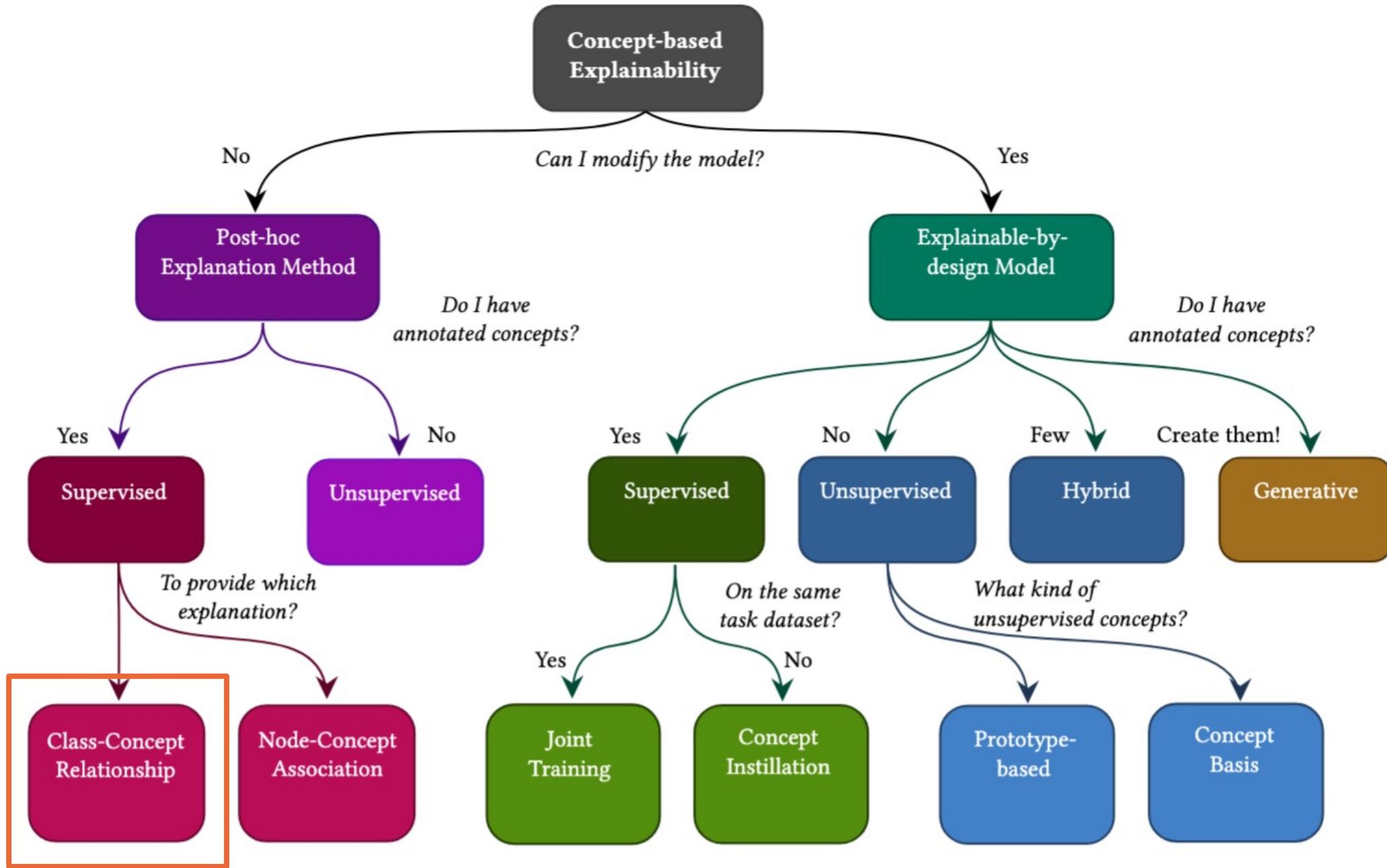
Generate!



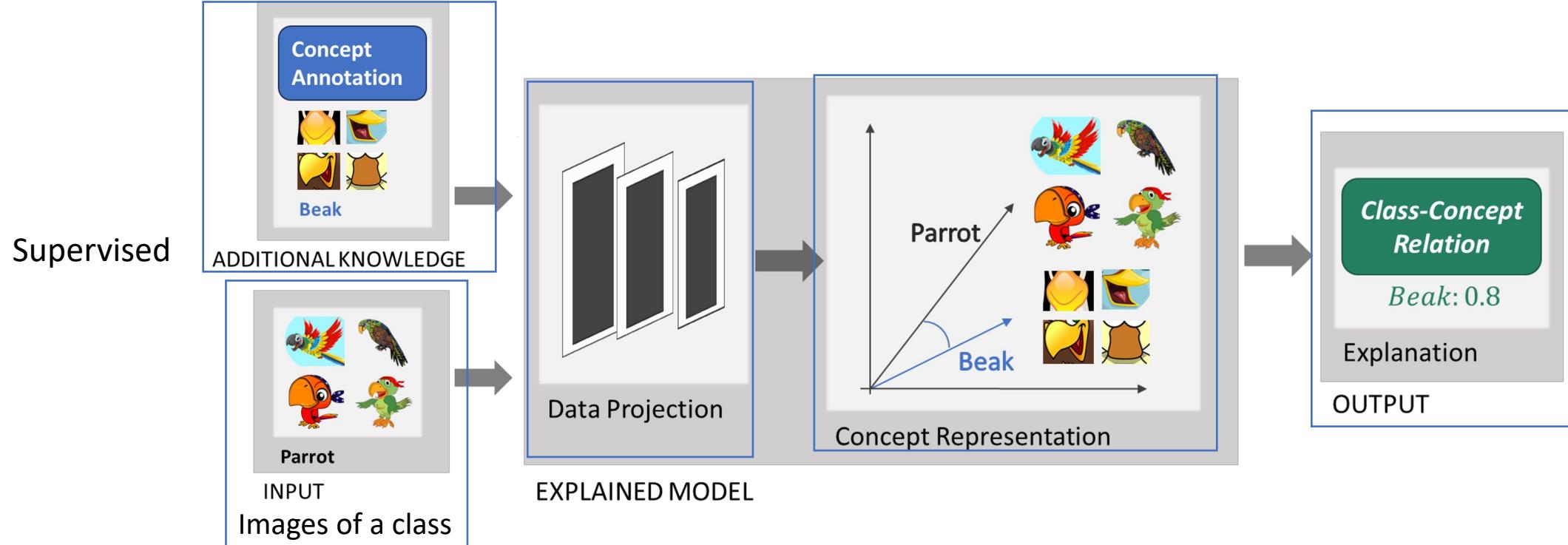
Taxonomy of C-XAI

- **Definition of concept**
- **Types of concepts**
 - Symbolic Concepts, Unsupervised Concept Basis, Prototypes, Textual Concepts
- **Types of concept-based explanations**
 - Class-Concept Relations , Node-Concept Association , Concept-Visualization
- **Stages of Explainability**
 - Post-hoc, By design
- **Annotation availability**
 - Supervised/annotated concepts, Unsupervised, Hybrid, Generated
- **Overview of C-XAI approaches**

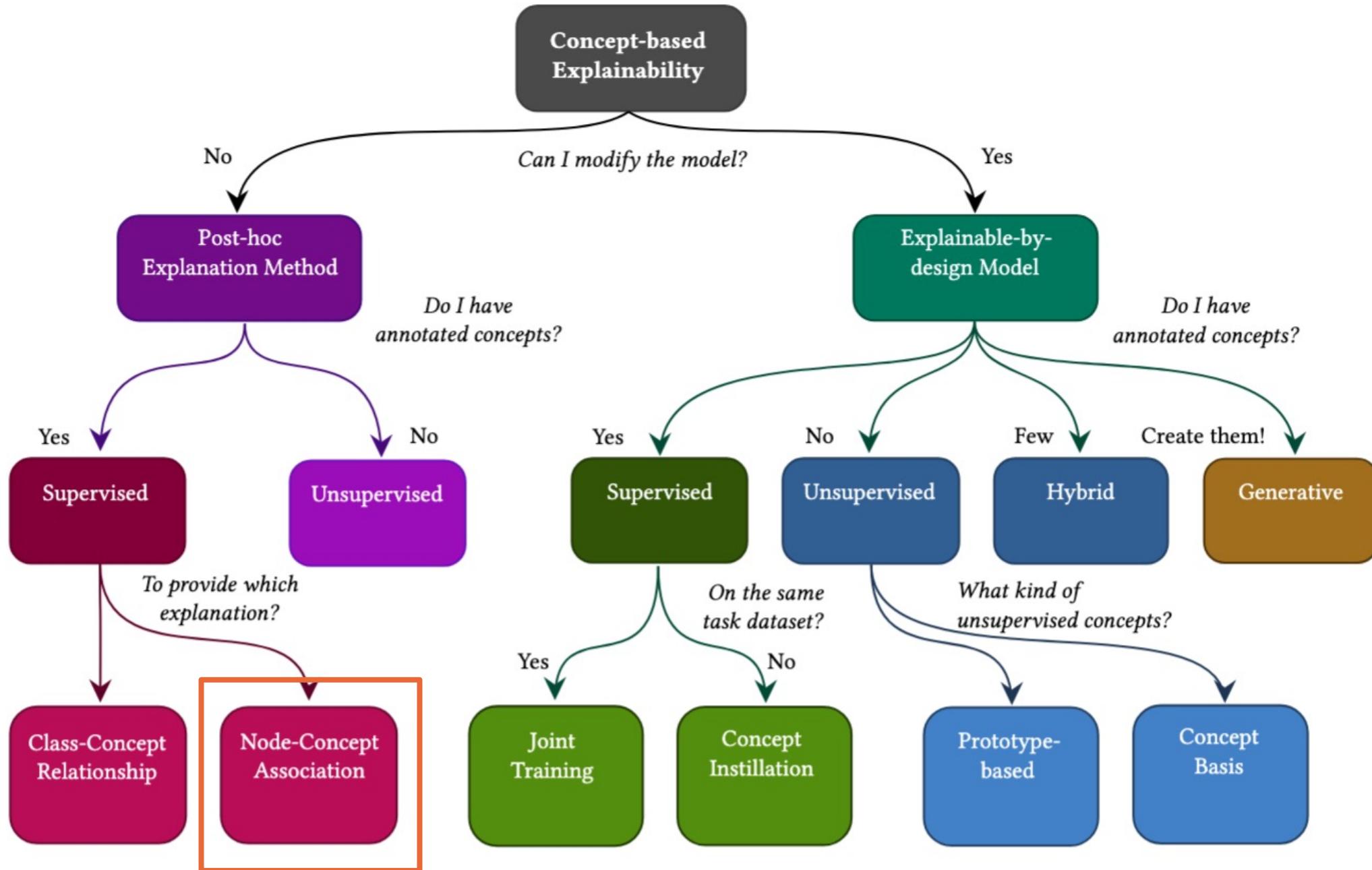




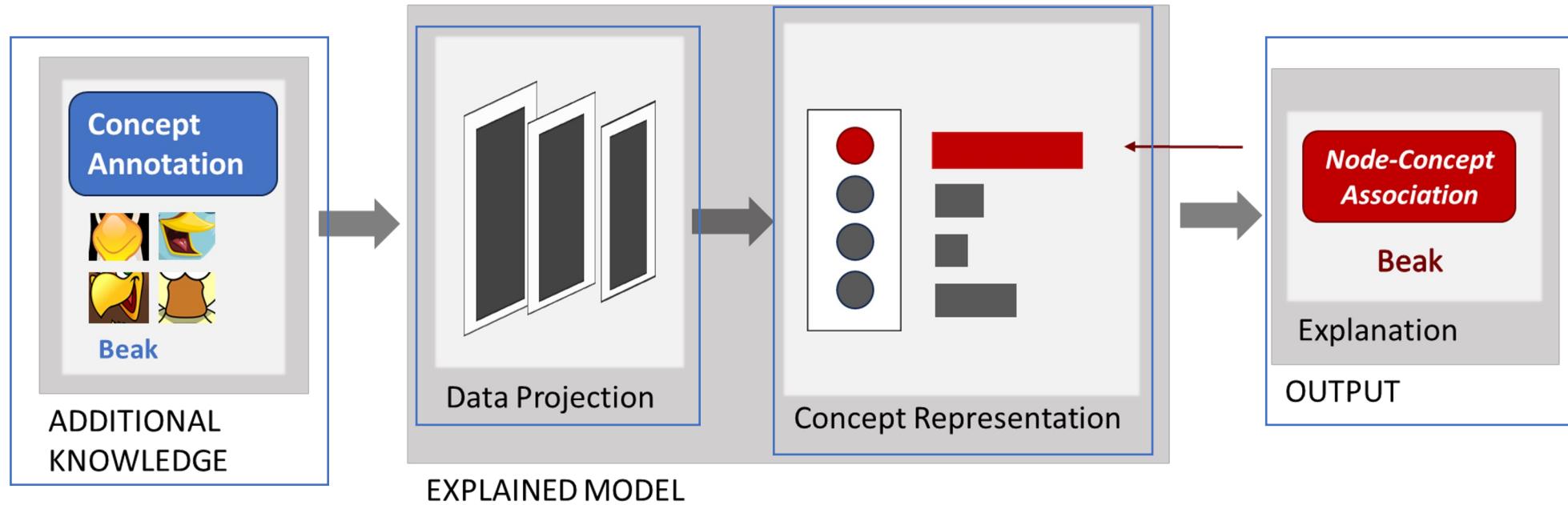
Post-hoc Supervised method providing class-concept relations



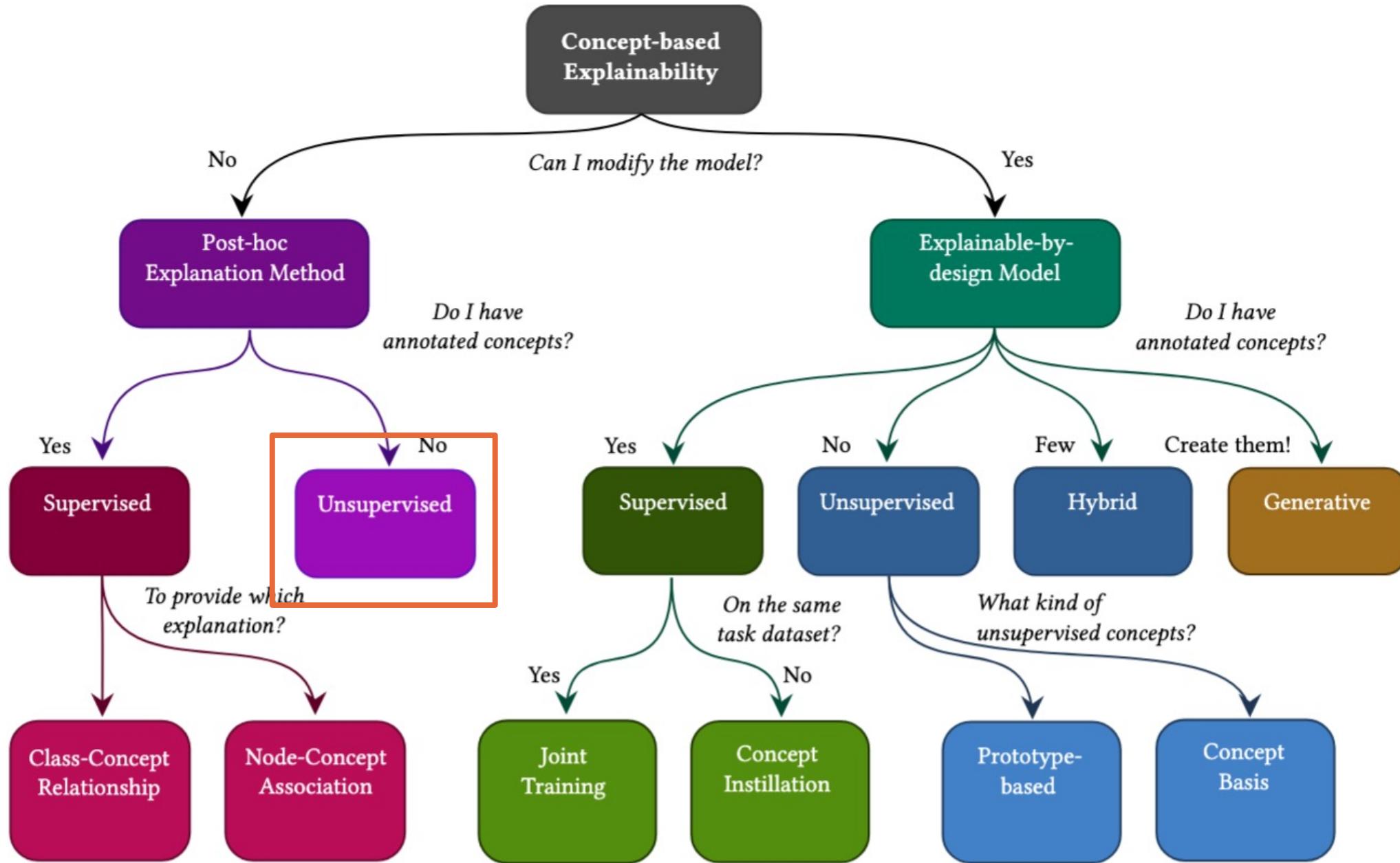
- Analyze the relationship of the representation of the concept data in the model latent space with the output class



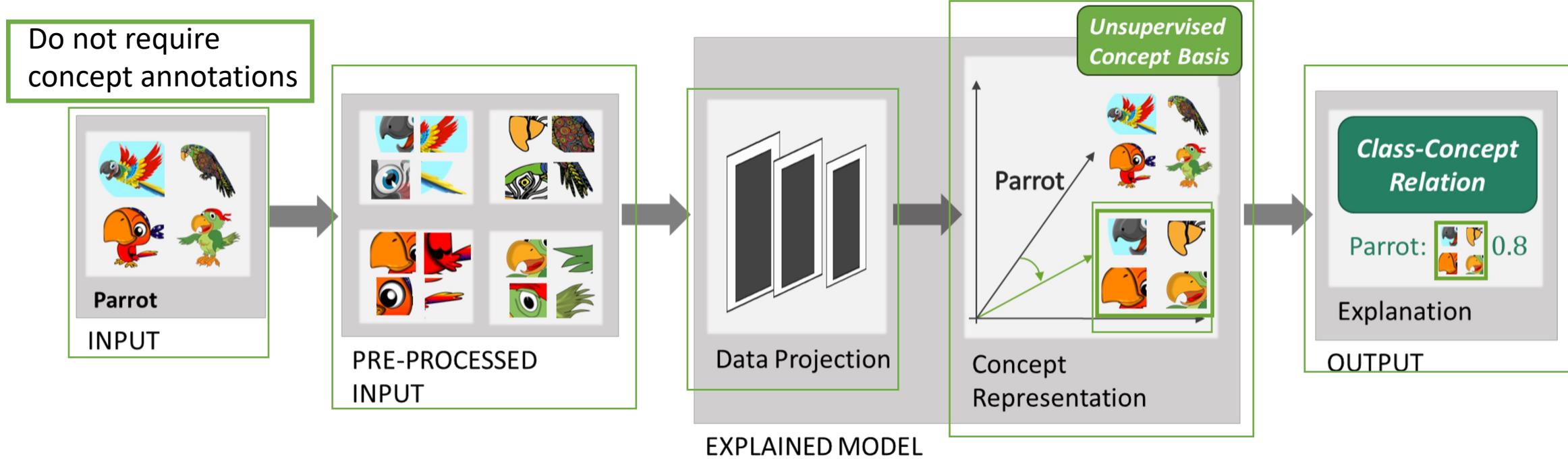
Post-hoc Supervised method providing node-concept association



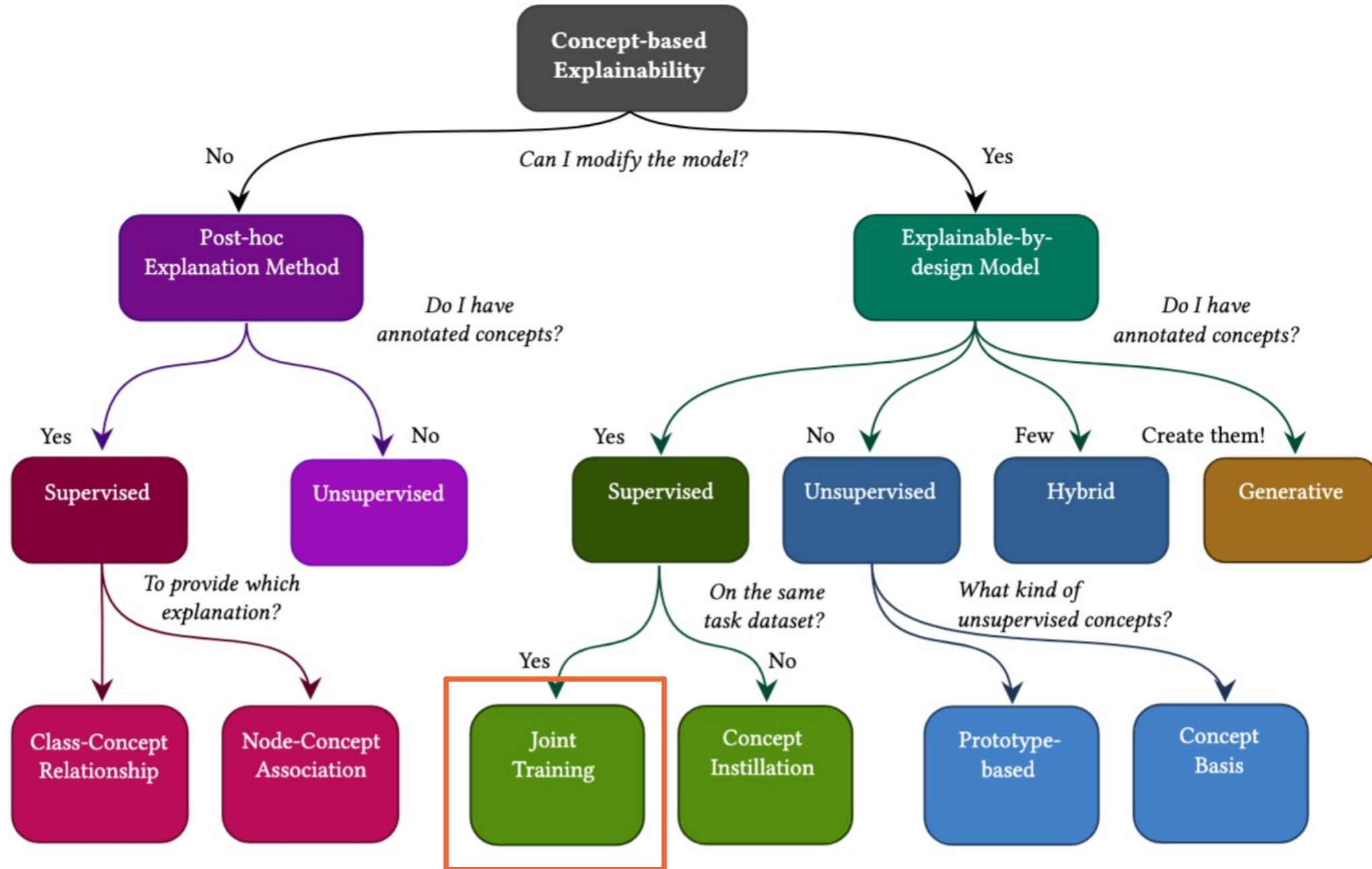
- **Analyze the activations** of the hidden nodes **when fed with the concept annotation data**
- They **associate to each node the concept** for which it is activated the most



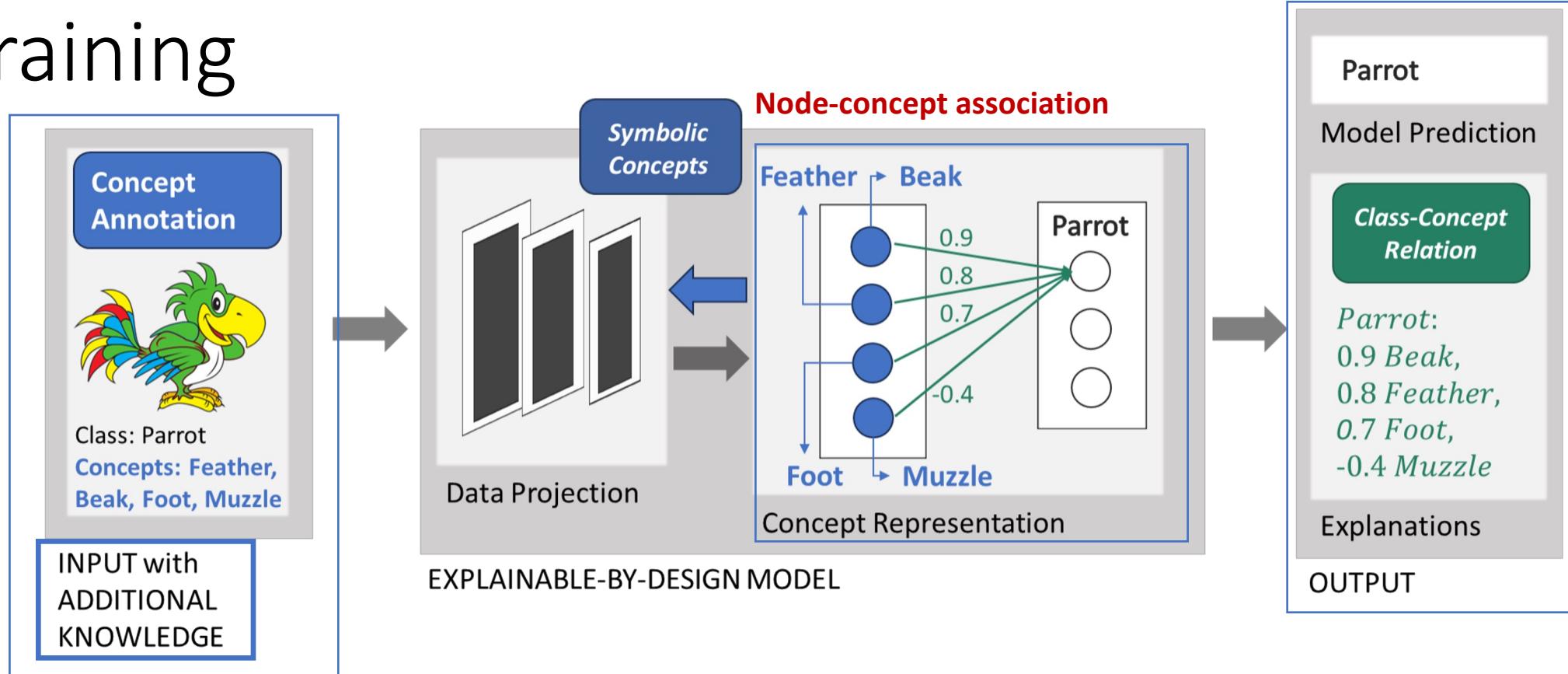
Post-hoc Unsupervised method providing class-concept relation association



