

CONCEPT-BASED INTERPRETABLE DEEP LEARNING

AAAI 2025 – Philadelphia, USA



WHO ARE WE?



Mateo Espinosa Zarlenga
Final-year PhD Student
University of Cambridge
me466@cam.ac.uk

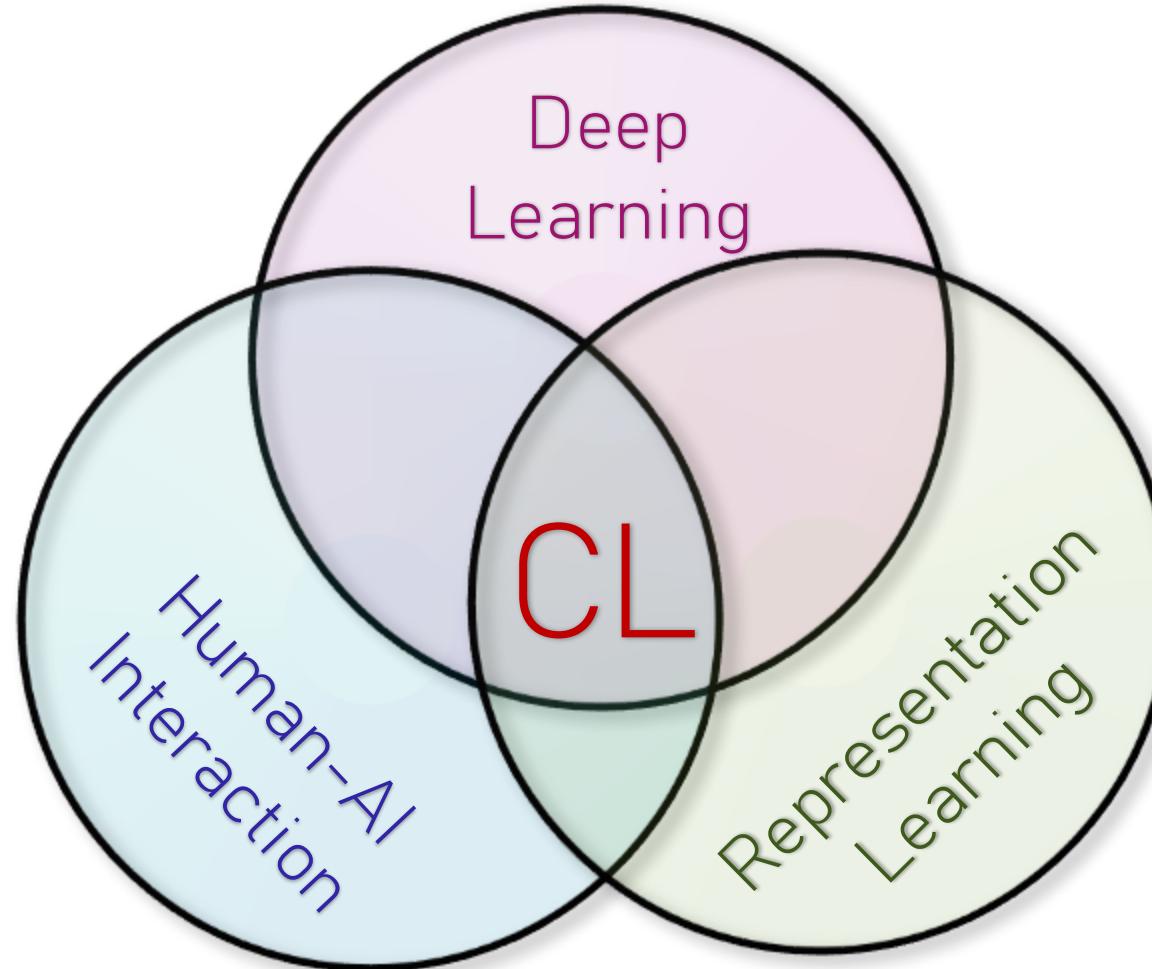


Pietro Barbiero
Swiss Postdoctoral Fellow
IBM Research (Switzerland)
pietro.barbiero@ibm.com



IBM Research

WHAT IS THIS TUTORIAL ABOUT?



We will look at how **Concept Learning (CL)** can be used to design **interpretable Deep Neural Networks**

TUTORIAL GOALS

Our **main goals** for this tutorial are threefold:

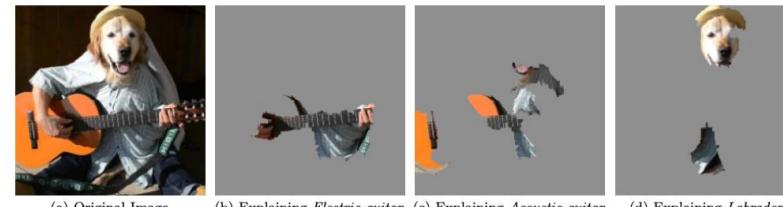
1. Provide a **non-exhaustive but well-rounded overview** of **concept learning (CL)**.
2. Convince you that **concept representations** can be very useful for **designing powerful but interpretable** neural models.
3. Bring together a variety of **resources** (surveys, method papers, libraries, etc.) to **facilitate access to the current state of CL**.

WHAT IS THIS TUTORIAL NOT ABOUT?

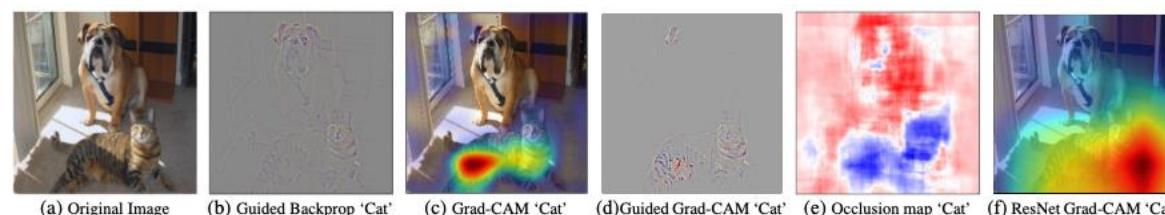
WHAT IS THIS TUTORIAL NOT ABOUT?

We will **not** have time to dive deep into:

1. “Traditional” explainable AI (XAI) methodologies



Example of *LIME* (taken from [1])



Example of *GradCAM* and other saliency methods (taken from [2])

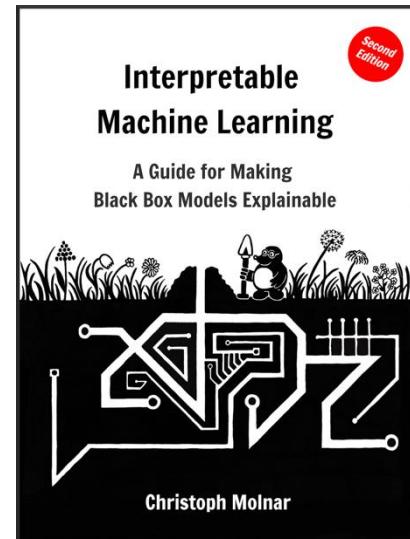
[1] Ribeiro et al. "Why should i trust you?" Explaining the predictions of any classifier." KDD (2016).

[2] Selvaraju, et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." ICCV (2017).

WHAT IS THIS TUTORIAL NOT ABOUT?

We will **not** have time to dive deep into:

1. “Traditional” explainable AI (XAI) methodologies



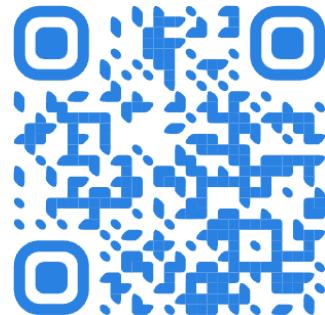
Link to Book

Interpretable Machine Learning
Christoph Molnar

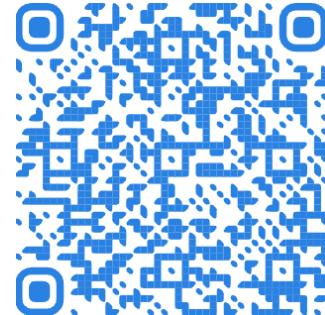
WHAT IS THIS TUTORIAL NOT ABOUT?

We will **not** have time to dive deep into:

1. “Traditional” explainable AI (XAI) methodologies
2. Deep **philosophical aspects** of explaining models



The Mythos of Interpretability
Lipton et al. (2018) [1]



To Explain or to Predict?
Galit (2018) [2]



Explanation Theory
Bromberger (1992) [3]

[1] Lipton, Zachary C. "The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery." *Queue* (2018).

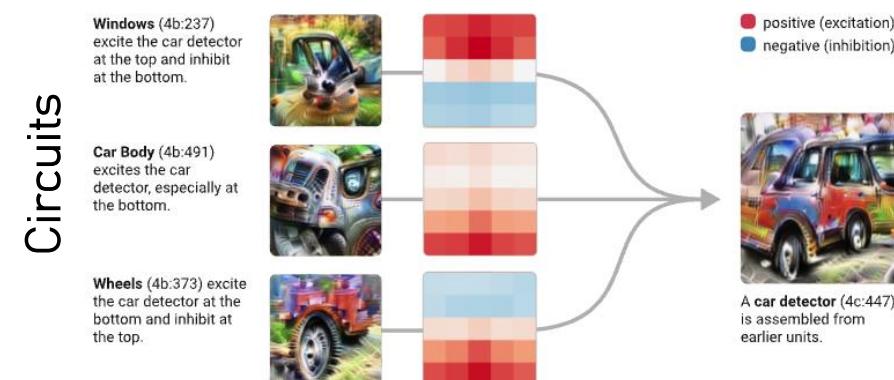
[2] Galit Shmueli. "To Explain or to Predict?." *Statist. Sci.* 25 (3) 289 - 310, August 2010. <https://doi.org/10.1214/10-STS330>

[3] Bromberger, Sylvain. *On what we know we don't know: Explanation, theory, linguistics, and how questions shape them*. University of Chicago Press, 1992.

WHAT IS THIS TUTORIAL NOT ABOUT?

We will **not** have time to dive deep into:

1. “Traditional” explainable AI (XAI) methodologies
2. Deep philosophical aspects of explaining models
3. Connections with **Mechanistic Interpretability**



Taken from [1]

[1] Olah, Chris, et al. "Zoom in: An introduction to circuits." *Distill* 5.3 (2020): e00024-001



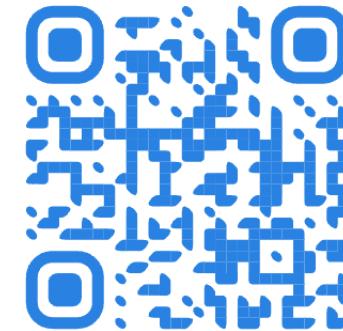
WHAT IS THIS TUTORIAL NOT ABOUT?

We will **not** have time to dive deep into:

1. “Traditional” explainable AI (XAI) methodologies
2. Deep philosophical aspects of explaining models
3. Connections with **Mechanistic Interpretability**



Distill Circuits Thread



Anthropic Circuits Thread

[1] Cammarata, Nick, et al. "Thread: circuits." *Distill* 5.3 (2020): e24.

[2] Anthropic "Transformer Circuits Thread" found at <https://transformer-circuits.pub/>

TUTORIAL OUTLINE

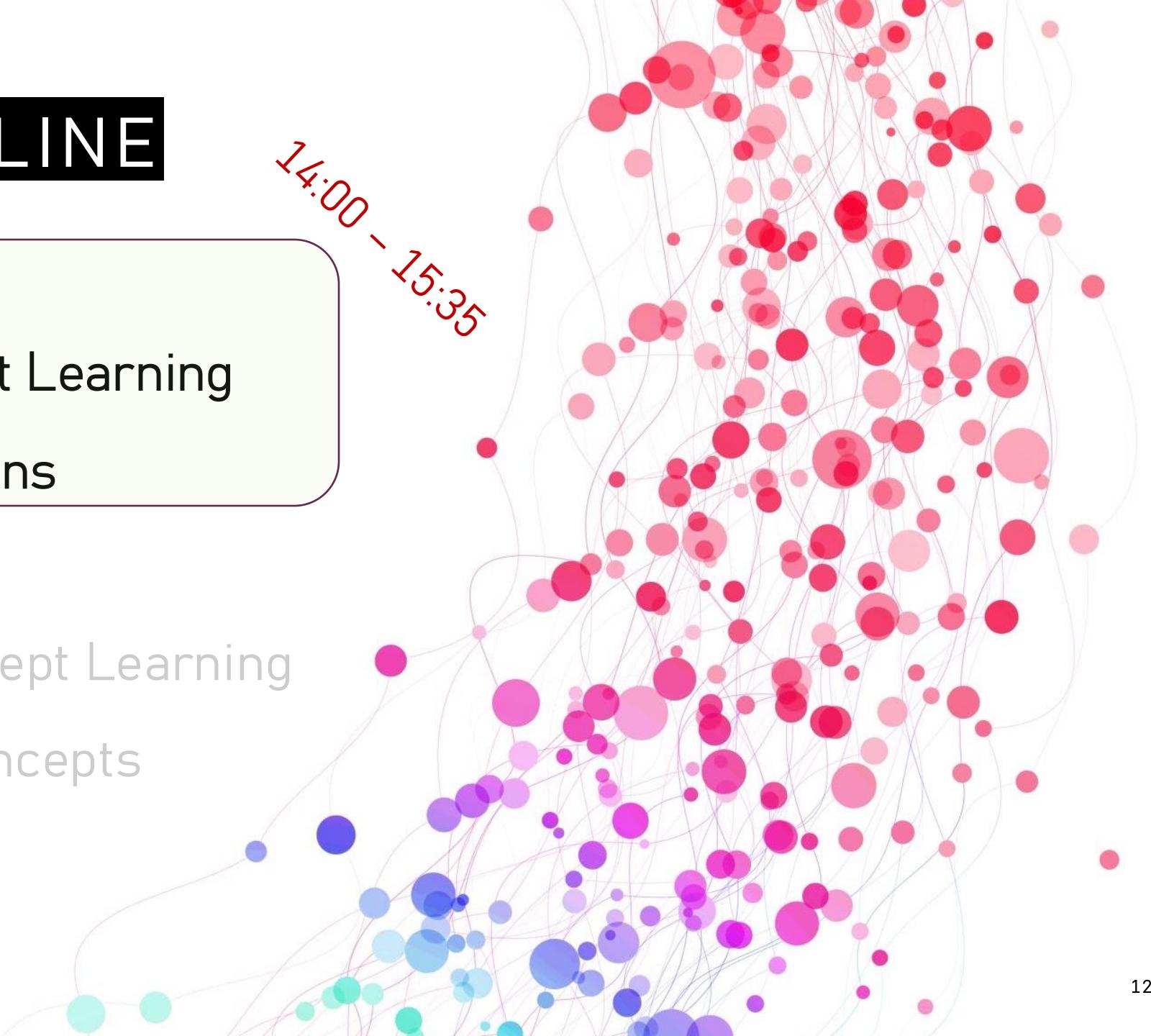
1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



TUTORIAL OUTLINE

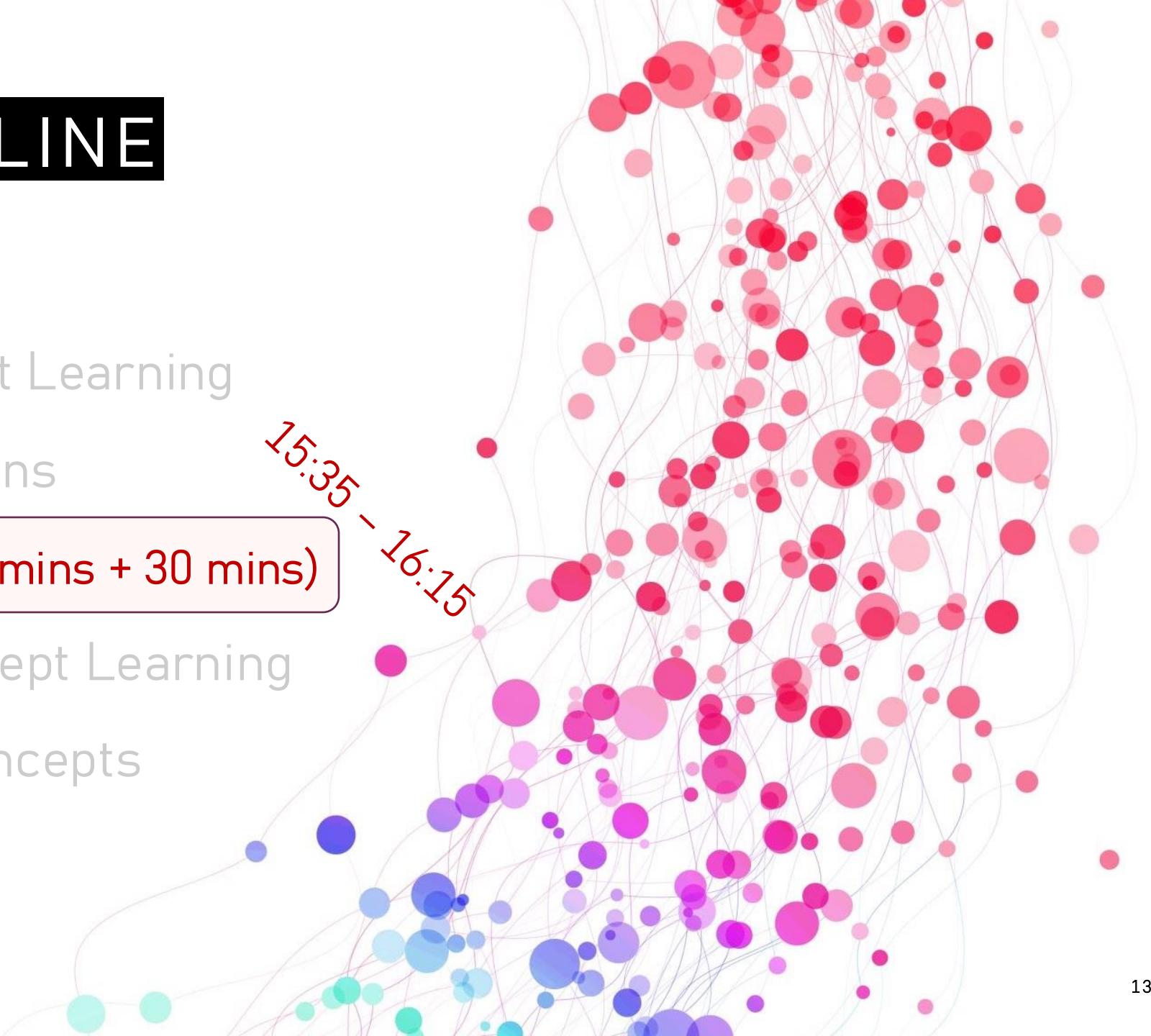
1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A

14:00 - 15:35



TUTORIAL OUTLINE

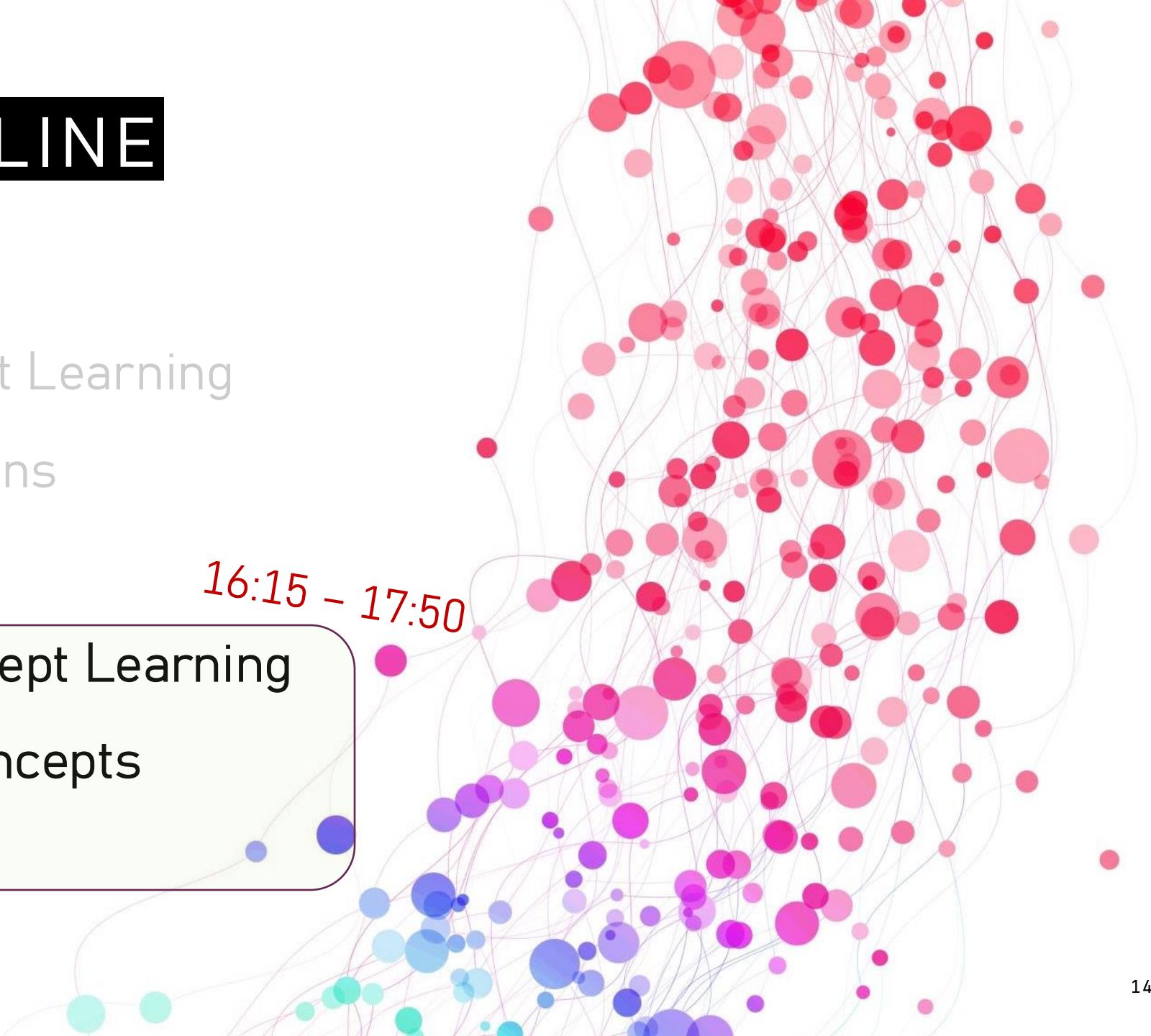
1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
- 4. Q&A + Break (10 mins + 30 mins)**
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



TUTORIAL OUTLINE

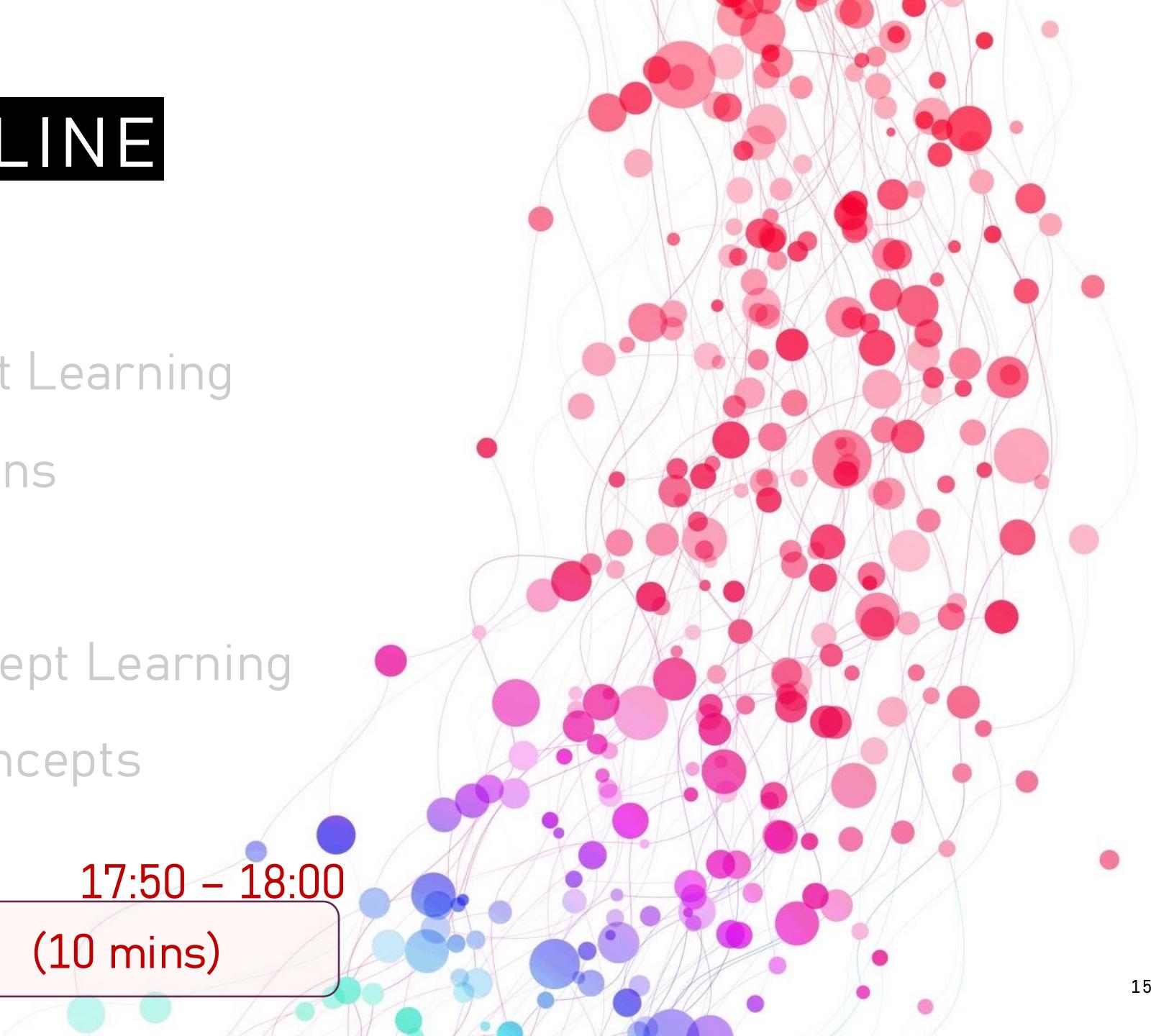
1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A

16:15 - 17:50



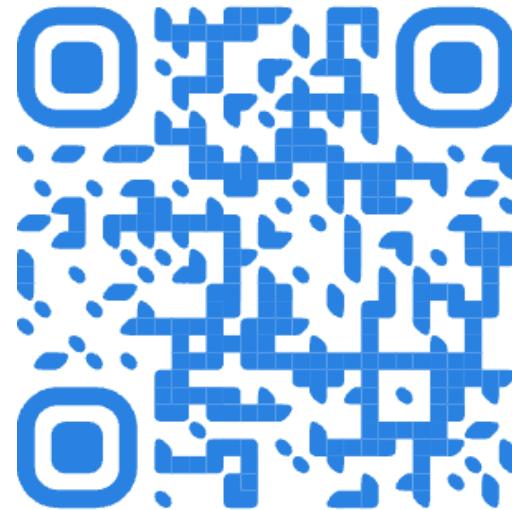
TUTORIAL OUTLINE

1. Introduction
 2. Supervised Concept Learning
 3. Concept Interventions
 4. Q&A + Break
 5. Unsupervised Concept Learning
 6. Reasoning With Concepts
 7. Future Directions
 8. Q&A
- 17:50 – 18:00
(10 mins)



TUTORIAL WEBSITE AND MATERIALS

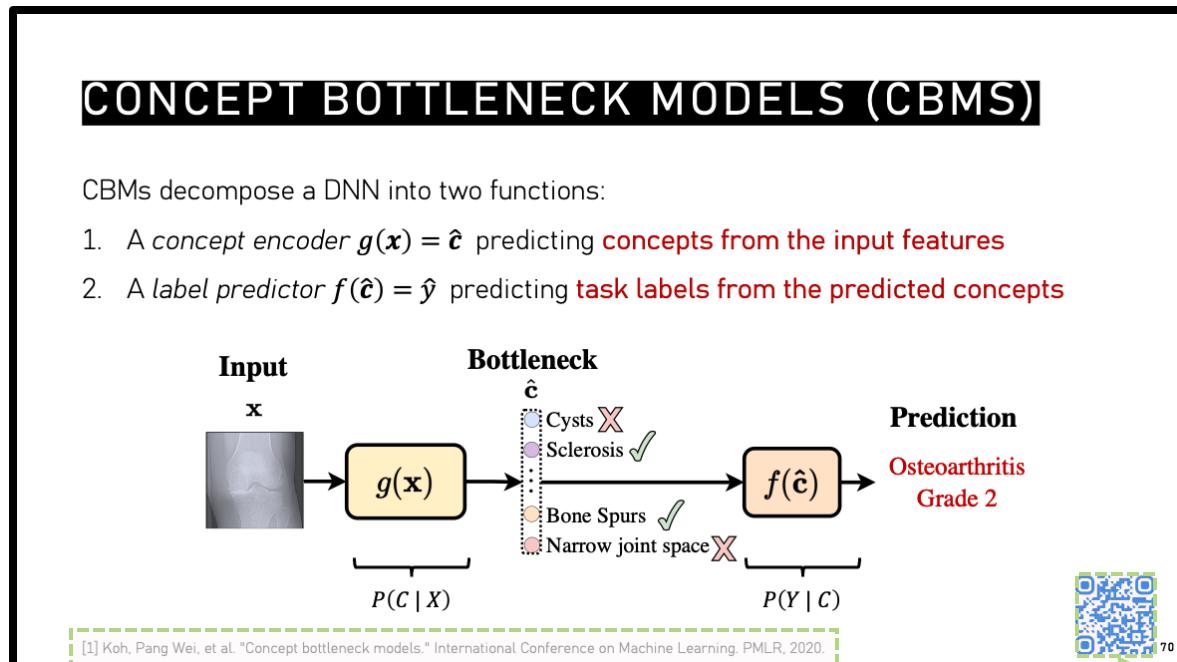
This tutorial's **slides**, **schedule**, and **resources** are in our website:



<https://conceptlearning.github.io/>

TUTORIAL WEBSITE AND MATERIALS

Throughout the tutorial, watch for **QR codes** to relevant references

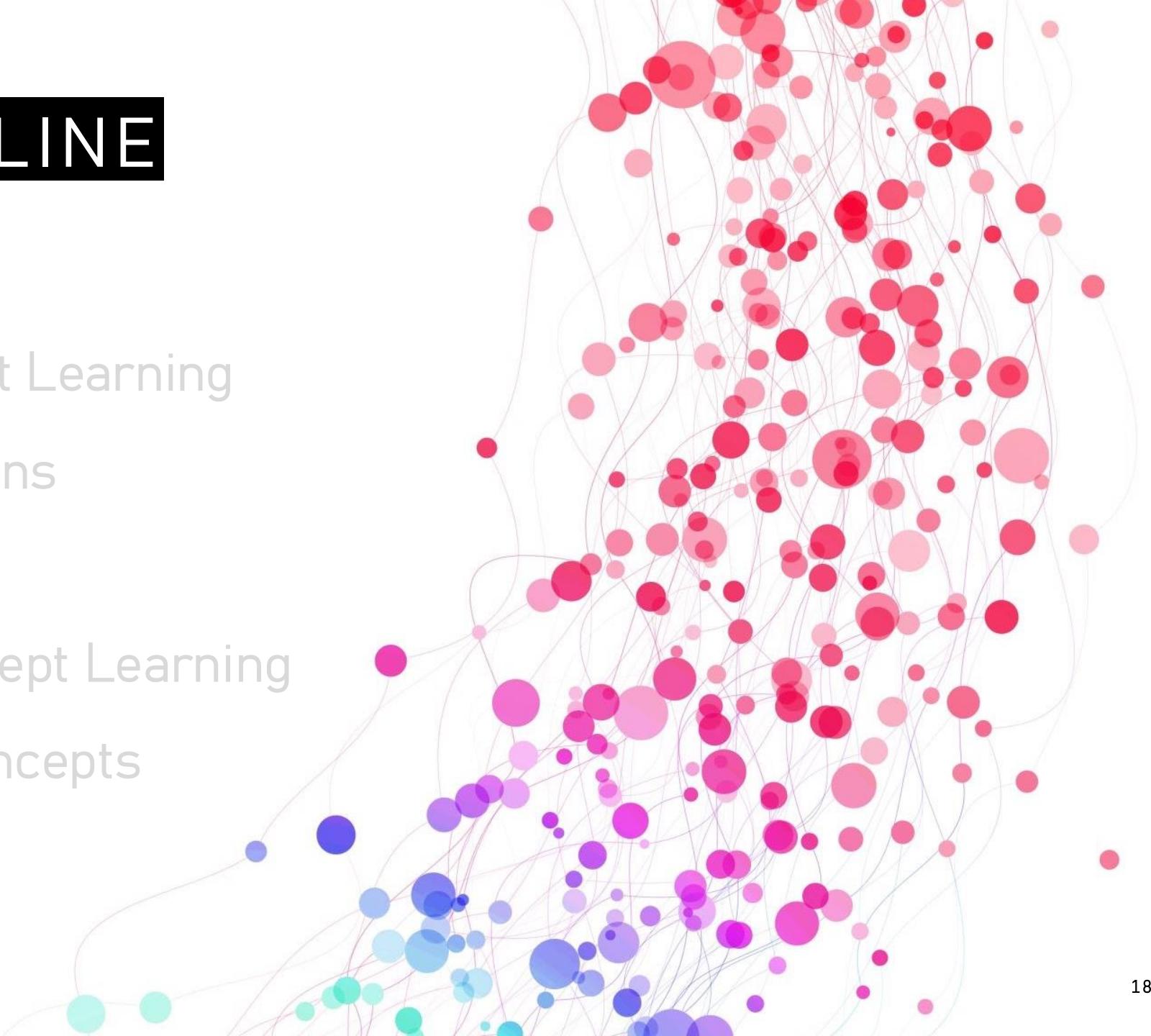


Citation + Hyperlink
(if you download slides)

QR code to paper/reference/extra material

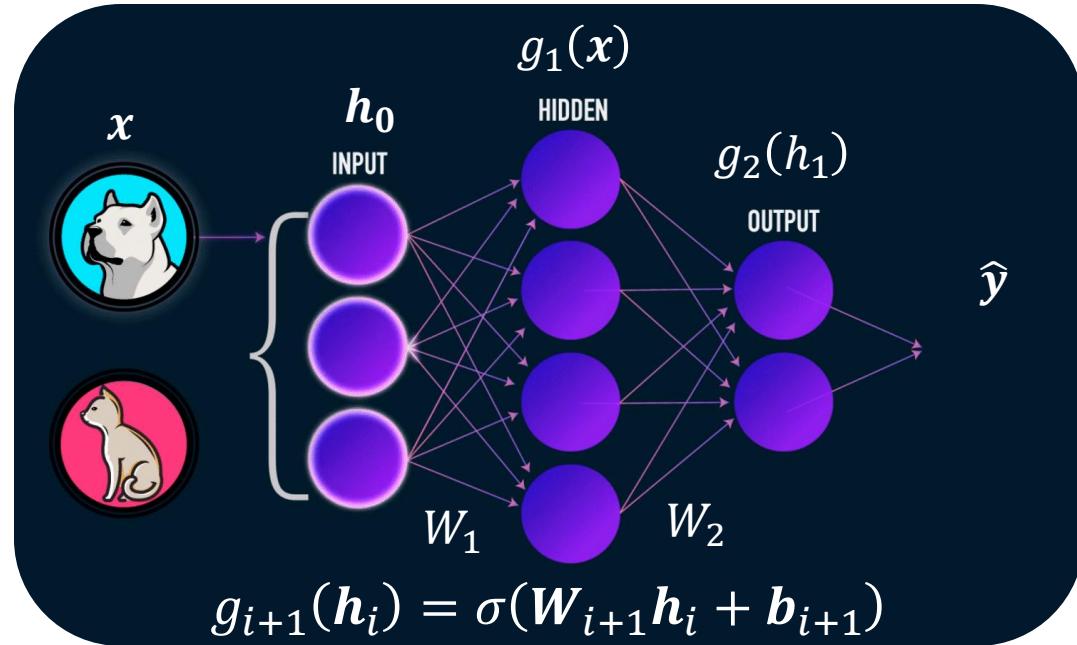
TUTORIAL OUTLINE

1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



A SWISS ARMY KNIFE FOR AI

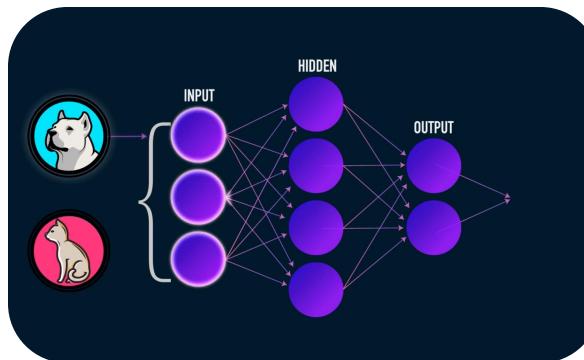
Artificial Intelligence (AI) has experienced **a boom** in the last decade driven by so-called **Deep Neural Networks (DNNs)**



Goal: learn $\{W_1, b_1, \dots, W_m, b_m\}$ s.t. $(g_m \circ \dots \circ g_2 \circ g_1)(x) = \hat{y} \approx y$

THE POWER OF SCALE

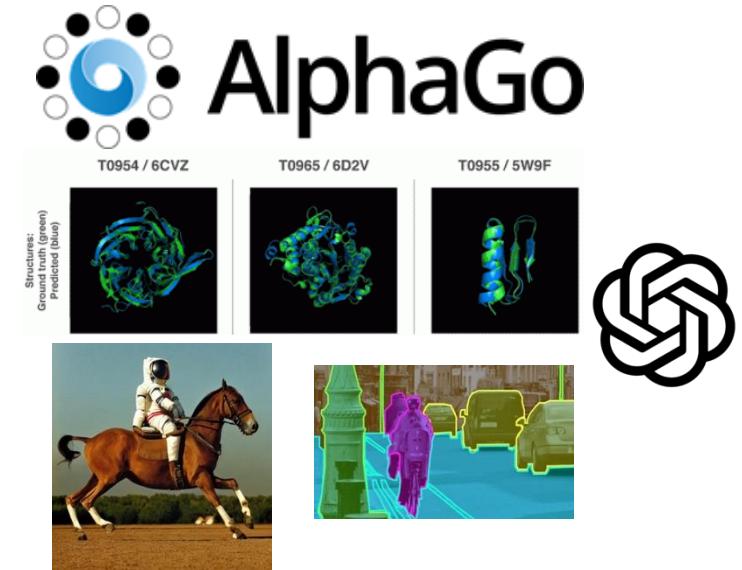
Scaling up DNNs can lead to **expressive** and **generalisable** models:



+

A LOT of **data**,
money, **time**,
and **sweat**

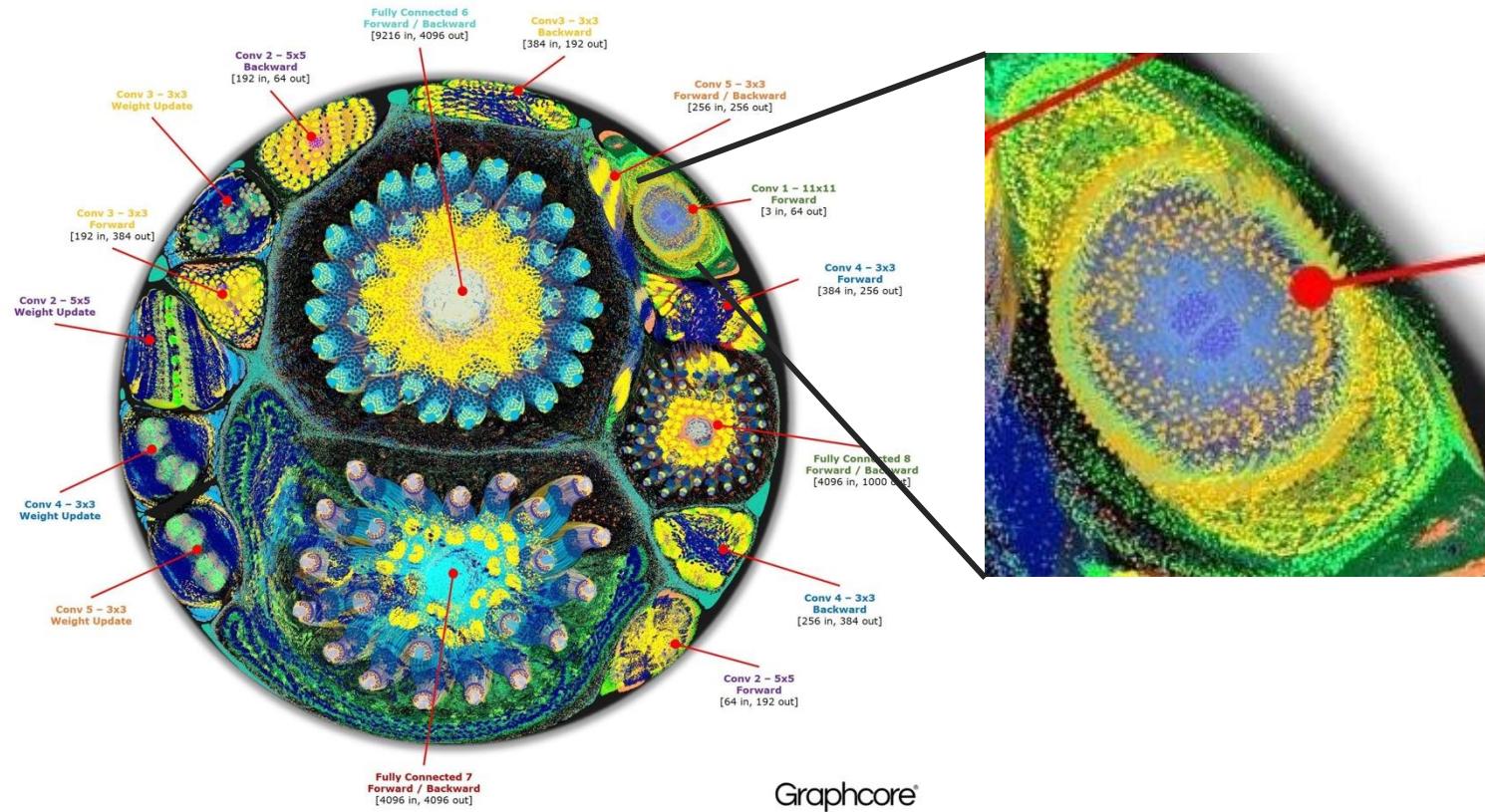
=



- [1] Silver, David, et al. "Mastering the game of Go with deep neural networks and tree search." *nature* 529.7587 (2016): 484-489.
- [2] OpenAI. "GPT-4 technical report." *arXiv* (2023).
- [3] Jumper, John, et al. "Highly accurate protein structure prediction with AlphaFold." *nature* 596.7873 (2021): 583-589.
- [4] Rombach, Robin, et al. "High-resolution image synthesis with latent diffusion models." *CVPR* (2022).

THE BLACK-BOX PROBLEM

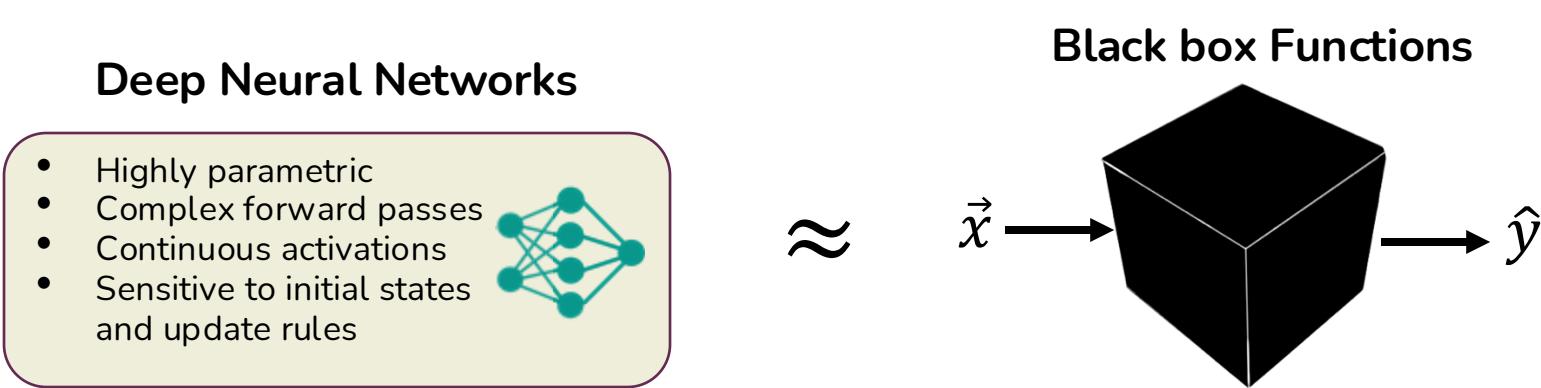
Scale, however, leads to **notoriously complex models!**



[1] Adapted from Graphcore, "Inside an AI 'Brain' – What does Machine Learning Look Like?" (2017)

THE BLACK-BOX PROBLEM

Scale, however, leads to **notoriously complex models!**



DNNs are "**black-box**" models

THE FLIP SIDE OF THE COIN

Blindly using black-box models can lead to all sorts of problems:

Wrongfully Accused by an Algorithm

In what may be the first known case of its kind, a faulty facial recognition match led to a Michigan man's arrest for a crime he did not commit.

Why Amazon's Automated Hiring Tool Discriminated Against Women

**Predictive policing algorithms are racist.
They need to be dismantled.**

[1] Kashmir Hill, "Wrongfully Accused by an Algorithm." The New York Times (2020).

[2] Rachel Goodman, "Why Amazon's Automated Hiring Tool Discriminated Against Women." ACLU (2018).

[3] Will Douglas Heaven, "Predictive policing algorithms are racist. They need to be dismantled." MIT Technology Review (2020).

THE FLIP SIDE OF THE COIN

Blindly using black-box models can lead to all sorts of problems:

It's not all bad news ☺

Why Amazon's Automated Hiring
Tool Discriminated Against Women

Predictive policing algorithms are racist.
They need to be dismantled.

[1] Natasha Bernal, "IBM Watson AI criticised after giving 'unsafe' cancer treatment advice." The Telegraph (2018).

[2] Kashmir Hill, "Wrongfully Accused by an Algorithm." The New York Times (2020).

[3] Rachel Goodman, "Why Amazon's Automated Hiring Tool Discriminated Against Women." ACLU (2018).

[4] Will Douglas Heaven, "Predictive policing algorithms are racist. They need to be dismantled." MIT Technology Review (2020).

EXPLAINING DNNs

Recent advances in AI came with **a rise in interest** in making these models "**interpretable**"

CIO BLOG

Companies Grapple With AI's Opaque Decision-Making Process

Uber, Xerox's PARC, Capital One among organizations investigating how AI solves problems

Opinion Artificial intelligence

Beware the rise of the black box algorithm

WHO calls for safe and ethical AI for health

16 May 2023 | Departmental update | Reading time: 2 min (507 words)

How to Build Artificial Intelligence We Can Trust

Computer systems need to understand time, space and causality.
Right now they don't.

Building Trust In AI: The Case For Transparency

Why businesses need explainable AI—and how to deliver it

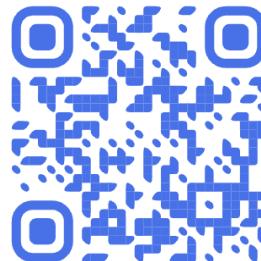
September 29, 2022 | Article

EXPLAINING DNNs

This interest has manifested itself at the **regulatory/legal level**

General Data Protection Regulations (GDPR, 2016):

- "The data subject shall have the right not to be subject to a **decision** based solely on **automated processing**, including profiling...." (Art. 22)
- The data subject has the right to "**meaningful information** about the **logic** involved" in the decision. (Art. 13 and 15)



GDPR

EU AI Act (2024):

- "Any affected person subject to a **decision** which is taken by.. a **high-risk AI system** ... shall have the right to obtain from the deployer **clear and meaningful explanations** (Art. 86)



EU AI Act

[1] GDPR. EU. "Automated individual decision-making, including profiling." (2022).

[2] Act. EU Artificial Intelligence. "The EU Artificial Intelligence Act." (2024).

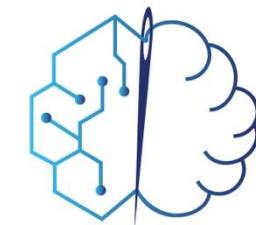
EXPLAINING DNNs

Researchers in Explainable Artificial Intelligence (XAI*) have developed a significant number of methods to explain DNNs



Explainable Artificial
Intelligence (XAI)

(DARPA 2016)



TAILOR

(EU Horizon Program)

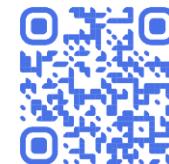
*Not to be confused with a certain bird-related company

FIRST THINGS FIRST: TERMINOLOGY

Welcome to the **Wild West** of XAI terminology



Goebel et al. (2018)



Freiesleben et al. (2023)



Gilpin et al. (2018)



Vilone et al. (2021)



Confalonieri et al. (2020)



Rudin et al. (2019)



Barredo Arrieta et al. (2020)

FIRST THINGS FIRST: TERMINOLOGY

Here we will use some the following definitions by Gilpin et al. [1]:

- **Explainability (why)**: the ability to answer questions of the form "*why does this particular input lead to that particular output?*"
- **Interpretability (how)**: the ability to describe "the internals of a system in a way that is understandable to humans."

[1] Gilpin et al. "Explaining explanations: An overview of interpretability of machine learning." DSAA (2018).



FEATURE ATTRIBUTION

XAI methods have traditionally explained a model's prediction by estimating how **important** each **input feature** is for the **output**



SHAP (Scott et al., 2017)

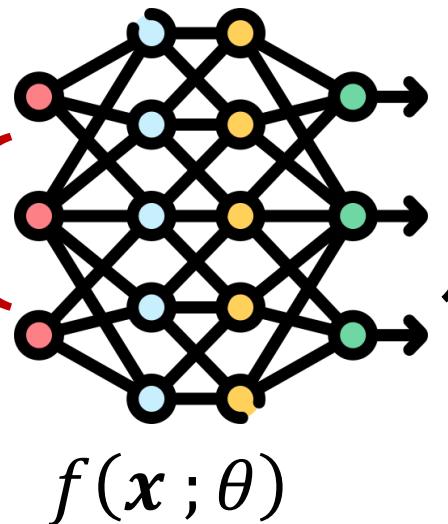
We call these **feature importance** or **feature attribution** methods

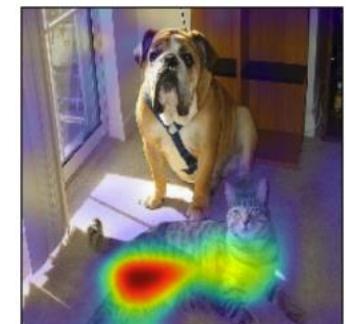
[1] Scott, M., and Lee Su-In. "A unified approach to interpreting model predictions." NeurIPS (2017).



SALIENCY: FEATURE ATTRIBUTION IN DNNS

DNN-specific attribution methods are called **saliency methods**

$$\psi(f(x; \theta), x, \hat{y}, \text{cat}) =$$




These are usually computed by measuring **model sensitivity via its gradient** $\frac{\partial f(x)_y}{\partial x_i}$

[1] Example taken from [Selvaraju et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization."](#) ICCV (2017).



WHAT'S WRONG WITH FEATURE ATTRIBUTION?

WHAT'S WRONG WITH FEATURE ATTRIBUTION?

1. Low-level features like individual pixels are **not always semantically meaningful**:



Can you guess what this is?

WHAT'S WRONG WITH FEATURE ATTRIBUTION?

1. Low-level features like individual pixels are **not always semantically meaningful**:



Limes!

WHAT'S WRONG WITH FEATURE ATTRIBUTION?

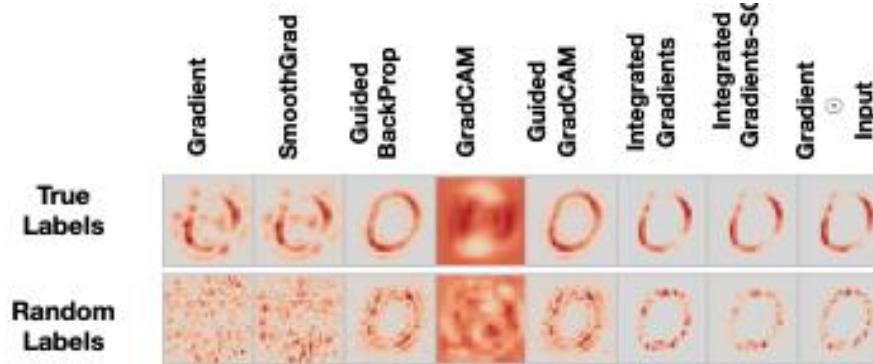
2. Saliency maps lack of **actionability!**



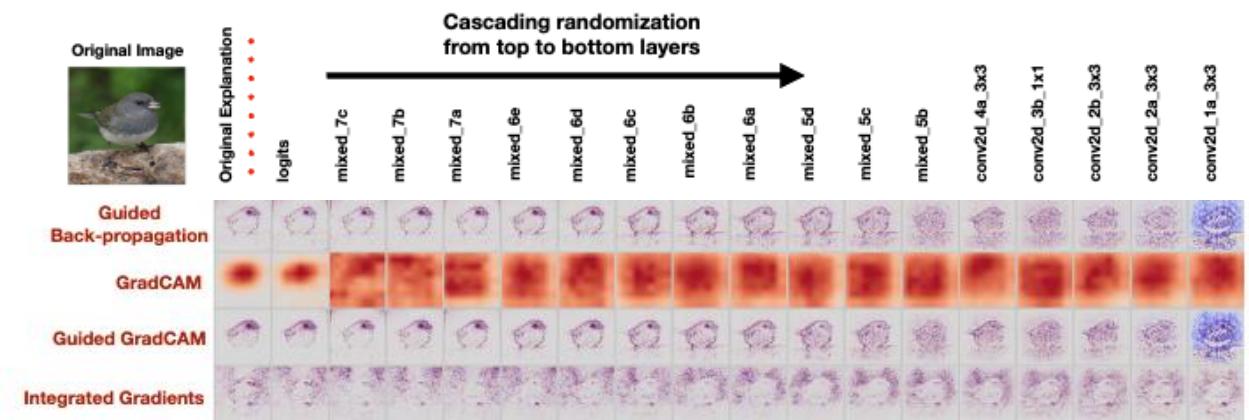
What does this really tell you about **how** the model made a prediction?

WHAT'S WRONG WITH FEATURE ATTRIBUTION?

3. Several saliency methods **fail very simple sanity checks**



Random training labels do not always lead to random maps [1]



Random weights do not always lead to random maps [1]

[1] Adebayo, Julius, et al. "Sanity checks for saliency maps." NeurIPS (2018).



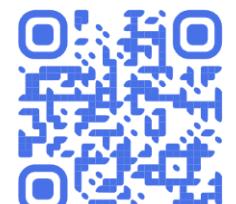
WHAT'S WRONG WITH FEATURE ATTRIBUTION?

4. Saliency methods are susceptible to **adversarial attacks** [1,2]



[1] Dombrowski et al. "Explanations can be manipulated and geometry is to blame." NeurIPS (2019).

[2] Also relevant Ghorbani et al. "Interpretation of neural networks is fragile." AAAI (2019).



WHAT'S WRONG WITH FEATURE ATTRIBUTION?

How can we go around the limitations of feature attribution?

Here, we will focus on using so-called
“*concepts*” to construct explanations

WHAT ARE CONCEPTS?

Concepts are **high-level** and semantically **meaningful** units of information

Task: bird species

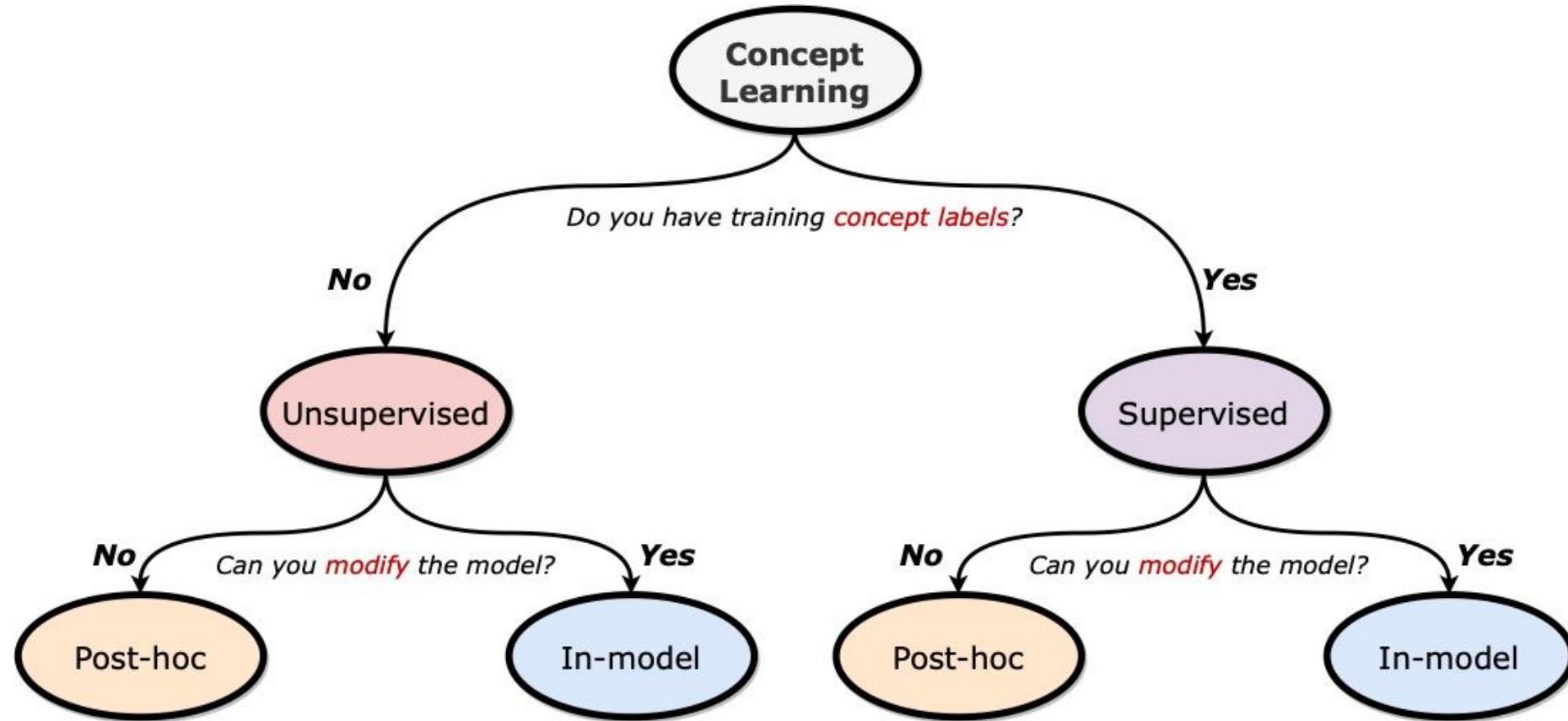


Explanation of the prediction:

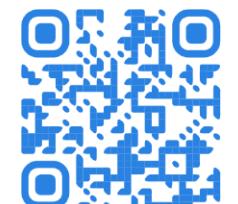
- wing color
- beak length
- tail shape

Concepts are terms or units of information **used by domain experts** to communicate or explain things to each other

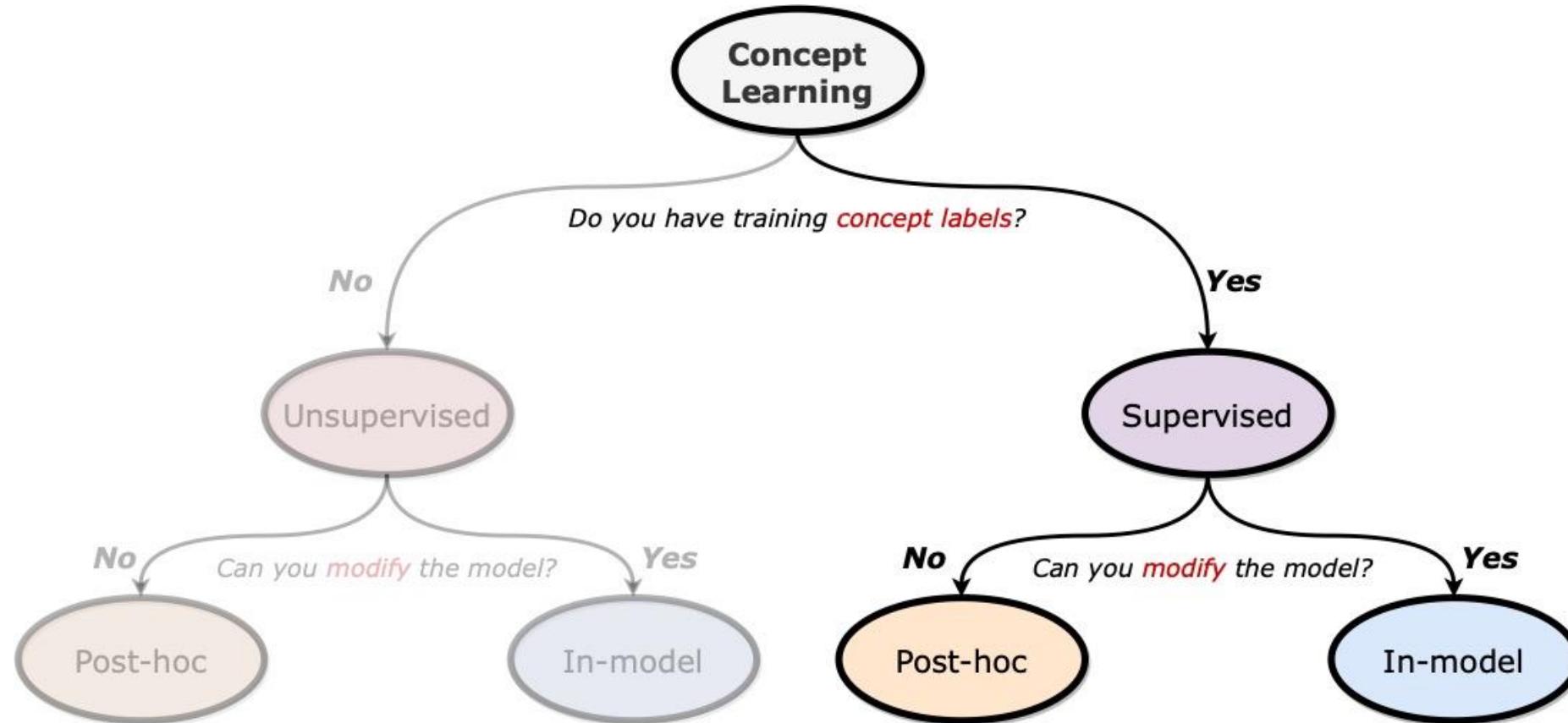
A MAP OF CONCEPT LEARNING



[1] For a more complete taxonomy, see Poeta et al. "Concept-based explainable artificial intelligence: A survey." arXiv (2023).