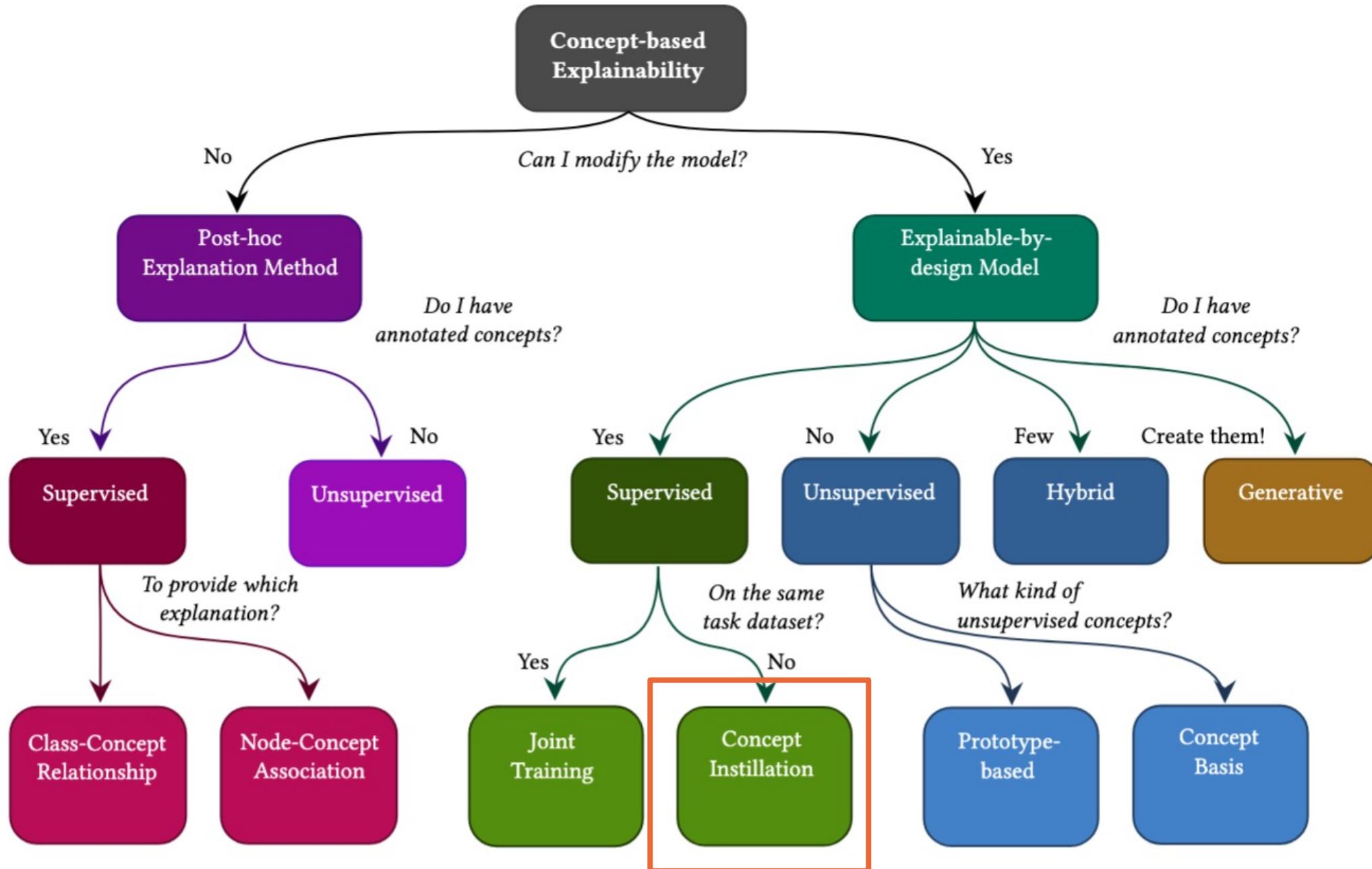
- Pre-process the inputs in parts of the samples
- Clusterize projections of the processed inputs in the latent space → clusters are **unsupervised concepts**
- **Analyze the relation of unsupervised concepts with output classes/predictions**



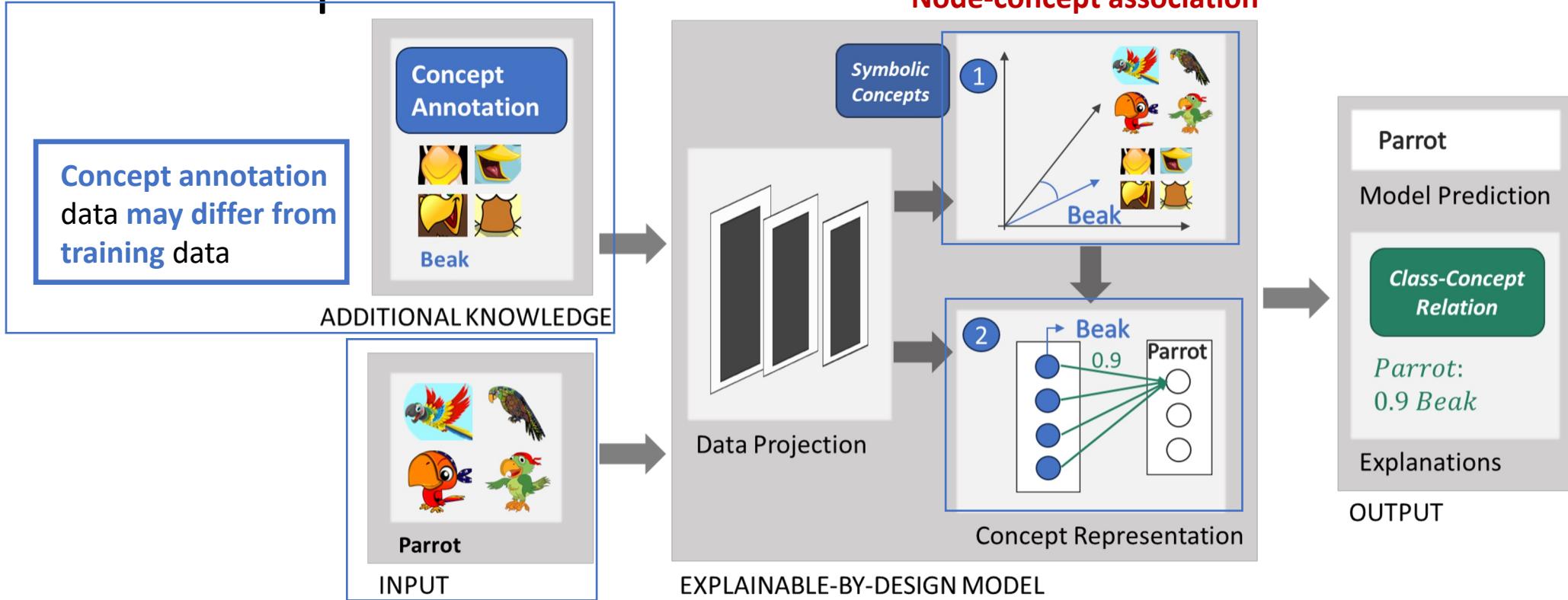
Supervised Concept-based models jointly training



- Train a model with a **hidden layer predicting the concepts**
- The **predicted concepts** are used **to make the final prediction**
 - If the task predictor is a white box model also **extract class-concept relations**

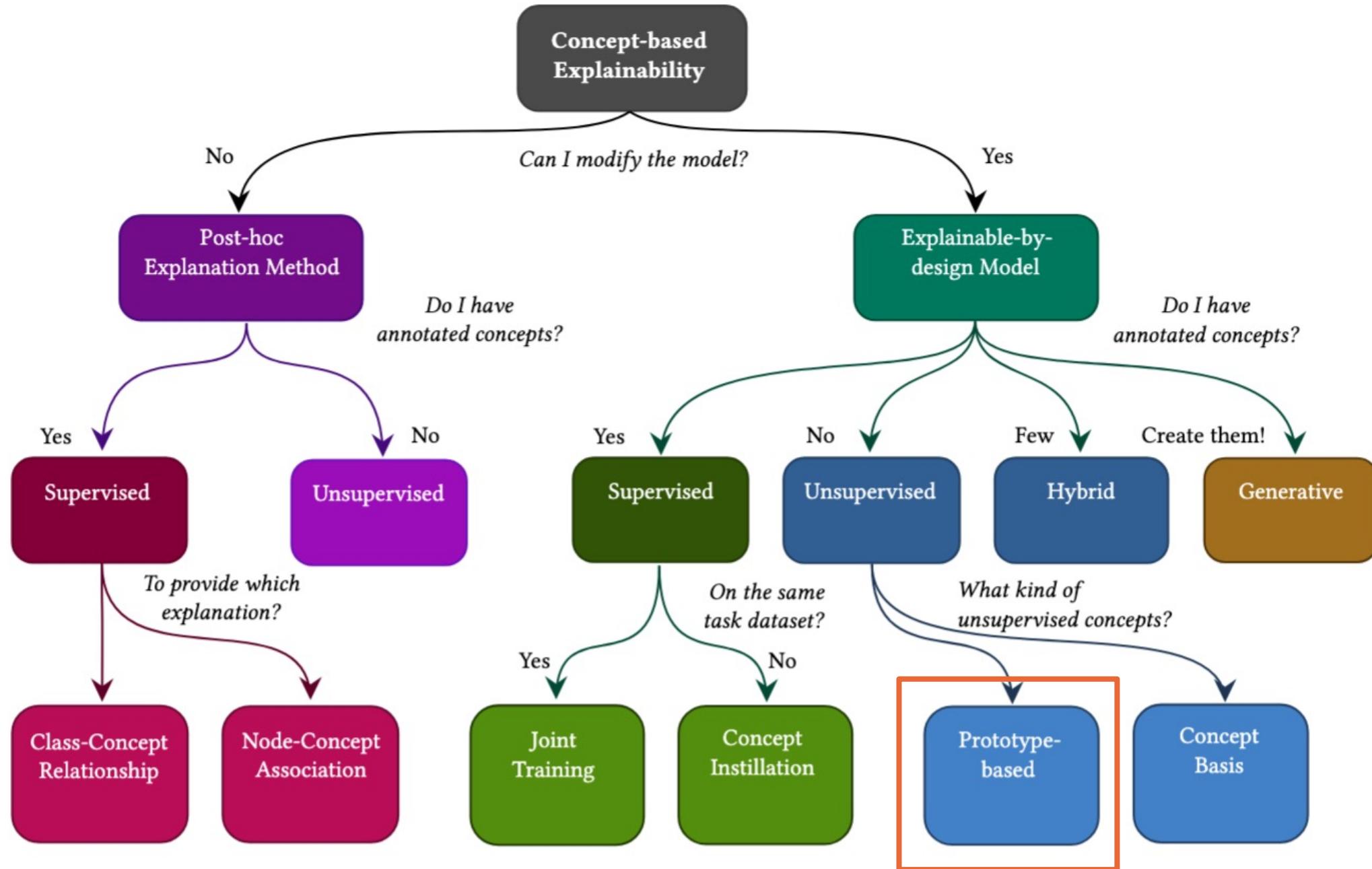


Supervised Concept-based models instilling concepts

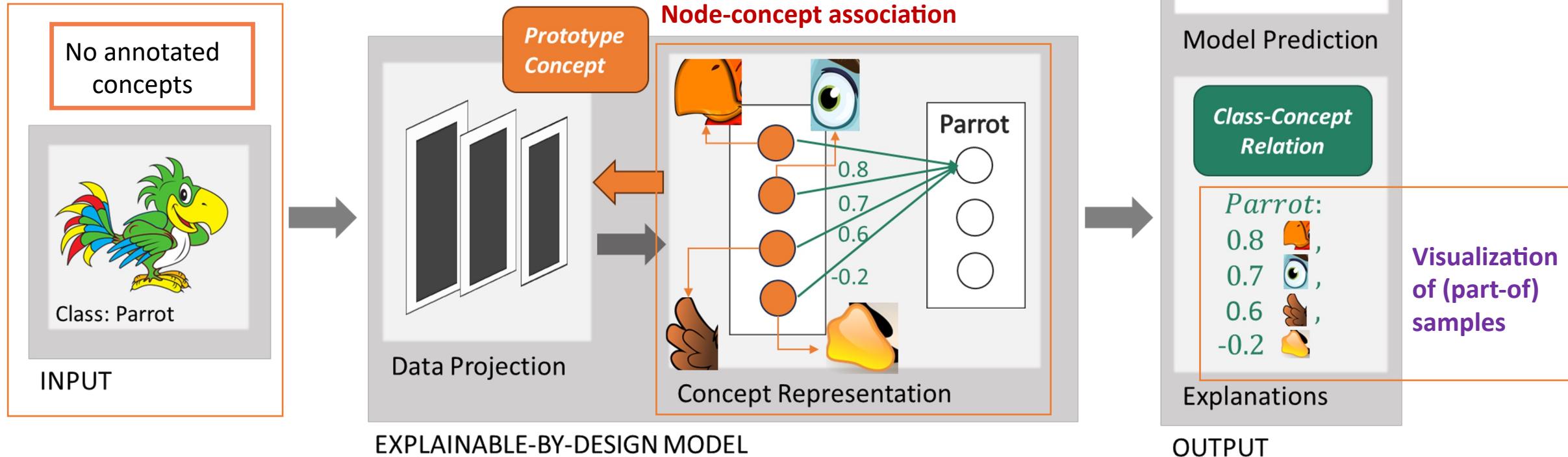


- Turn a **black-box model** into an **explainable-by-design** one:
 - **Concept instillation:** fine-tune a layer to predict the concepts
 - Keep **training the top of the network to predict the classes**

e.g., Chen, Zhi, et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence (2020).

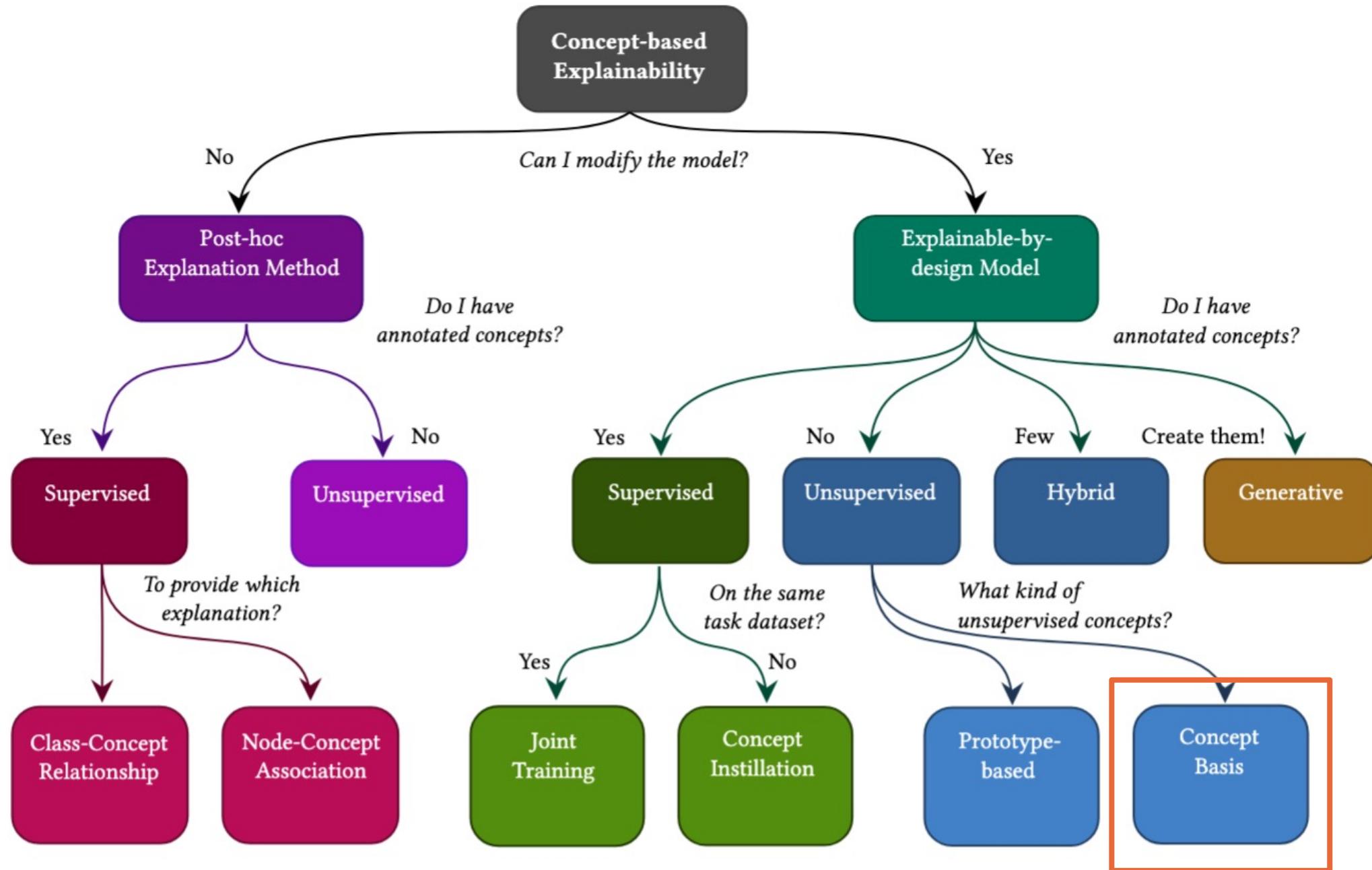


Unsupervised Concept-based models employing prototype concepts

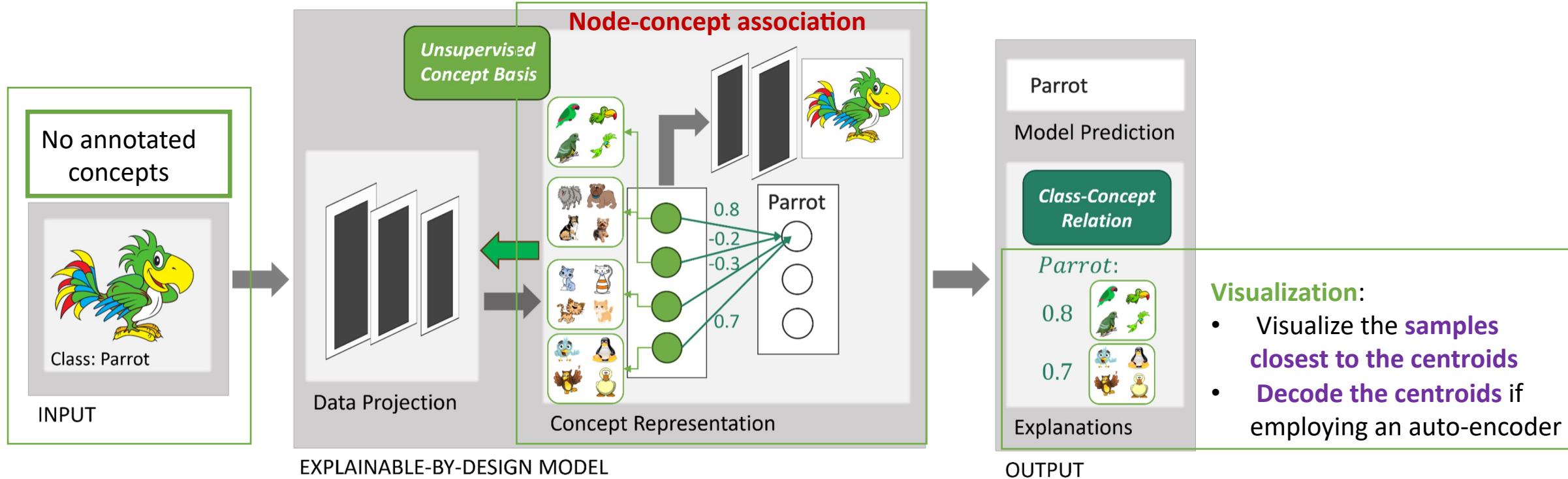


Train the network to:

- Learn to predict the output classes
- **Encode** in the hidden layers **representative (part-of) train examples**

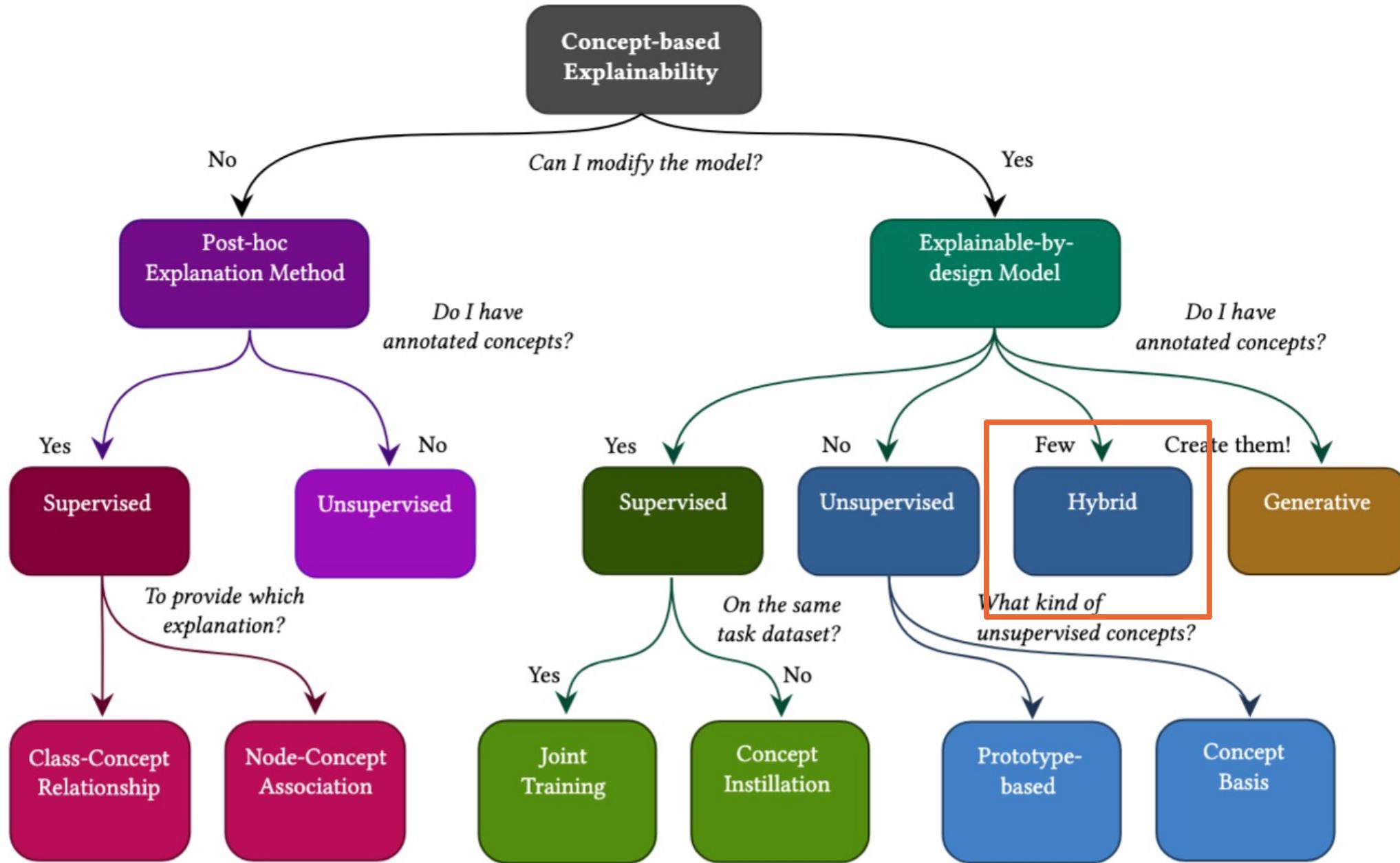


Unsupervised Concept-based models employing unsupervised concept basis

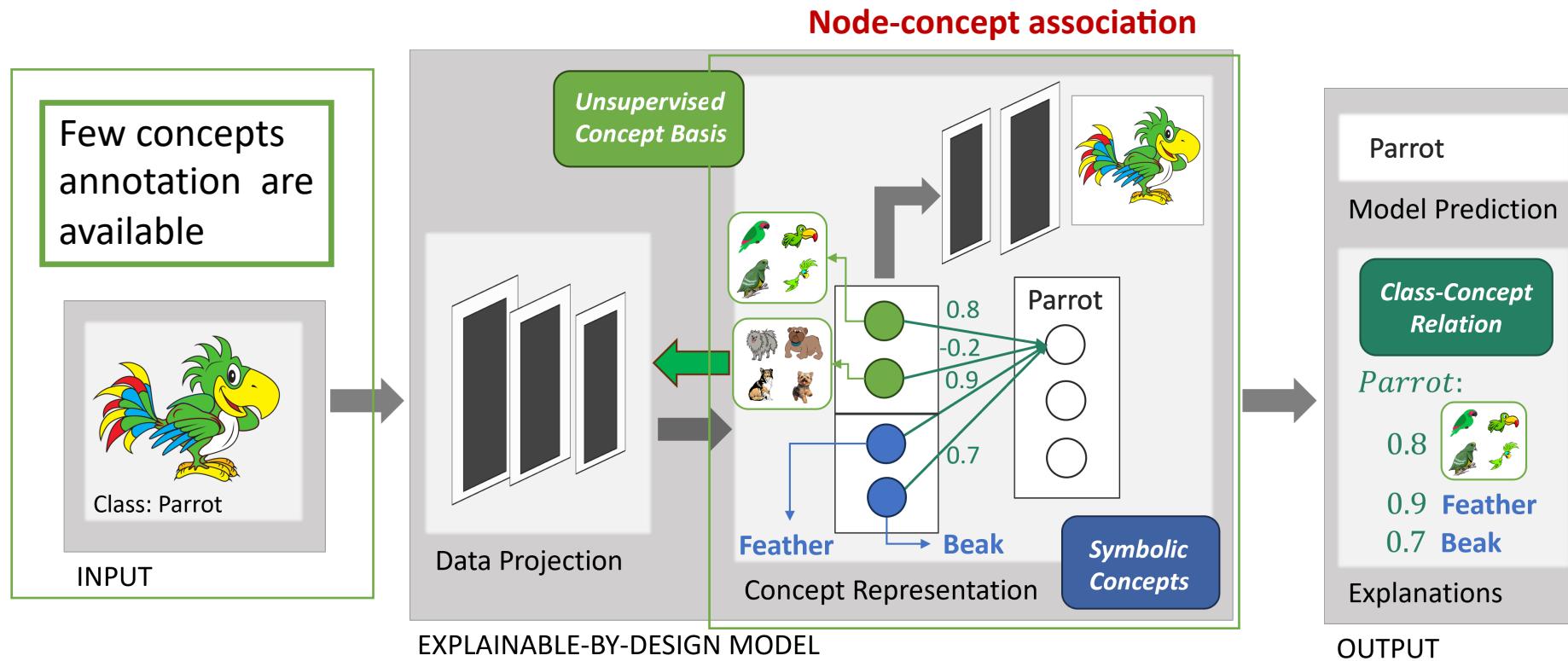


They train the network to:

- Learn to predict the output class
- **Create clusters of samples in the latent representation**, associating clusters to nodes

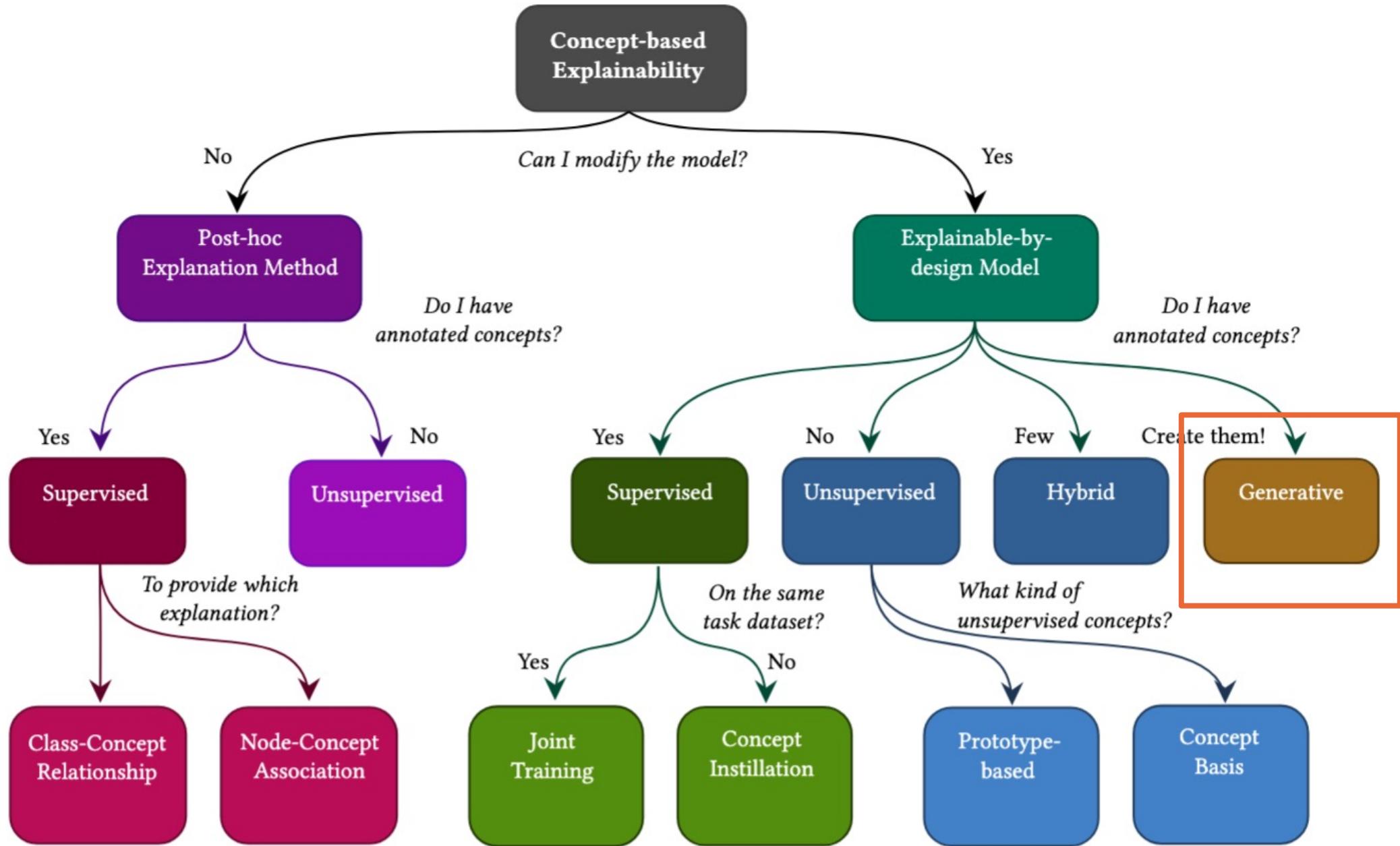


Hybrid Concept based Models

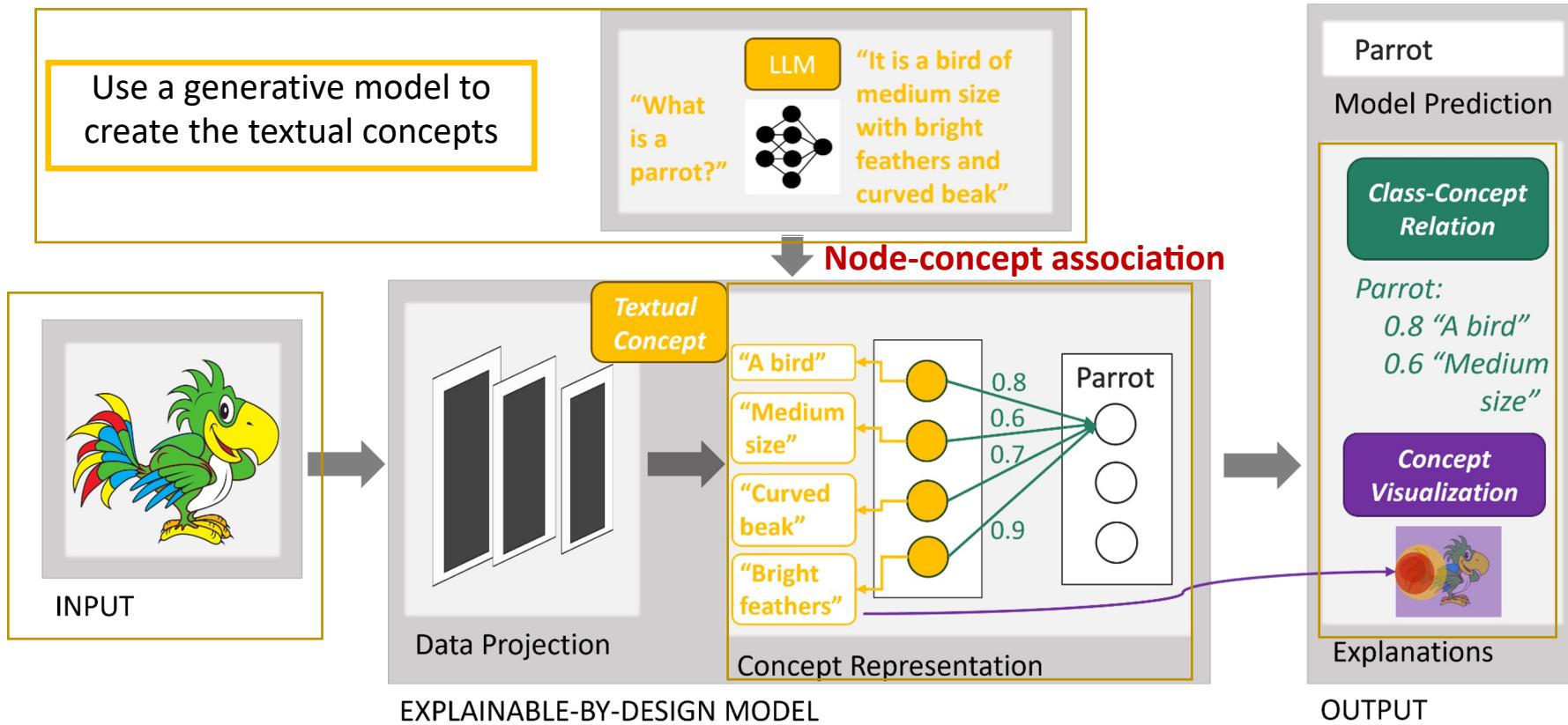


Train the network to:

- Learn to **predict a given set of concepts** with a **subset of neurons**
- Create a **clusterized representation** in the remaining neurons



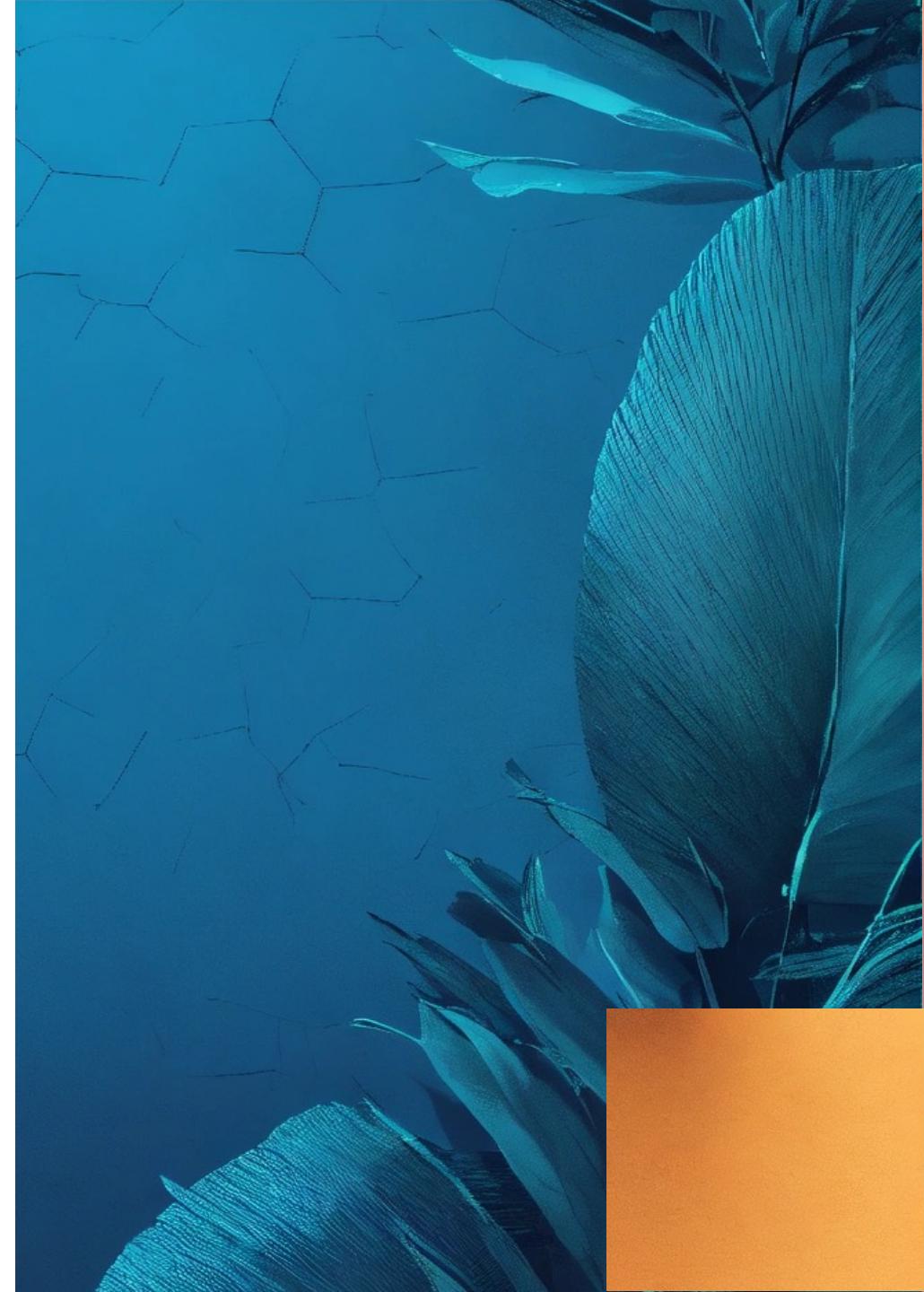
Generative concept-based models



- The **embeddings of textual concepts** are **aligned to the latent input representation** to produce **concept scores**
- The **concept scores** are used to **provide the final classification**

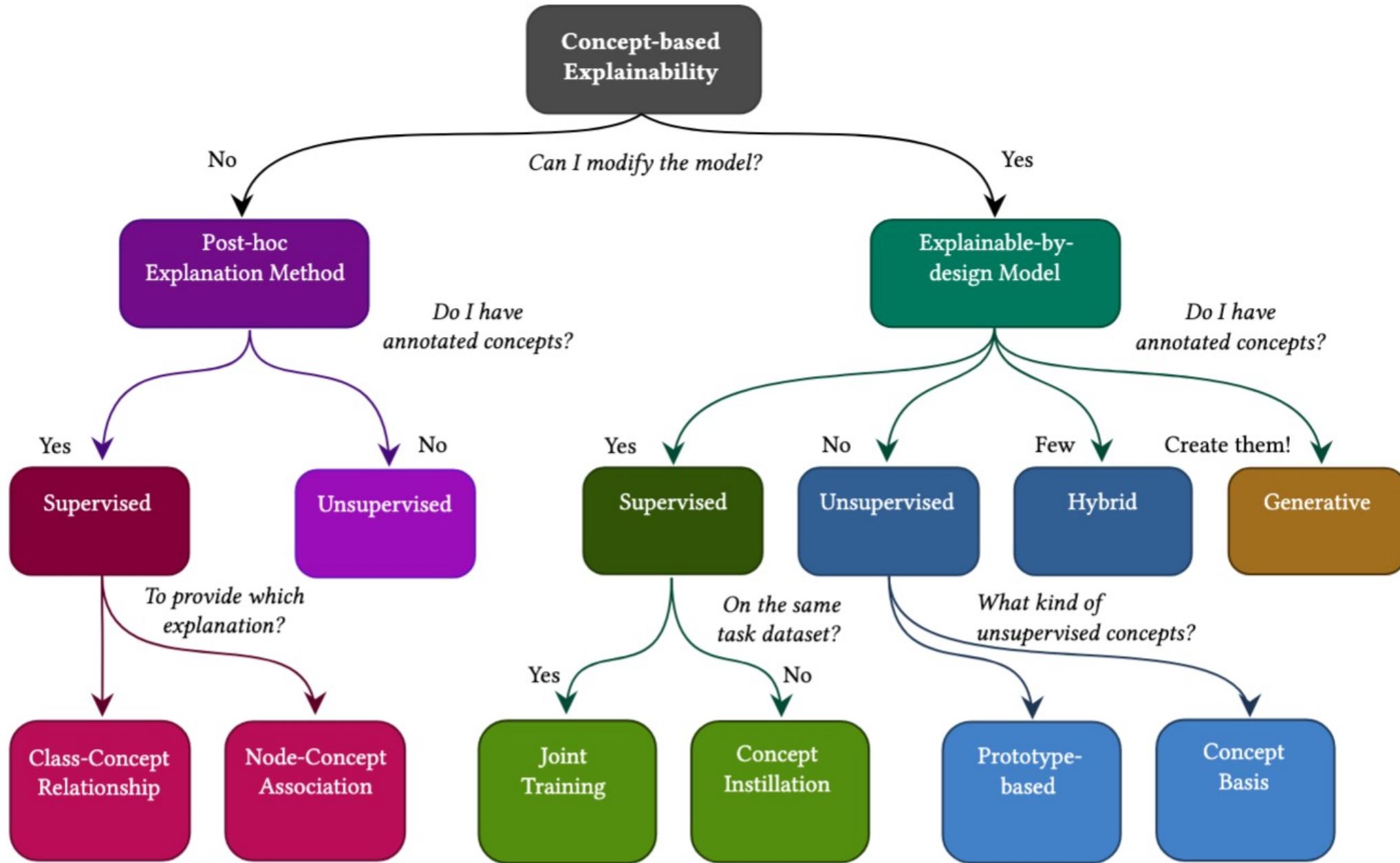
Tutorial structure

- Introduction to XAI techniques
- Taxonomy of C-XAI
- **Post-hoc vs By-design C-XAI approaches**
- Evaluation, Resources and applications
- Hands-on session



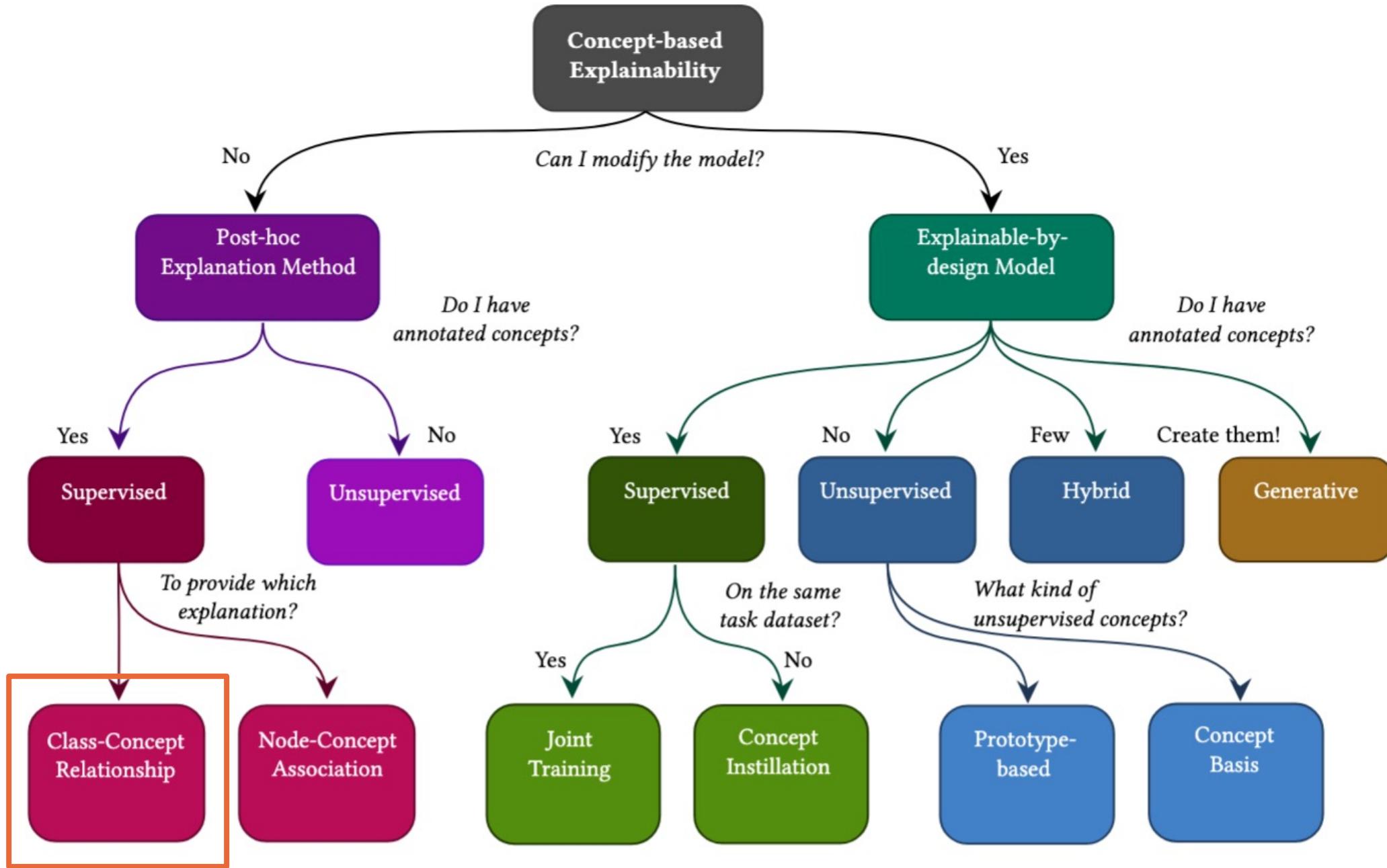
OUTLINE

1. Testing with Concept Activation Vector (T-CAV)
2. Concept Bottleneck Models (CBM)
3. Concept Embedding Models (CEM)
4. Self-Explainable Neural Network (SENN)
5. Label-free Concept Bottleneck Models (LF-CBM)

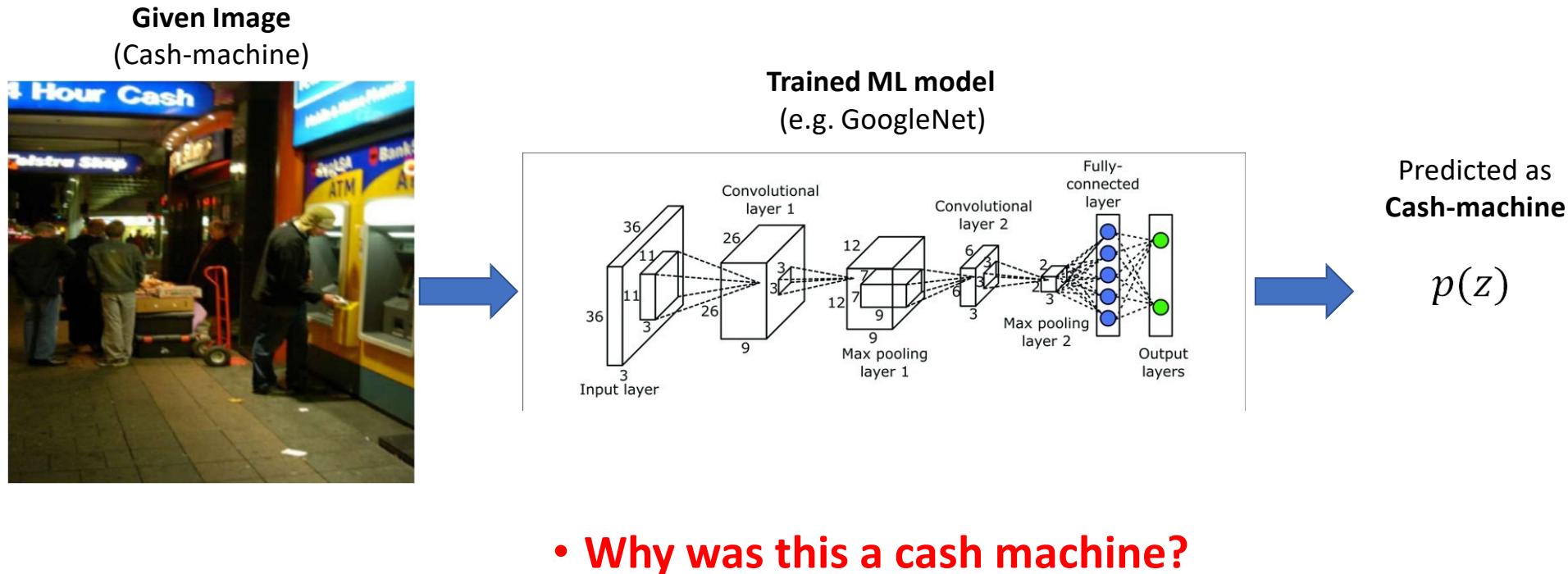


Credits for some of these slides goes to Kim et al.