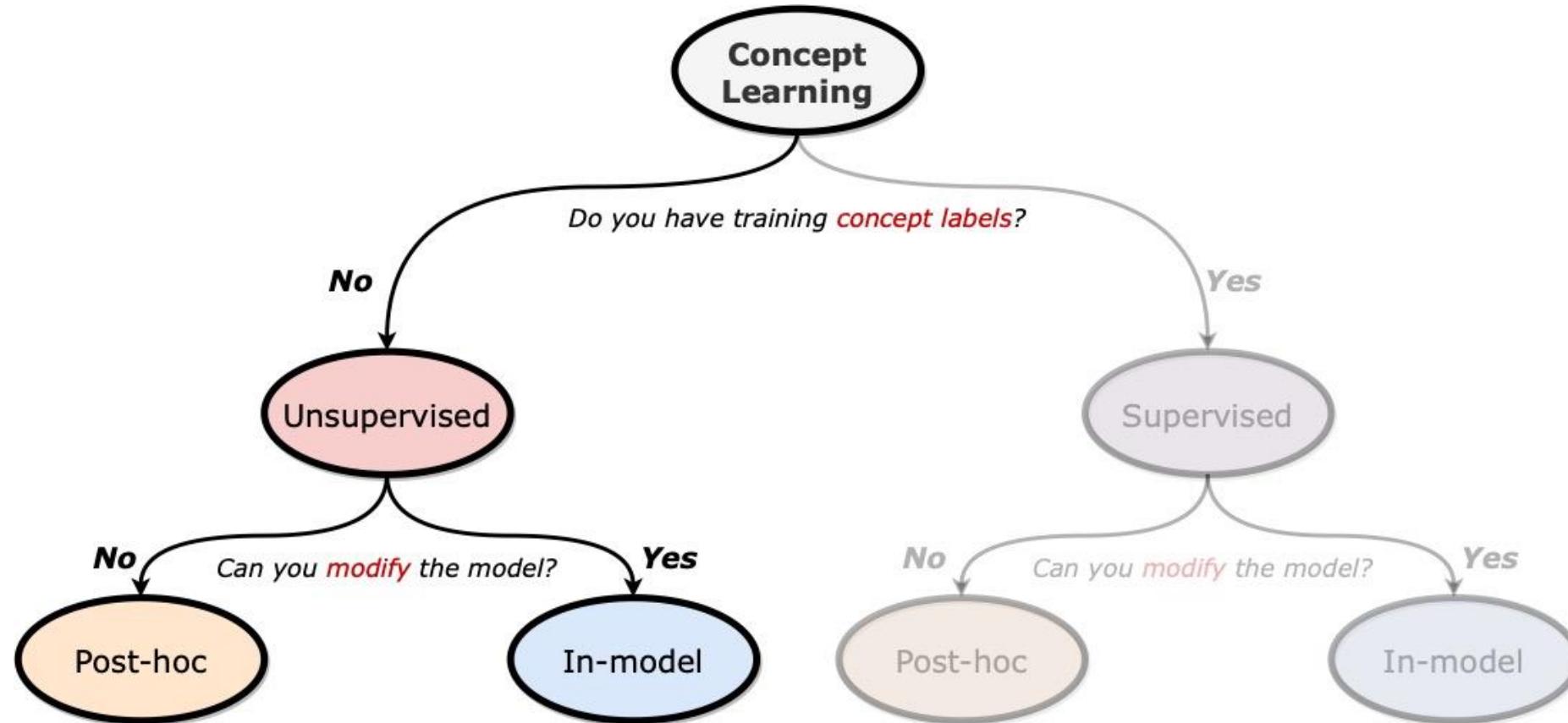
A MAP OF CONCEPT LEARNING



1/3

In the first third of this tutorial, we will discuss supervised concept learning

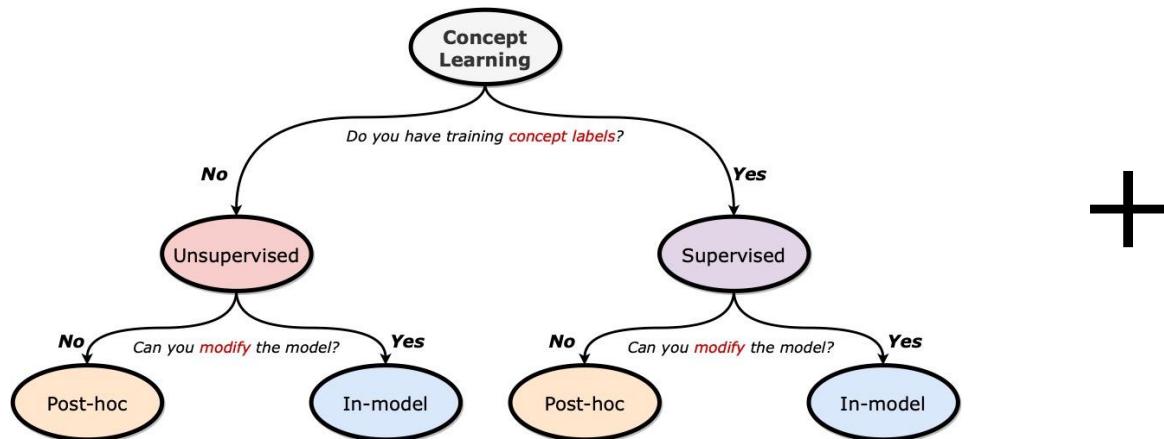
A MAP OF CONCEPT LEARNING



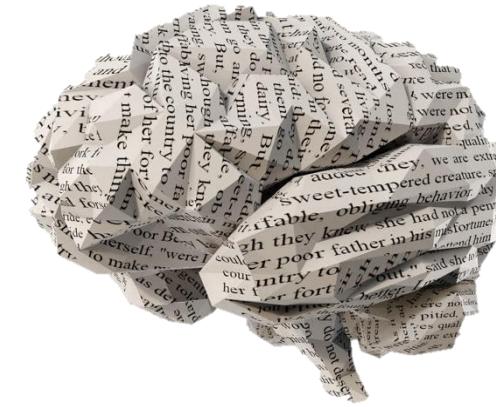
2/3

The second third discusses unsupervised concept learning approaches

A MAP OF CONCEPT LEARNING



Concepts



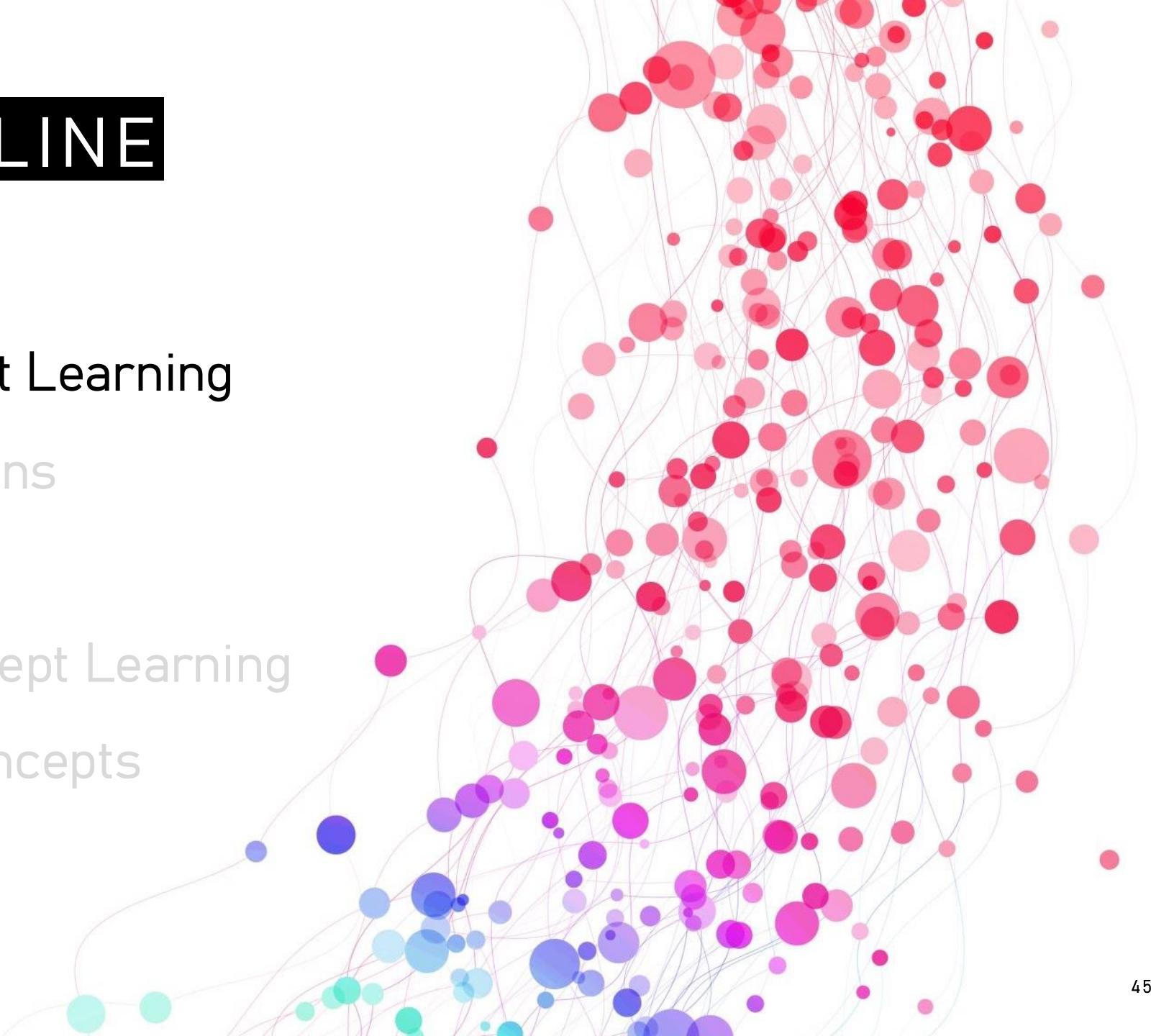
Reasoning

3/3

Finally, in the last third we discuss applications of CL to symbolic reasoning

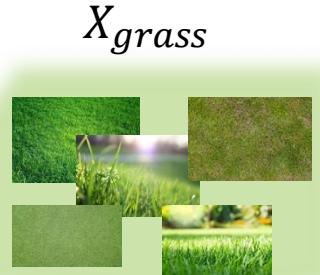
TUTORIAL OUTLINE

1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



DIFFERENT LEVELS OF SUPERVISION

"Supervised" is a **loaded term**. In this tutorial's context,
"**supervised**" means a method has **access to "concept" labels**



OR



$$\begin{pmatrix} \mathbf{c}^{(i)} \\ \text{lays eggs} \\ \text{has scales} \\ \text{has wings} \\ \text{eats only plants} \\ \text{eats meat} \\ \text{black wings} \\ \text{is colorful} \\ \text{has teeth} \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

Sparse sets of images containing a concept

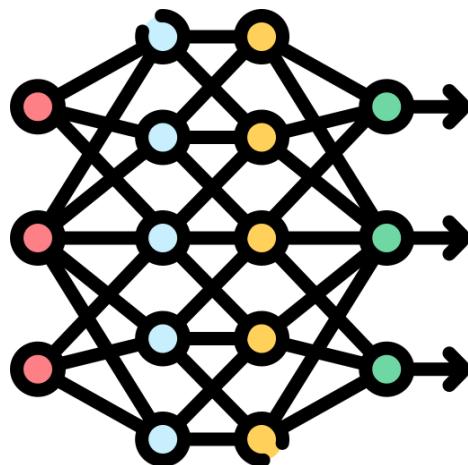
Dense binary vector annotations

These labels could come **besides** other downstream **task "labels"**

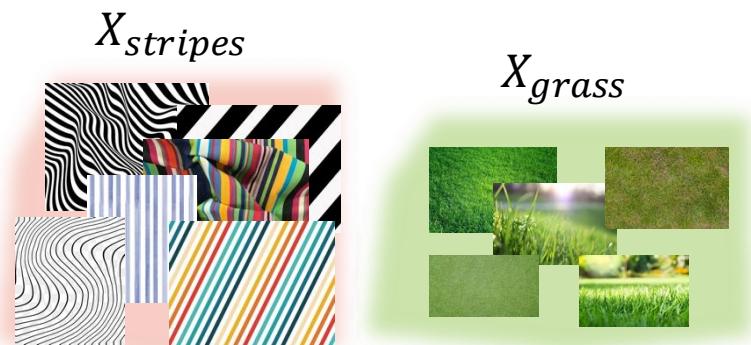
POST-HOC CONCEPT LEARNING

We will start by looking into supervised **post-hoc** concept learning:

A **trained** DNN $f(\mathbf{x}; \theta)$



Sparse concept annotations



Concept-based Explanations

Local

How important is "stripes" for a prediction?

Global

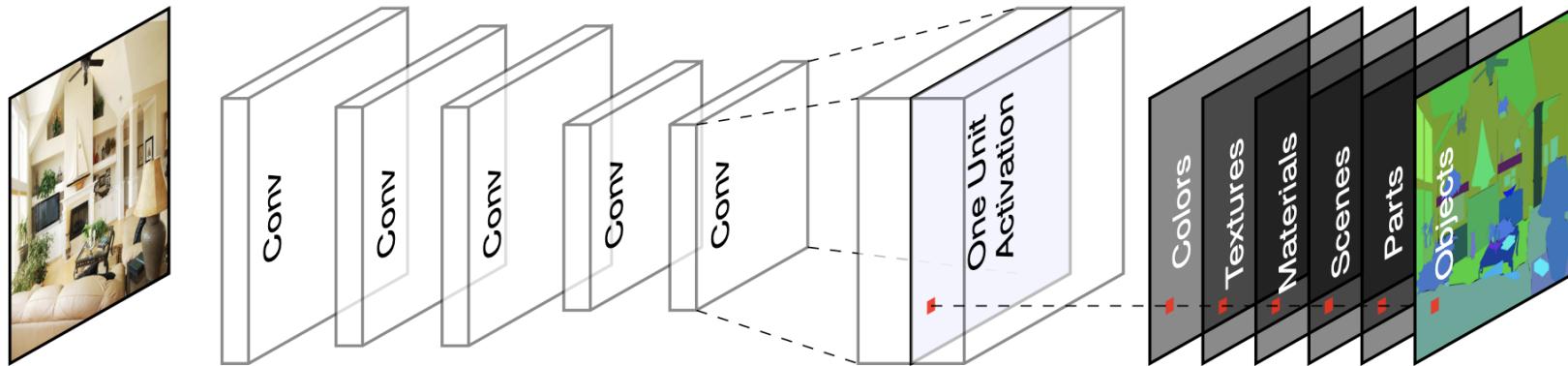
how important is "grass" for the class "cow"?

Given

We Want

WHY SHOULD WE EVEN ATTEMPT THIS?

Evidence suggests DNNs **may** predict based on **concepts**

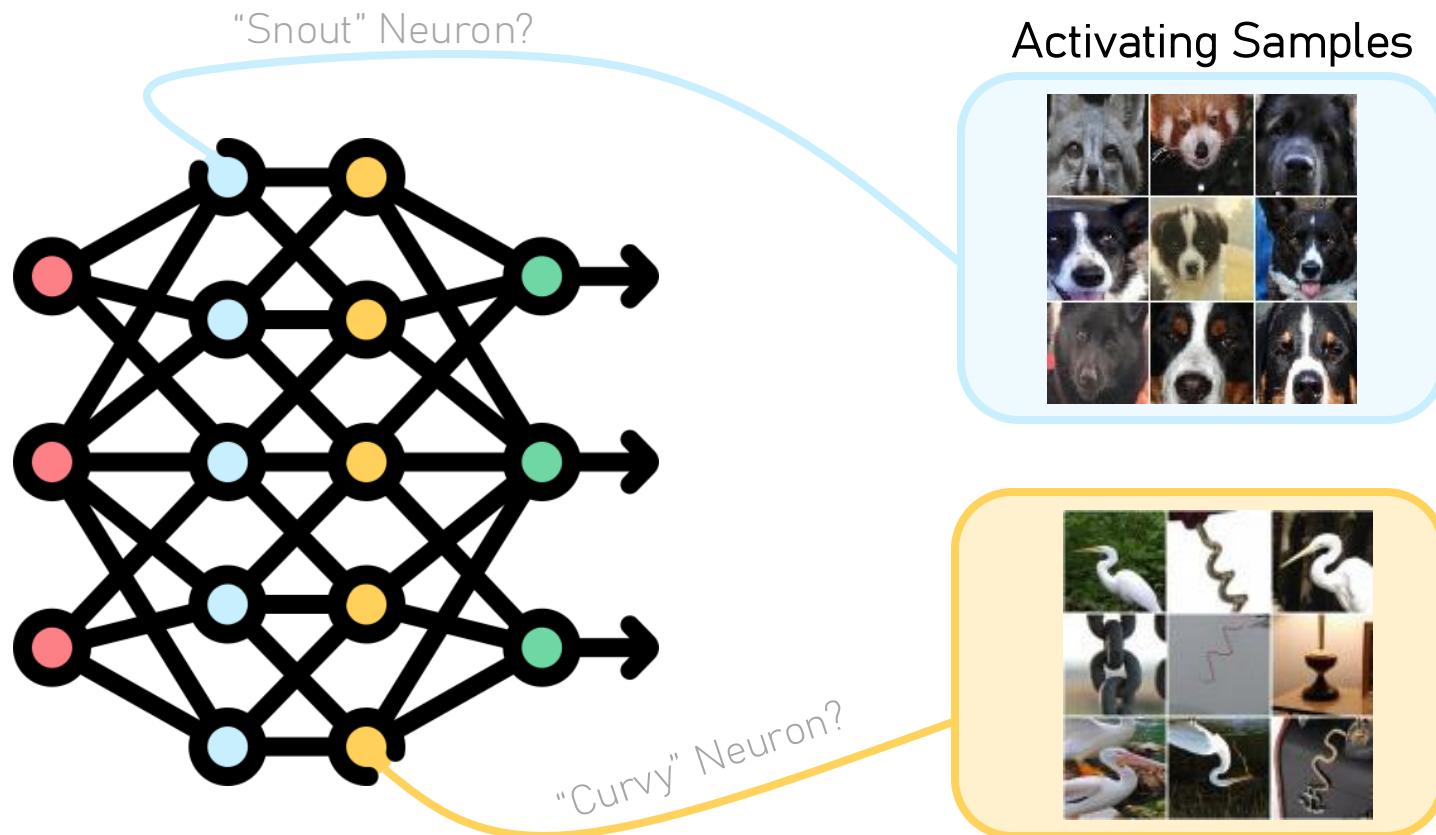


[1] Bau, David, et al. "Network dissection: Quantifying interpretability of deep visual representations." CVPR (2017).



WHY SHOULD WE EVEN ATTEMPT THIS?

Evidence suggests DNNs **may** predict based on **concepts**

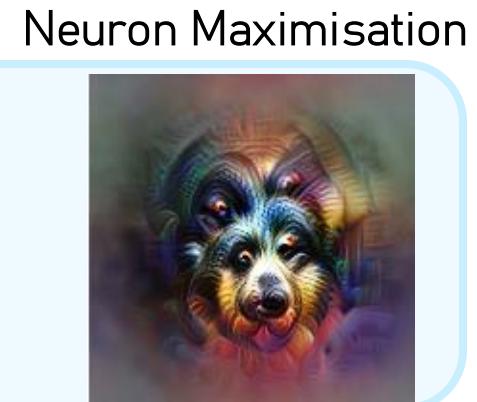
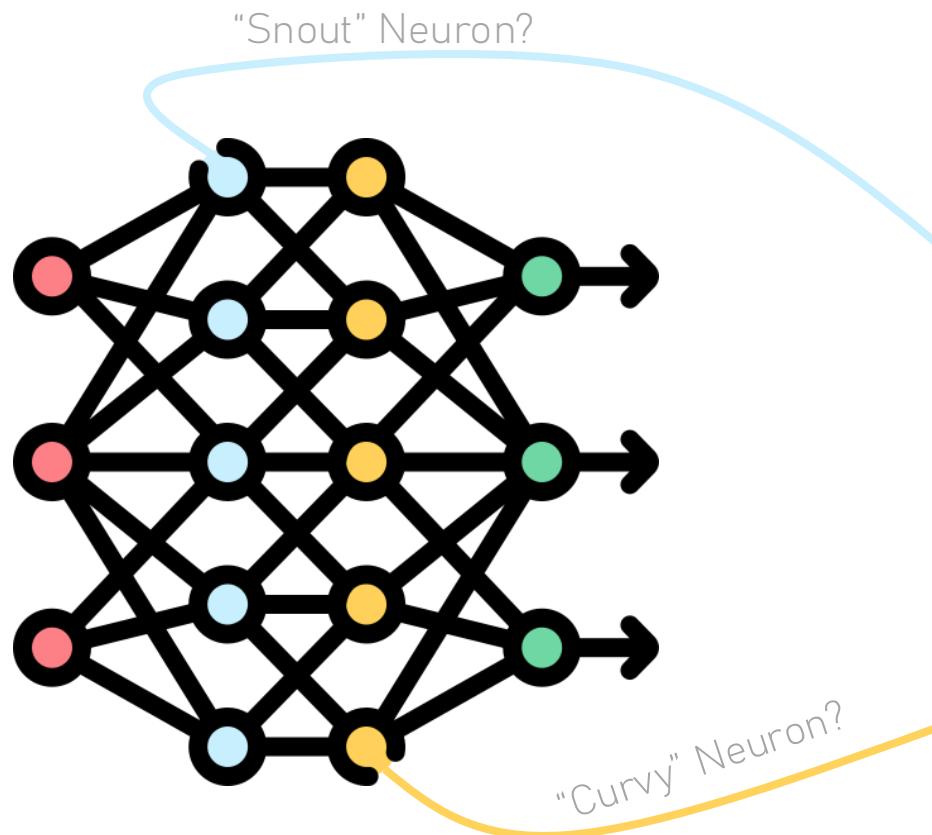


[1] Examples taken from Olah et al. "Feature visualization." Distill 2.11 (2017).



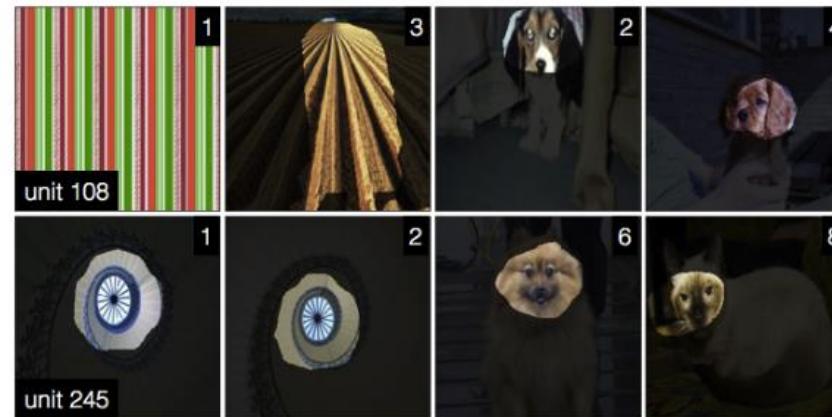
WHY SHOULD WE EVEN ATTEMPT THIS?

Evidence suggests DNNs **may** predict based on **concepts**



CONCEPTS ARE NOT ALWAYS LOCALISED

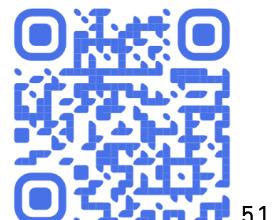
Concepts may **not be always localised** to specific neurons/maps
but they may be **distributed across the DNN's latent space**



The same units appear to represent different concepts

This is sometimes called **Polysemy** (Olah et al., 2020)

[1] Fong et al. "Net2vec: Quantifying and explaining how concepts are encoded by filters in deep neural networks." CVPR (2018).



CONCEPTS ARE NOT ALWAYS LOCALISED

Concepts may **not be always localised** to specific neurons/maps

Could we then try and capture **directions in the latent space** that are **associated with known concepts?**

unit 245

The same units appear to represent different concepts

This is sometimes called **Polysemy** (Olah et al., 2020)

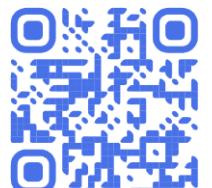
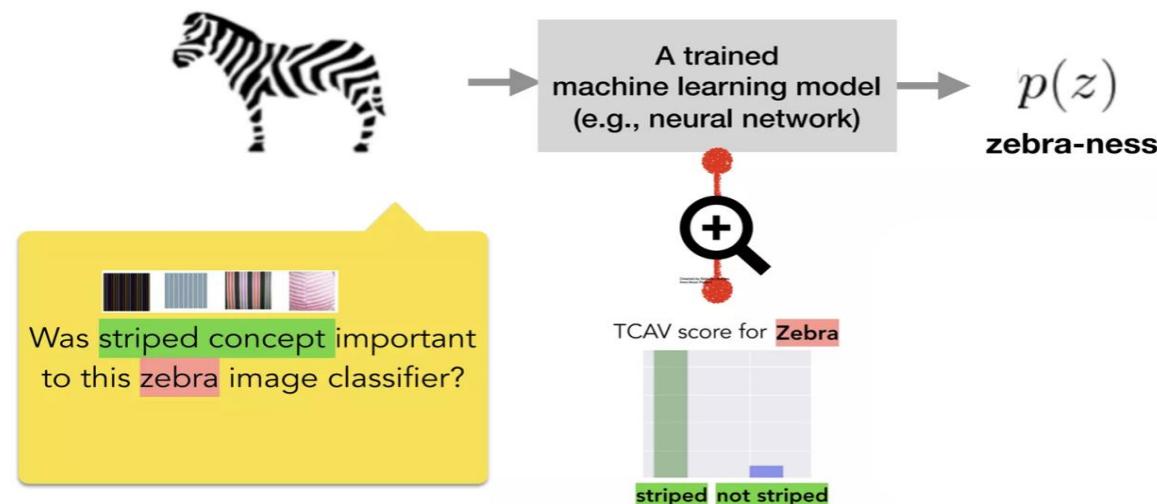
[1] Fong et al. "Net2vec: Quantifying and explaining how concepts are encoded by filters in deep neural networks." CVPR (2018).



TESTING WITH CONCEPT ACTIVATION VECTORS

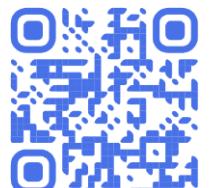
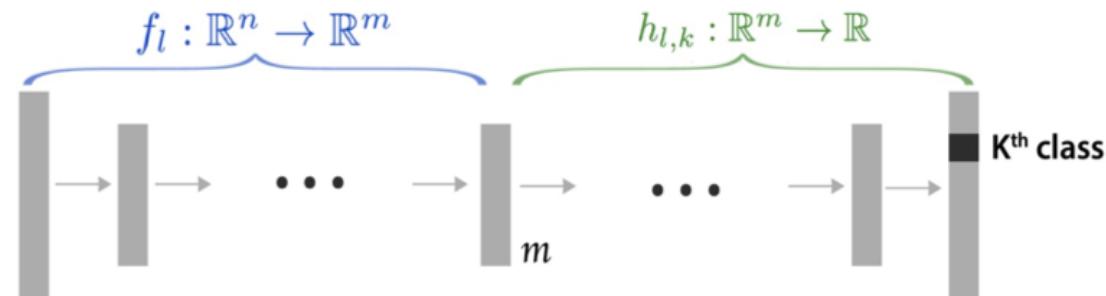
This is the idea behind **T-CAV** (Testing with concept activation vectors)

How **sensitive** is the prediction of zebra is to the **presence of the concept** of "stripes"?



PARTITIONING THE NETWORK

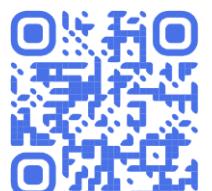
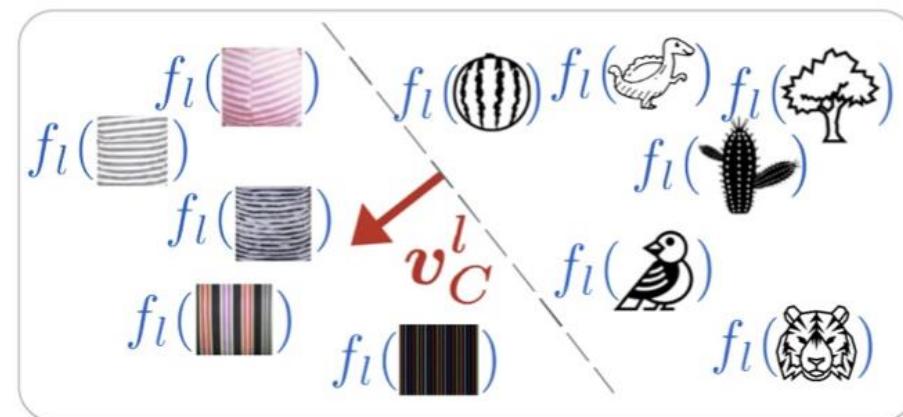
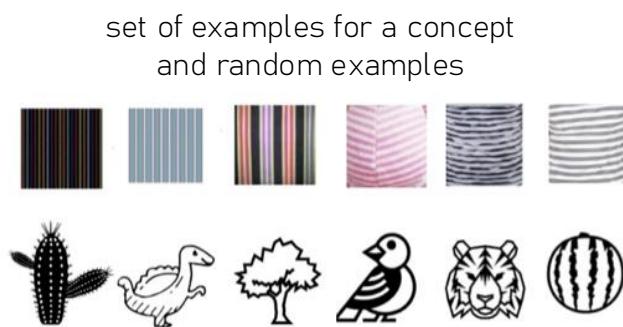
Step 1: Choose an intermediate layer $f_l: \mathbb{R}^n \rightarrow \mathbb{R}^m$ with m neurons



LEARNING CONCEPT ACTIVATION VECTORS

Step 2: Learn the *Concept Activations Vectors* (CAVs)

- Train a linear classifier to distinguish between the activations of concept's examples and random ones
- The CAV is the **vector orthogonal** to the classification boundary v_C^l



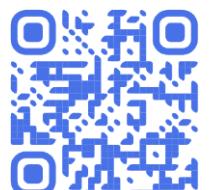
TESTING WITH CAVS (T-CAV)

Step 3: Given a sample \mathbf{x} , construct a **local importance score** $S_{C,k,l}(x)$ indicating how important concept C is for the k -th output label.

We want $S_{C,k,l}(x)$ to capture “**how much would the prediction of class k change if I “increase” concept C in sample \mathbf{x} ?**”

$$S_{C,k,l}(x) = \lim_{\epsilon \rightarrow 0} \frac{h_{l,k}(f_l(\mathbf{x}) + \epsilon \mathbf{v}_C^l) - h_{l,k}(f_l(\mathbf{x}))}{\epsilon}$$

(**Read as:** how much would the prediction of label k change if I take a small step in the direction of concept C ?)

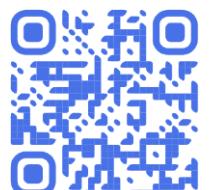


TESTING WITH CAVS (T-CAV)

Step 3: Given a sample \mathbf{x} , construct a **local importance score** $S_{C,k,l}(\mathbf{x})$ indicating how important concept C is for the k -th output label.

We want $S_{C,k,l}(\mathbf{x})$ to capture “**how much would the prediction of class k change if I “increase” concept C in sample \mathbf{x} ?**”

This is the same as a directional derivative!



TESTING WITH CAVS (T-CAV)

Step 3: Given a sample x , construct a **local importance score** $S_{C,k,l}(x)$ indicating how important concept C is for the k -th output label.