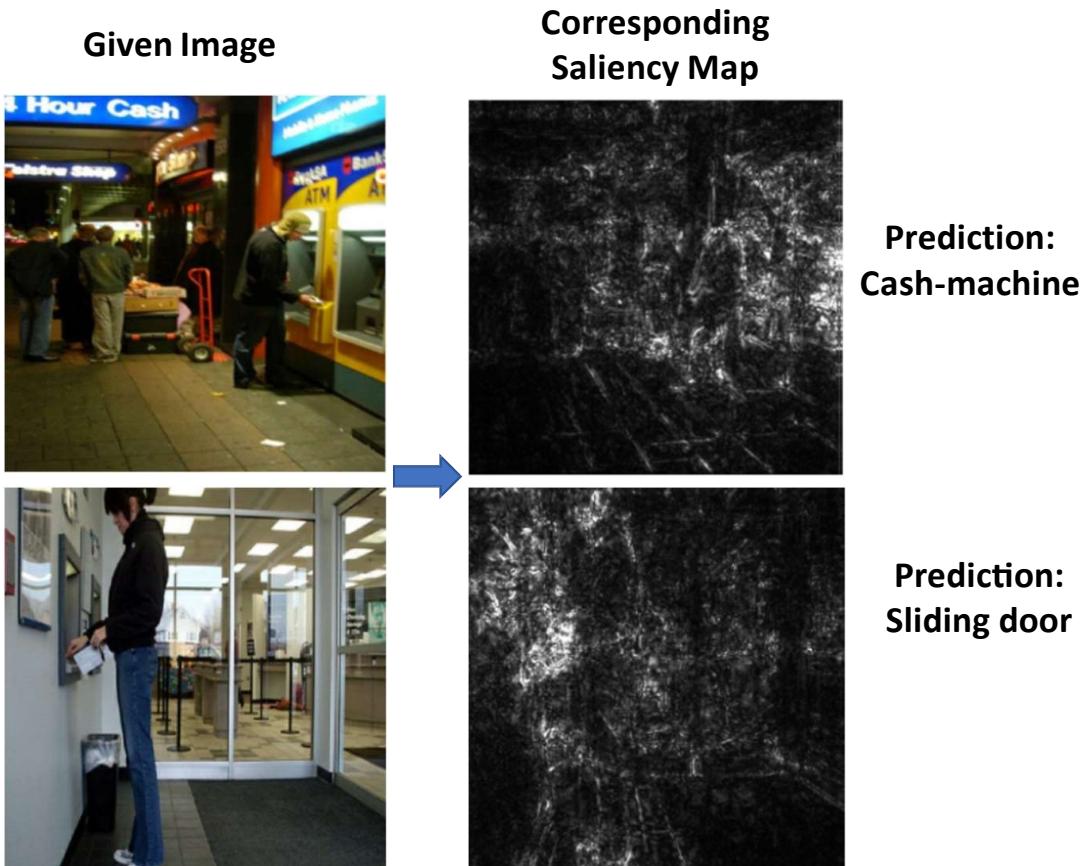
1. Testing with Concept Activation Vectors (T-CAV)



Example: Post-training explanation



Problem Objective

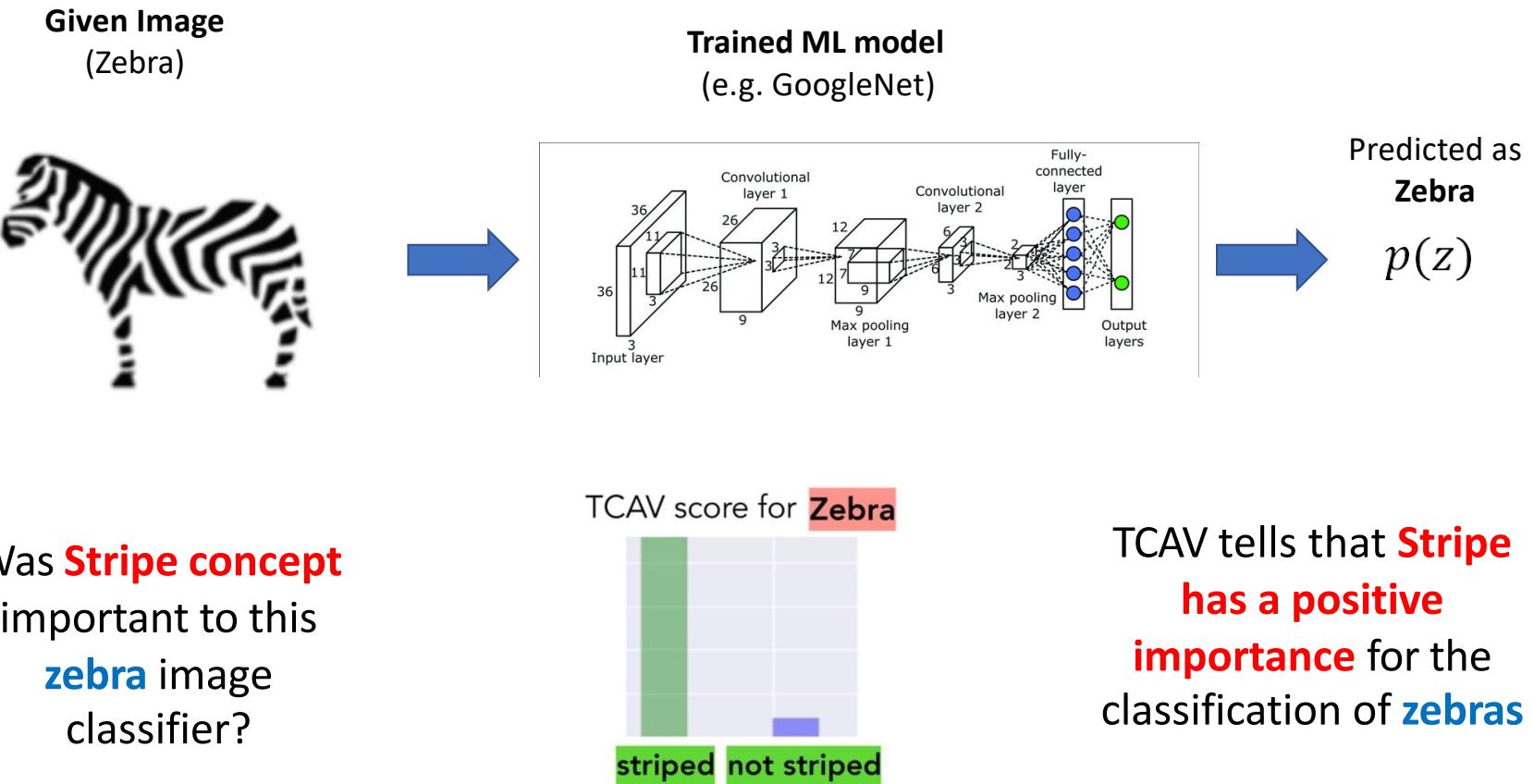


- Did the '**human**' concept matters?
- Did the '**paper**' concept matters?
- Did the '**ATM**' or '**Cash**' concept matters?

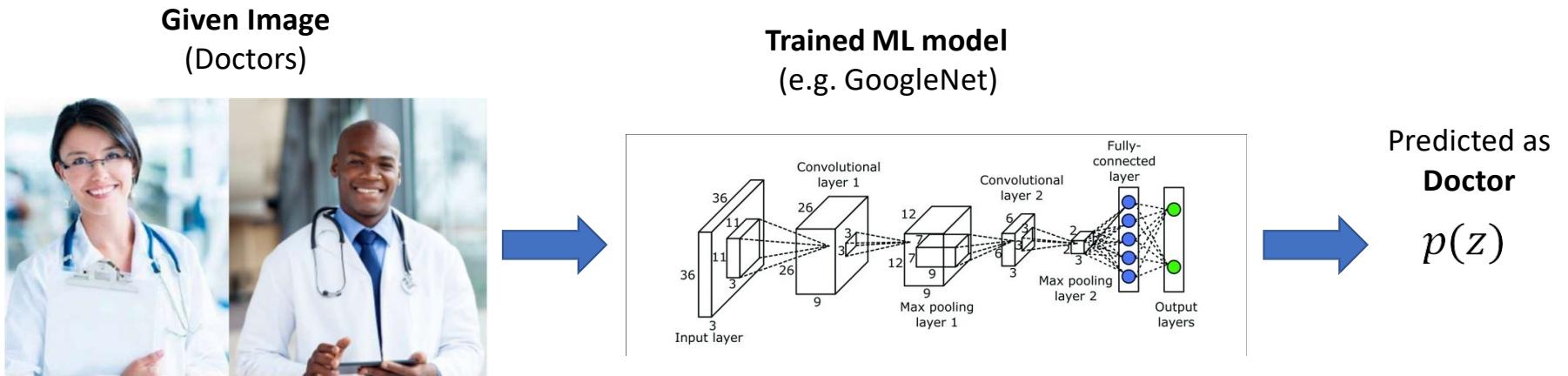
TCAV objective:

Quantitatively measure how important are "**user- chosen concepts**"

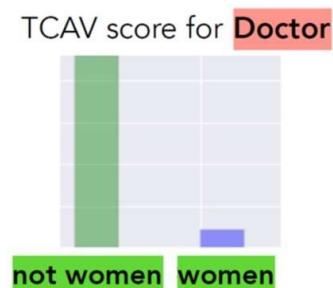
TCAV Example 1 (Zebra)



TCAV Example 2 (Doctor)



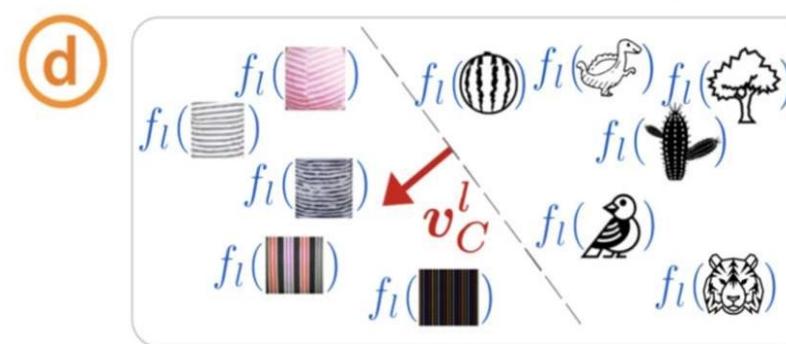
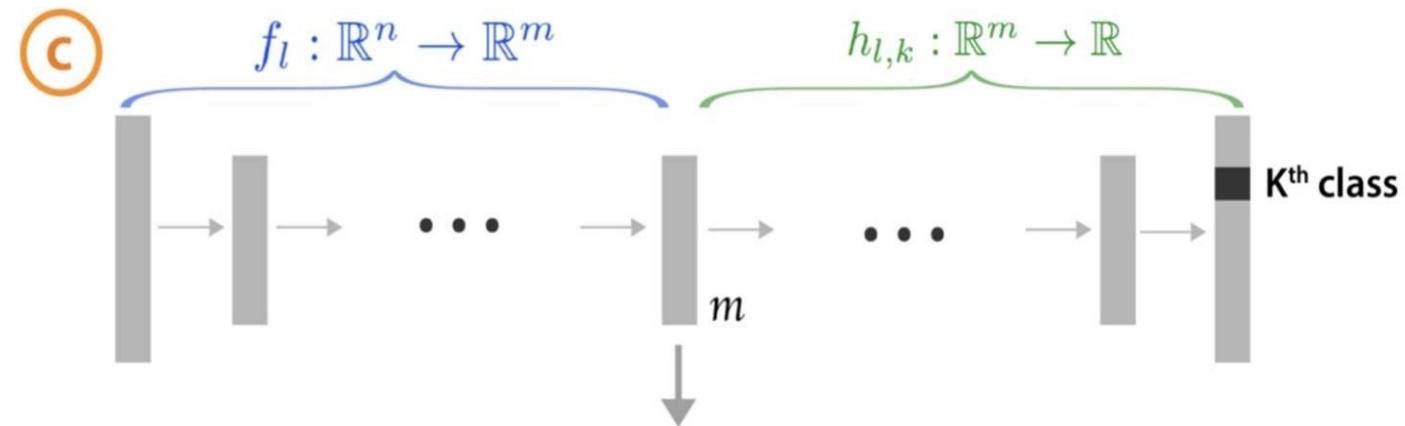
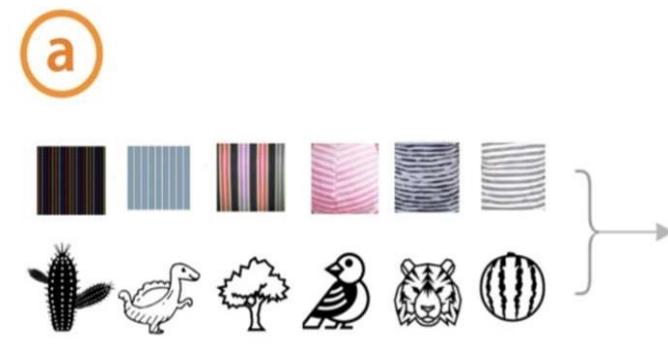
Was **Woman** concept important to this **doctor** image classifier?



TCAV tells that **Woman** has a negative importance for the classification of doctors

BIAS IDENTIFICATION!

TCAV: Overview



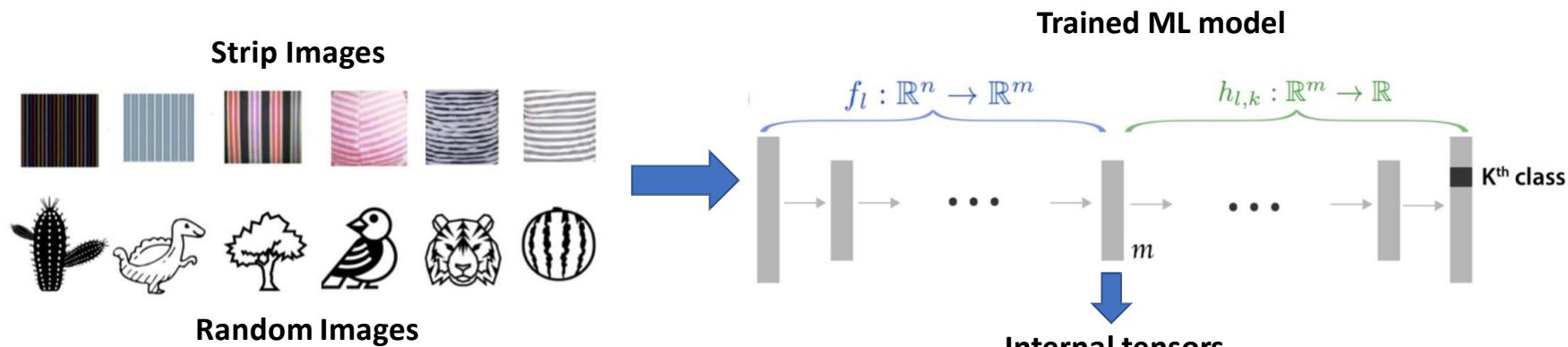
e

$$S_{C,k,l}(\text{zebra}) \\ = \nabla h_{l,k}(f_l(\text{zebra})) \cdot v_C^l$$

TCAV components

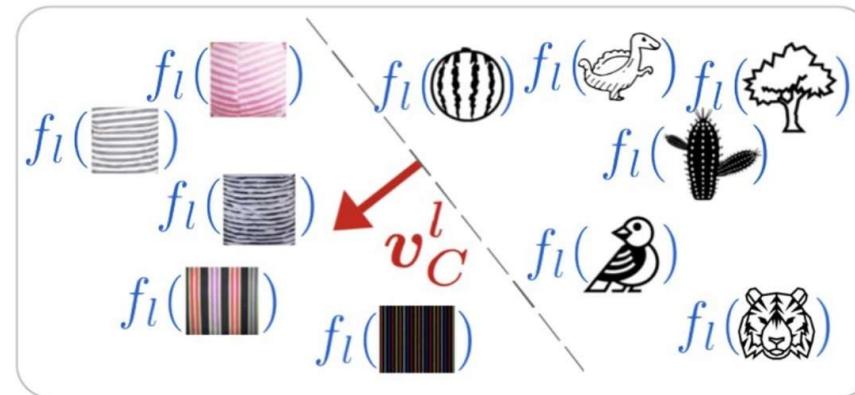
- a) A dataset annotated with both **examples of concepts** and **random images**
- b) The dataset with the **original classes**
- c) The **model** to explain
- d) The Concept Activation Vectors (CAV)
- e) The TCAV score showing the **influence** of a concept on a given class

TCAV: (1) How to define CAV?



Train a **linear classifier** to separate the projection of the concepts from the random images

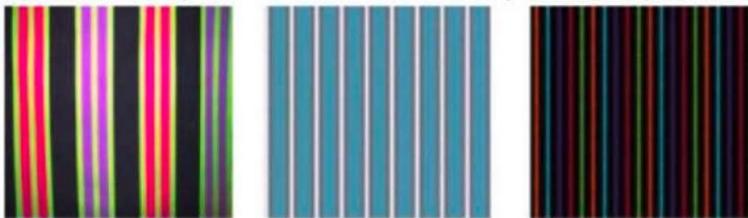
CAV (v_C^l) is the vector **orthogonal** to the decision boundary



Sorting Images with CAVs

- Given a set of images (e.g., belonging to the same class)
- Compute the cosine similarity between
 - the latent representation of an image $f_l(x)$
 - the CAV v_C^l of the selected concept

CEO concept: most similar striped images



Model Women concept: most similar necktie images



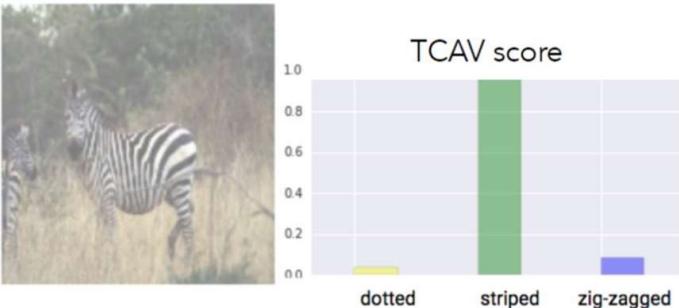
CEO concept: least similar striped images



Model Women concept: least similar necktie images



TCAV: (2) How to compute TCAV scores?



$$\begin{aligned} \text{zebra-ness} &\rightarrow \frac{\partial p(z)}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x}) \\ \text{striped CAV} &\rightarrow \frac{\partial}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x}) \end{aligned}$$

$$\begin{aligned} S_{C,k,l}(\mathbf{x}) &= \lim_{\epsilon \rightarrow 0} \frac{h_{l,k}(f_l(\mathbf{x}) + \epsilon \mathbf{v}_C^l) - h_{l,k}(f_l(\mathbf{x}))}{\epsilon} \\ &= \nabla h_{l,k}(f_l(\mathbf{x})) \cdot \mathbf{v}_C^l, \end{aligned} \tag{1}$$

Directional derivative with CAV:

- $S_{C,k,l}(\mathbf{x}) > 0$: positive influence
- $S_{C,k,l}(\mathbf{x}) < 0$: negative influence

$$S_{C,k,l}(\text{zebra})$$

$$S_{C,k,l}(\text{zebra})$$

$$S_{C,k,l}(\text{zebra})$$

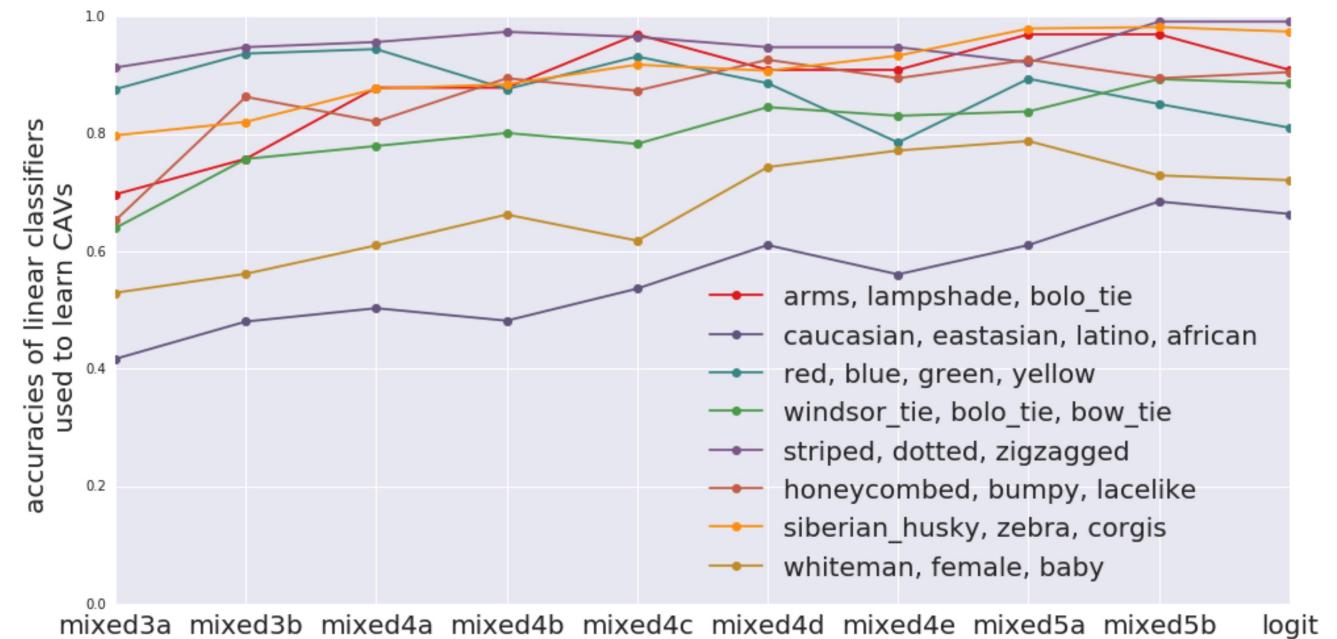
$$S_{C,k,l}(\text{zebra})$$

$$\text{TCAV}_{Q_{C,k,l}} = \frac{|\{\mathbf{x} \in X_k : S_{C,k,l}(\mathbf{x}) > 0\}|}{|X_k|}$$

The **TCAV score** is the number of class samples having a positive directional derivative w.r.t. the **CAV**

When and where can concept be learnt?

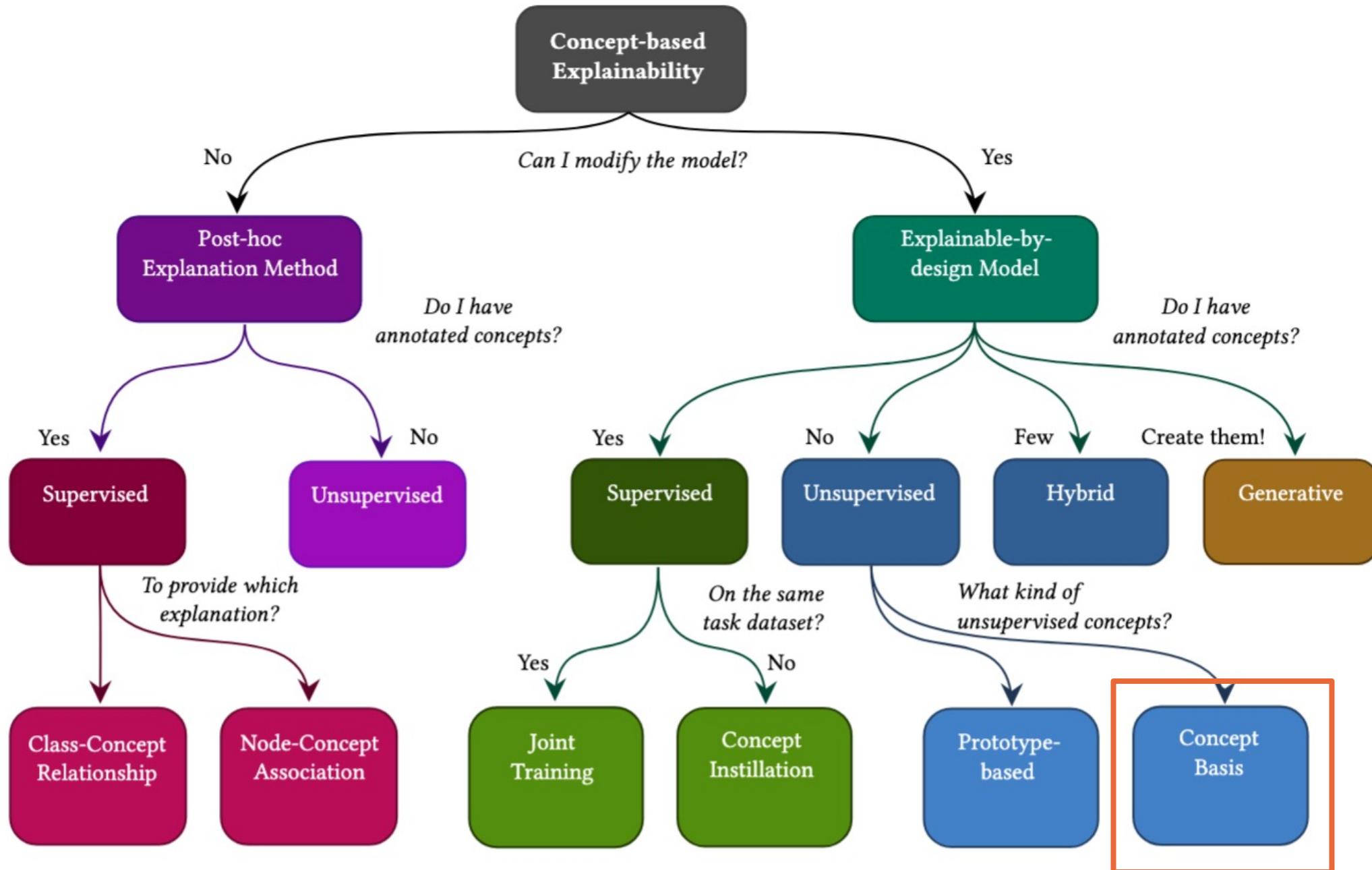
- Accuracy of the «linear probe»
 - *High* implies the network **has automatically learnt** a concept
 - *Low* implies the network **does not use** that concept for predicting the final class



- Simpler concepts have high accuracy throughout the NN
- High-level concepts can be detected better at higher layers

Credits for some of these slides goes to Alvarez Melis & Jaakkola

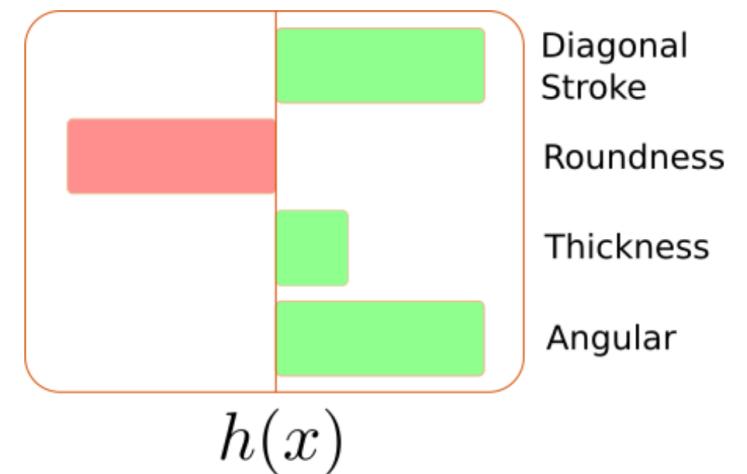
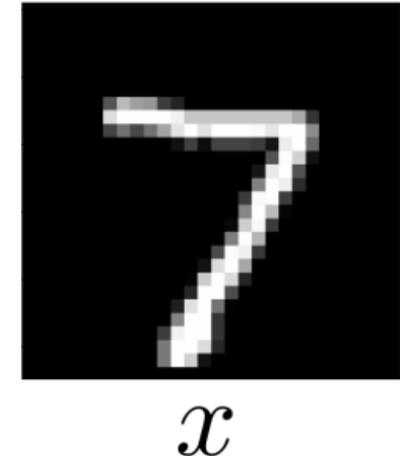
2. Self Explainable Neural Networks (SENNs)



Let's Learn Concepts by Design

- Instead of explaining a network in terms of concepts
- Let's use a linear model working on top of interpretable features:

- Linear model: $f(x) = \sum_i^k \theta_i x_i$
- Linear model working on top of Basis Concepts: $f(x) = \sum_i^k \theta_i h_i(x)$

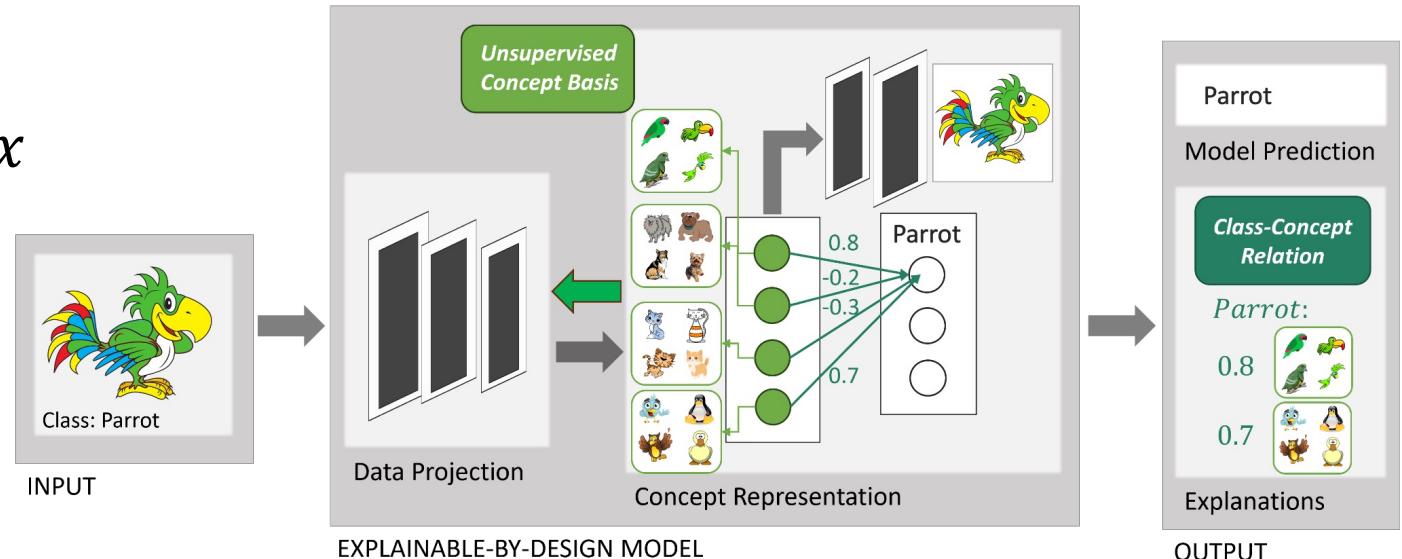


How are Concepts learnt?

- **Fidelity:** the representation of x should preserve relevant information of the input

- Autoencoder:

$$x' = h_{dec}(h(x))$$



- **Sparsity:** the classes should be represented by few concepts:

- Sparse Auto Encoder:

$$L_{SAE} = \|x - x'\|^2 + \lambda KL(h(x)||\rho)$$

Beyond Linear Models

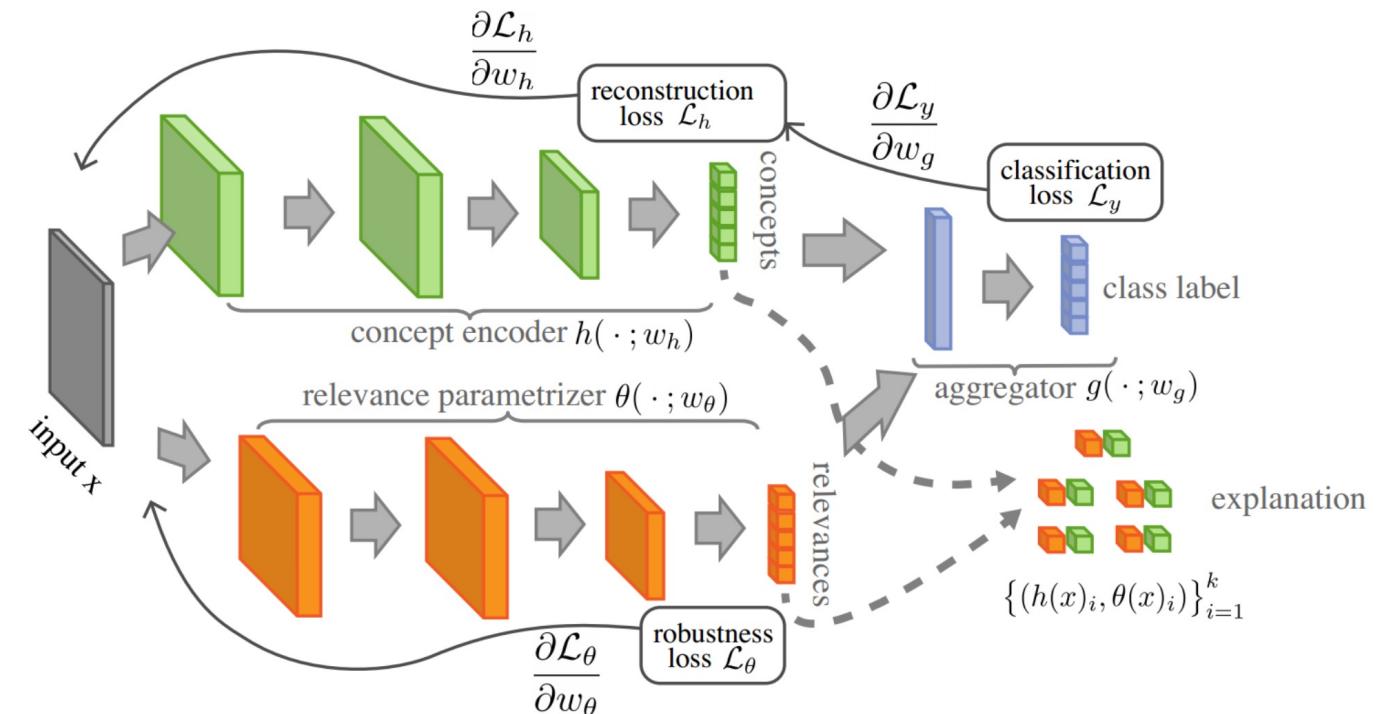
- Relevance scores θ_i are fixed globally for all samples.
Let's predict them!

- Linear model working on top of Basis Concepts:

$$f(x) = \sum_i^k \theta_i h_i(x)$$

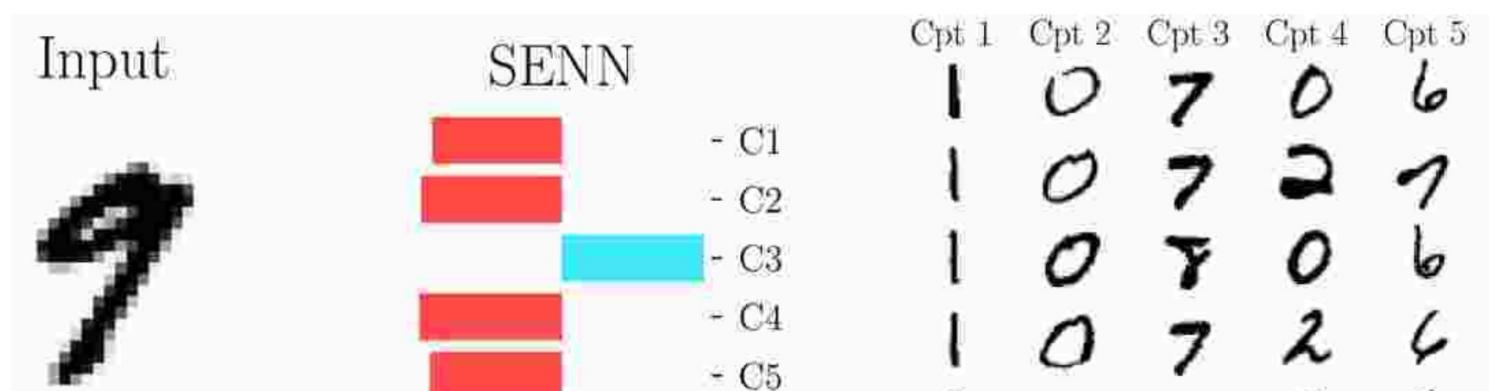
- Linear model Predicting local weights:

$$f(x) = \sum_i^k \theta_i(x) h_i(x)$$



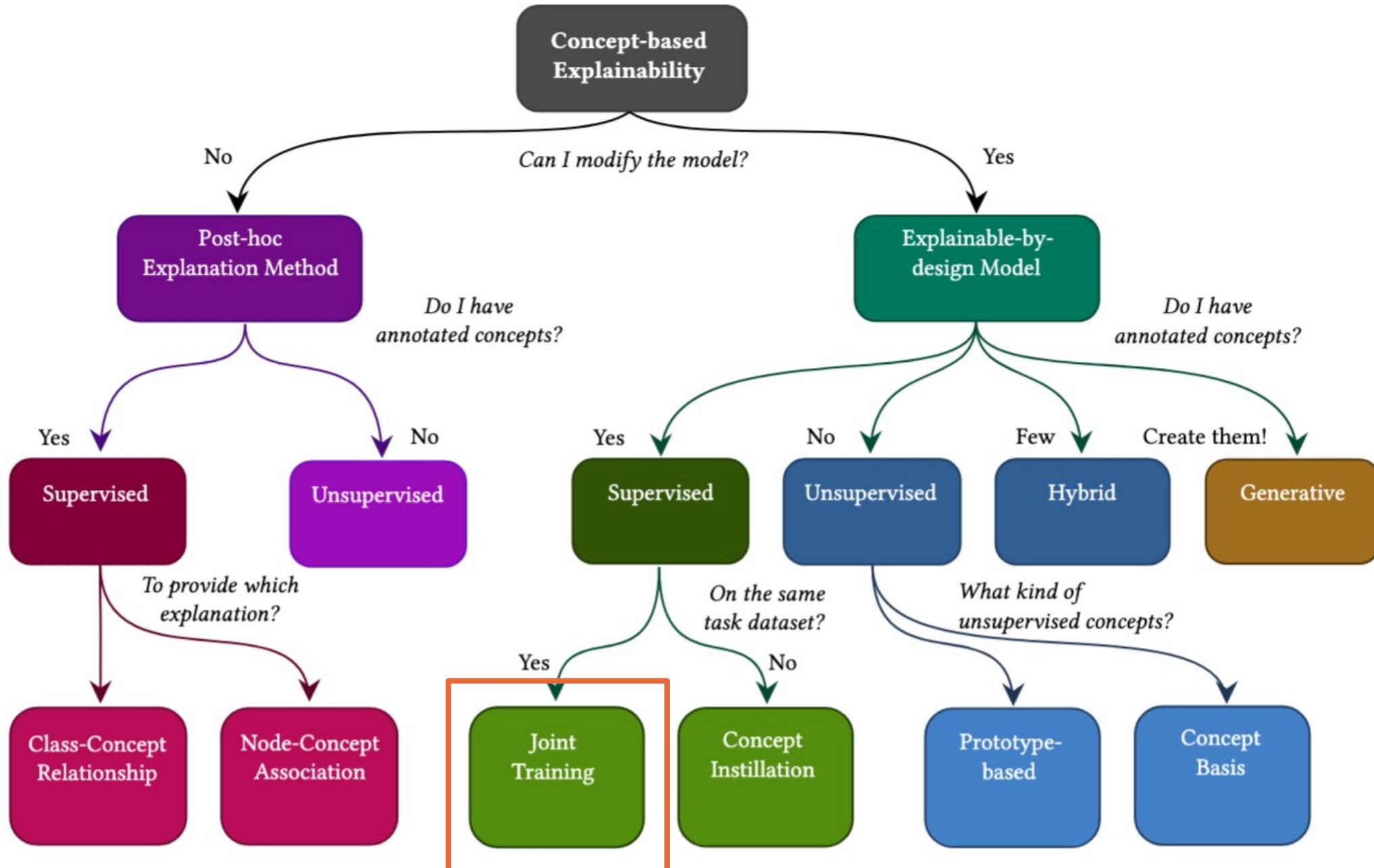
Interpreting Unsupervised Basis Concept

- $h(x)$ is faithful and sparse
- But it still requires to be interpreted (similarly to mechint)
- Let's use its «concept dictionary»

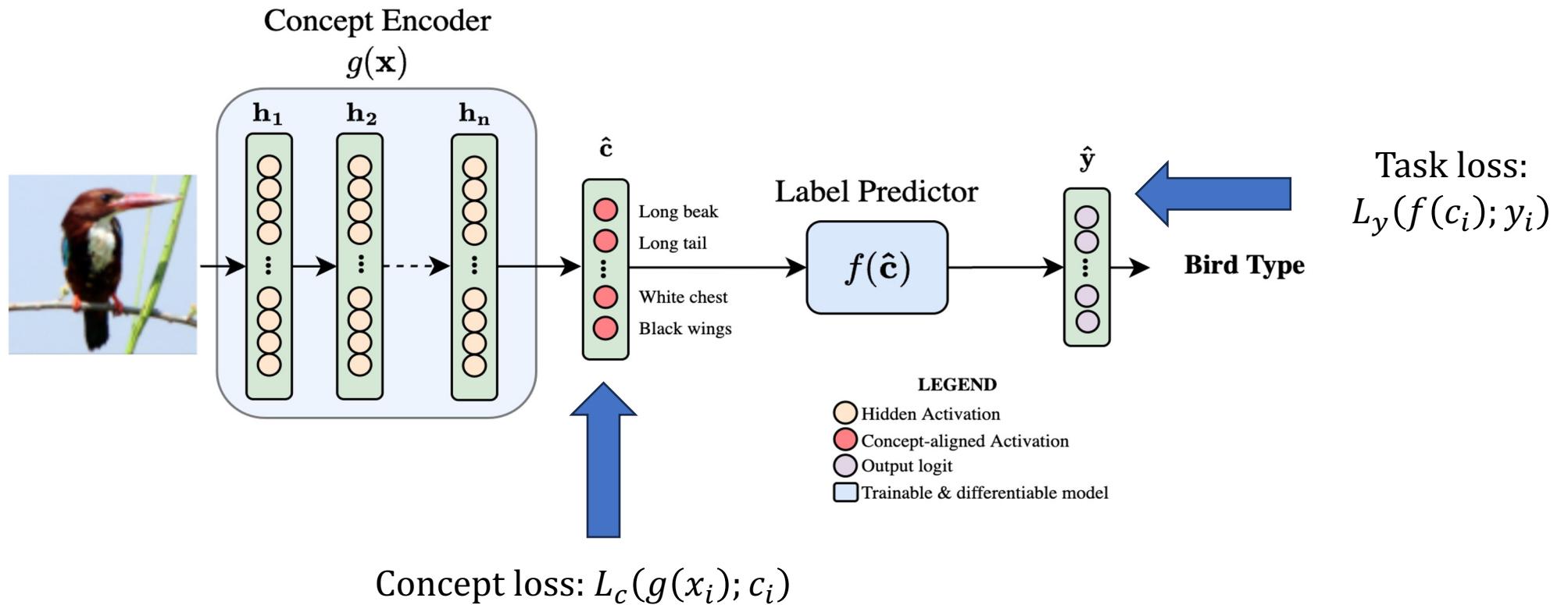


Credits for some of these slides goes to Koh et al.

3. Concept Bottleneck Models (CBMs)



Concept bottleneck models architecture

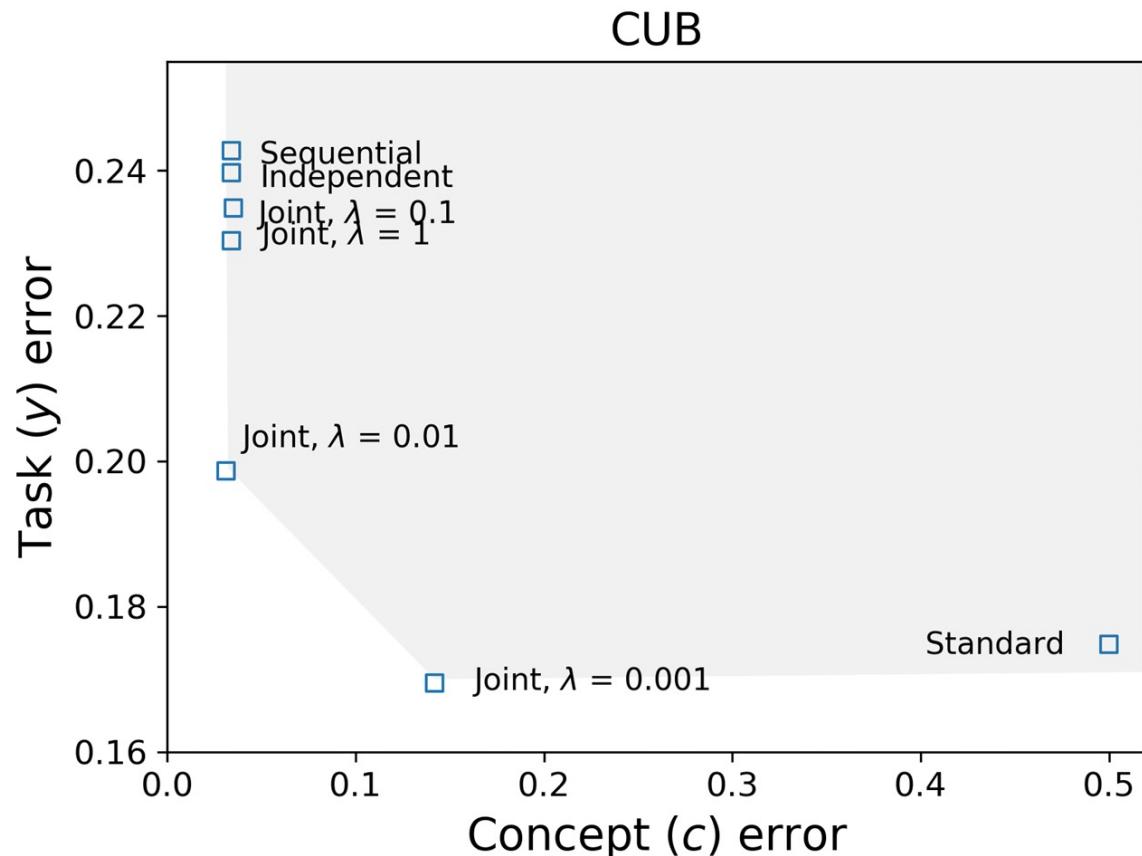


Different training strategy

- Indipendent: $\hat{f} = \arg \min_f \sum_i L_y(f(c_i), y_i)$ $\hat{g} = \arg \min_g \sum_i L_c(g(x_i), c_i)$

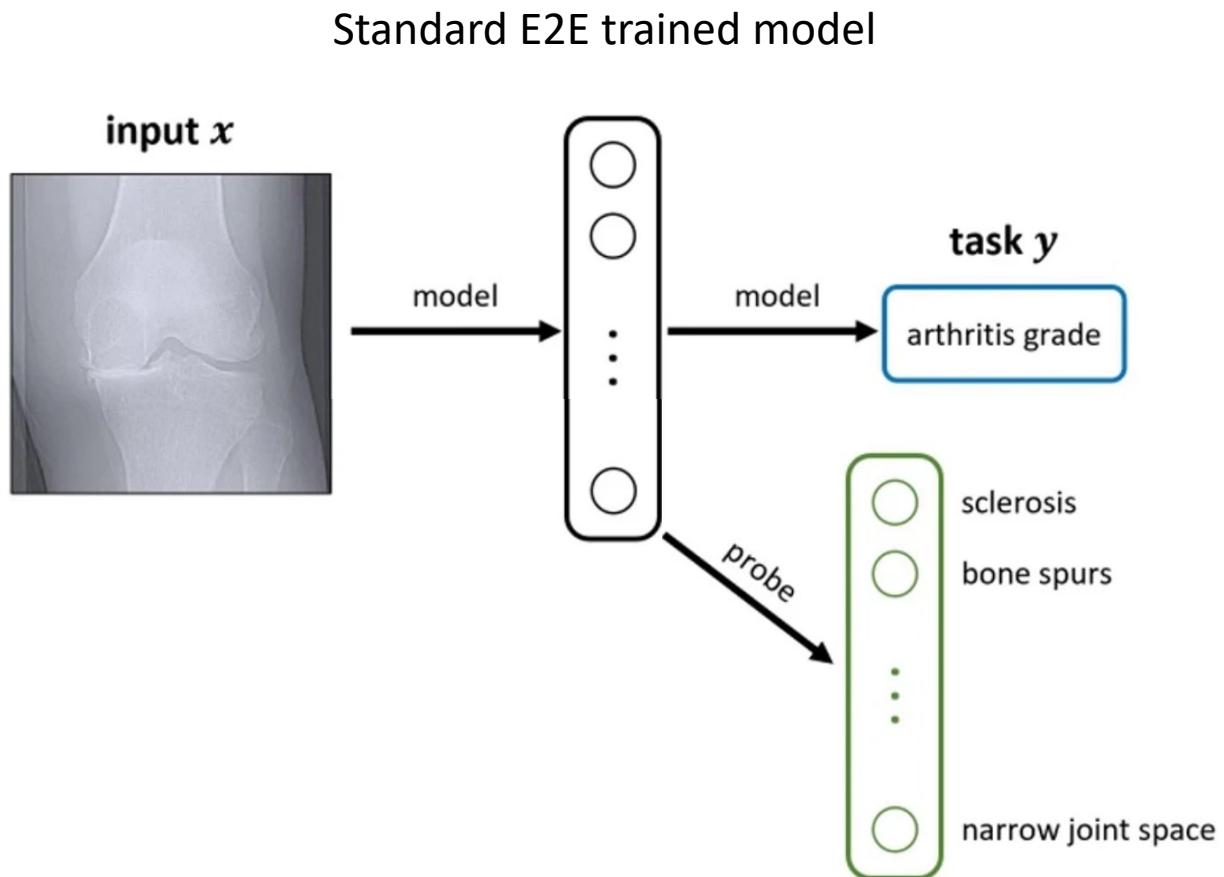
f is trained using the truth concepts
- Sequential: $\hat{f} = \arg \min_f \sum_i L_y(f(g(x_i)), y_i)$ g is trained first as above, then freezed
- Joint: $\hat{f}, \hat{g} = \arg \min_f \sum_i L_y(f(c_i), y_i) + \lambda \arg \min_g \sum_i L_c(g(x_i), c_i)$ f, g trained together for some $\lambda > 0$
- Standard: $\hat{f}, \hat{g} = \arg \min_f \sum_i L_y(f(c_i), y_i)$ It ignores the concepts loss