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

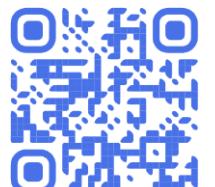
Intermediate representation of
at layer l

Output
function

CAV for concept
 C (e.g., stripes)

The rate of change of output function
as we move in the direction of a
concept from data point 

Intuition: "high directional derivative" = "large positive change in class label if we 'increase' C in input x "



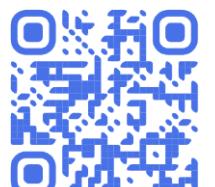
TESTING WITH CAVS (T-CAV)

Step 4: Get a **global importance score** (T-CAV) for each concept by **combining the local sensitivities** of samples in an **evaluation set**:

The **T-CAV score** is the fraction of samples with label \mathbf{k} that are positively influenced by concept C

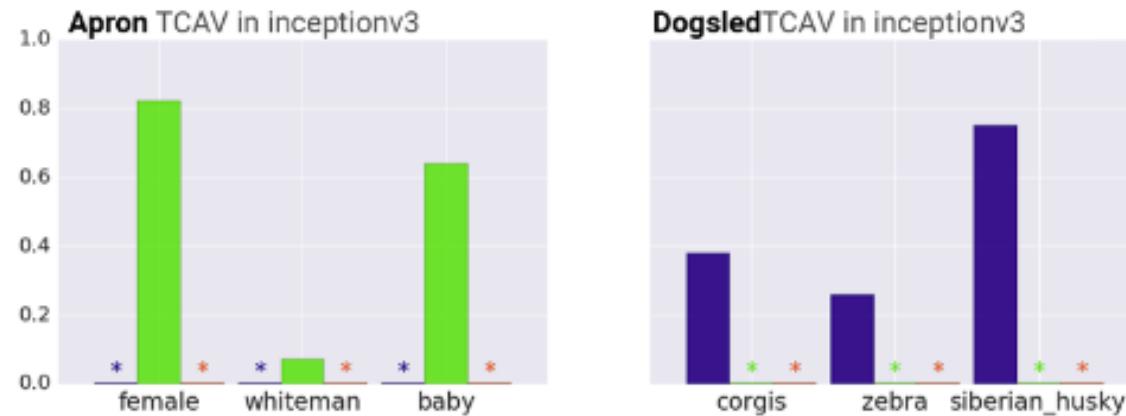
$$\text{TCAV}_{Q_{C,k,l}} = \frac{|\{x \in X_k : S_{C,k,l}(x) > 0\}|}{|X_k|}$$

X_k : inputs
with label k



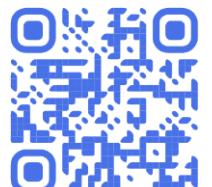
EXAMPLE: DETECTING BIASES WITH T-CAV

You can use T-CAV scores to **explore/identify model biases**



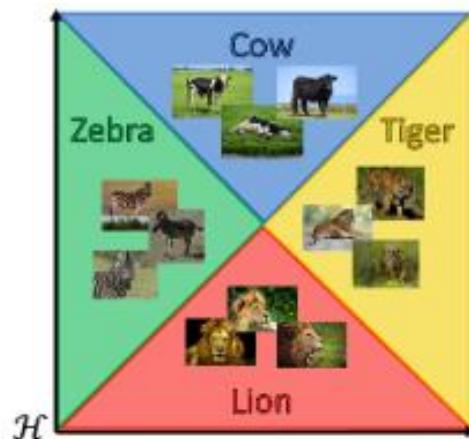
The concept of “female” was found to be significant for predicting the class “Apron”

The concept of “Siberian husky” was found to be significant for predicting the class “Dogsled”



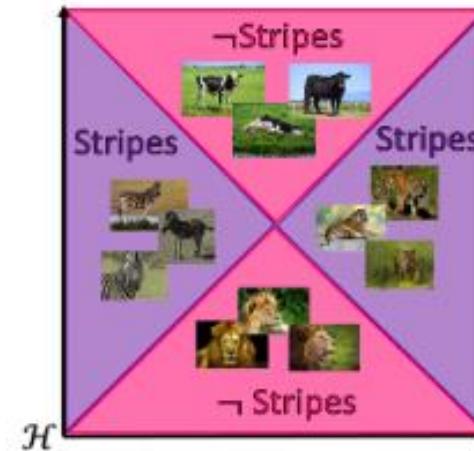
T-CAV LIMITATIONS

Assuming concepts are **linearly separable** is a **strong and unrealistic assumption**



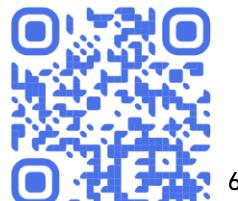
Classes can be linearly separable

VS



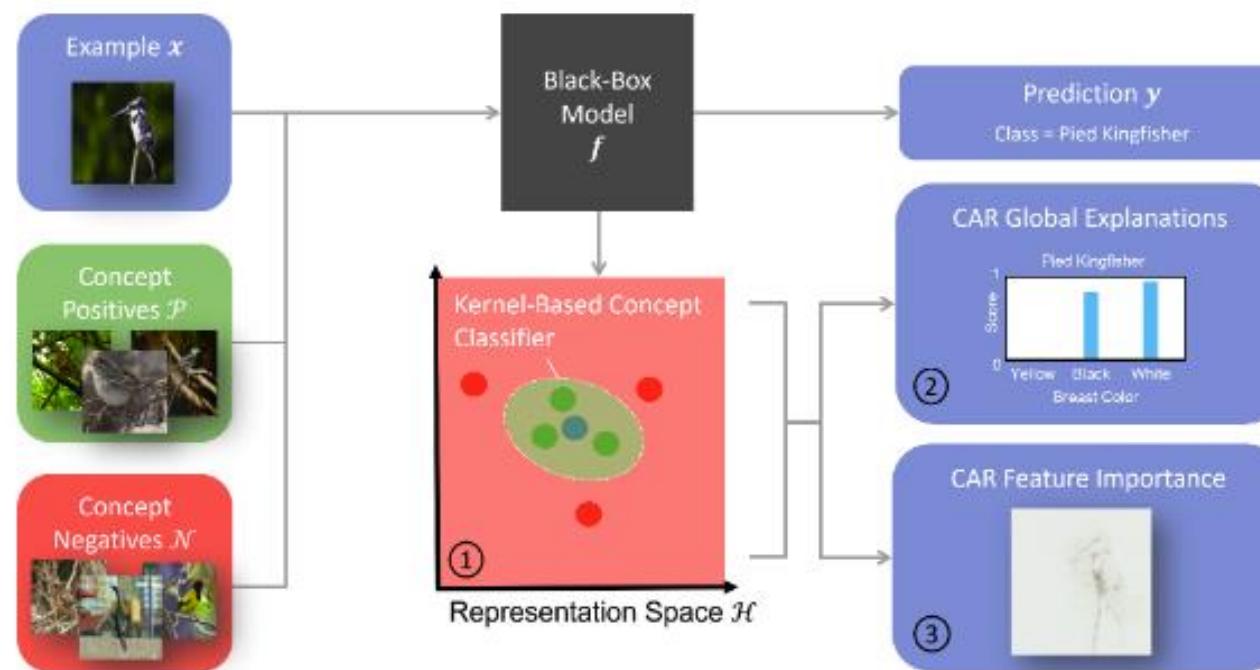
While concepts may not be separable

[1] Crabbé et al. "Concept activation regions: A generalized framework for concept-based explanations." NeurIPS (2022).

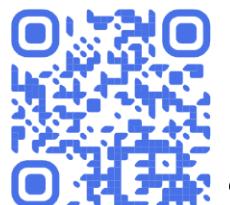


CONCEPT ACTIVATION REGIONS

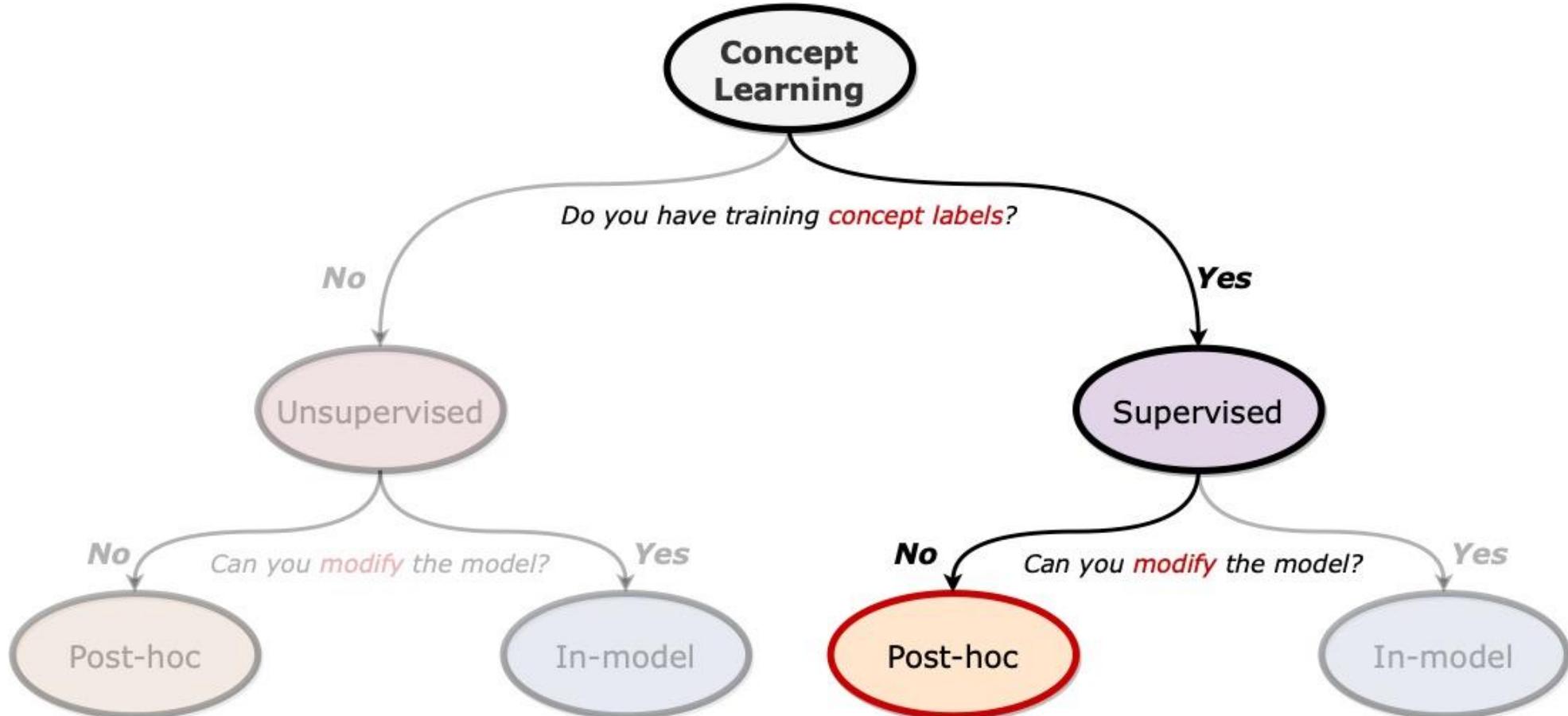
This can be solved by using **kernel methods** to perform our concept probing **on higher-dimensional space where concepts may be separable**



[1] Crabbé et al. "Concept activation regions: A generalized framework for concept-based explanations." NeurIPS (2022).



THE POST-HOC STORY SO FAR



THE POST-HOC STORY SO FAR

Post-hoc methods have a clear set of **important limitations**:

1. They may fail to properly explain a model → potentially **doubling the source of error!**



THE POST-HOC STORY SO FAR

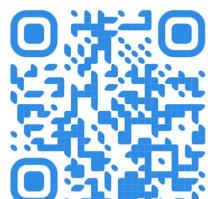
Post-hoc methods have a clear set of **important limitations**:

1. They may fail to properly explain a model → potentially **doubling the source of error!**



In fact, **these methods often disagree** with each other [1]

[1] Krishna et al. "The disagreement problem in explainable machine learning: A practitioner's perspective." TMLR (2022).



THE POST-HOC STORY SO FAR

Post-hoc methods have a clear set of **important limitations**:

2. Explanations are prone to **confirmation bias** [1]

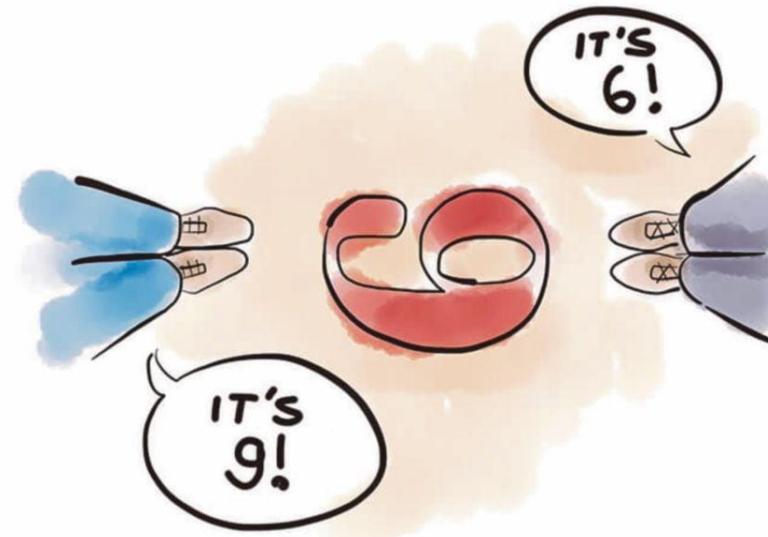


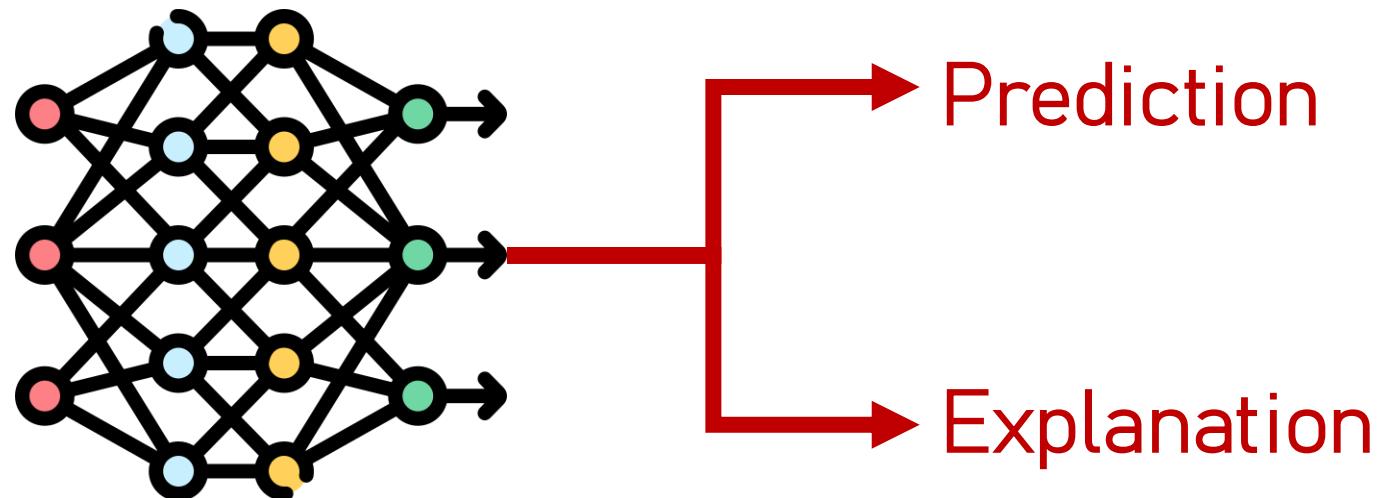
Image taken from "Confirmation Bias and the new Malaysia" by Datuk Steven Wong (New Straits Times)

[1] Bertrand et al. "How cognitive biases affect XAI-assisted decision-making: A systematic review." AIES (2022)

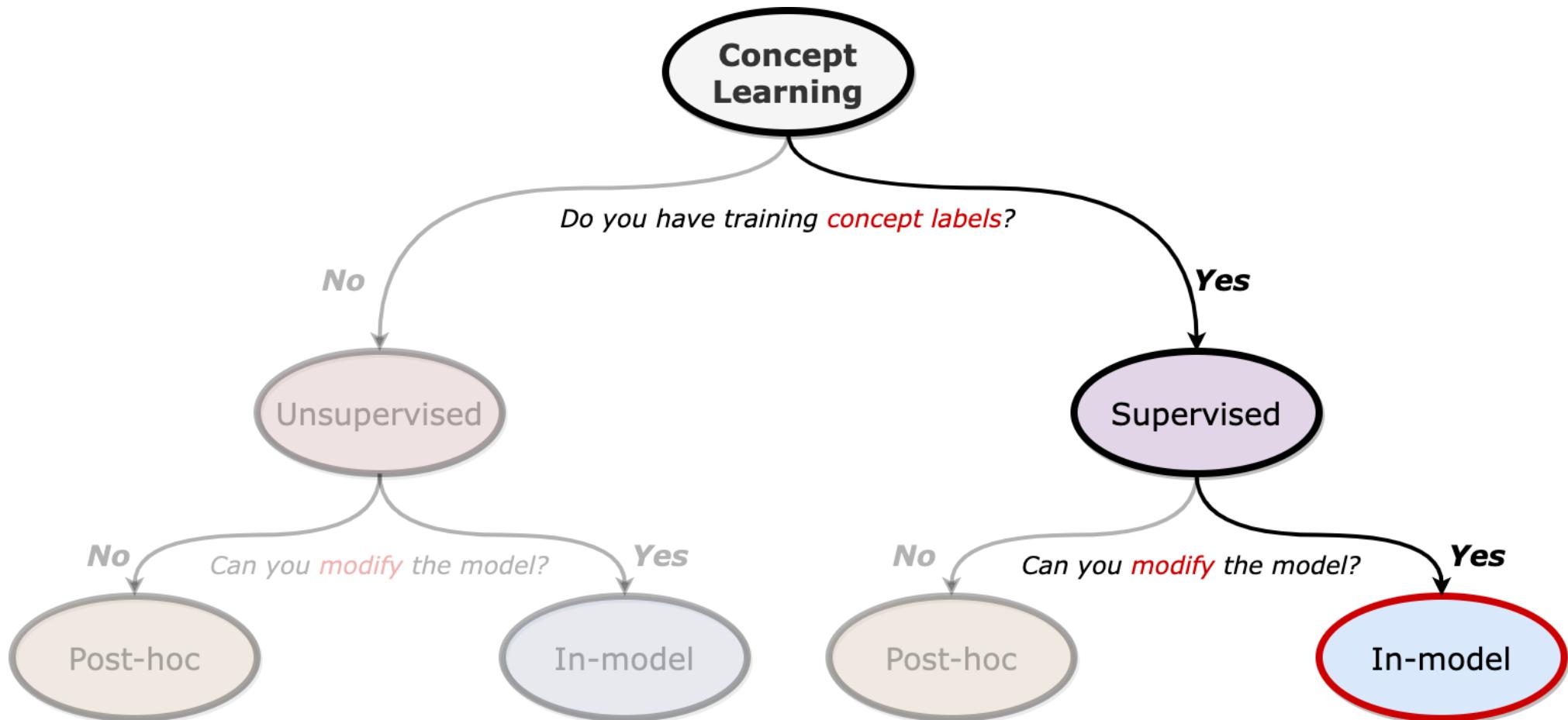


GOING IN-MODEL

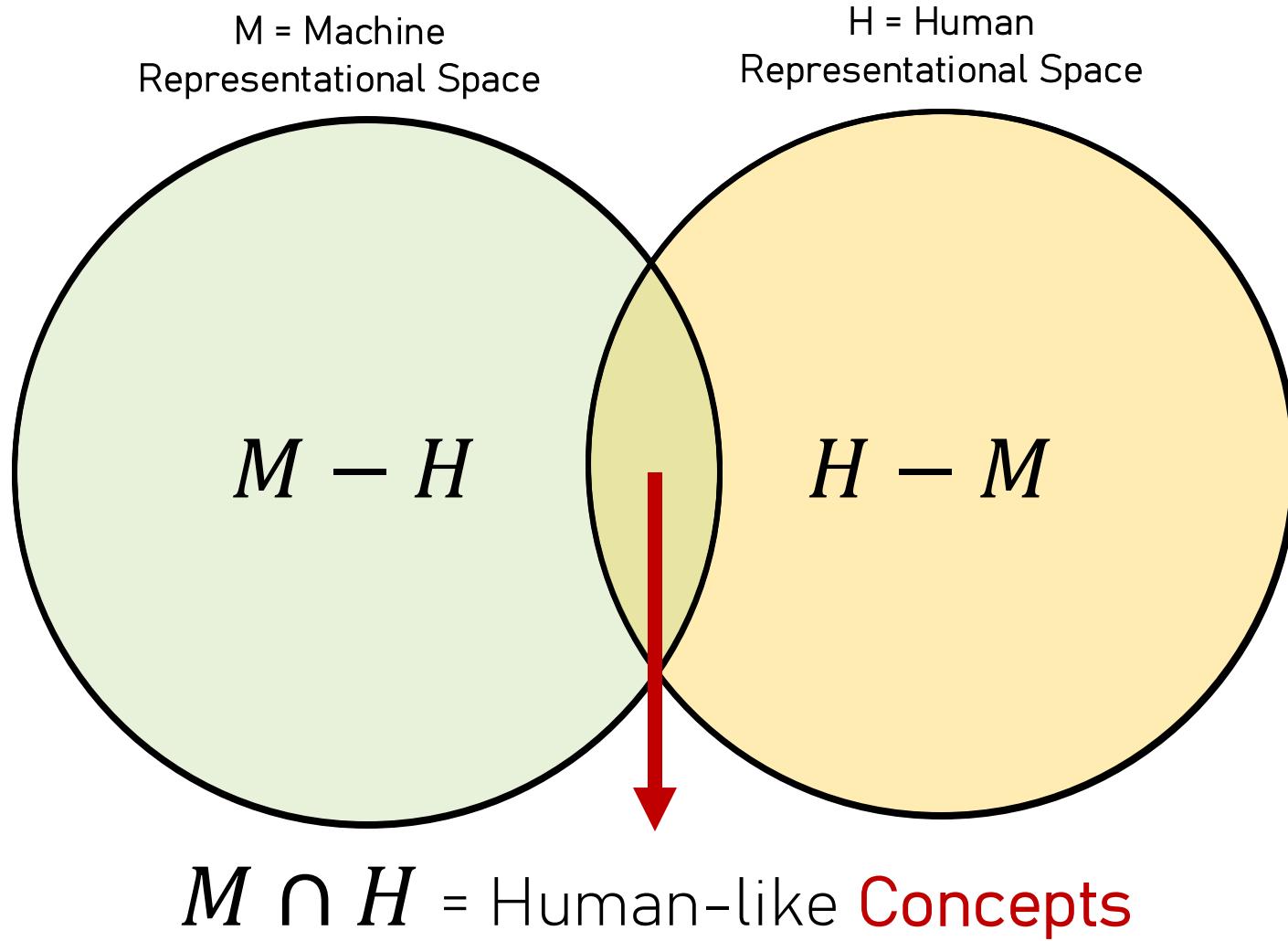
Rather than explaining an already trained model, **let the model explain itself!**



GOING IN-MODEL



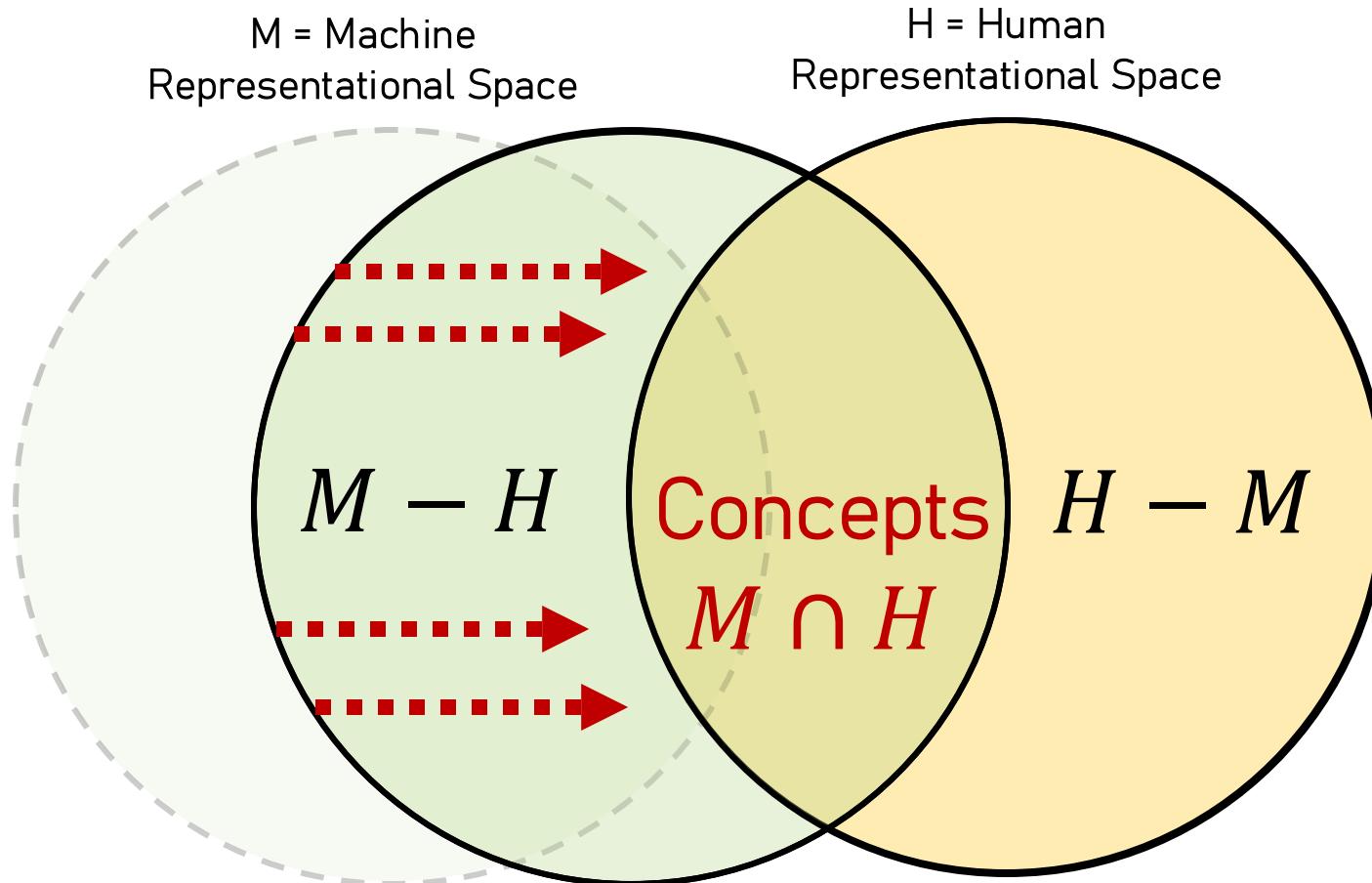
ALIGNING MACHINES AND HUMANS



[1] Inspired by Schut et al. "Bridging the human-ai knowledge gap: Concept discovery and transfer in alphazero." arXiv (2023).



ALIGNING MACHINES AND HUMANS

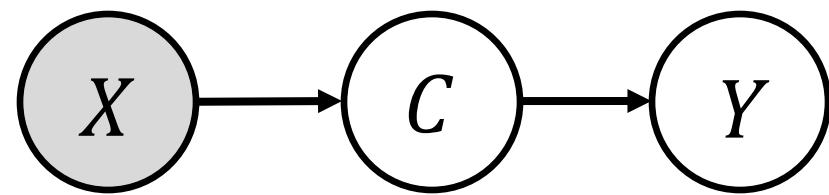


[1] Inspired by Schut et al. "Bridging the human-ai knowledge gap: Concept discovery and transfer in alphazero." arXiv (2023).



CONCEPT-BASED REASONING

Concept-based reasoning can be framed as a **Concept Bottleneck Model** [1]



$$P(C, Y | X) = P(C | X)P(Y | C)$$

X = Sample Features

C = Human-interpretable “Concepts”

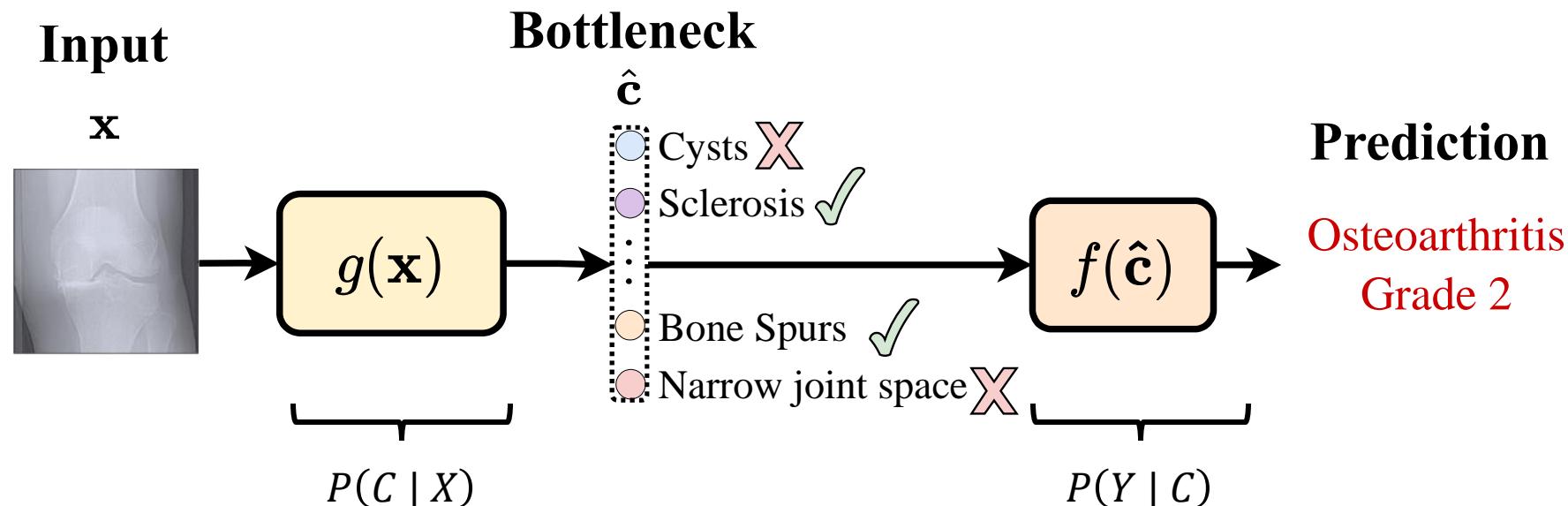
Y = Target Task Labels



CONCEPT BOTTLENECK MODELS (CBMs)

CBMs decompose a DNN into two functions:

1. A *concept encoder* $g(\mathbf{x}) = \hat{\mathbf{c}}$ predicting **concepts from the input features**
2. A *label predictor* $f(\hat{\mathbf{c}}) = \hat{y}$ predicting **task labels from the predicted concepts**



[1] Koh et al. "Concept bottleneck models." International Conference on Machine Learning, PMLR, 2020.