Different interpretability/performance trade-offs



- **Sequential** and **independent** are the more «trustworthy» because they ensure no concept leakage
- **Joint** strategy provides better task accuracy
 - Different trade-offs according to the λ value
- **Standard** model still has higher accuracy on average

Explicitly concept training ensure model learns the concepts

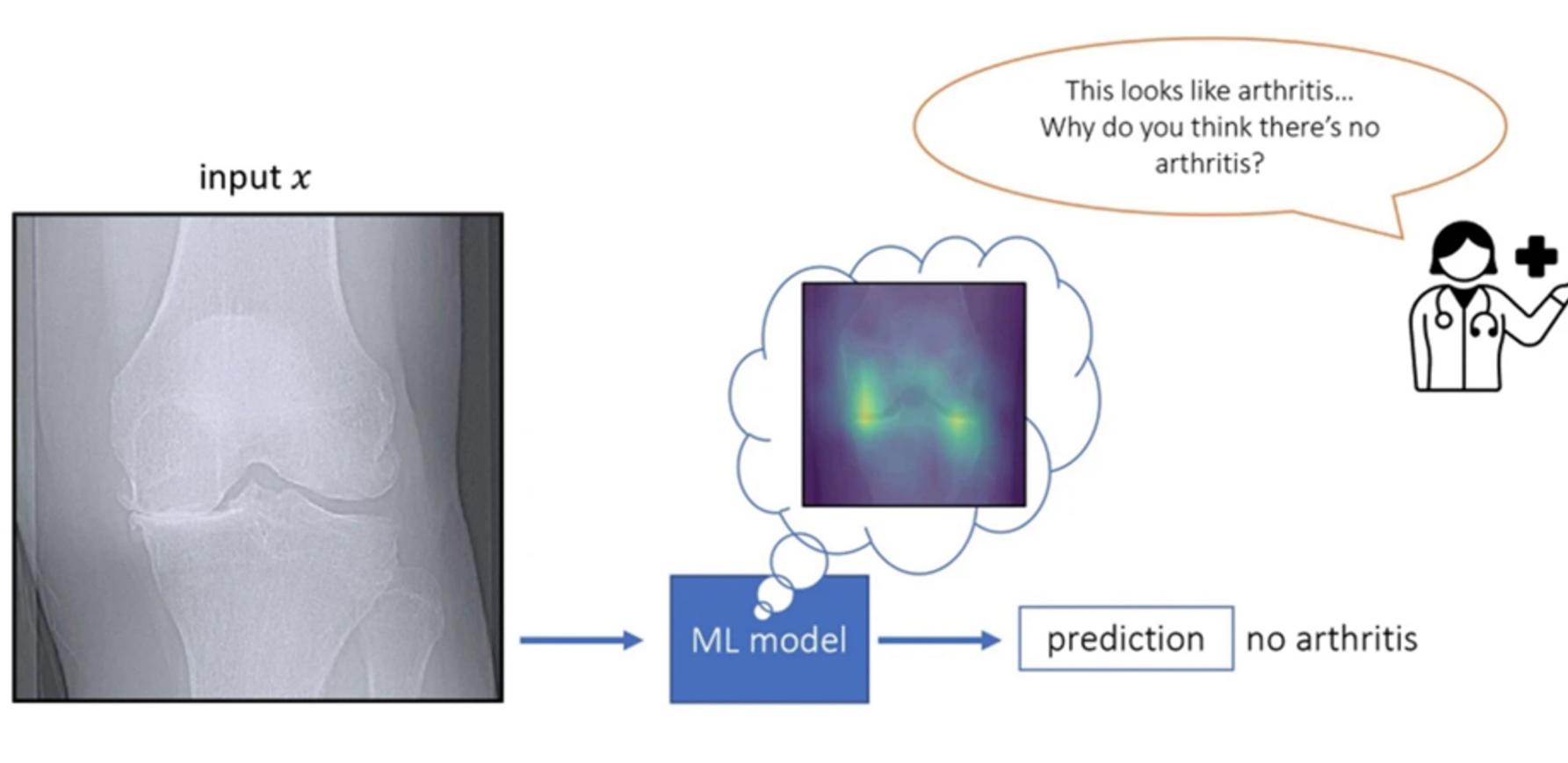


Method	X-Ray Concept Error (\downarrow)
Independent	0.53
Sequential	0.53
Joint	0.54
TCAV [Probe]	0.68

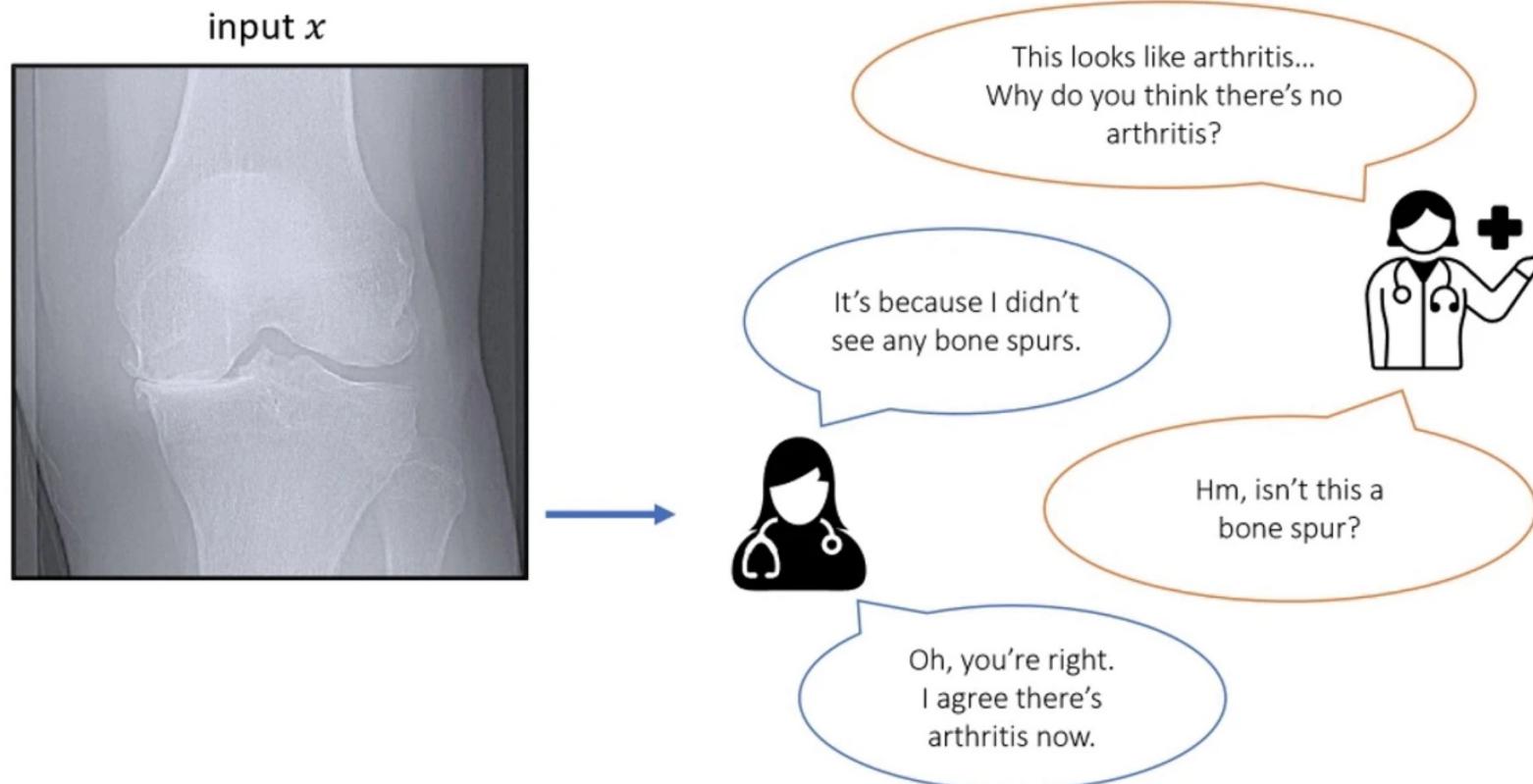
In a trained model, identifying some concepts may not be possible, because it might not have learnt them automatically

→ Only by explicitly training a model we can ensure it represents all concepts!

End-2-End models are difficult to interact with



Ideal: Interact through high-level concepts



Importance of Concept Intervention



*bone
spurs*

→ { 1.48 ... **0.15** ... 0.09

I don't think
there are
bone spurs

wrong
no arthritis

Intervention

{ 1.48 ... 1.0 ... 0.09 }

**correct
moderate
arthritis**

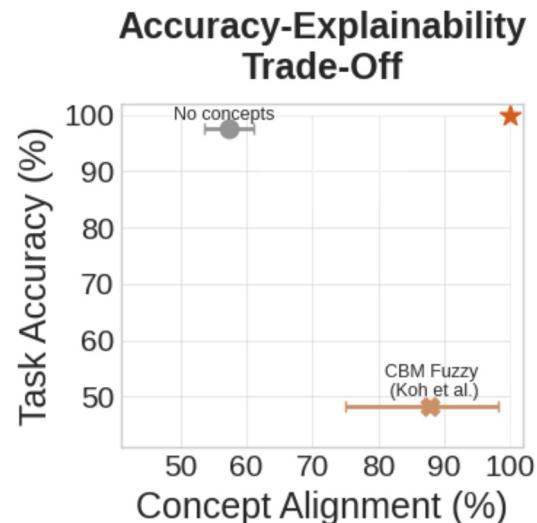
Actually,
there is a
bone spur in
this x-ray



CBM Drawbacks

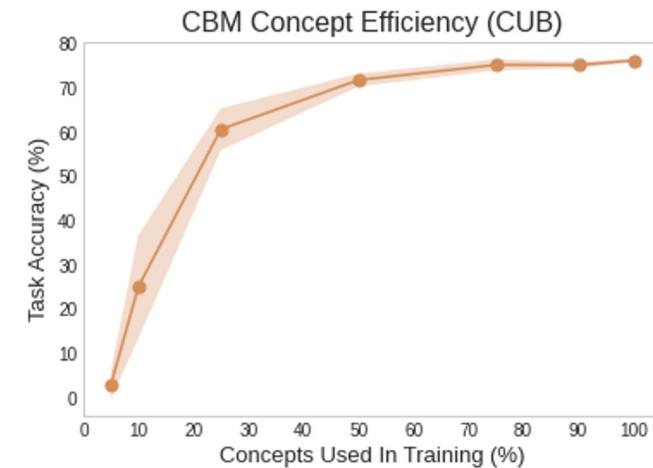
Poor Trade-offs

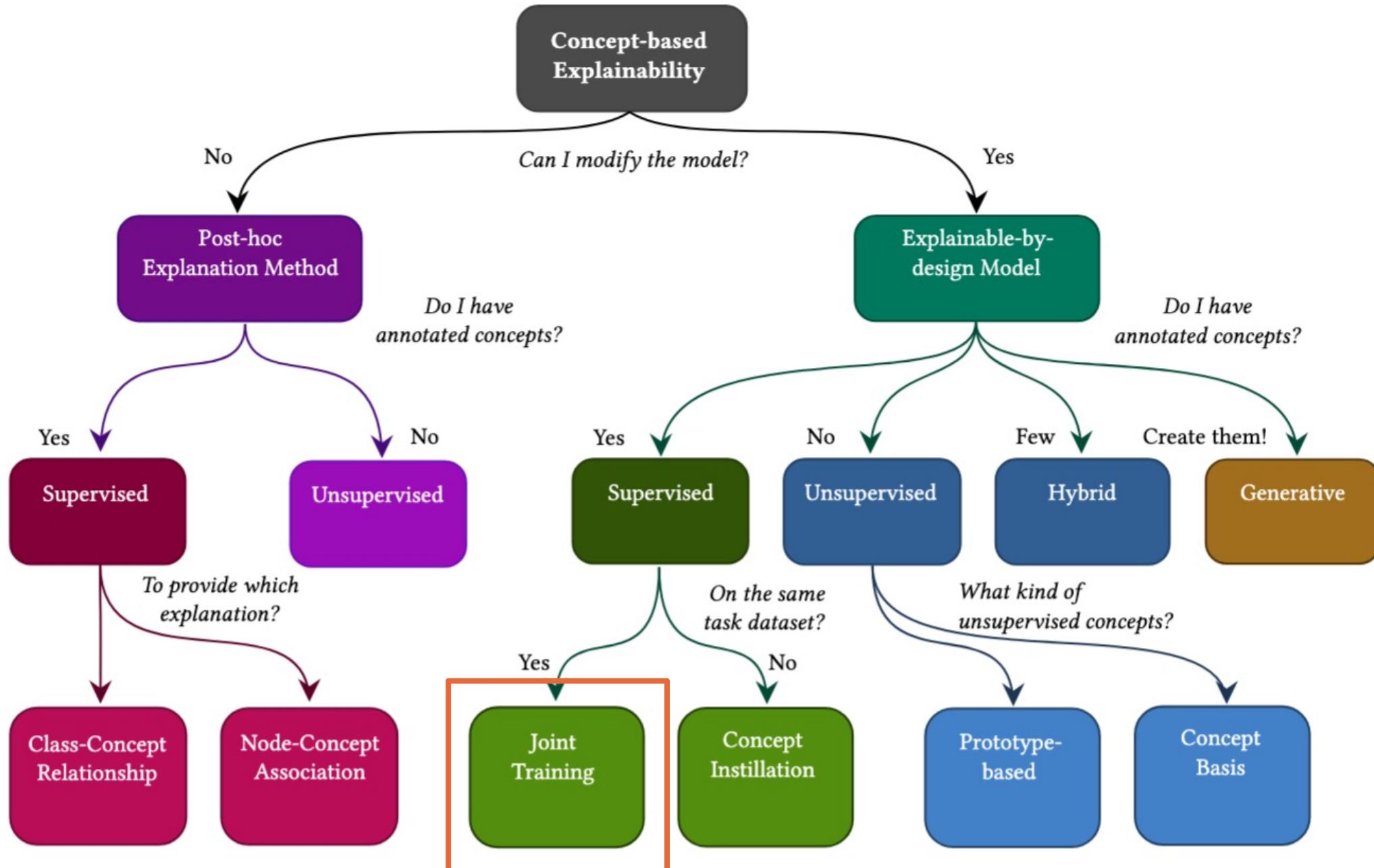
Struggle to compromise between accuracy and explainability



Low Concept Efficiency

CBMs do not scale in real-world conditions

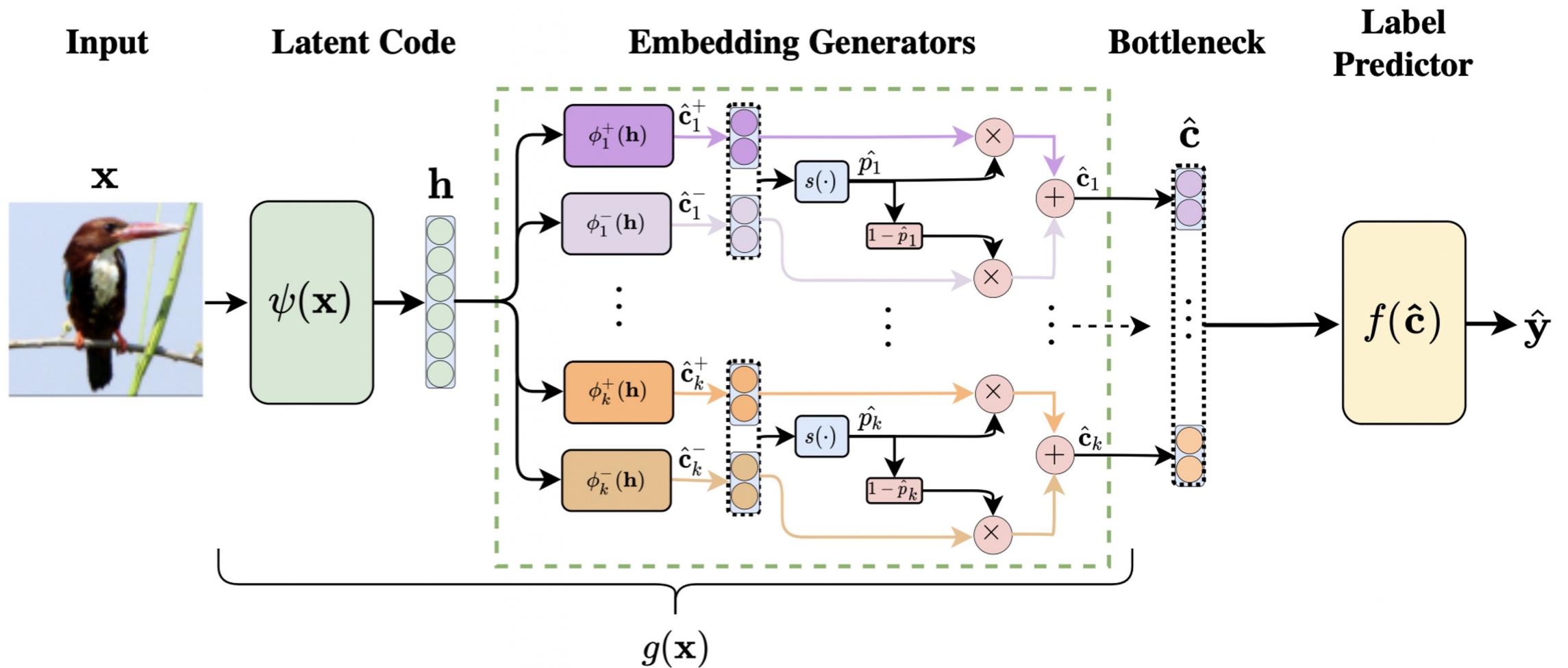




Credits for some of these slides goes to Barbiero et al.

4. Concept Embedding Models (CEM)

Concept Embedding Models: overview



Concept Embedding workflow

1. $h = \psi(x)$: the latent space of the model
2. $\mathbf{c}_i^+ = \phi_i^+(x)$: neural model dedicated to represent the i-th **positive** concept embedding
3. $p_i = s([\mathbf{c}_i^+, \mathbf{c}_i^-])$: the *concept score* (i.e., probability of presence of the ith concept) is a function shared among concepts working on the concatenations of the concept embeddings
4. $\hat{\mathbf{c}}_1 = p_1 \mathbf{c}_1^+ + (1 - p_1) \mathbf{c}_1^-$: the *concept embedding* is represented by the weighted combination of the positive and negative concept embeddings according to its presence
5. $f([\hat{\mathbf{c}}_1, \dots, \hat{\mathbf{c}}_i, \dots, \hat{\mathbf{c}}_k])$: the task predictor works on the concatenation of all the concept embeddings

CEM: A neural-symbolic approach

Neural

—
Concepts are represented
with: unsupervised
embeddings

$$c_i \in \mathbb{R}^k$$

Symbolic (CBM)

—
Concepts are represented
with: **supervised** scalars

$$c_i \in [0,1]$$

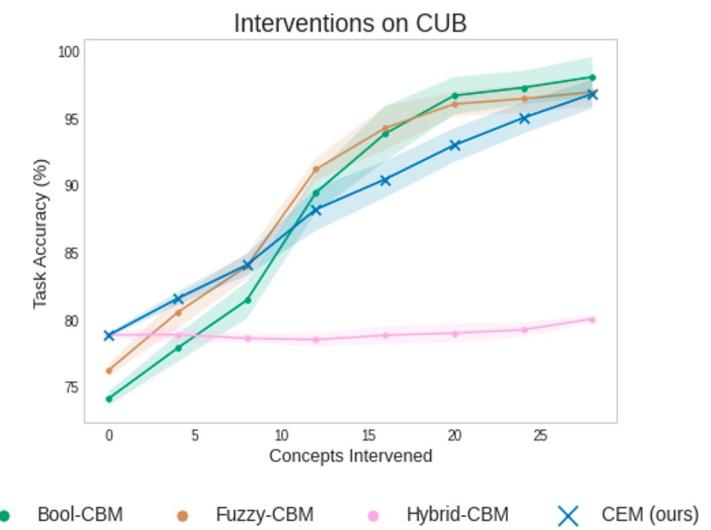
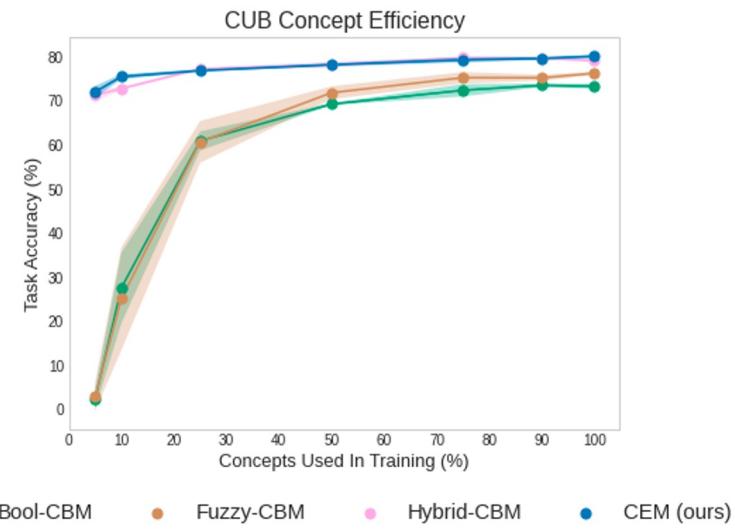
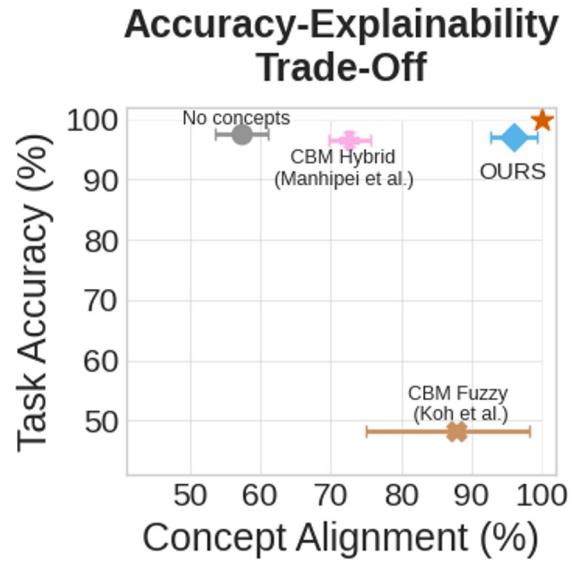
Neural Symbolic (CEM)

—
Concepts are represented
with: pairs of **supervised**
embeddings

$$c_i \in \mathbb{R}^k$$

$$c_i = \text{agg}(c_i^+, c_i^-)$$

CEM Advantages



Beyond Trade-offs

CEMs overcome the current accuracy-explainability trade-off

High Concept Efficiency

CEMs scale to real-world conditions where concept supervisions are scarce

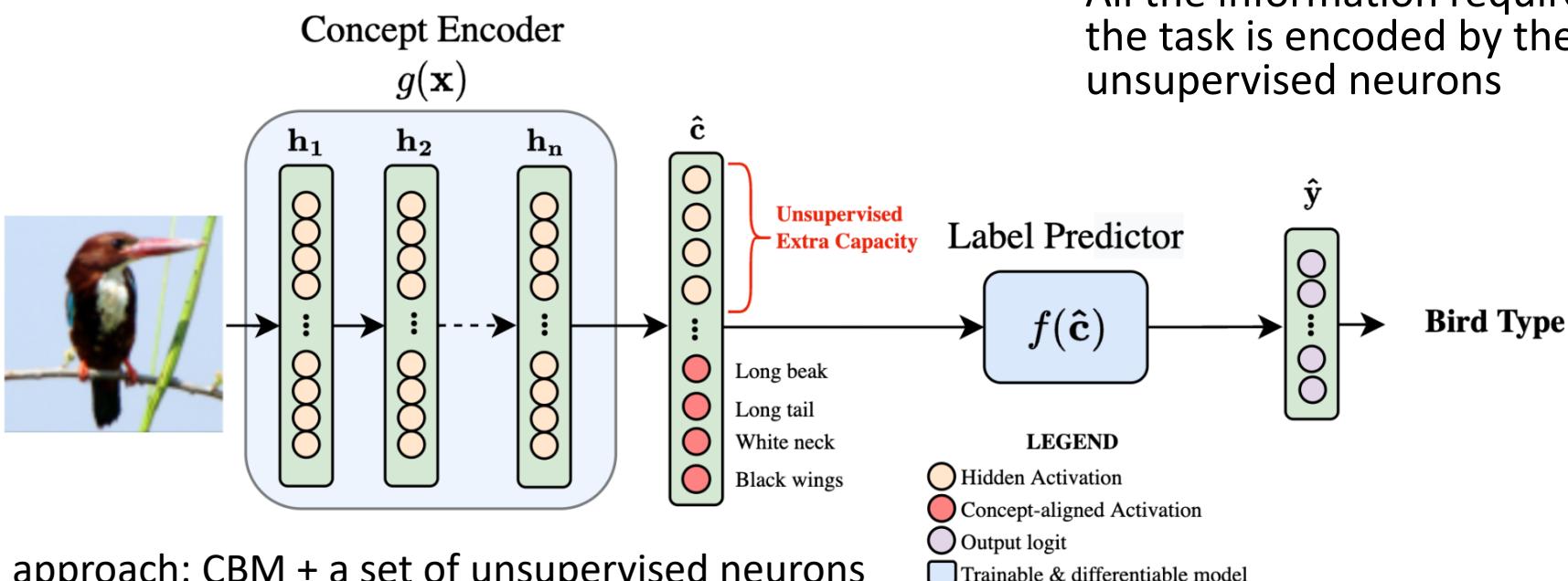
Effective interventions

CEMs are responsive to concept interventions

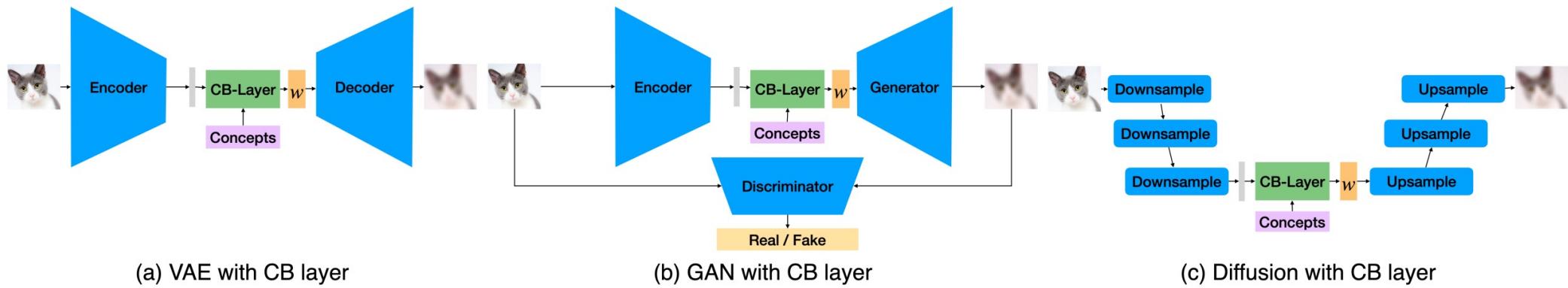
CEM vs Hybrid approach

- PROS:
 - Retain high accuracy
 - Has high concept efficiency like CEM

- CONS:
 - Prevent any effect of concept intervention
 - Changing the predicted scores has no effect on the task prediction
 - All the information required to predict the task is encoded by the unsupervised neurons



Generative models based on CEM

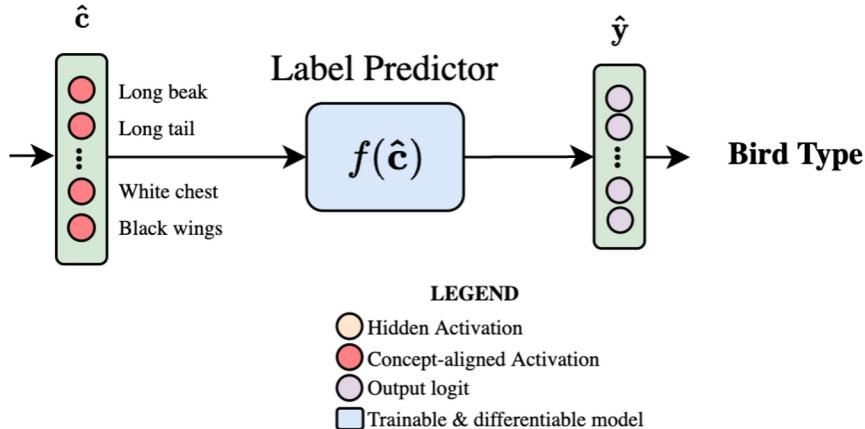


Ismail, Aya Abdelsalam, et al. "Concept bottleneck generative models." ICLR 2024.