TRAINING A CBM

Given a **concept-annotated dataset** $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{c}^{(i)}, y^{(i)})\}_{i=1}^N$ we can **train** a CBM in three different forms:

(1) Independently

$$\begin{cases} \mathbb{E}_{(\mathbf{x}, \mathbf{c}) \sim \mathcal{D}}[\text{BCE}(g(\mathbf{x}), \mathbf{c})] \\ \mathbb{E}_{(\mathbf{c}, y) \sim \mathcal{D}}[\text{CE}(f(\mathbf{c}), y)] \end{cases}$$

(2) Sequentially

(a) $\mathbb{E}_{(\mathbf{x}, \mathbf{c}) \sim \mathcal{D}}[\text{BCE}(g(\mathbf{x}), \mathbf{c})]$

↓
Freeze g

(b) $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\text{CE}(f(g(\mathbf{x})), y)]$

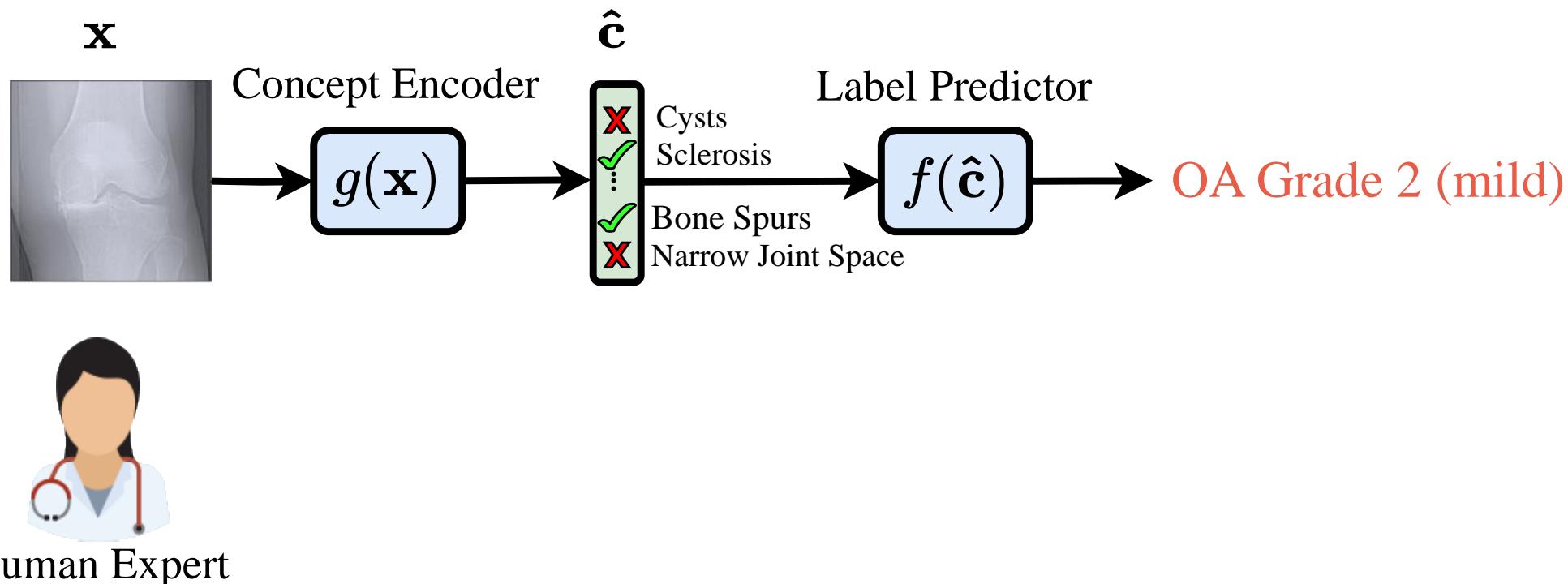
(3) Jointly

$$\mathbb{E}_{(\mathbf{x}, \mathbf{c}, y) \sim \mathcal{D}}[\text{CE}(f(g(\mathbf{x})), y) + \lambda \cdot \text{BCE}(g(\mathbf{x}), \mathbf{c})]$$



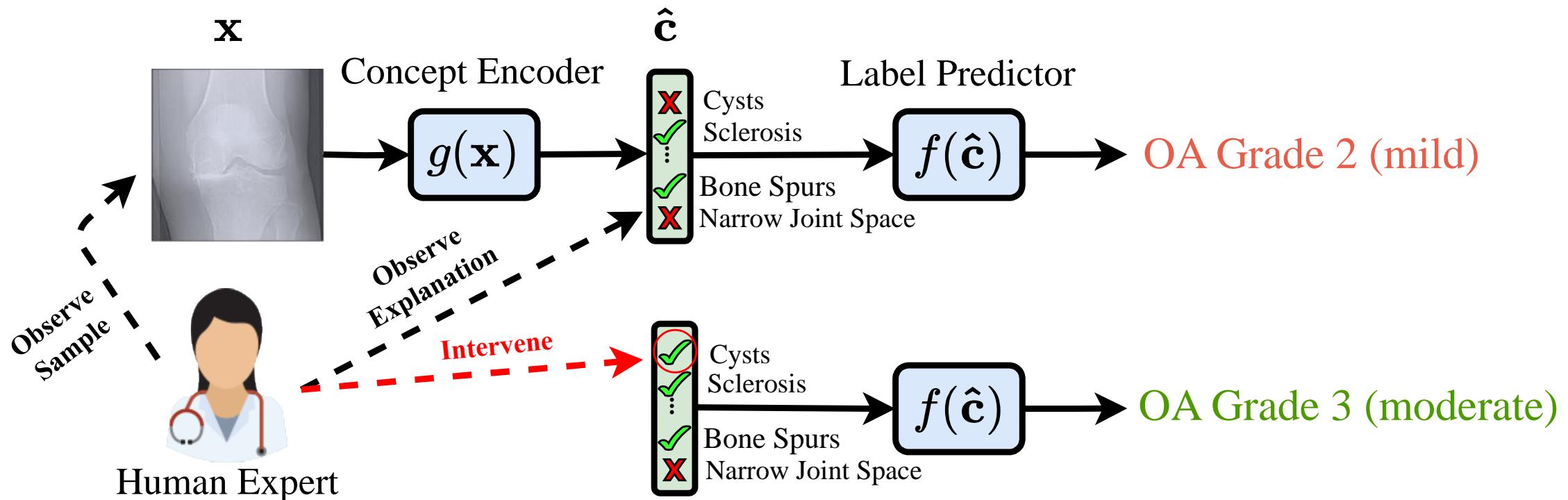
CONCEPT-LEVEL INTERVENTIONS

Concept-based reasoning enables powerful **human-AI interactions**



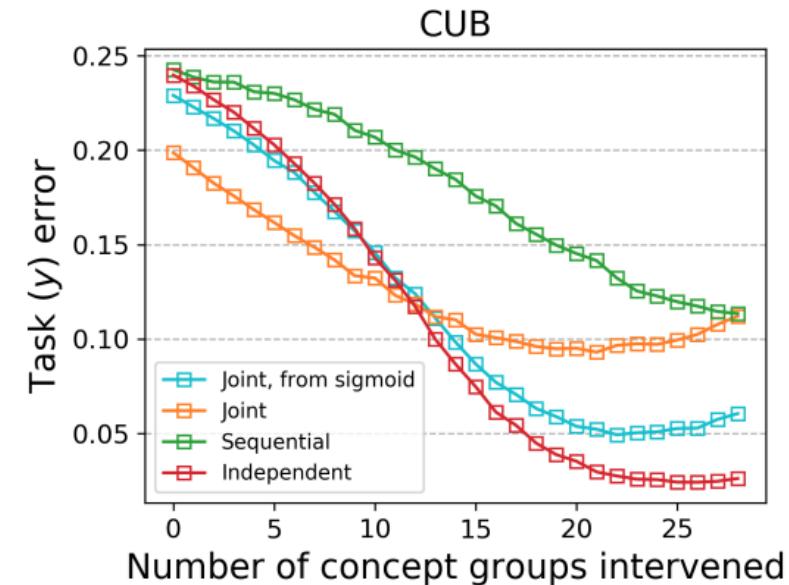
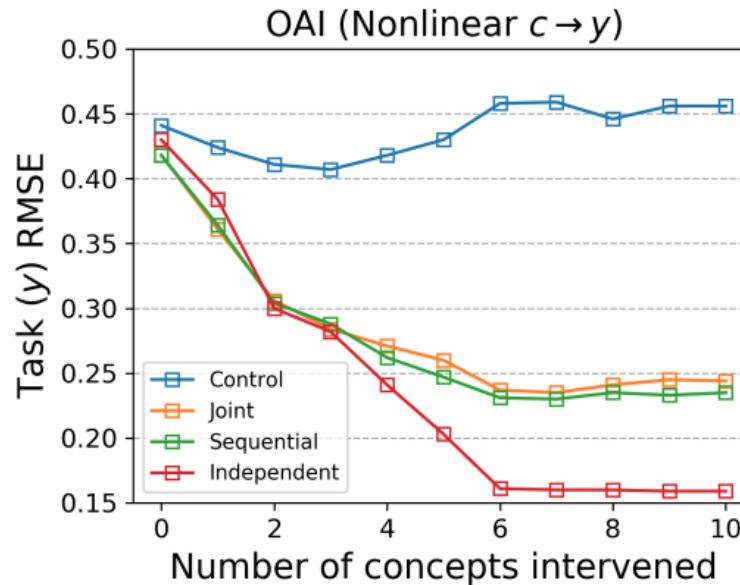
CONCEPT-LEVEL INTERVENTIONS

Concept-based reasoning enables powerful **human-AI interactions**



CONCEPT INTERVENTIONS

As we intervene on more concepts, CBM's **test error goes down!**



ARE CBMS ALL WE NEED?

CBMs are great in a lot of ways:

1. They are simple to understand and provide **high-level explanations**.
2. They enable **test-time interventions** that improve their accuracy.
3. They are very **stable**, expressive and **easy to train**.

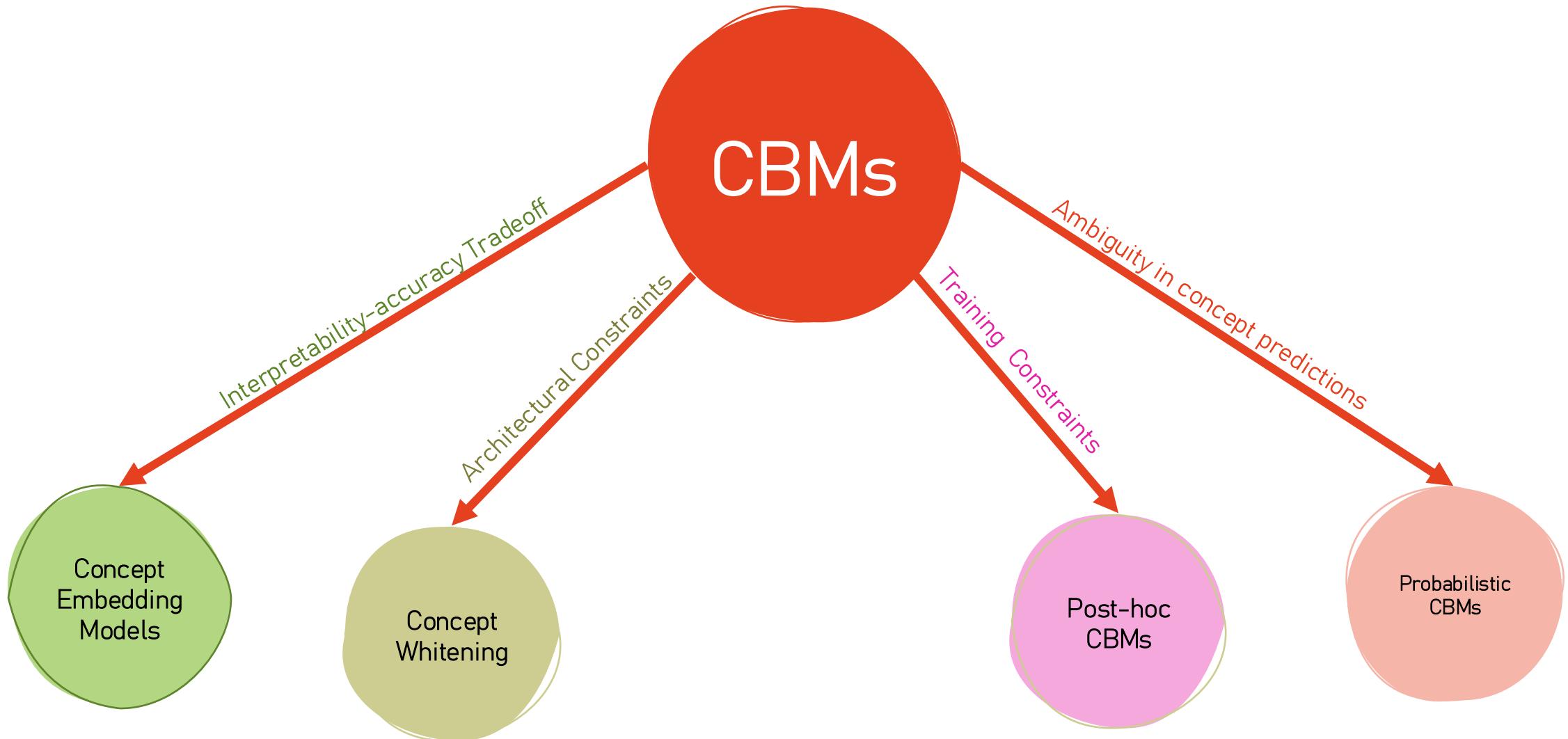
So, are we done?

Short Answer: No

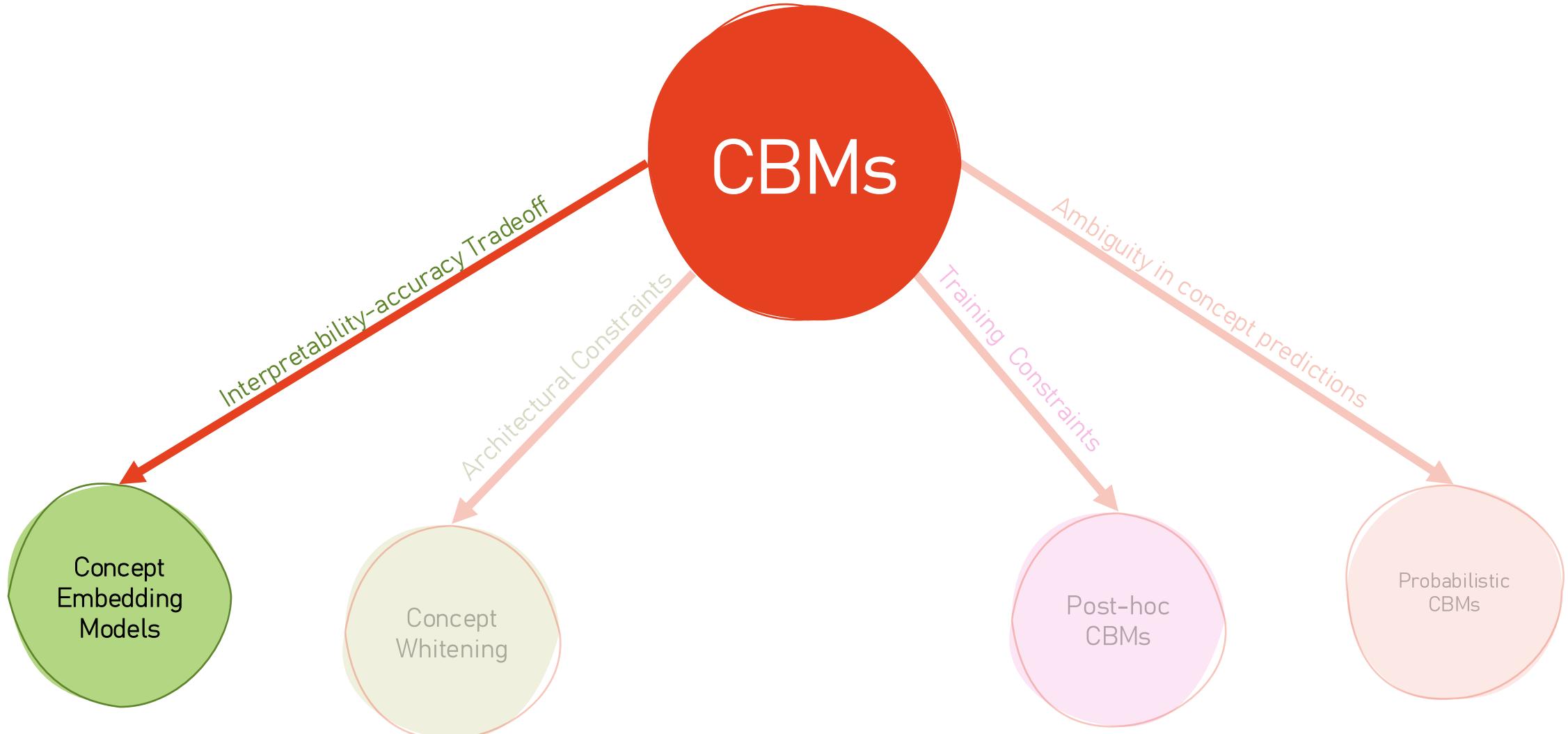
Long Answer:



INTRODUCING CBM'S FRIENDS



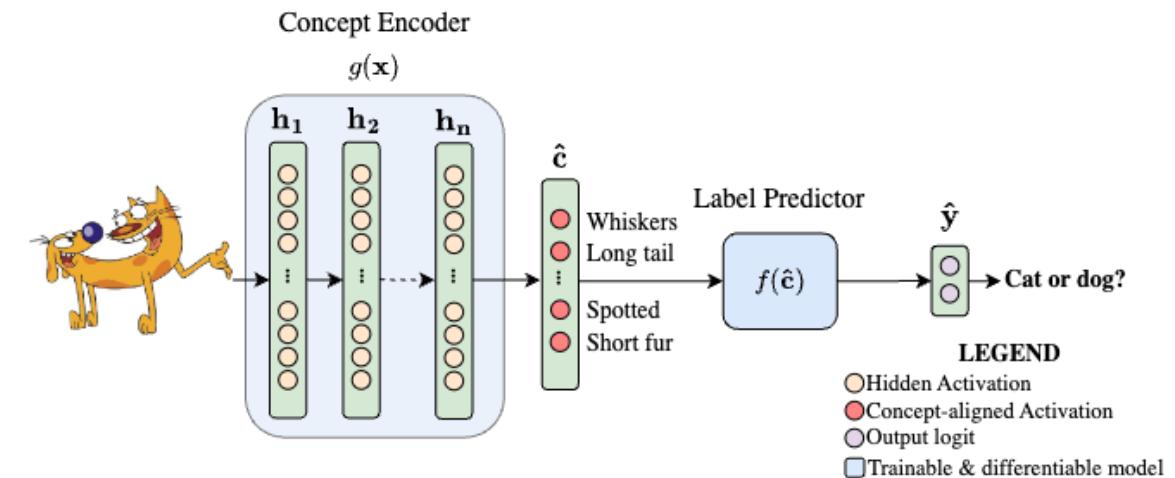
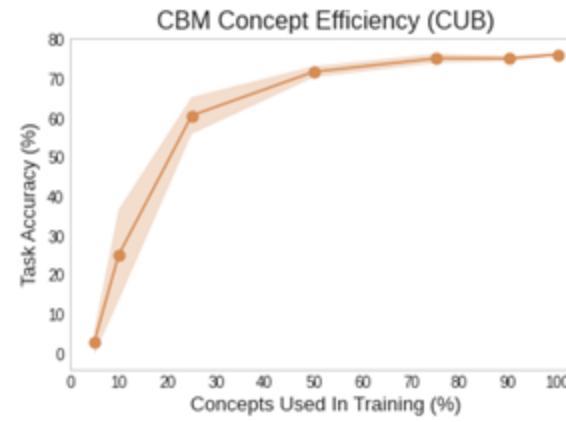
SPEED-DATING WITH CBM'S FRIENDS



CONCEPT EMBEDDING MODELS

Limitation Being Addressed

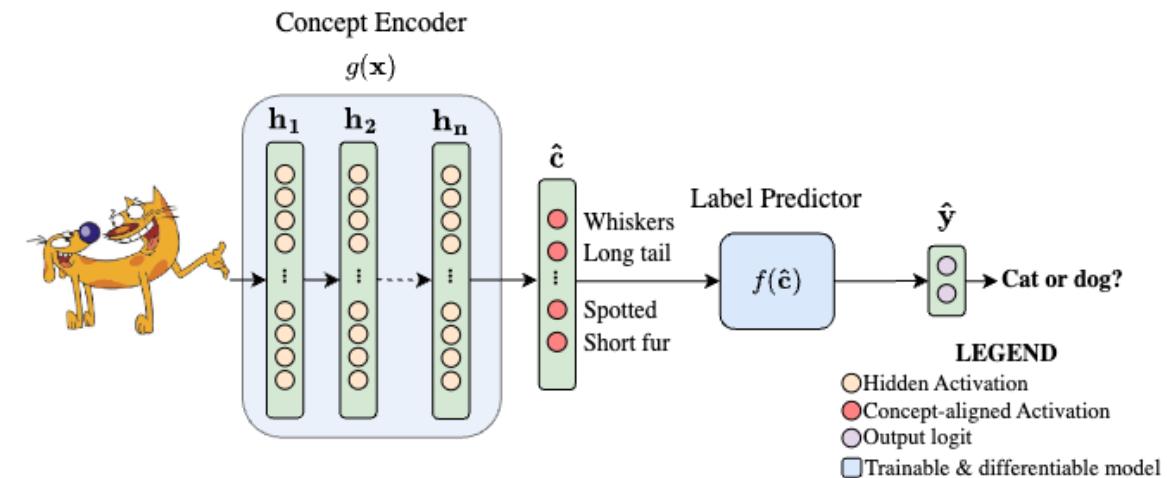
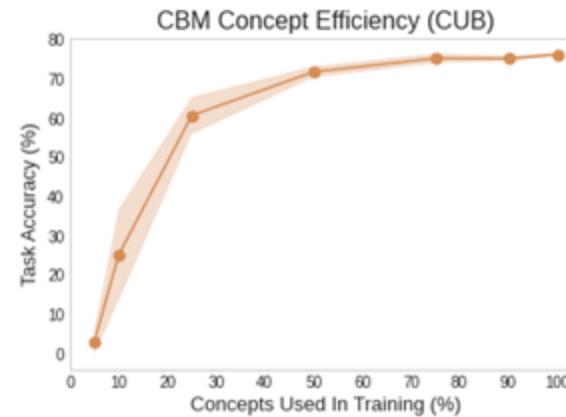
Provided concepts need to be “**complete**” or else we observe a **trade-off**!



CONCEPT EMBEDDING MODELS

Limitation Being Addressed

Provided concepts need to be “**complete**” or else we observe a **trade-off**!



Why can't we just add a bypass from the input to the output?

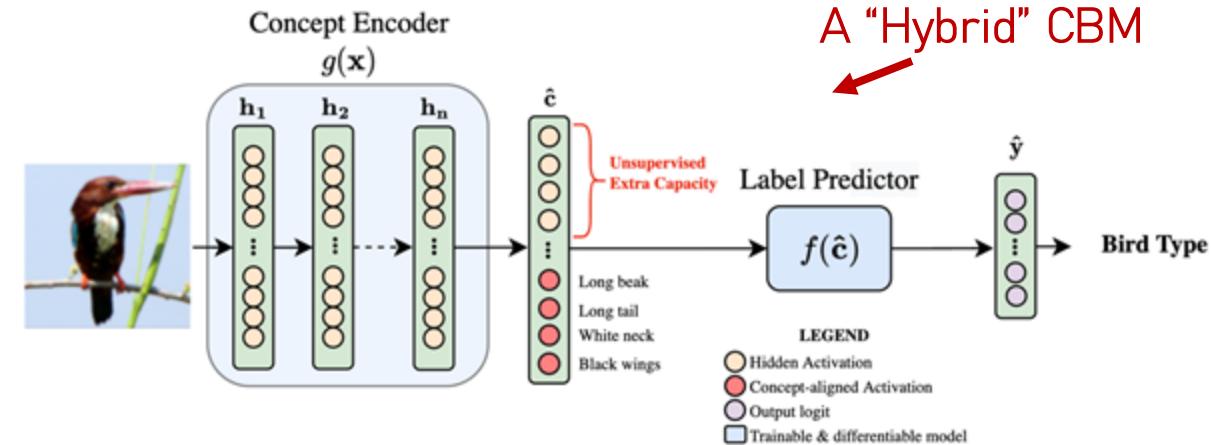
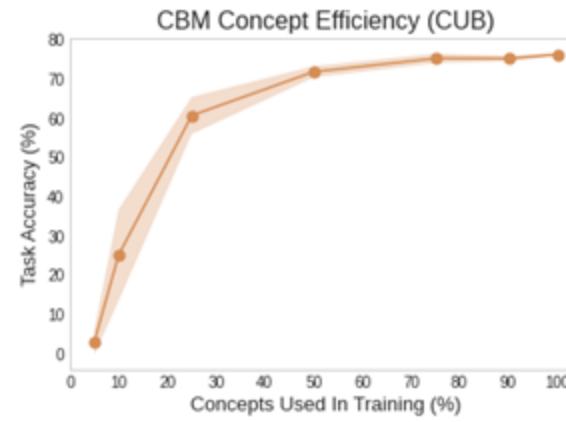
[1] Espinosa Zarlenga, Barbiero et al. "Concept embedding models: Beyond the accuracy-explainability trade-off." NeurIPS (2022)



CONCEPT EMBEDDING MODELS

Limitation Being Addressed

Provided concepts need to be “**complete**” or else we observe a **trade-off**!



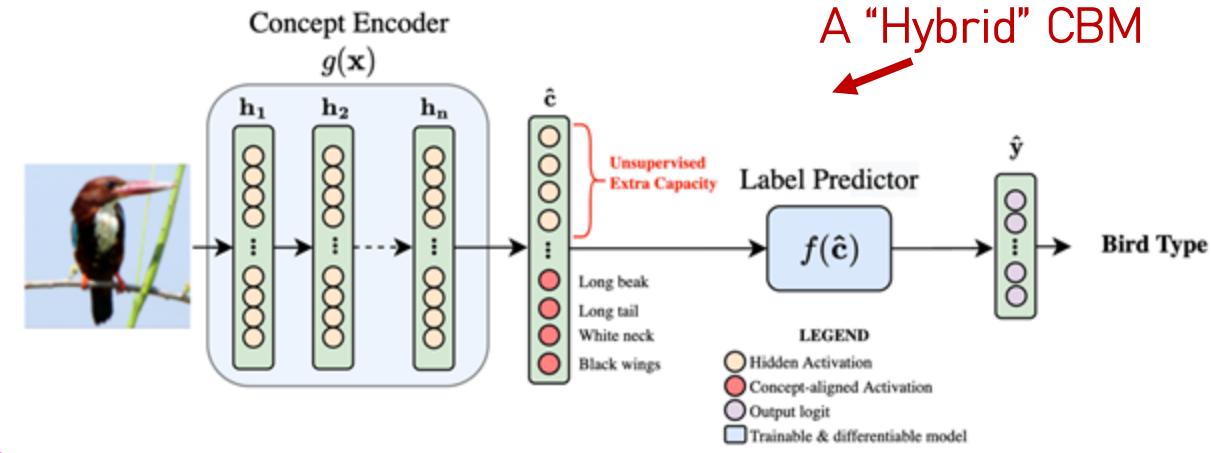
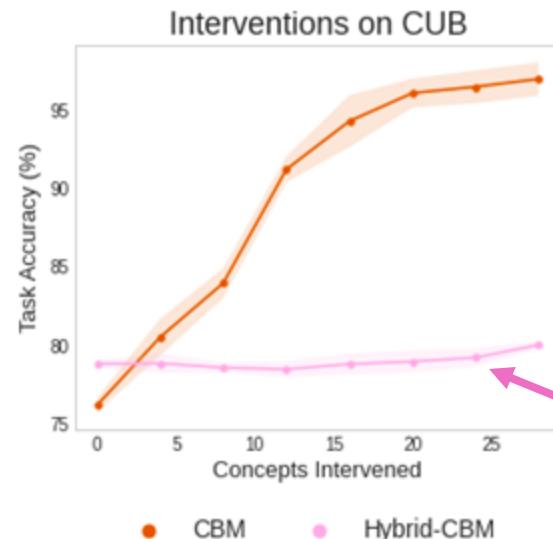
Why can't we just add a bypass from the input to the output?



CONCEPT EMBEDDING MODELS

Limitation Being Addressed

Provided concepts need to be “**complete**” or else we observe a **trade-off**!



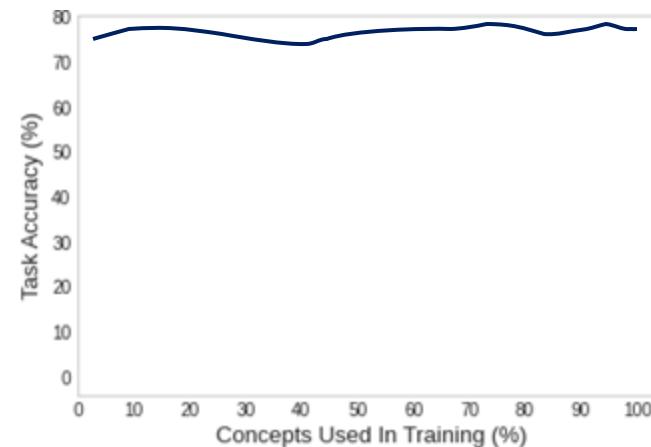
Interventions do not necessarily work with Hybrid CBMs!



CONCEPT EMBEDDING MODELS

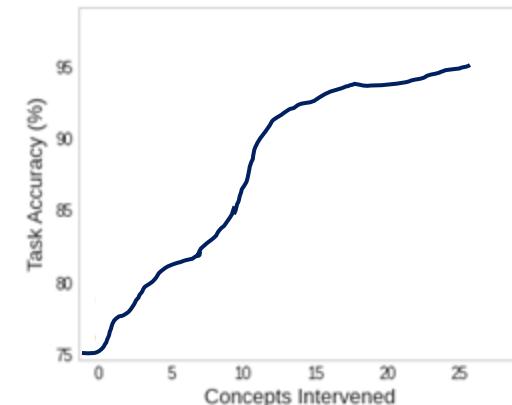
What we want

Similar performance regardless of the number of concepts used during training



"Completeness Agnostic"

Better performance as we intervene in more concepts



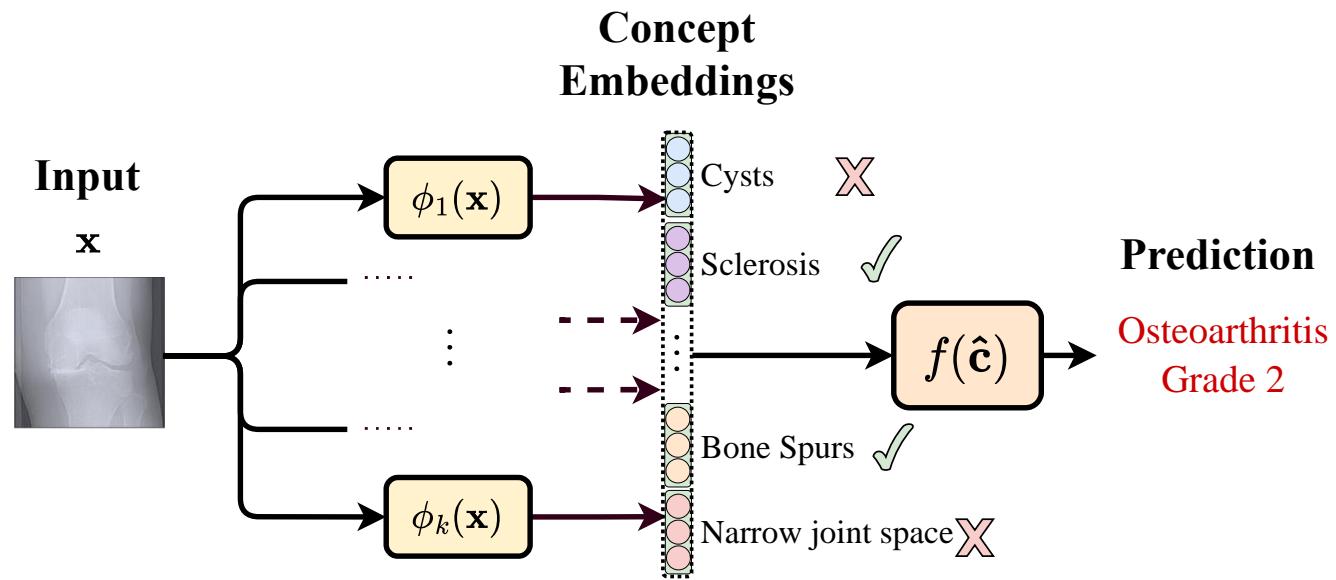
"Intervenable"



CONCEPT EMBEDDING MODELS

Proposed Solution

We can achieve **completeness agnosticism** by extending the **concept representations** to **higher-dimensions**



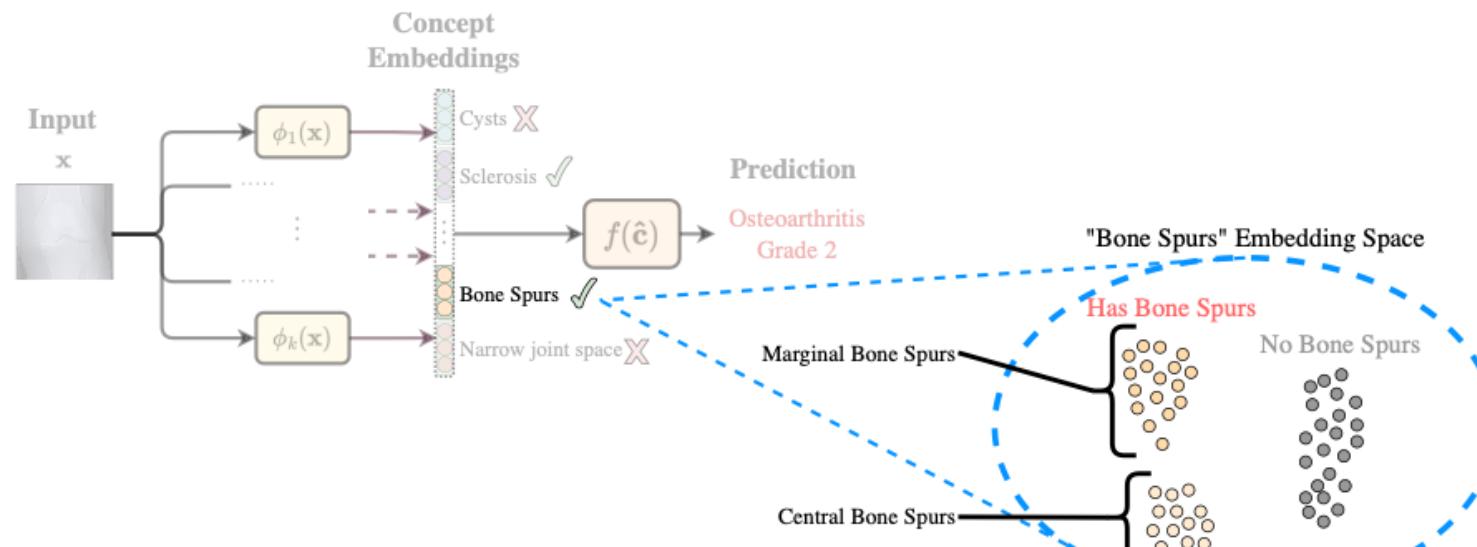
[1] Espinosa Zarlenga, Barbiero et al. "Concept embedding models: Beyond the accuracy-explainability trade-off." NeurIPS (2022)



CONCEPT EMBEDDING MODELS

Proposed Solution

We can achieve **completeness agnosticism** by extending the **concept representations** to **higher-dimensions**



[1] Espinosa Zarlenga, Barbiero et al. "Concept embedding models: Beyond the accuracy-explainability trade-off." NeurIPS (2022)



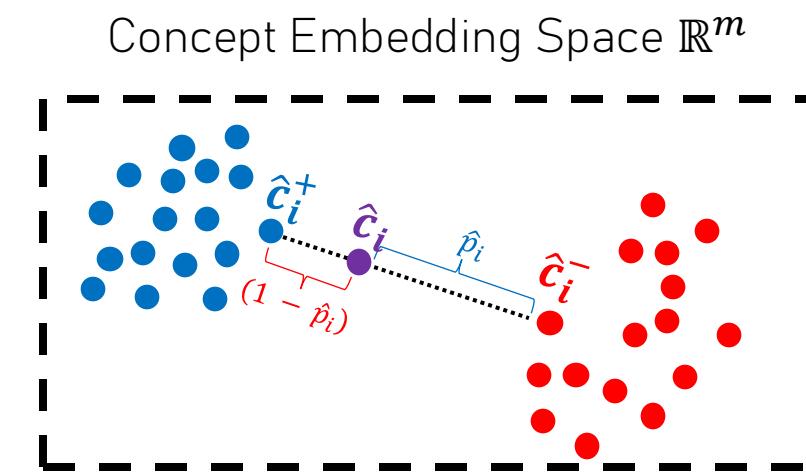
CONCEPT EMBEDDING MODELS

Proposed Solution

We can achieve **intervenability** by decomposing $\hat{\mathbf{c}}_i$ as the **mixture** between two representations $\{\hat{\mathbf{c}}_i^+, \hat{\mathbf{c}}_i^-\}$:

$$\hat{\mathbf{c}}_i := \hat{p}_i \hat{\mathbf{c}}_i^+ + (1 - \hat{p}_i) \hat{\mathbf{c}}_i^-$$