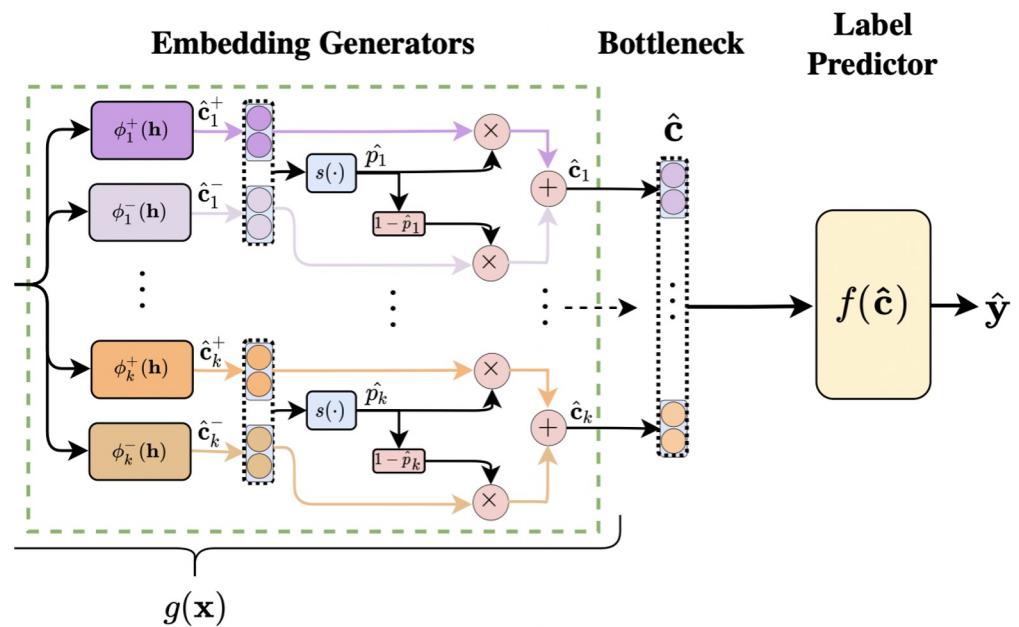
Have we lost something?

Interpretability

CBM: Interpretable



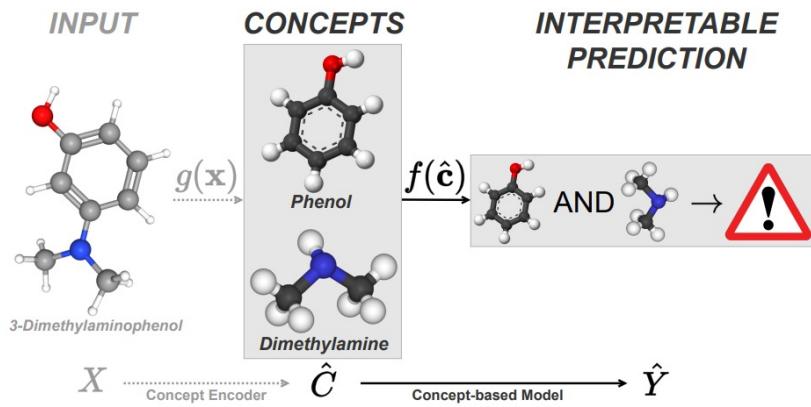
CEM: NON-Interpretable



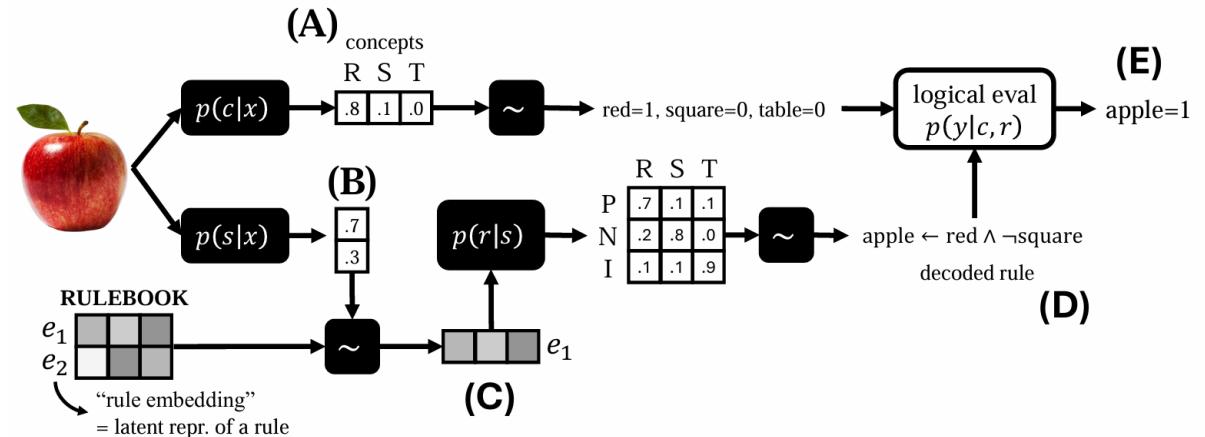
$$\hat{\mathbf{c}}_{\text{yellow}} = [2.3, 0.3, -3.5, \dots]^T$$

Task-Interpretable CEMs

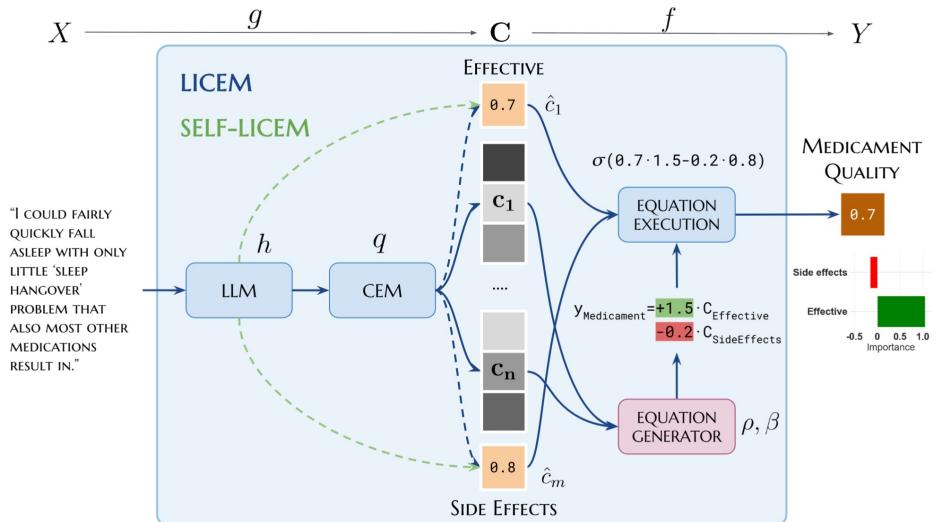
DCR¹: Local, Rule-based Predictions



CMR²: Global, Rule-based Predictions



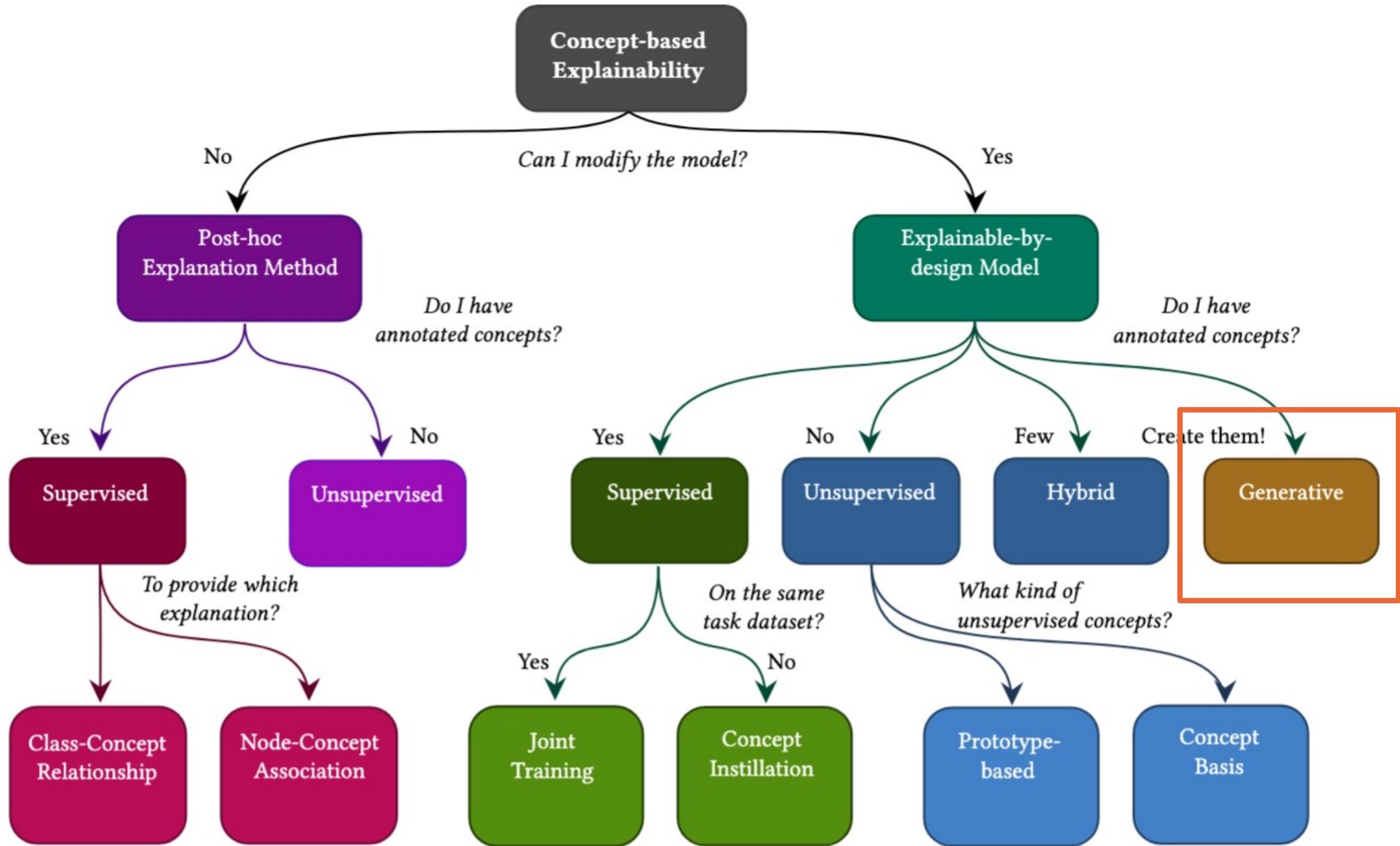
LICEM³: Local, Equation-based Predictions



1. Barbiero, Pietro, et al. "Interpretable Neural-Symbolic Concept Reasoning." International Conference on Machine Learning (2023).
2. Debot, David, et al. "Interpretable Concept-Based Memory Reasoning." The 38th Annual Conference on Neural Information Processing Systems (2024).
3. De Santis, Francesco, et al. "Linearly Interpretable Concept Embedding Model for Text Classification." Joint European Conference on Machine Learning and Knowledge Discovery in Databases (2025).

Credits for some of these slides goes to Oikarinen et al.

5. Label-free CBM (LF-CBMs)

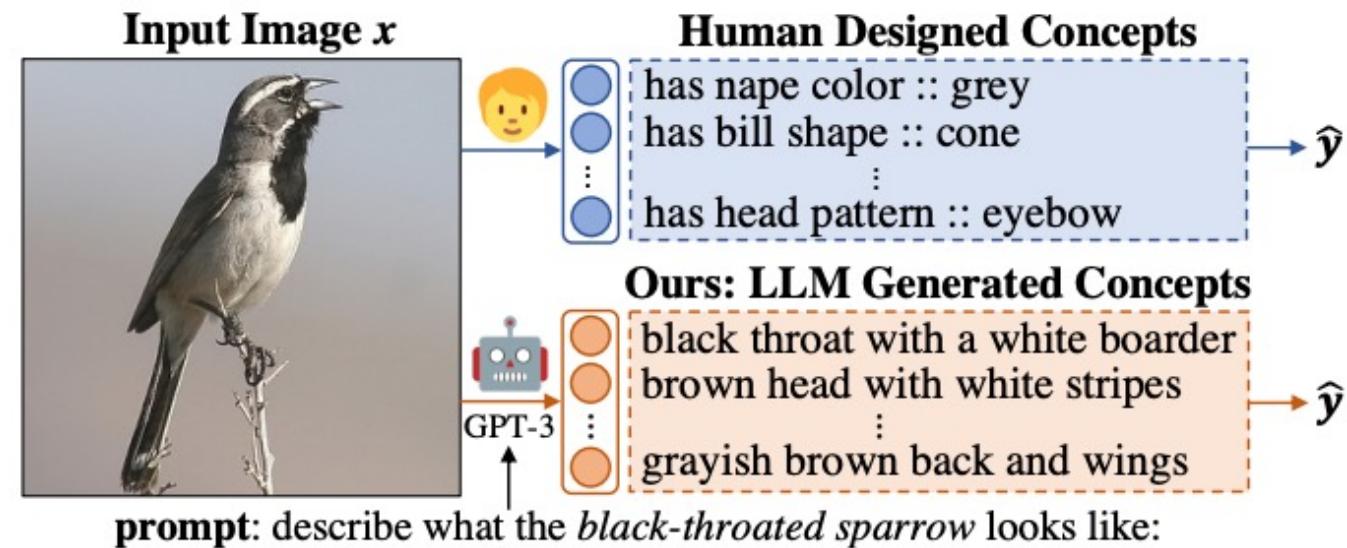


CBM Issues

- Strong extra annotation effort required from humans: Up to $|C| * |N|$
 - $|C|$ # concepts
 - $|N|$ # samples
- Low performance when using few concepts without embeddings
 - Unable to distinguish classes characterized by the same concepts
 - E.g. Plane vs Helicopter cannot be distinguished with only the concept fly

Let's use a pretrained VLM!

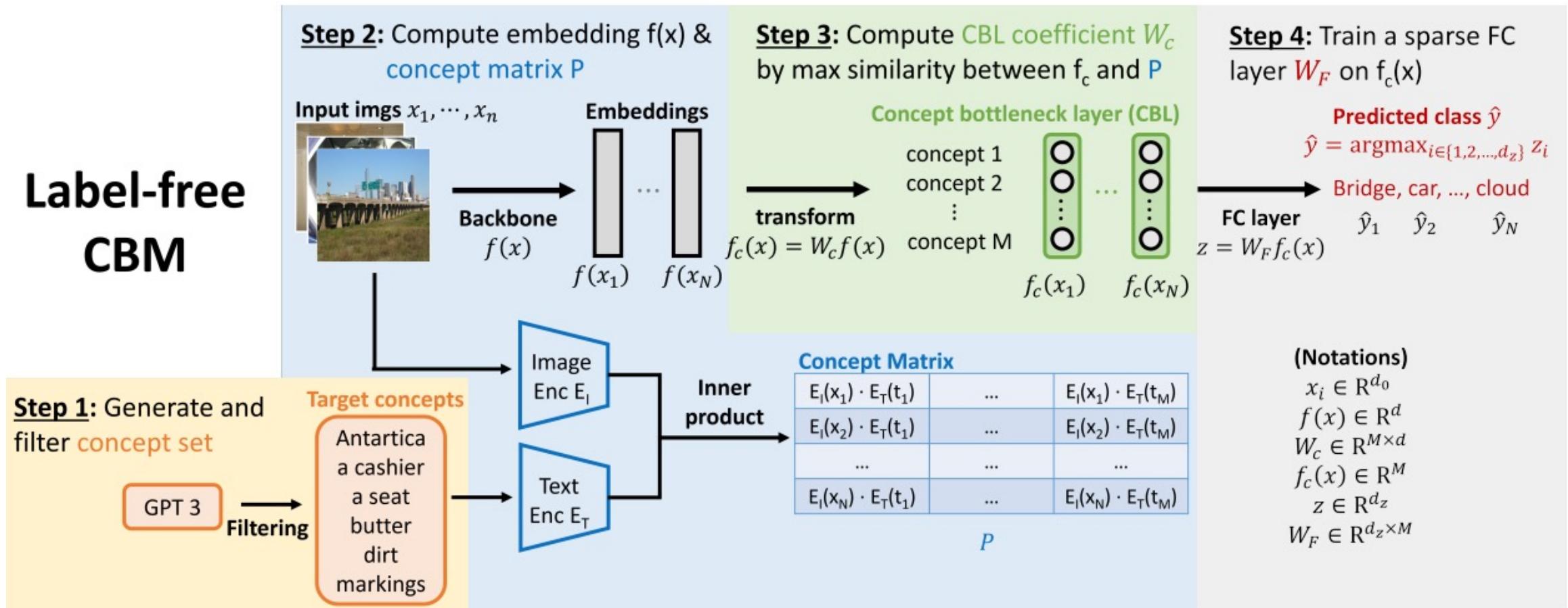
- No extra effort required by human annotators



- Potentially infinite number of concepts
--> no bottleneck in the representation

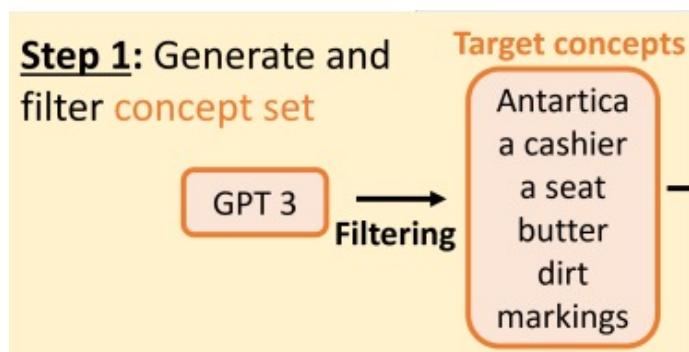
LF-CBM architecture

Label-free CBM



LF-CBM architecture

Label-free CBM

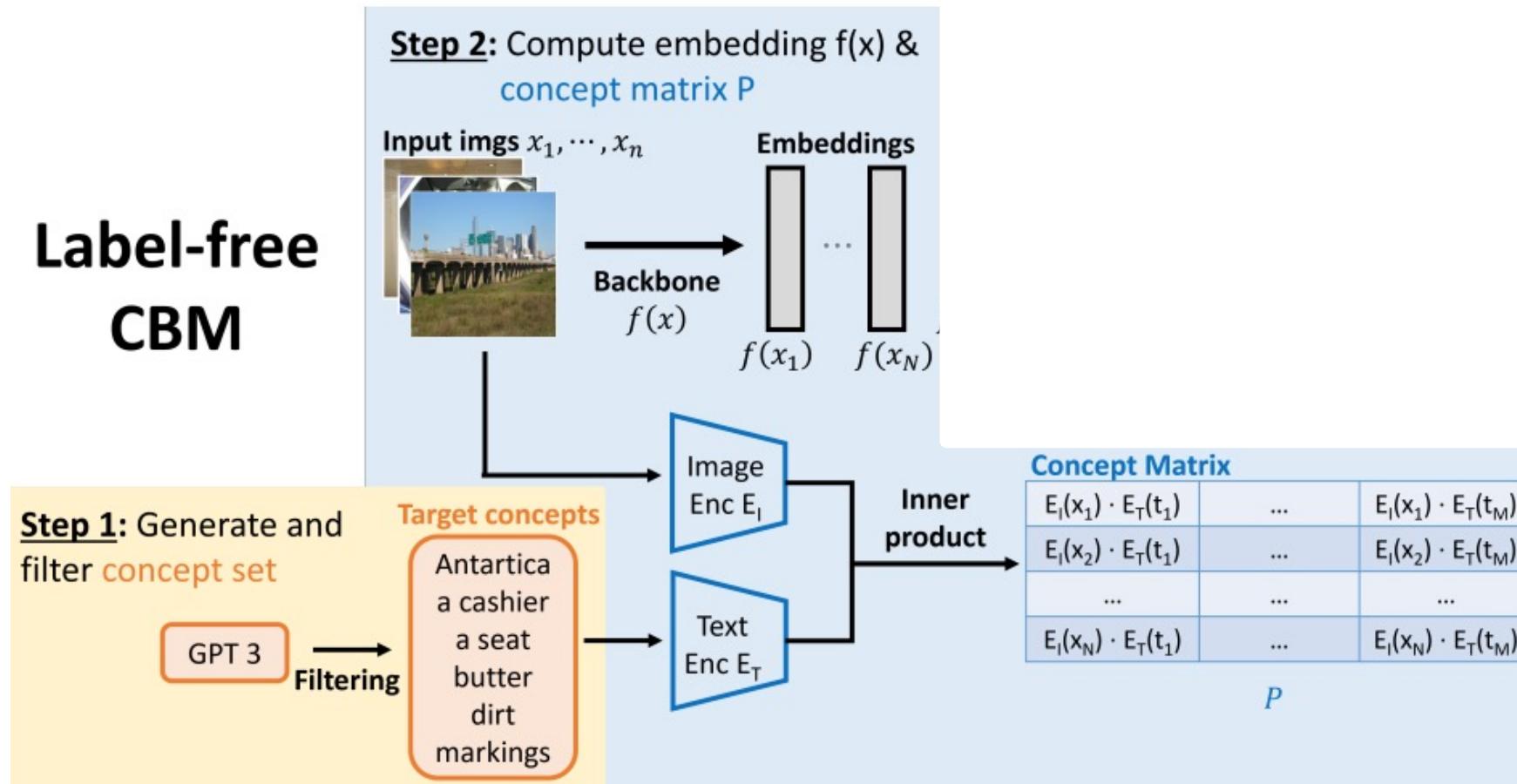


Step 1: Concept set Generation and Filtering

- Generation from GPT-3 prompted with:
 - «*List the most important features for recognizing something as a {class}:»*
- Filtering concepts with:
 - Excessively long descriptions
 - Too similar to classes
 - Too similar to each other
 - Not present in training data
 - Unable to learn with external model