- “Positive” concept embeddings
- “Negative” concept embeddings



CONCEPT EMBEDDING MODELS

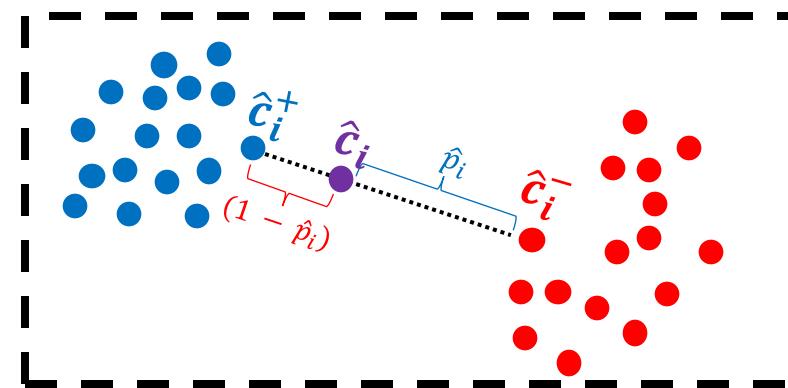
Proposed Solution

We can achieve **intervenability** by decomposing $\hat{\mathbf{c}}_i$ as the **mixture** between two representations $\{\hat{\mathbf{c}}_i^+, \hat{\mathbf{c}}_i^-\}$:

$$\hat{\mathbf{c}}_i := \hat{p}_i \hat{\mathbf{c}}_i^+ + (1 - \hat{p}_i) \hat{\mathbf{c}}_i^-$$

- “Positive” concept embeddings
- “Negative” concept embeddings

Concept Embedding Space \mathbb{R}^m



Determining a concept's activation given $\hat{\mathbf{c}}_i$ then comes down to determining whether $\hat{\mathbf{c}}_i$ comes from $P(\hat{\mathbf{c}}_i^+ | x)$ or $P(\hat{\mathbf{c}}_i^- | x)$

[1] Espinosa Zarlenga, Barbiero et al. "Concept embedding models: Beyond the accuracy-explainability trade-off." NeurIPS (2022)



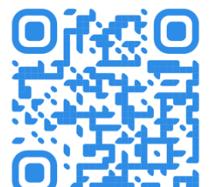
CONCEPT EMBEDDING MODELS

Proposed Solution

You can then **intervene on a concept** by fixing its representation to the embedding corresponding to the ground-truth concept label:

$$\hat{c}_i := \begin{cases} \hat{c}_i^+ & \text{if } c_i = 1 \\ \hat{c}_i^-, & \text{otherwise} \end{cases}$$

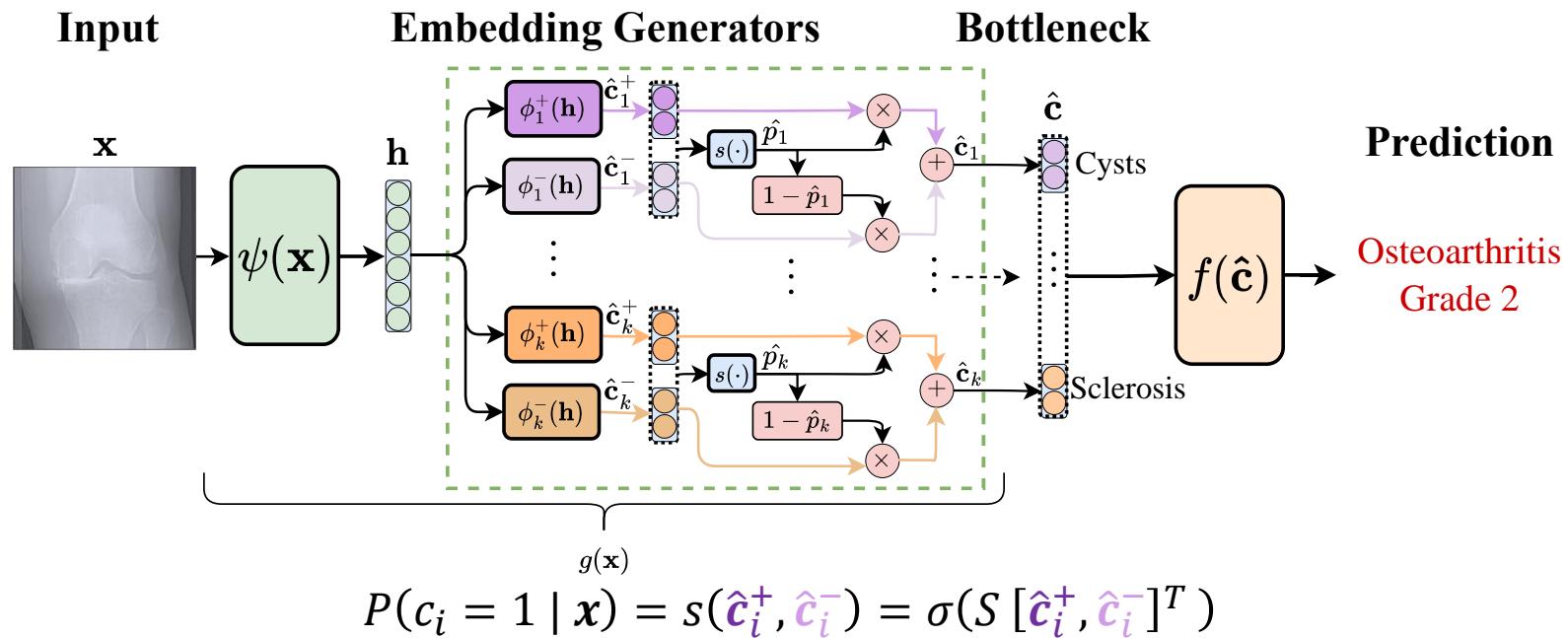
We can randomly do these **interventions at training time** to learn more useful representations!



CONCEPT EMBEDDING MODELS

Proposed Solution

Learn two functions (ϕ_i^+ , ϕ_i^-) mapping \mathbf{x} to a positive $\hat{\mathbf{c}}_i^+$ and a negative embedding $\hat{\mathbf{c}}_i^-$



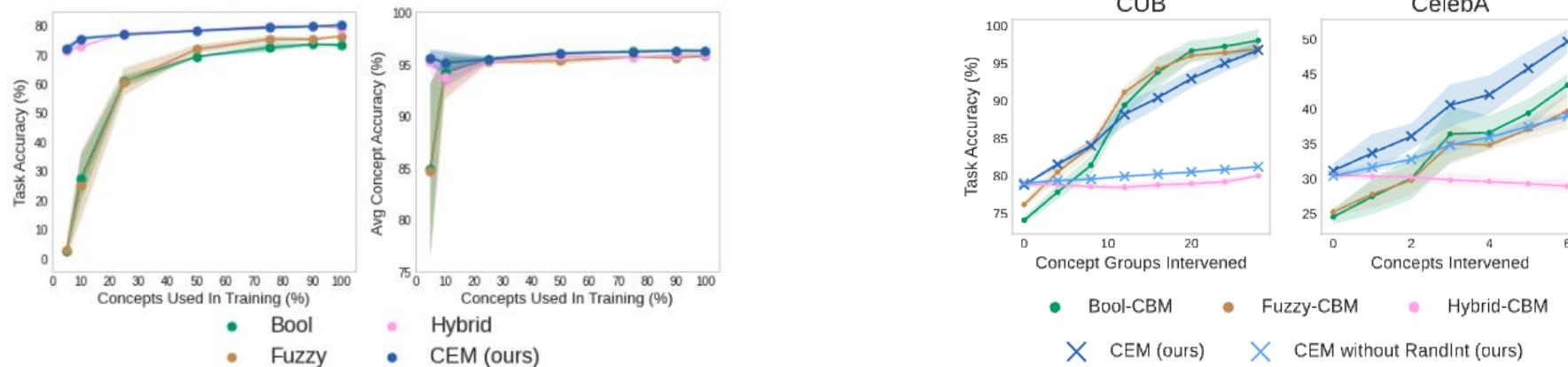
[1] Espinosa Zarlenga, Barbiero et al. "Concept embedding models: Beyond the accuracy-explainability trade-off." NeurIPS (2022)



CONCEPT EMBEDDING MODELS

Proposed Solution

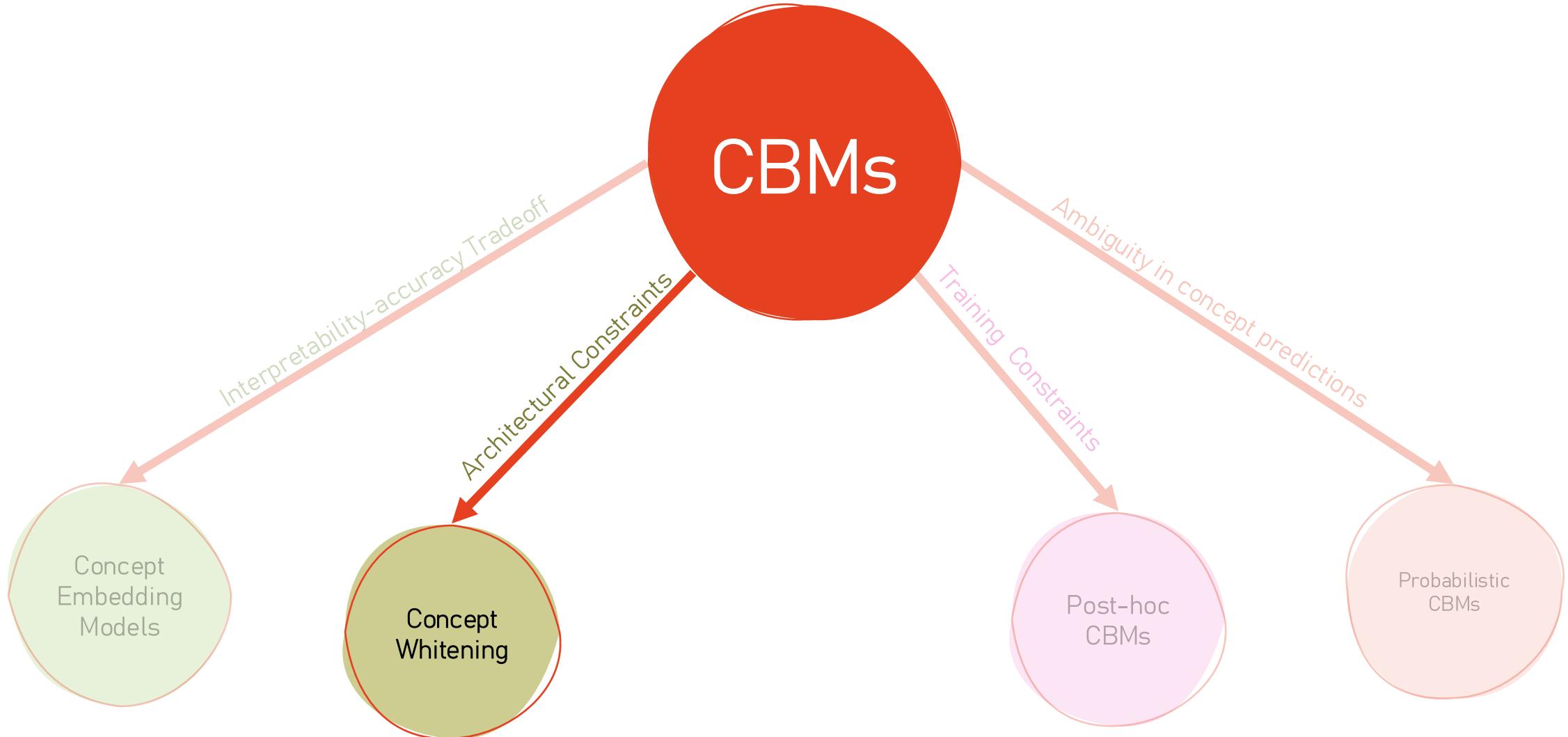
This gives you models that are **completeness agnostic** and **intervenable***



*This is particularly true when, during training, you **randomly intervene**** on a concept with probability p_{int} .

** These sorts of training-time interventions are useful here **only** because by using embeddings, **we can backpropagate gradients to the concept encoder** even when a concept is intervened on (a CBM wouldn't).

SPEED-DATING WITH CBM'S FRIENDS



CONCEPT WHITENING (CW)

Limitation Being Addressed

Training a CBM has **impractical architectural constraints** and **requires all training samples to be concept annotated!**



This limits our ability to exploit powerful pre-trained models

[1] Chen et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence 2.12 (2020).



INTERPRETABLE-BY-DESIGN NEURAL LAYER

Proposed Solution

Design an **interpretable-by-design layer** which we can use to replace an equivalent component in a **pre-trained model** and **quickly fine-tune it** to make it interpretable

We will target the commonly used Batch Normalization (BN) layer

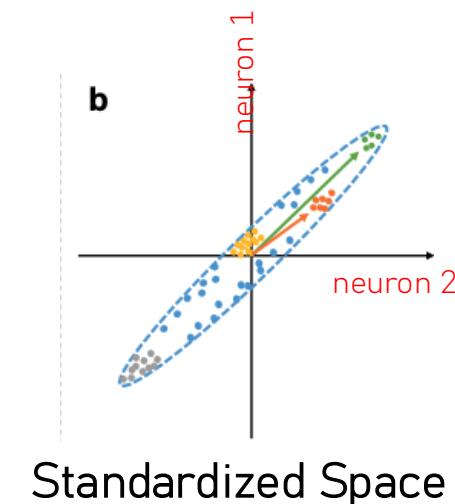
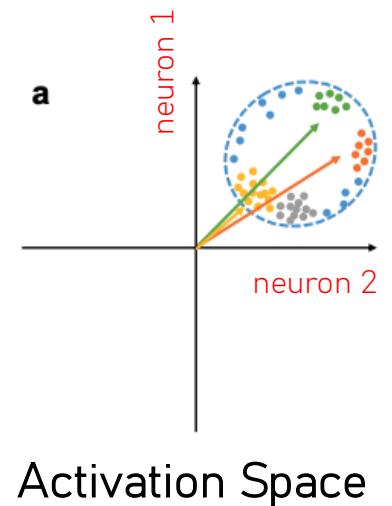
[1] Chen et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence 2.12 (2020).



WHITENING FOR DISENTANGLING CONCEPTS

Intuition

Normalization can somewhat **help disentangle concepts** in a DNN's latent space



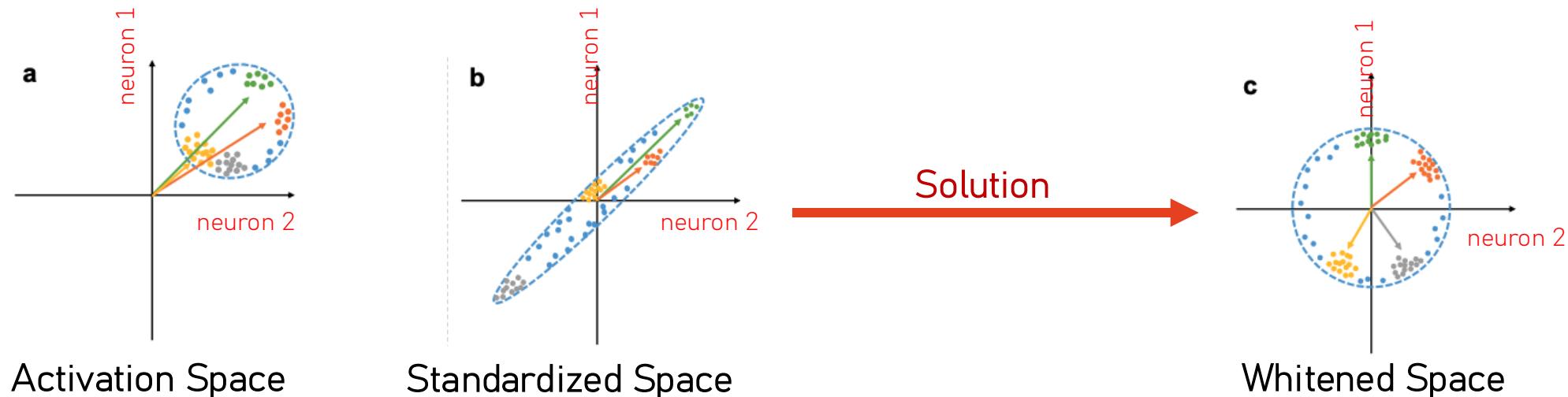
[1] Chen et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence 2.12 (2020).



WHITENING FOR DISENTANGLING CONCEPTS

Intuition

Whitening a latent space can allow us to properly separate concepts in the latent space



[1] Chen et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence 2.12 (2020).

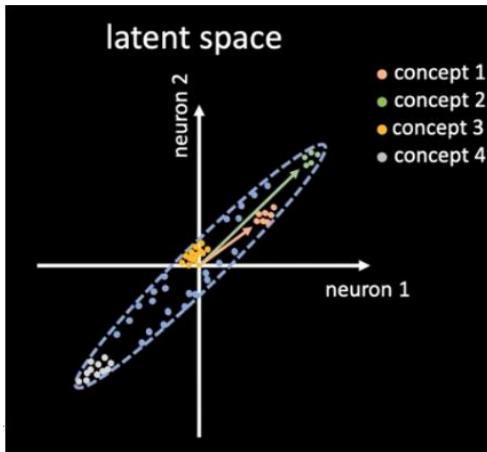


ROTATING TO ENSURE CONCEPT ALIGNMENT

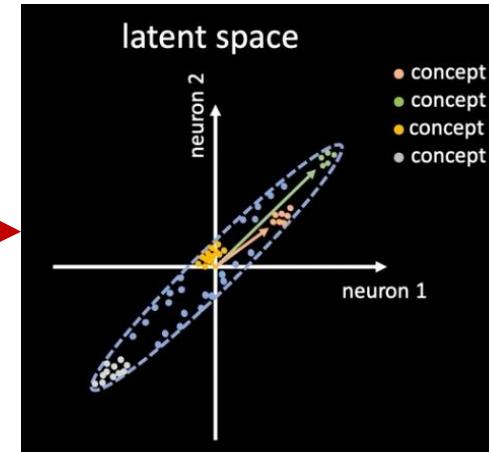
Approach

More importantly, once an input is whitened, **we can apply a rotation to align a specific concept to a specific axis!**

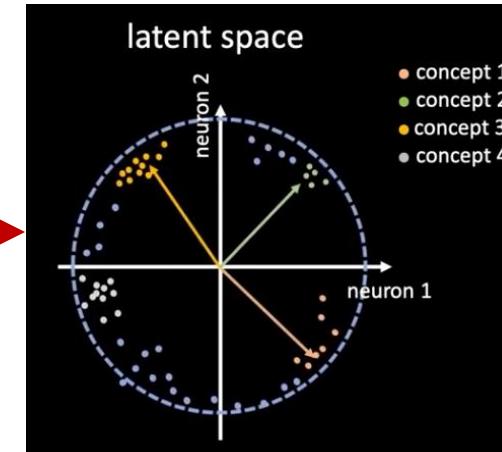
Normalize



Whiten



Rotate



[1] Chen et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence 2.12 (2020).



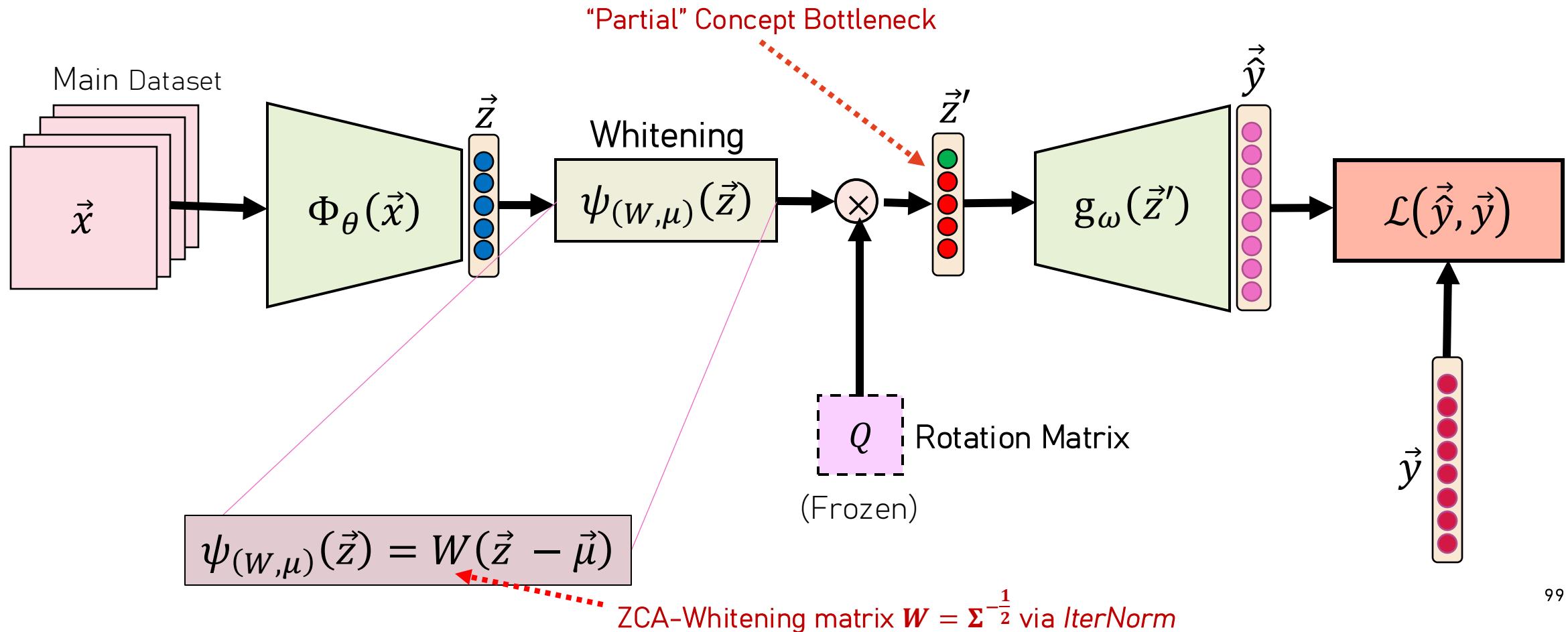
CONCEPT WHITENING (CW)

Approach

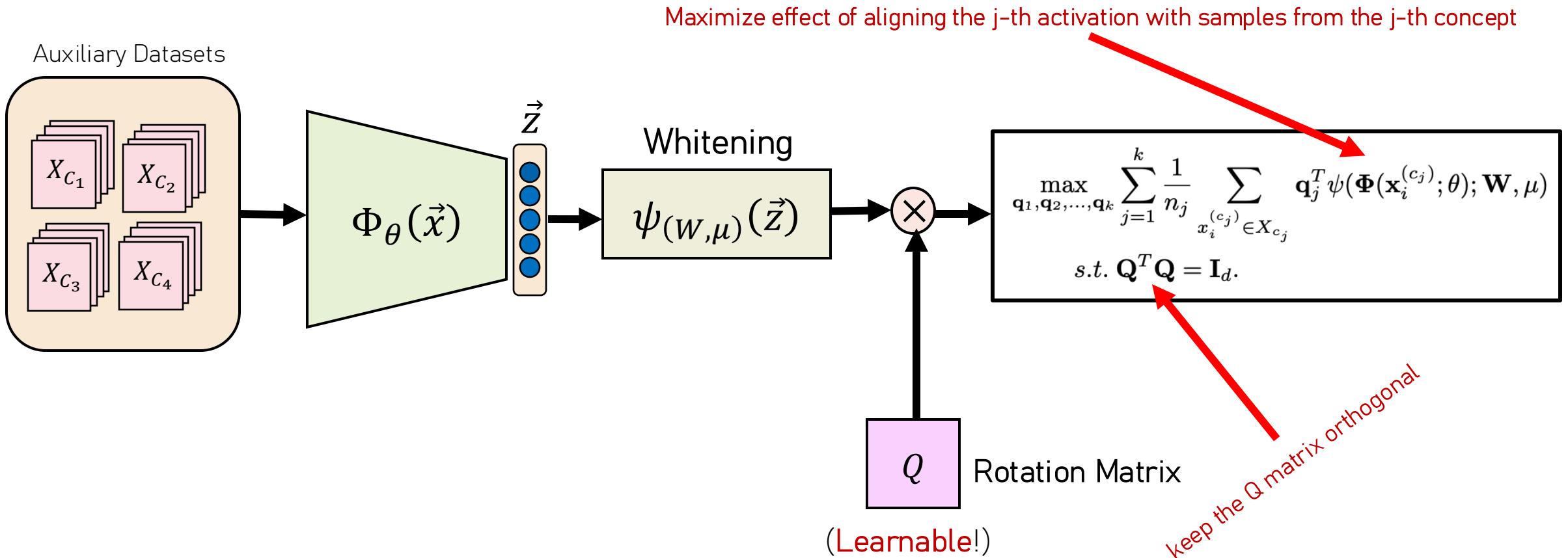
Given a **fine-tuning training set** $\mathcal{D}_t = \{(x^{(i)}, y^{(i)})\}_i$ and **k concept sets** $\mathcal{D}_c = \{X_{C_1}, X_{C_2}, \dots, X_{C_k}\}$, we will **learn a rotation matrix** $Q \in \mathbb{R}^{m \times m}$ by **iterating between**:

1. **Task training step** → make sure the downstream task prediction is accurate
2. **Concept alignment step** → make sure each concept is aligned to a latent activation

TRAINING CW: TASK TRAINING STEP

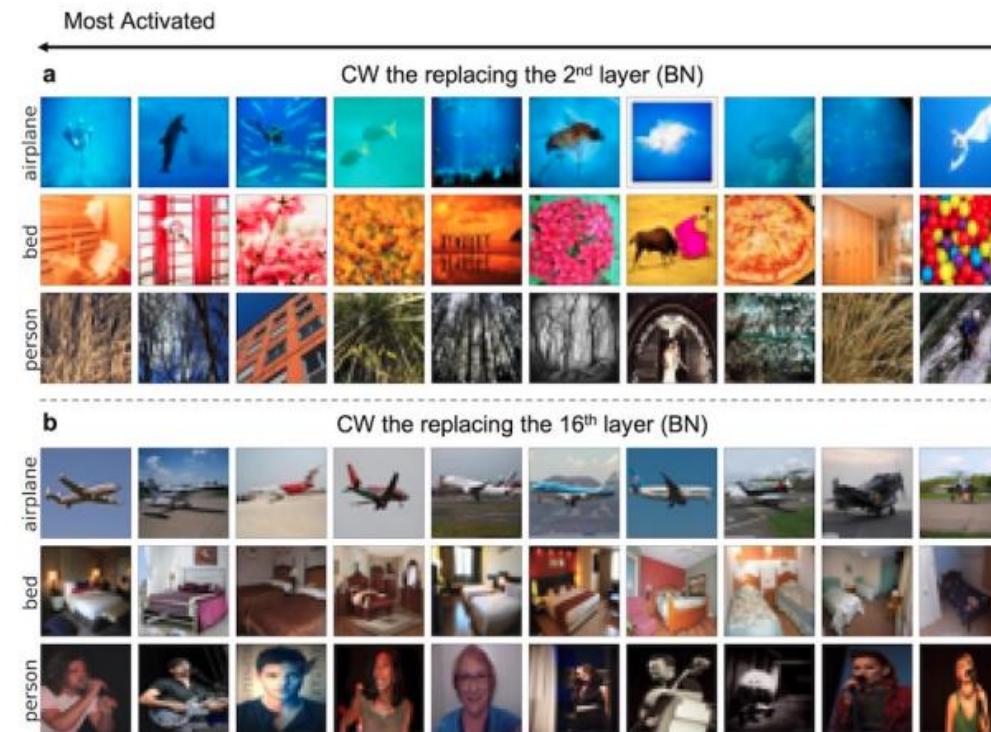


TRAINING CW: CONCEPT ALIGNMENT STEP



FOUND PROTOTYPES

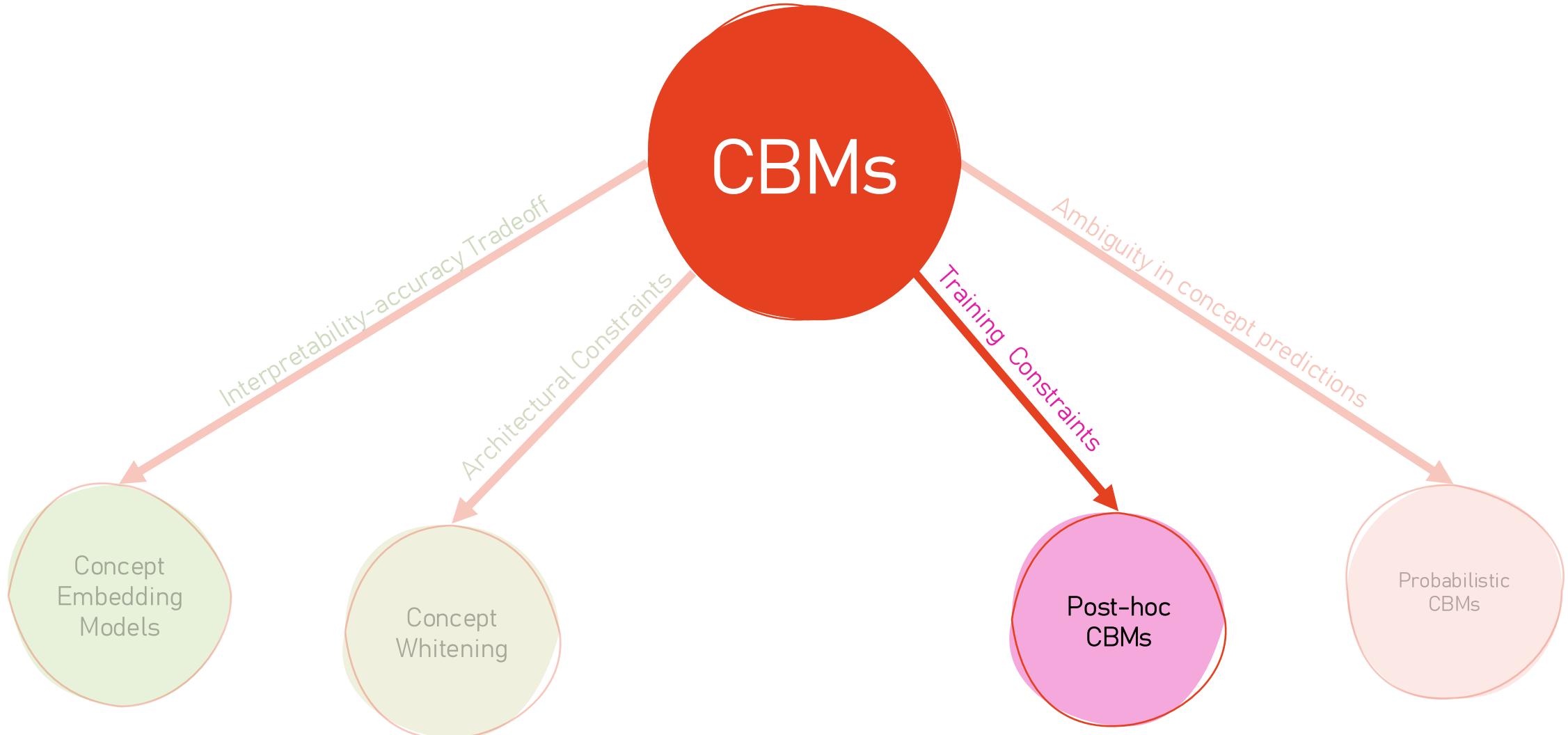
CW supports the hypothesis that **DNNs learn more complex concepts in later layers**:



[1] Chen et al. "Concept whitening for interpretable image recognition." Nature Machine Intelligence 2.12 (2020).