LF-CBM architecture

Label-free CBM



Step 2: Concept matrix definition

Concept Matrix		
$E_I(x_1) \cdot E_T(t_1)$...	$E_I(x_1) \cdot E_T(t_M)$
$E_I(x_2) \cdot E_T(t_1)$...	$E_I(x_2) \cdot E_T(t_M)$
...
$E_I(x_N) \cdot E_T(t_1)$..	$E_I(x_N) \cdot E_T(t_M)$

P

→ CUB



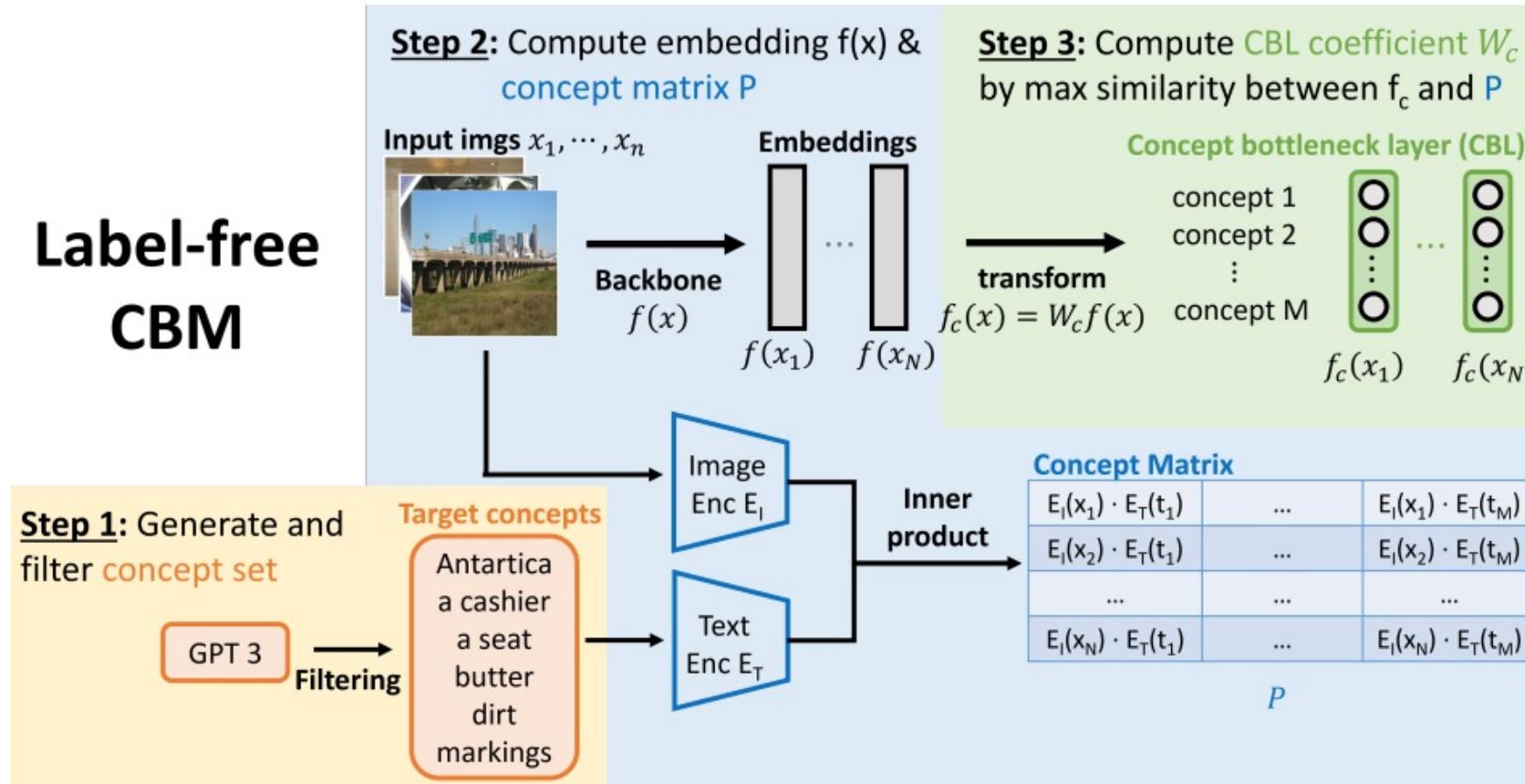
$E_I(x_i)$: CLIP embedding of the ith image

$E_T(t_j)$: CLIP embedding of the jth concept textual description



LF-CBM architecture

Label-free CBM



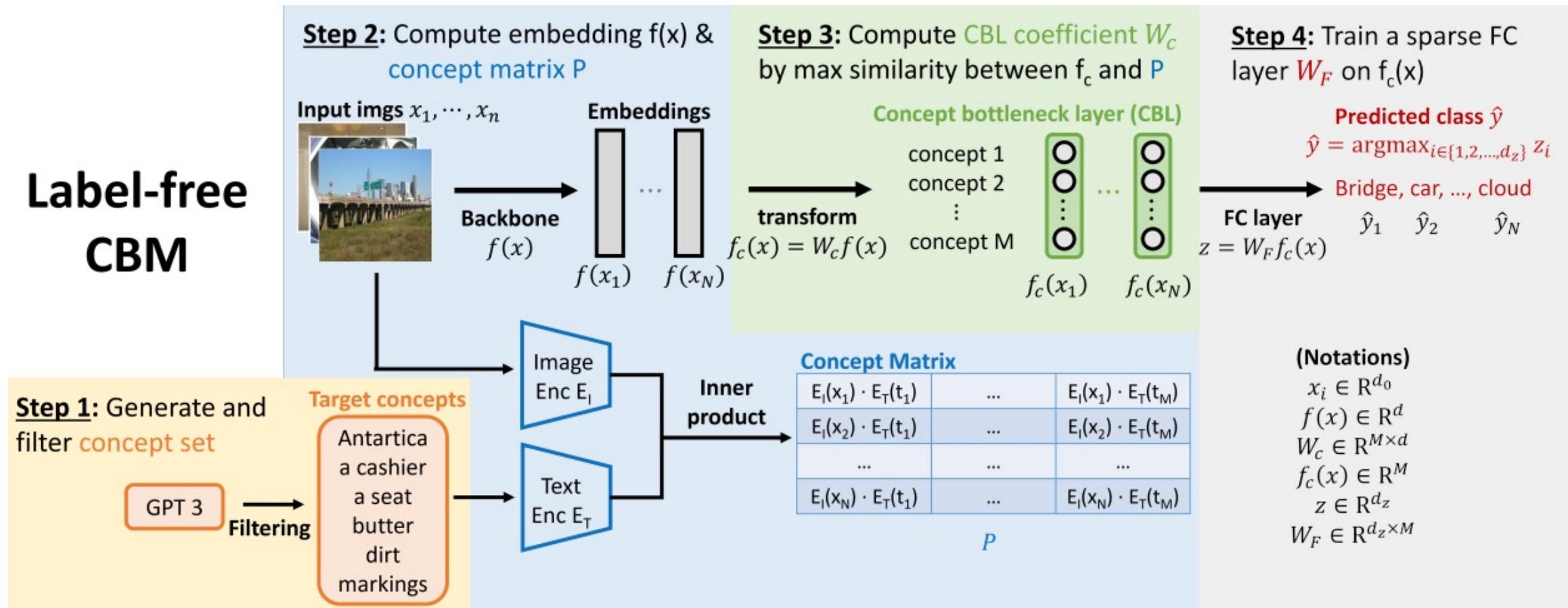
Step 3: Learn to predict the Concept Matrix

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

- $t_i, P_{:,i}$: all CLIP concepts representation of a given concepts for all samples in the training set
- q_i : all predicted representation for a given concepts
 - $q_i = [W_c f(x_1), W_c f(x_2), \dots, W_c f(x_N)]$
 - On top of the backbone image representation $f(x_N)$ that is also trained

LF-CBM architecture

Label-free CBM



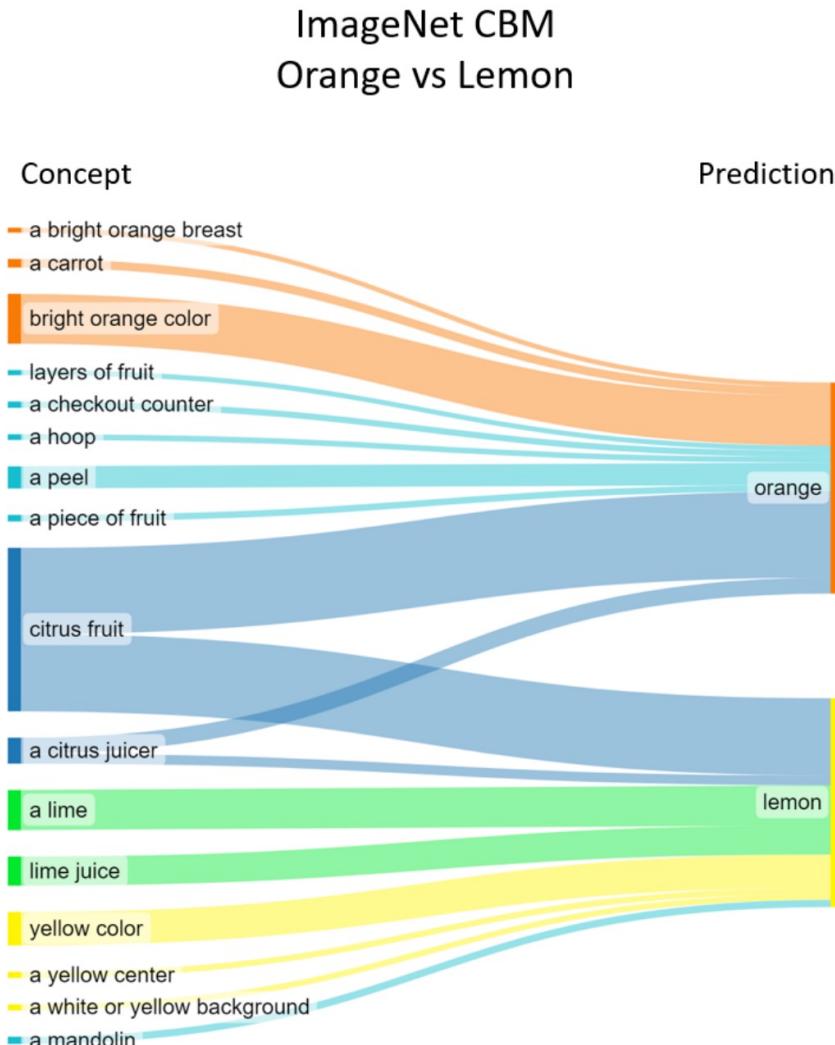
Step 4: Sparse FC layer on top of concept predictions

Final Prediction:

$$\hat{y} = W_F f_c(x)$$

Training:

- sequentially, after having trained the CBL
- $\min \sum_i^N L_{CE} (\hat{y}, y) + ||W_F||$



Performance

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

Open Issues

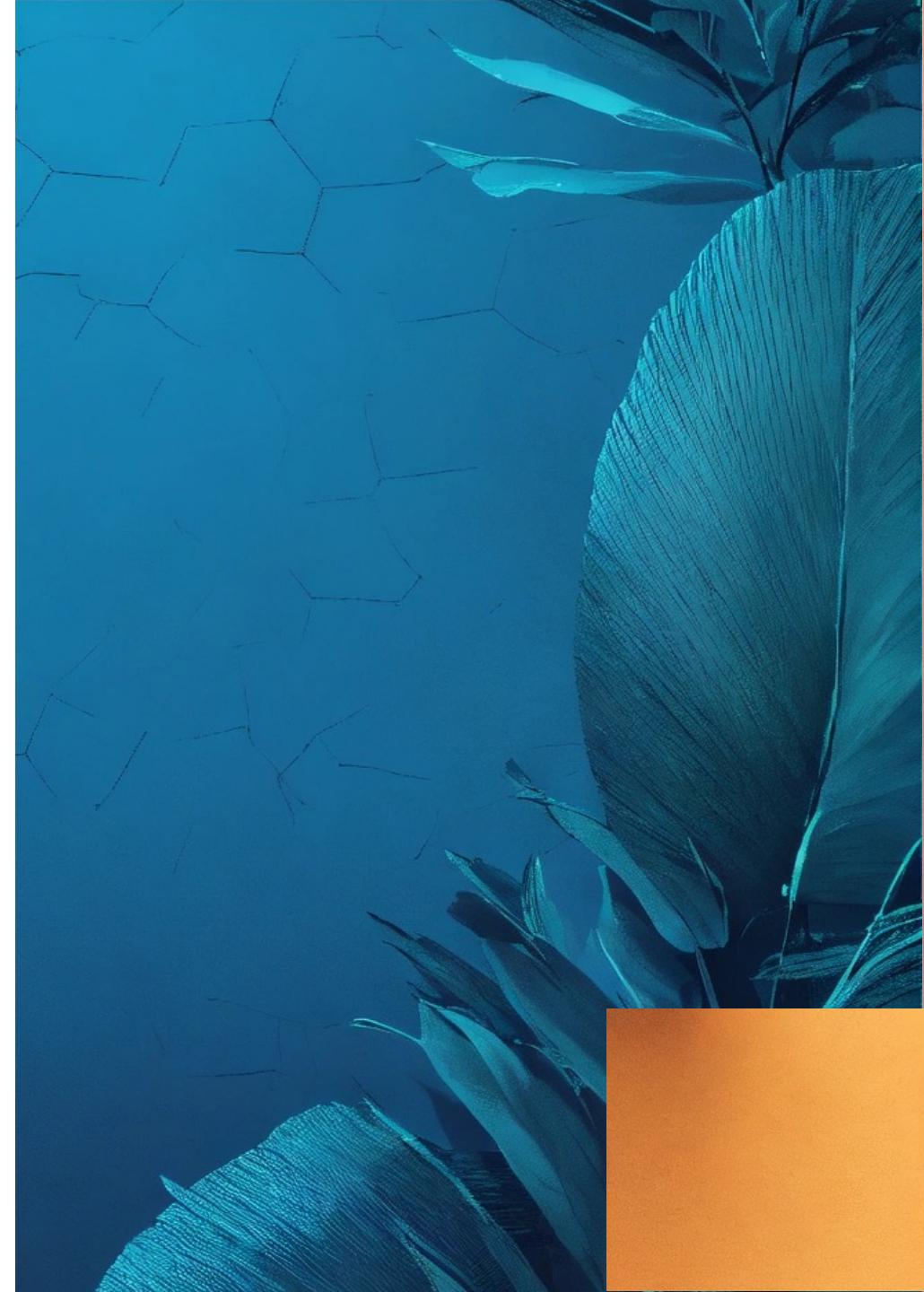
1. The employment of 3 Models to make one prediction is redundant
 1. GPT for concept list generation, CLIP for matrix description, Resnet/ViT for final prediction
2. LF-CBM concept accuracy is very low¹

	MNIST E/O MNIST	Add.	CIFAR100	Skin	CUB200
BotCL (recon)	0.47 ± 0.01	0.41 ± 0.01	0.38 ± 0.03	0.47 ± 0.02	0.34 ± 0.01
BotCL (contr)	0.47 ± 0.02	0.45 ± 0.02	0.40 ± 0.04	0.44 ± 0.03	0.37 ± 0.02
SENN	0.61 ± 0.02	0.58 ± 0.01	0.44 ± 0.02	0.52 ± 0.02	0.41 ± 0.02
ProtoPNet	0.26 ± 0.01	0.24 ± 0.01	0.31 ± 0.01	0.16 ± 0.02	0.28 ± 0.03
LF-CBM	0.52 ± 0.01	0.50 ± 0.03	0.45 ± 0.01	0.58 ± 0.01	0.45 ± 0.01
LCBM (ours)	0.88 ± 0.08	0.81 ± 0.04	0.60 ± 0.01	0.58 ± 0.01	0.55 ± 0.01

1. De Santis, Francesco, et al. "Towards Better Generalization and Interpretability in Unsupervised Concept-Based Models." ECML-PKDD 2025

Tutorial structure

- Introduction to XAI techniques
- Taxonomy of C-XAI
- **Evaluation, Resources and applications**
- Hands-on session



Metrics & Evaluation of Concept-based Models

- **Quantitative Evaluation**

Assesses how much **concepts contribute** to **class predictions** and **task performance** and their **quality**.

Automated, Replicable.

- **Qualitative Evaluation**

Examines the ***quality*** of concept-based explanations ***from a human perspective***.

Reflects **human understanding, hard to scale**.

Metrics & Evaluation of Concept-based Models

- **Quantitative Evaluation**

Two main categories

1. **Contribution to Prediction & Performance**
2. **Quality of Concepts**

Quantitative Evaluation

1. Contribution to Prediction & Performance

*How well do **concepts contribute** to class predictions and task performance?*

Concept effect on class prediction



Provide the **concept relevance** by means of the concept weights, connecting the concept to the final classes

Concept effect on task performance



Measure the **concepts' contribution** to the predictive capacity of the model

Quantitative Evaluation

2. Quality of Concepts

*What is the **quality of the learned concepts**?*

Properties

**Network
relationships**

**Concept Prediction
Error**

Intrinsic characteristics
of the concepts
*e.g., purity,
distinctiveness*

Relationship
between a **node** of
the model and **the**
concept it represents

Alignment between
learned concepts and
concept ground truth

Qualitative Evaluation

Human-centered assessment

Understandability & Trust

Plausibility,
Trustability,
Reasonability &
Factuality

Coherence with Human Concepts

Assess if learned concepts group similar examples

Utility for decision-making

Practical usefulness and effectiveness in providing insights into the model

Quantitative

Strengths.

- Automated
- Scalable
- Reproducible

Limitations.

- Focuses on model-side importance, not human meaning
- Limited when concept labels are not available.

Qualitative

Strengths.

- Reflects real human understanding
- Captures trust and usefulness

Limitations.

- Costly, time-consuming, hard to scale.
- Lack of reproducibility or comparability due to subjectivity.

They are needed together!

Datasets for Concept-based xAI

Concept annotation data

Ground-truth annotation for (symbolic) concepts (e.g., color, texture, object parts)

Why concept datasets matter?

- **Enable supervised learning**

Models can be trained to recognize concepts explicitly

- **Allow evaluation**

Ground-truth concepts allow quantitative evaluation

- **Bridge human ↔ model reasoning**

Ensure concepts are meaningful to users, not only to the network

Datasets for Concept-based xAI

Available for **different modalities**



Image

CUB (200 bird species, 100+ attributes), BRODEN (1000+ visual concepts), OAI.



Text

CEBaB, SST/IMDB.



Video

BDD-OIA (driving actions, 21 concept types).



Tabular

COMPAS (recidivism), MIMIC-II (patient survival).

Datasets for Concept-based xAI

Limitations.

- **Scarcity** of explicit concept labels
Few have explicit annotation, most infer concepts indirectly
- **Limited coverage beyond images**
- **Lack of standardization** across modalities

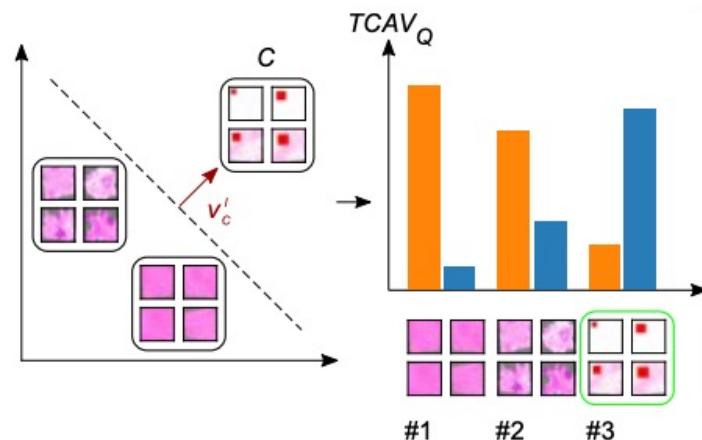
Examples of application of C-XAI methods



Medical.

Stronger adoption of concept-based XAI

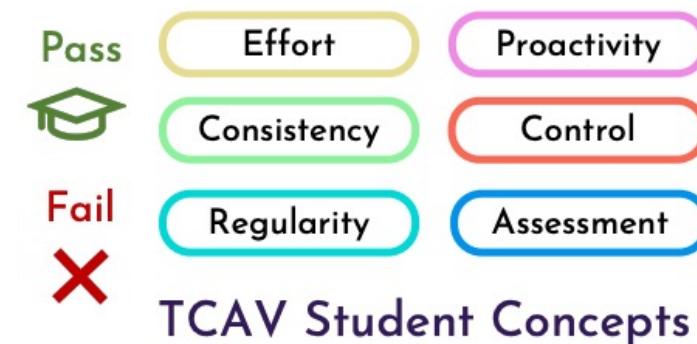
e.g., *TCAV* has been applied to skin cancer diagnosis



Education.

Makes the model's reasoning much more *interpretable* for educators.

e.g., *TCAV* adapted to GNN to interpret student interaction patterns

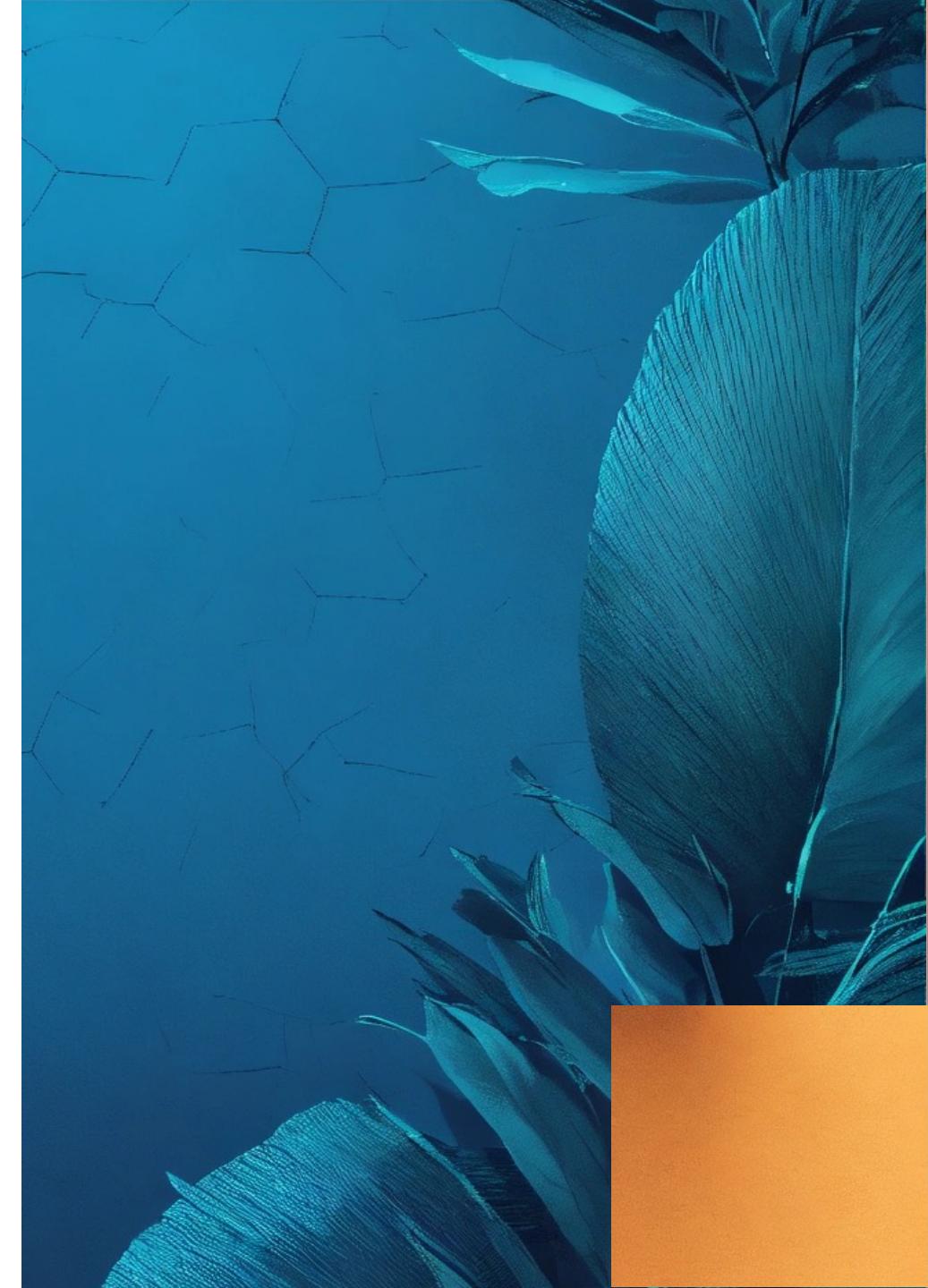


Trends and future directions

- Integration of **foundation and generative models** into C-XAI
 - Increasingly used to **support explanation tasks** and for **concept generation**
- Need for **standardize evaluation metrics**
- Need for **concept-based datasets**
 - Explicit, standardized, multi-modal concept datasets

Tutorial structure

- Introduction to XAI techniques
- Taxonomy of C-XAI
- Evaluation, Resources and applications
- **Hands-on session**



Hands-on time!

<https://github.com/overview-cxai-tutorial-ecml-pkdd-2025/hands-on-cxai-ecml-pkdd>



All notebooks will run in Google Colab  -

no local setup required 😊

- 01 : Testing with Concept Activation Vectors (TCAV)
- 02 : Concept Bottleneck Model (CBM)
- 03 : Label Free Concept Bottleneck Model (LF-CBM)



Github Repository