SPEED-DATING WITH CBM'S FRIENDS

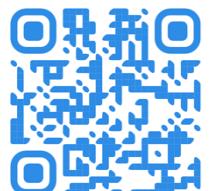


POST-HOC CBMS

Limitation Being Addressed

CW is great but has some key limitations:

1. It requires **a batch norm layer** in the pretrained model (read: **architecture-specific**)
2. It requires fine-tuning of **all the of the model's weights** (read: **could be expensive**)

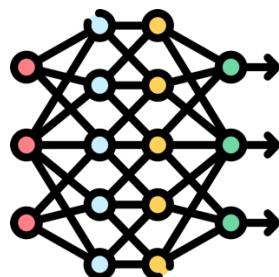


POST-HOC CBMS

Proposed Solution

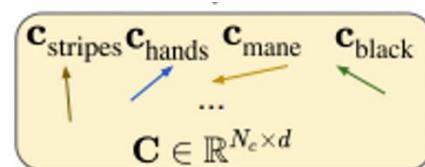
Given a bank of concept activation vectors in a pre-trained model, we **learn an interpretable mapping projected concept scores** and **a downstream task** of interest

A *trained* DNN



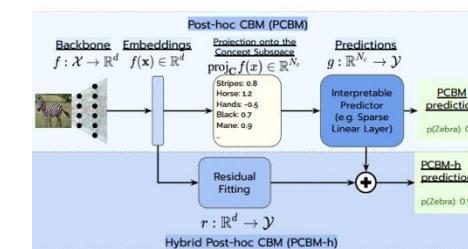
+

Concept Activation Vector Bank

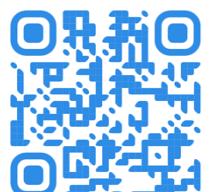


=

Post-hoc CBMs!



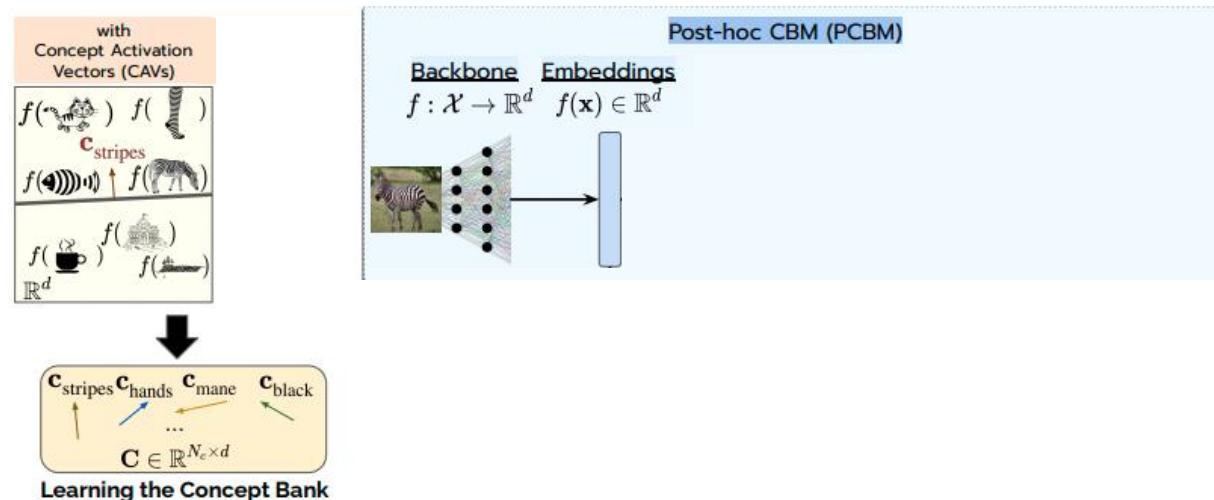
[1] Yuksekgonul et al. "Post-hoc concept bottleneck models." ICLR (2023).



POST-HOC CBMS

Proposed Solution

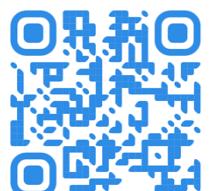
Goal: learn an **interpretable model** mapping **concept similarity scores** to **task labels**



Make the final prediction with an **interpretable predictor**

Step 1: learn CAVs from the frozen latent space of the pre-trained DNN

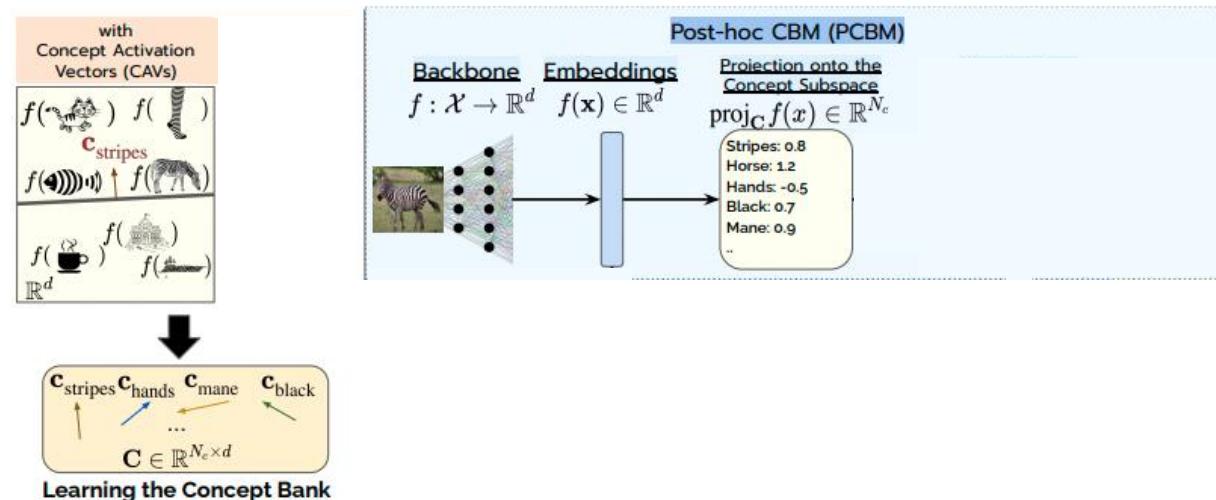
[1] Yuksekgonul et al. "Post-hoc concept bottleneck models." ICLR (2023).



POST-HOC CBMS

Proposed Solution

Goal: learn an **interpretable model** mapping **concept similarity scores** to **task labels**



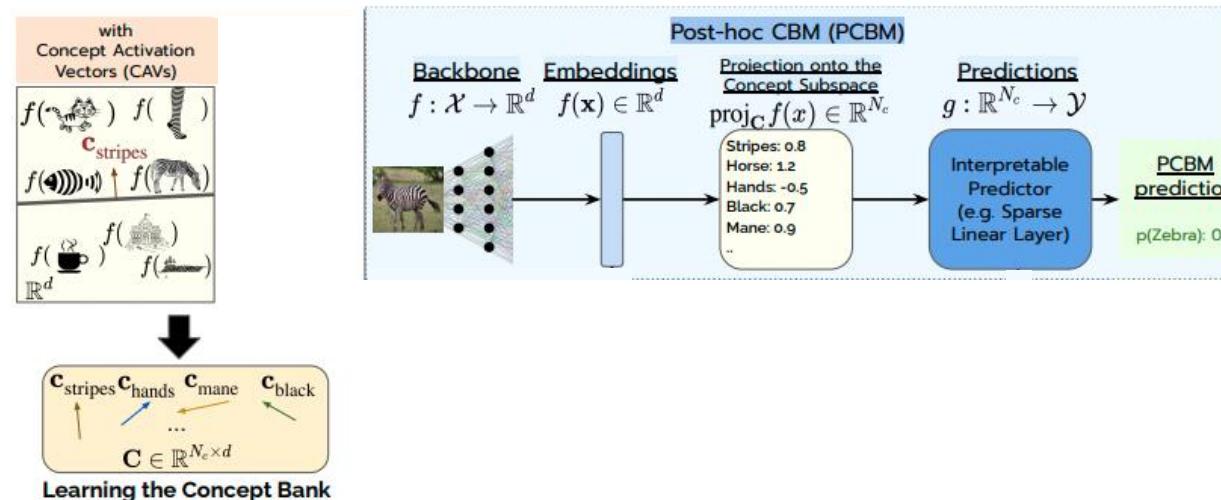
Step 2: project all training samples to the concept activation space using the cavs



POST-HOC CBMS

Proposed Solution

Goal: learn an **interpretable model** mapping **concept similarity scores** to **task labels**



Step 3: learn an interpretable predictor from the concept scores to the task labels

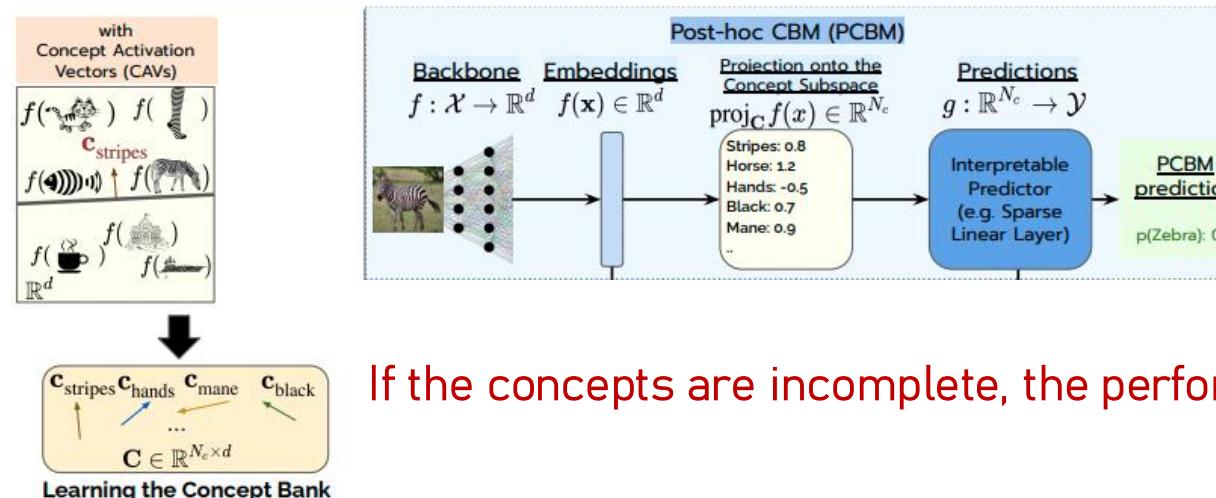
[1] Yuksekgonul et al. "Post-hoc concept bottleneck models." ICLR (2023).



POST-HOC CBMS

Proposed Solution

Goal: learn an **interpretable model** mapping **concept similarity scores** to **task labels**



If the concepts are incomplete, the performance will drop significantly!

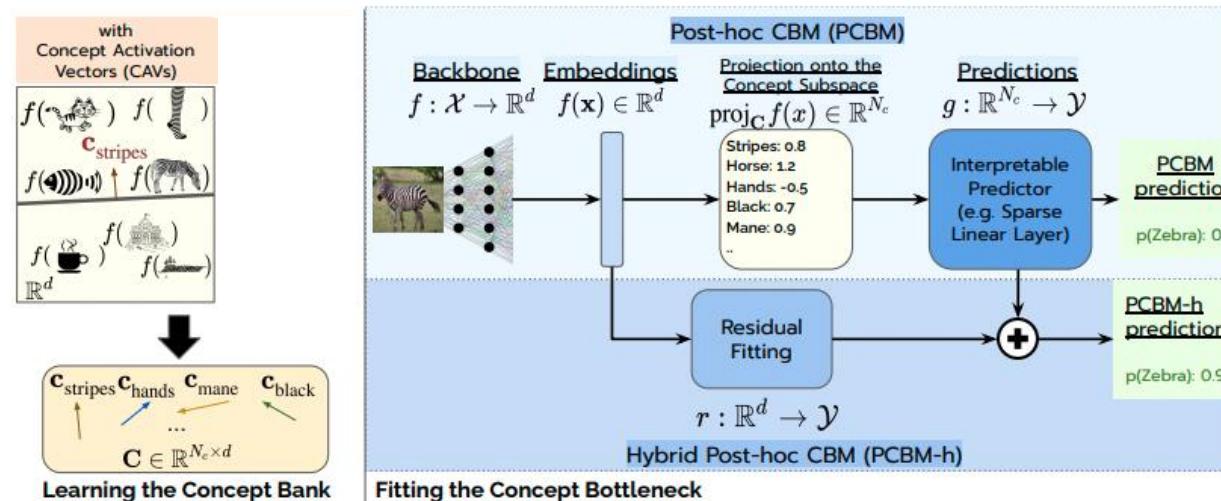
Step 3: learn an interpretable predictor from the concept scores to the task labels



POST-HOC CBMS

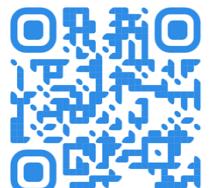
Proposed Solution

Goal: learn an **interpretable model** mapping **concept similarity scores** to **task labels**



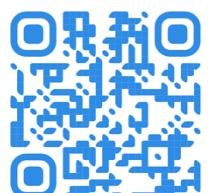
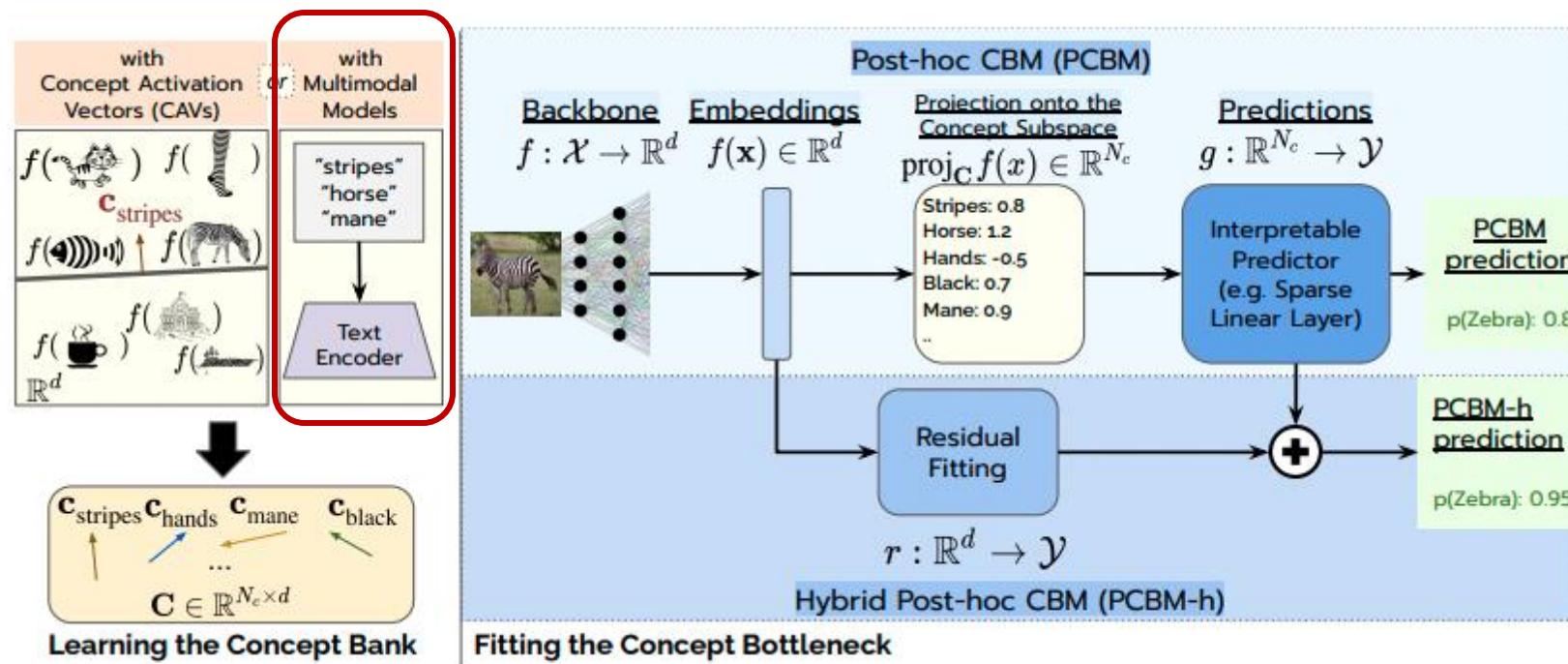
Step 4 (optional): fit a residual model if the concepts are incomplete

[1] Yuksekgonul et al. "Post-hoc concept bottleneck models." ICLR (2023).

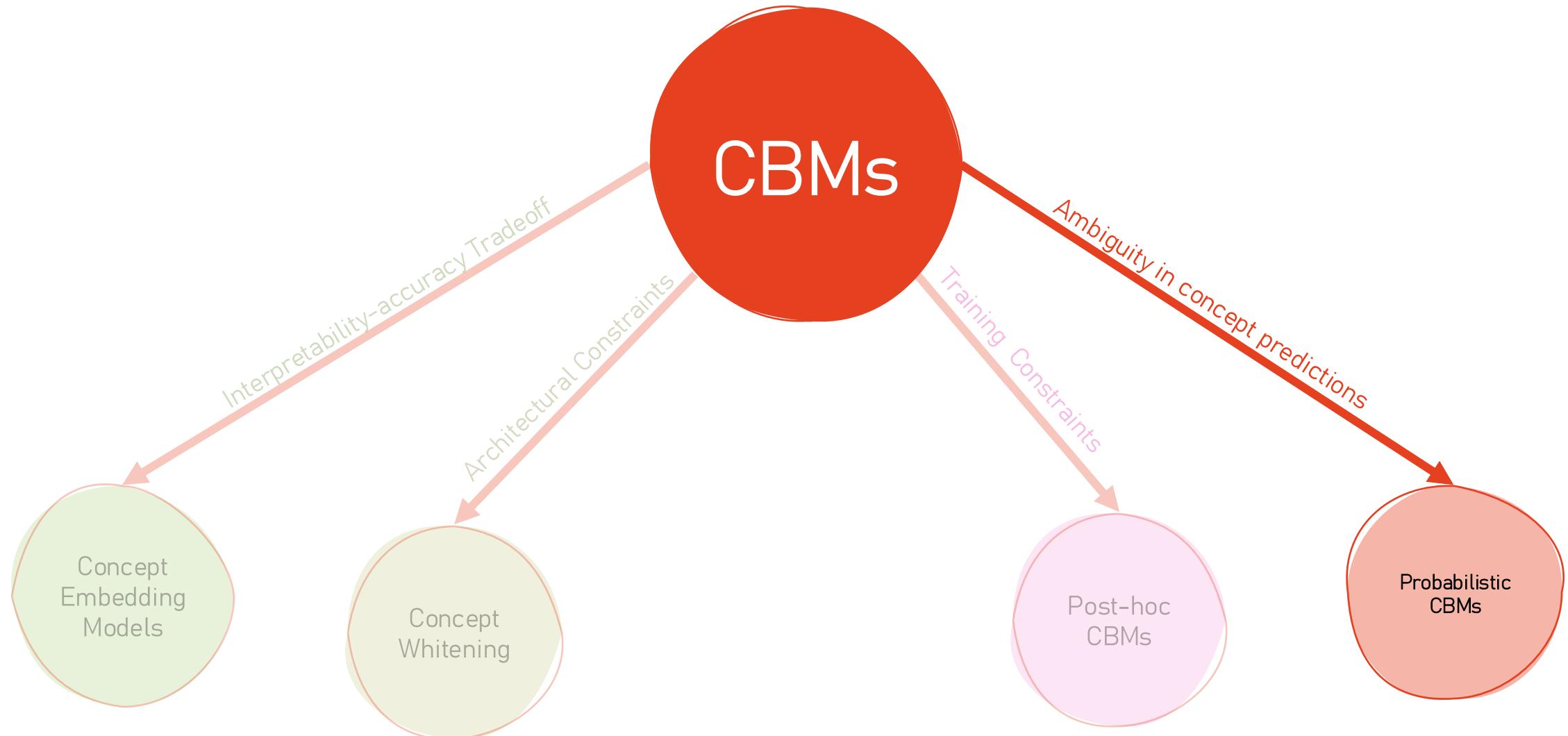


POST-HOC CBMS WITHOUT CONCEPT SETS

Post-hoc CBMs can be learnt **without concept sets** if we have access to **language-based concepts** together with a **multimodal model**



SPEED-DATING WITH CBM'S FRIENDS



PROBABILISTIC CBMS

Limitation being addressed

CBMs **must predict concepts** for all samples even they are **ambiguous**

Class: Green Jay

Concepts:
forehead color: blue
throat color: black
belly color: yellow
tail pattern: solid



ambiguity in tail



ambiguity in belly



ambiguity in color

The cross-entropy loss **does not encourage the concept predictor to be uncertain**

[1] Kim et al. "Probabilistic Concept Bottleneck Models." ICML (2023).

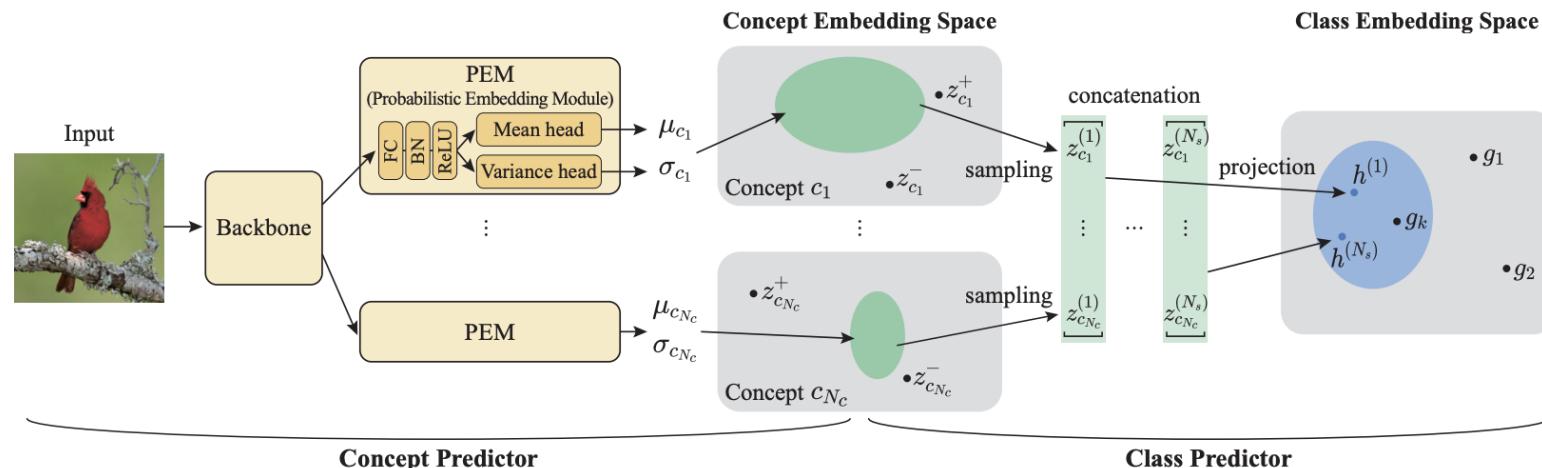


PROBABILISTIC CBMS

Proposed Solution

Use **probabilistic embeddings** that enable **uncertainty estimation** of each concept!

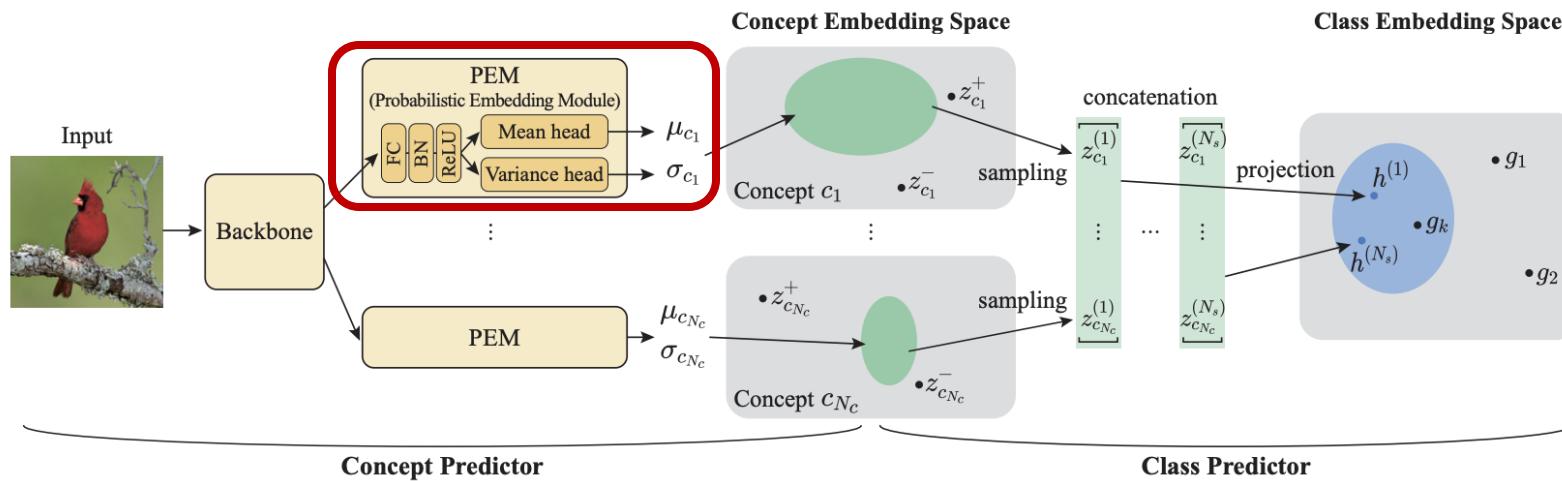
Learn a **distribution over concept embeddings** and use its **variance** to estimate **uncertainty**



PROBABILISTIC CBMS

Proposed Solution

Each **Probabilistic Embedding Module (PEM)** generates a **mean** μ_{c_i} and a **variance** σ_{c_i} for the concept embedding



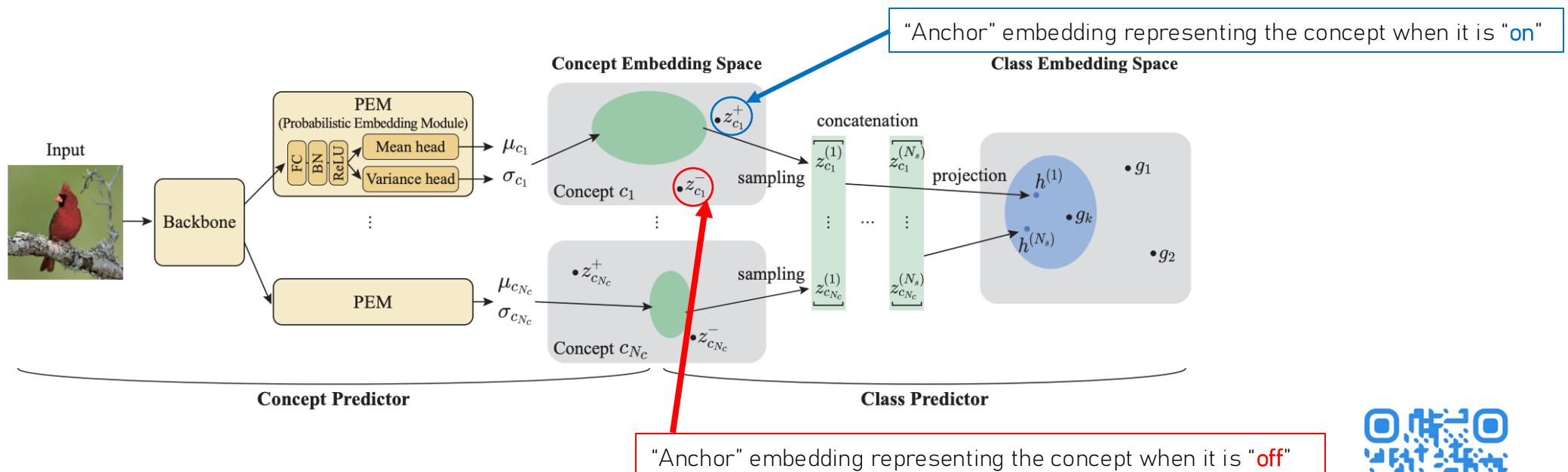
$$p(z_c|x) \sim \mathcal{N}(\mu_c, \text{diag}(\sigma_c))$$



PROBABILISTIC CBMS

Proposed Solution

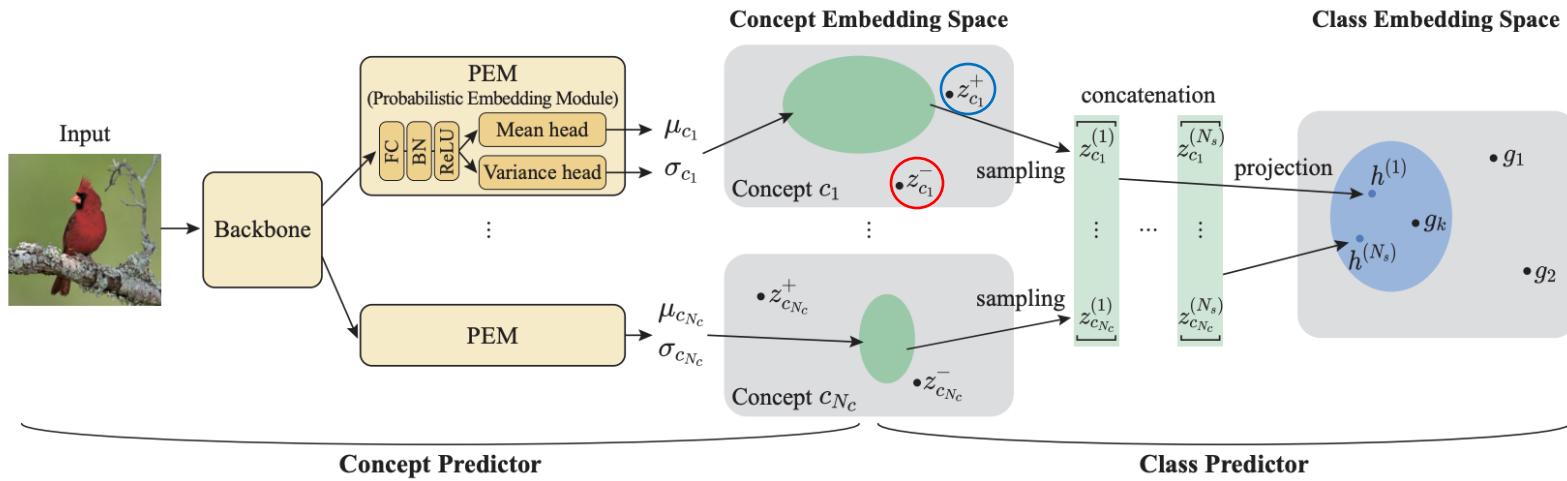
We learn a set of fixed **anchor embeddings** representing the concept when it is **on** vs **off**



PROBABILISTIC CBMS

Proposed Solution

The **distance** from the sampled embedding to each anchor can be used to **predict a concept**



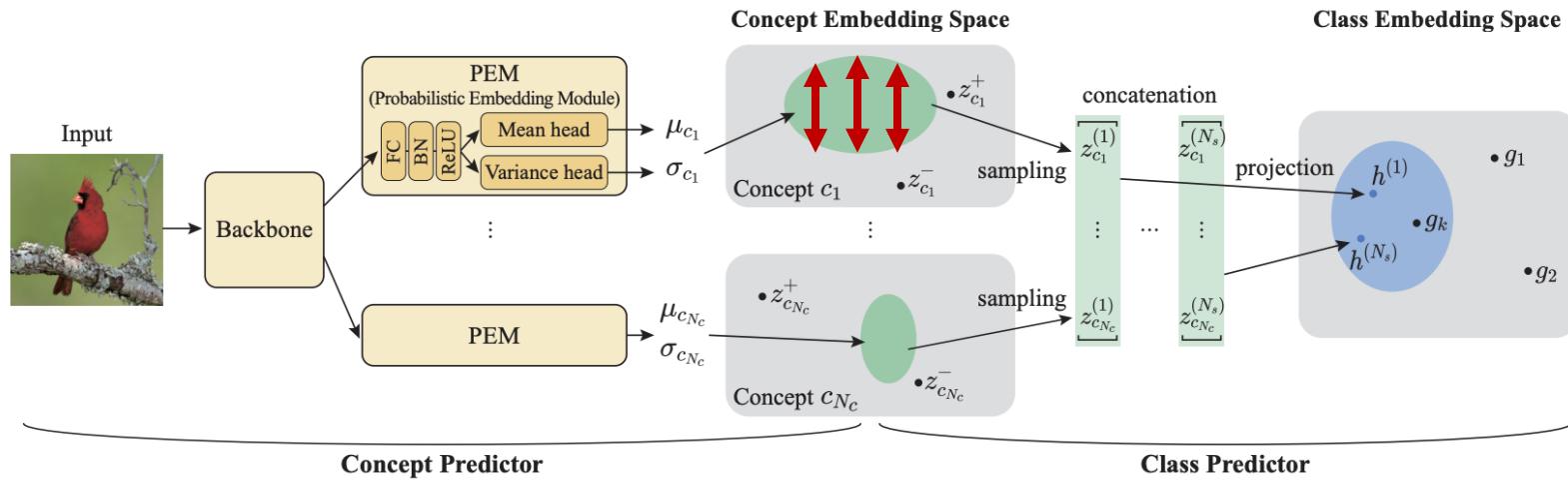
$$p(c = 1 | z_c) = \sigma \left(a \left(\|z_c - z_c^- \|_2 - \|z_c - z_c^+ \|_2 \right) \right)$$



PROBABILISTIC CBMS

Proposed Solution

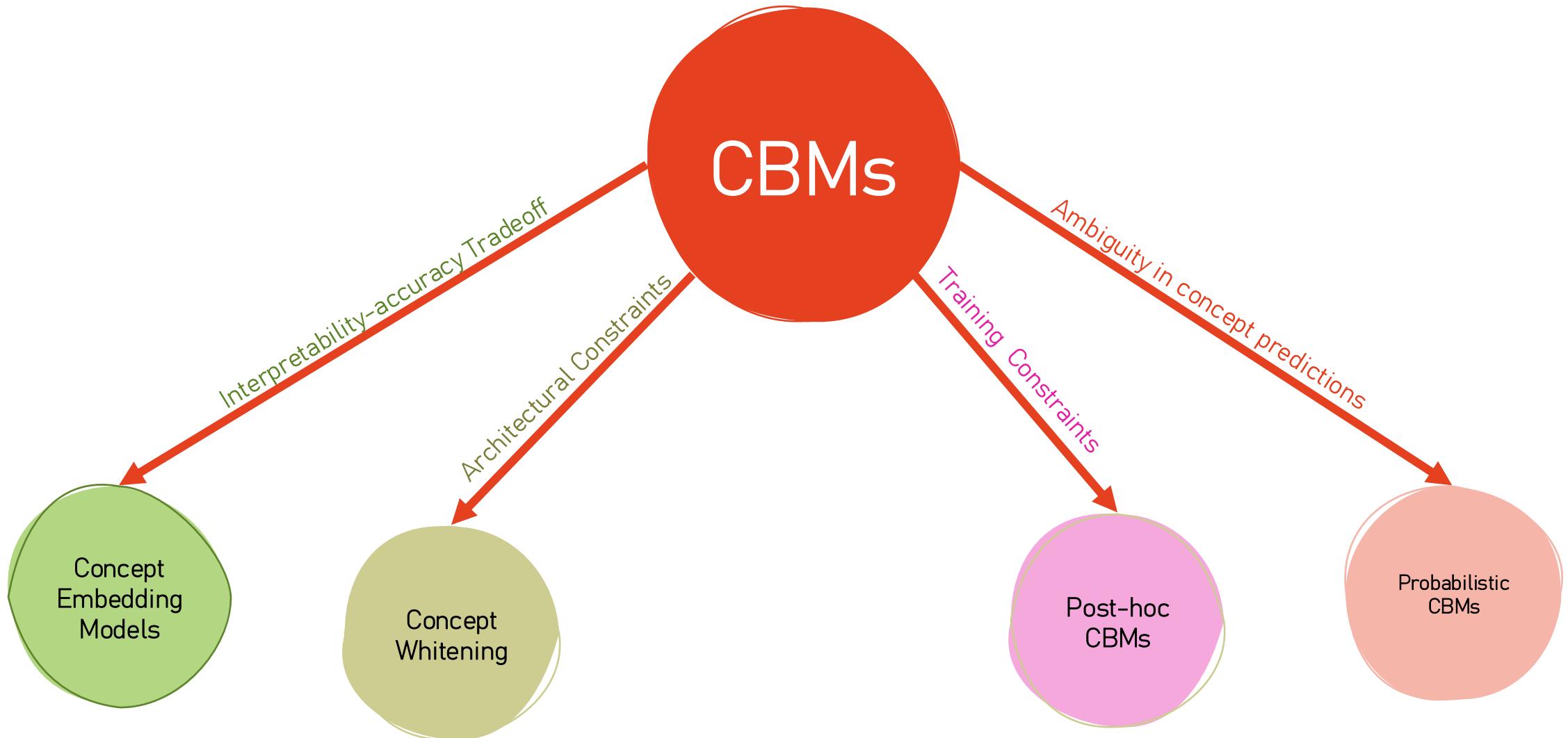
A **concept's distribution's volume** can be used to **quantify its uncertainty**



As **embeddings are modelled as Gaussians**, this is the **determinant of the covariance!**



END OF THE SPEED DATES!



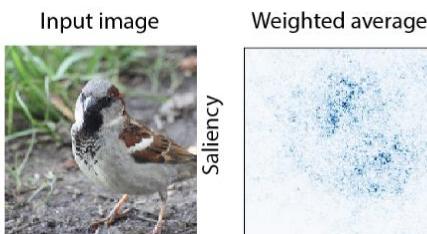
DO CBMS PROPERLY LEARN TO EXPLAIN?

DO CBMS PROPERLY LEARN TO EXPLAIN?

Several recent works suggest CBMs may have issues with **unwanted leakage**

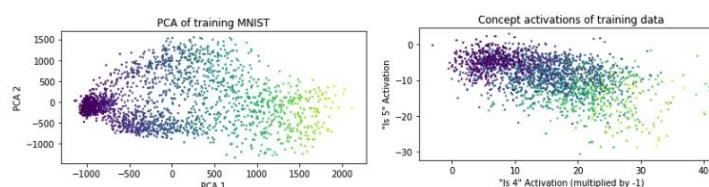
Attending spurious features

Joint model on concept: leg color



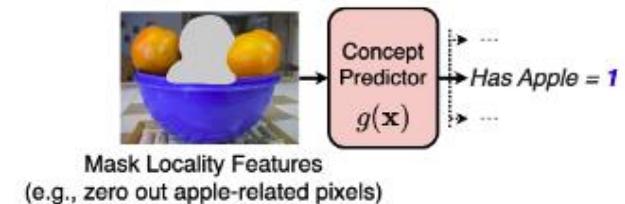
Saliency maps seem to suggest
concepts are not properly attended
(Margeloiu et al.) [1]

Leaking unwanted information



CBMs may have incentives to **encode the entire data representation** in the concepts' soft predictions
(Mahinpei et al.) [2]

Failing to capture concept locality



CBMs may **fail to capture a concept's locality** (e.g., physical location) even if it is only found on a fixed feature subset
(Raman et al.) [3]

[1] Margeloiu et al. "Do concept bottleneck models learn as intended?" ICLR Workshop on Responsible AI (2021).

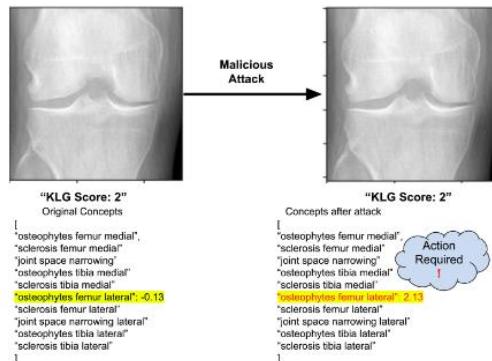
[2] Mahinpei et al. "Promises and pitfalls of black-box concept learning models." ICML Workshop on Theoretic Foundation, Criticism, and Application of XAI (2021).

[3] Raman et al. "Do Concept Bottleneck Models Respect Localities?" NeurIPS Workshop on XAI in Action (2024).

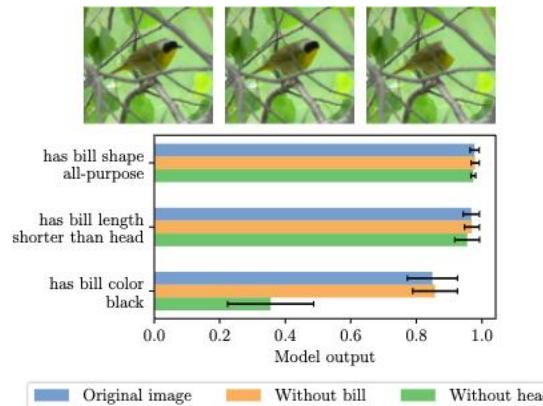
DO CBMS PROPERLY LEARN TO EXPLAIN?

Many more works have dived deeper into these issues!

Adversarial attacks and defences



Studying concept correlations



CBMs concepts predictions can be changed
without affecting the final prediction
(Sinha et al.) [1]

Simple changes to the loss, like **loss weighting**, can help **avoid CBMs exploiting unwanted correlations**
(Heidenmann et al.) [2]

Formalisms and metrics for leakage



Metrics for
unwanted
leakage [3]

Benchmark
suite for
reasoning
robustness [4]

Formalisation
of leakage [5]

Several works proposed ways to **formalise** or
measure concept leakage [3, 4, 5]

[1] Sinha et al. "Understanding and enhancing robustness of concept-based models." AAAI (2023).

[2] Heidemann et al. "Concept correlation and its effects on concept-based models." WACV (2023).

[3] Espinosa Zarlenga, Barbiero, Shams et al. "Towards robust metrics for concept representation evaluation." AAAI (2023).

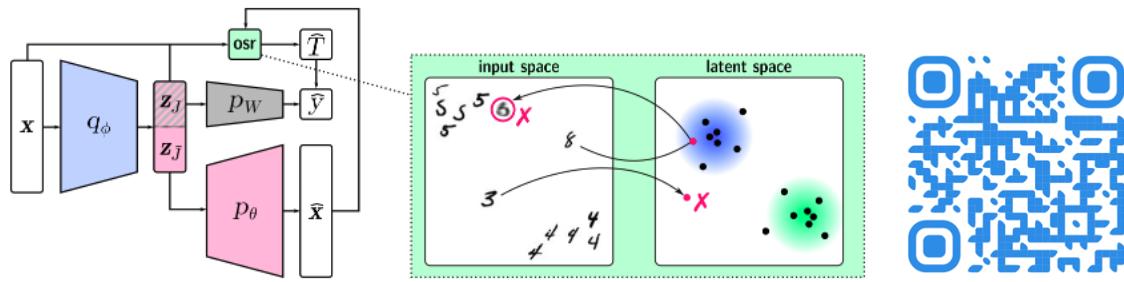
[4] Bortolotti, Marconato, et al. "A Neuro-Symbolic Benchmark Suite for Concept Quality and Reasoning Shortcuts." NeurIPS (2024).

[5] Marconato et al. "Interpretability is in the mind of the beholder: A causal framework for human-interpretable representation learning." Entropy (2023).

STEPS TOWARDS ADDRESSING LEAKAGE

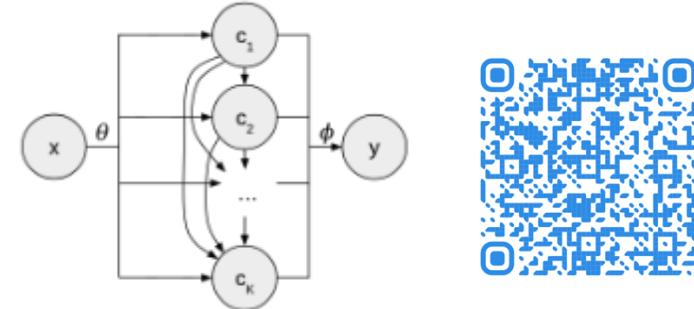
This have brought forth attempts to **address or mitigate the effects of leakage**:

GlanceNets



[Main Idea] Frame leakage in terms of disentanglement learning and use an open-set recognition to detect it at inference

Autoregressive CBMs



[Main Idea] Reduce leakage between concepts by modeling cross-concept relationships using an autoregressive architecture

- [1] Marconato et al. "Glancenets: Interpretable, leak-proof concept-based models." NeurIPS (2022).
[2] Havasi et al. "Addressing leakage in concept bottleneck models." NeurIPS (2022).

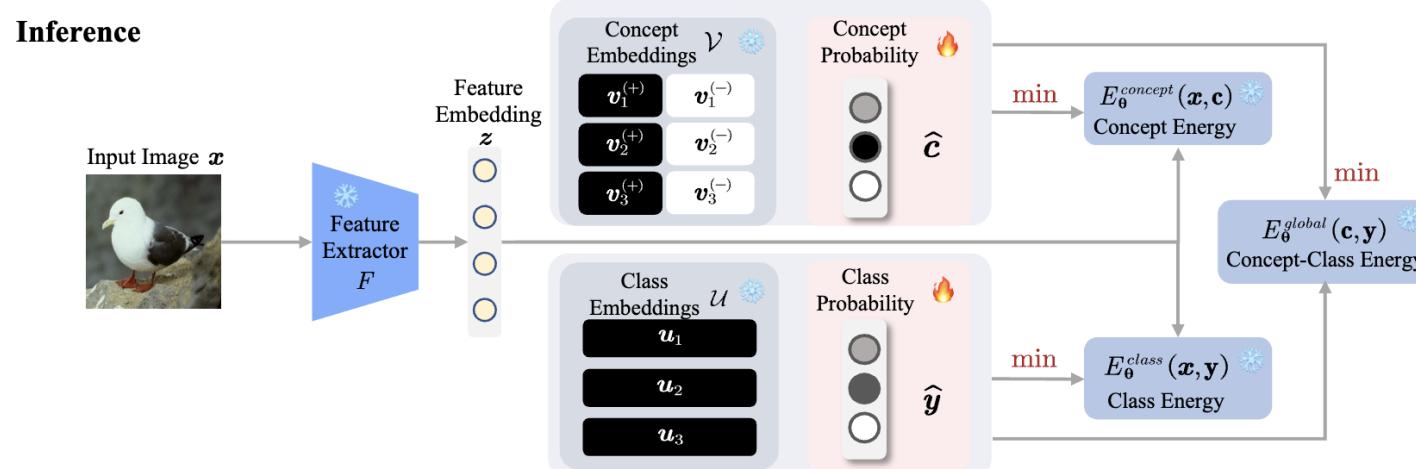
RECENT DIRECTIONS

CBMs have become **very popular** in XAI with several active **areas of research**:

RECENT DIRECTIONS

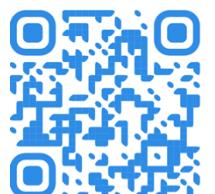
CBMs have become **very popular in XAI** with several active **areas of research**:

1. Capturing more **complex relationships** between concepts and tasks labels



Energy-based Concept Bottleneck Models (Xu et al., 2024)

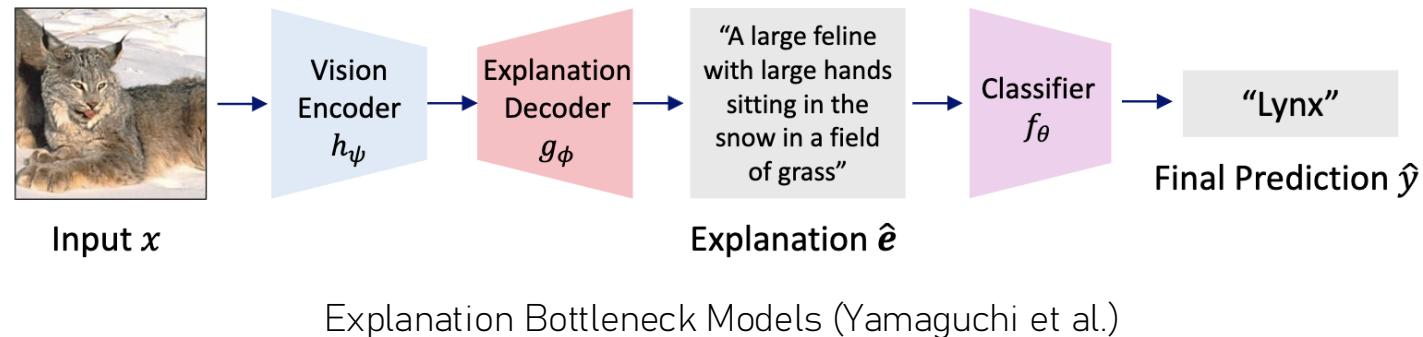
[1] Xu et al. "Energy-based concept bottleneck models: unifying prediction, concept intervention, and conditional interpretations." ICLR (2024).



RECENT DIRECTIONS

CBMs have become **very popular in XAI** with several active **areas of research**:

1. Capturing more **complex relationships** between concepts and tasks labels
2. Producing entirely **language-based bottlenecks** (accepted to this AAAI!)



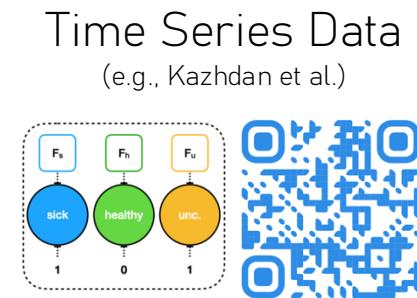
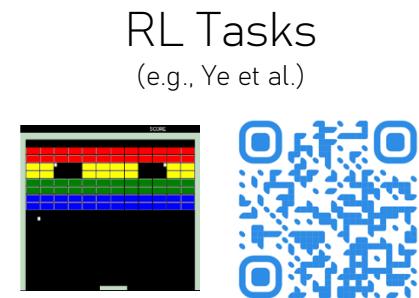
[1] Yamaguchi et al. "Explanation Bottleneck Models." AAAI (2025).



RECENT DIRECTIONS

CBMs have become **very popular in XAI** with several active **areas of research**:

1. Capturing more **complex relationships** between concepts and tasks labels
2. Producing entirely **language-based bottlenecks** (accepted to this AAAI!)
3. Exploring concepts in **modalities and tasks** other than **supervised visual tasks**



[1] Xuanyuan al. "Global concept-based interpretability for graph neural networks via neuron analysis." AAAI (2023).

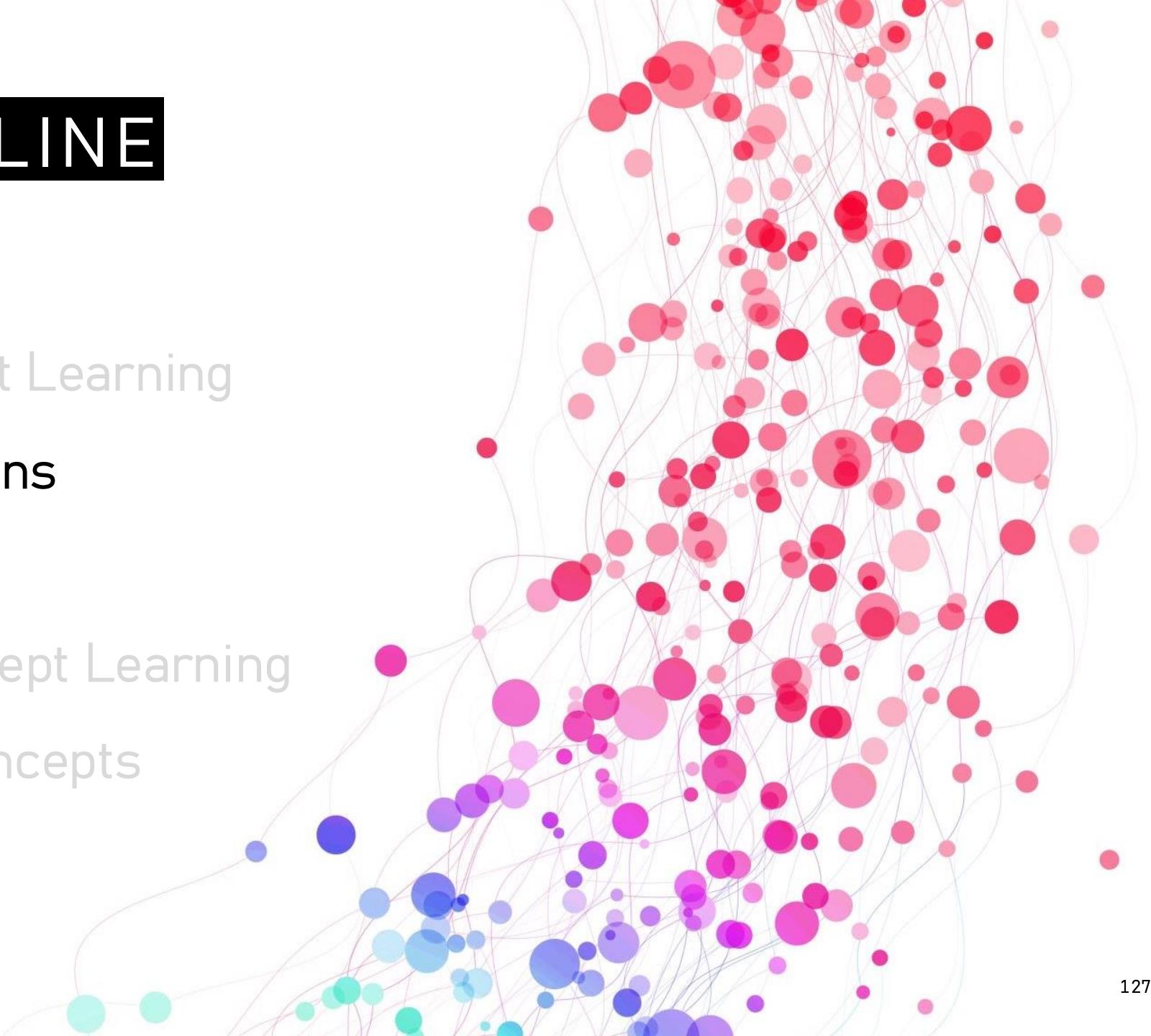
[2] Espinosa Zarlenga et al. "Tabcbm: Concept-based interpretable neural networks for tabular data." TMLR (2024).

[3] Ye et al. "Concept-based interpretable reinforcement learning with limited to no human labels." ICML (2024).

[4] Kazhdan et al. "MFME: generating RNN model explanations via model extraction." arXiv (2020).

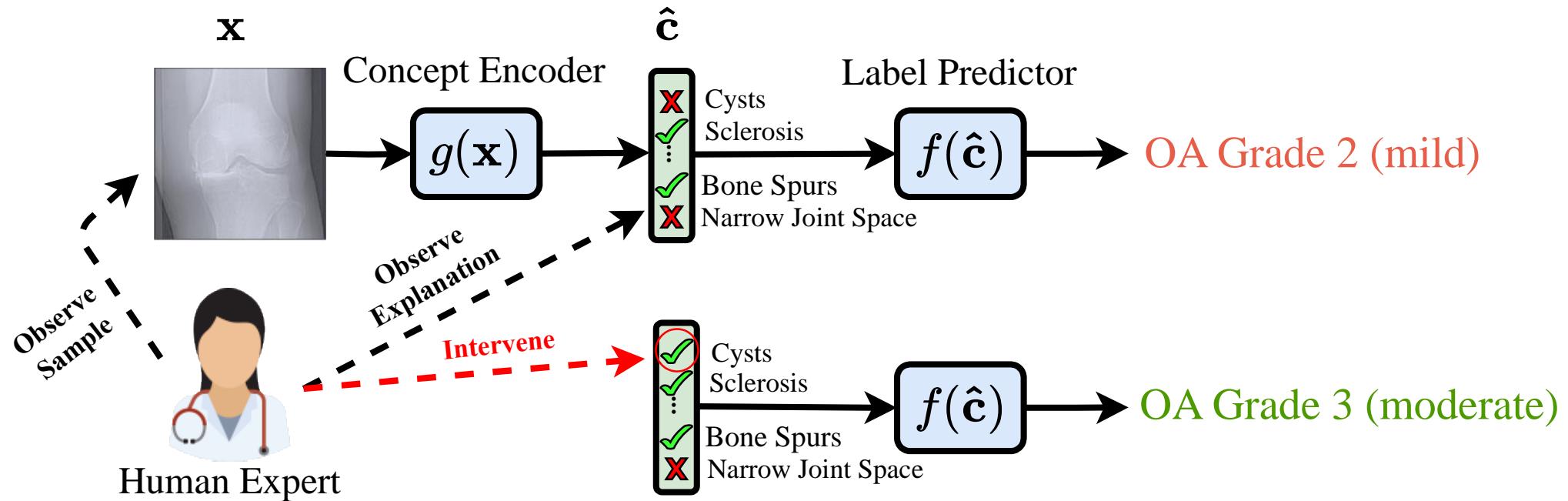
TUTORIAL OUTLINE

1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



RECALL CONCEPT INTERVENTIONS

Concept interventions enable experts to “inject” knowledge during inference



SOME CONCEPTS ARE BETTER THAN OTHERS

When intervening on a CBM, it is important to realise that some concepts are:

1. Less **informative** than others (e.g., redundant w.r.t. other concepts)
2. Less **certain** than others (e.g., due to occlusions or inherent difficulties)



Concept "Belly Color" is partially occluded

To identify a Raven from a Crow, "tail shape" is more informative than "wing color"

SOME CONCEPTS ARE BETTER THAN OTHERS

When intervening on a CBM, it is important to realise that some concepts are:

1. Less **informative** than others (e.g., redundant w.r.t. other concepts)
2. Less **certain** than others (e.g., due to occlusions or inherent difficulties)

Hence, an intervention's effectiveness **depends on the intervened concept!**

SELECTING MEANINGFUL CONCEPTS

Intervention policies select which concept to intervene on next by assigning each concept c_i a score s_i and selecting concepts in decreasing score order:

Given \mathbf{x} and concept predictions $\hat{\mathbf{c}}$, what concept should I intervene on next to minimize my model's task uncertainty?



SELECTING MEANINGFUL CONCEPTS

Intervention policies select which concept to intervene on next by assigning each concept c_i a score s_i and selecting concepts in decreasing score order:

Uncertainty of concept prediction (UCP)

Select the concept c_i with the highest predicted entropy $s_i = \mathcal{H}(\hat{c}_i)$



SELECTING MEANINGFUL CONCEPTS

Intervention policies select which concept to intervene on next by assigning each concept c_i a score s_i and selecting concepts in decreasing score order:

Uncertainty of concept prediction (UCP)

Select the concept c_i with the highest predicted entropy $s_i = \mathcal{H}(\hat{c}_i)$

Contribution of concept on target prediction (CCTP)

Select the concept c_i with the highest contribution on target prediction $s_i = \sum_{j=1}^L \left| \hat{c}_i \frac{\partial f_j(x)}{\partial \hat{c}_i} \right|$.



SELECTING MEANINGFUL CONCEPTS

Intervention policies select which concept to intervene on next by assigning each concept c_i a score s_i and selecting concepts in decreasing score order:

Uncertainty of concept prediction (UCP)

Select the concept c_i with the highest predicted entropy $s_i = \mathcal{H}(\hat{c}_i)$

Contribution of concept on target prediction (CCTP)

Select the concept c_i with the highest contribution on target prediction $s_i = \sum_{j=1}^L \left| \hat{c}_i \frac{\partial f_j(x)}{\partial \hat{c}_i} \right|$.

Expected change in target prediction (ECTP)

Select the concept c_i with the highest expected change in the target predictive distribution
$$s_i = (1 - \hat{c}_i) D_{KL}(\hat{y}_{\hat{c}_i=0} || \hat{y}) + \hat{c}_i D_{KL}(\hat{y}_{\hat{c}_i=1} || \hat{y})$$



SELECTING MEANINGFUL CONCEPTS

Intervention policies select which concept to intervene on next by assigning each concept c_i a score s_i and selecting concepts in decreasing score order:

Uncertainty of concept prediction (UCP)

Select the concept c_i with the highest predicted entropy $s_i = \mathcal{H}(\hat{c}_i)$

Contribution of concept on target prediction (CCTP)

Select the concept c_i with the highest contribution on target prediction $s_i = \sum_{j=1}^L \left| \hat{c}_i \frac{\partial f_j(x)}{\partial \hat{c}_i} \right|$.

Expected change in target prediction (ECTP)

Select the concept c_i with the highest expected change in the target predictive distribution
$$s_i = (1 - \hat{c}_i) D_{KL}(\hat{y}_{\hat{c}_i=0} || \hat{y}) + \hat{c}_i D_{KL}(\hat{y}_{\hat{c}_i=1} || \hat{y})$$

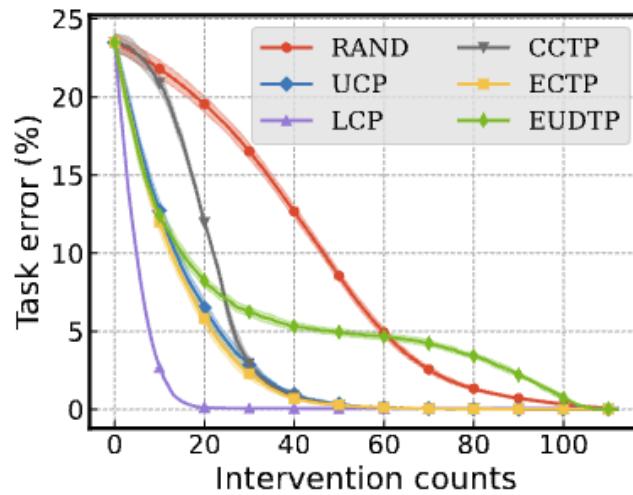
One can think of these policies as **proxies** for a concept's **information content** and **certainty**

[1] Shin et al. "A Closer Look at the Intervention Procedure of Concept Bottleneck Models." ICML 2023.

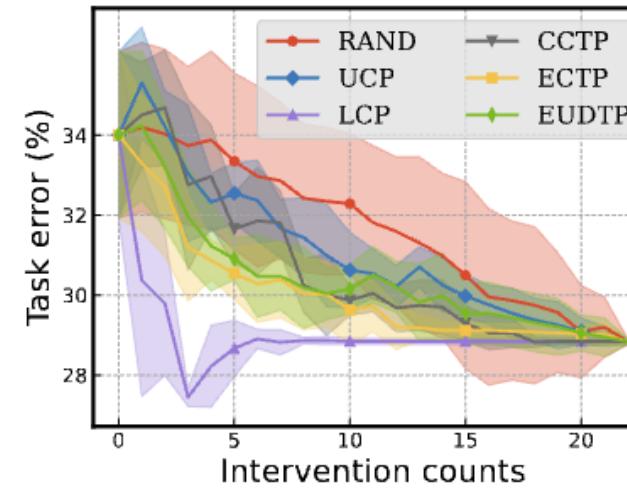


INTERVENTION POLICIES RESULTS

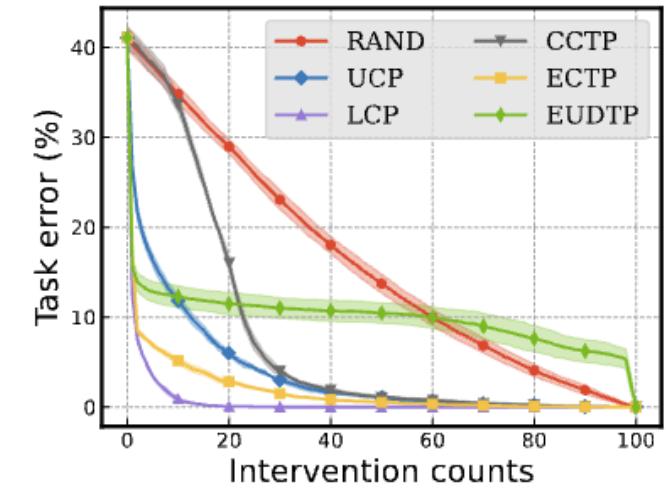
This leads to significantly different intervention curves:



(a) CUB



(b) SkinCon



(c) Synthetic

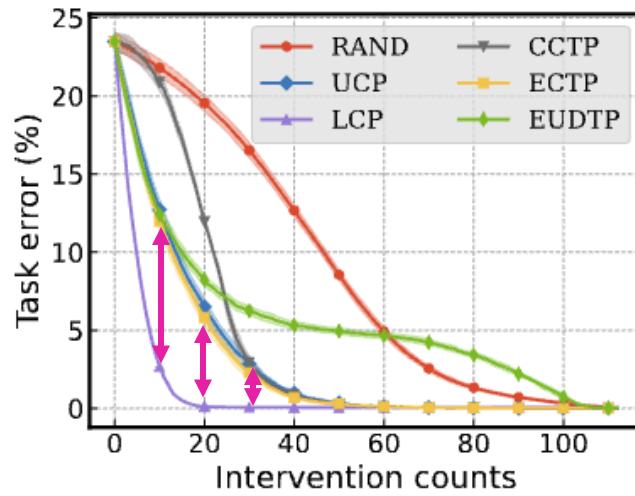
Best performing non-oracle policy is **Expected change in target prediction (ECTP)** but even the simple **Uncertainty of concept prediction (UCP)** is significantly better than the **random policy (RAND)**

(**Intuition:** one should select the concept leading to the highest expected change in the task's distribution)

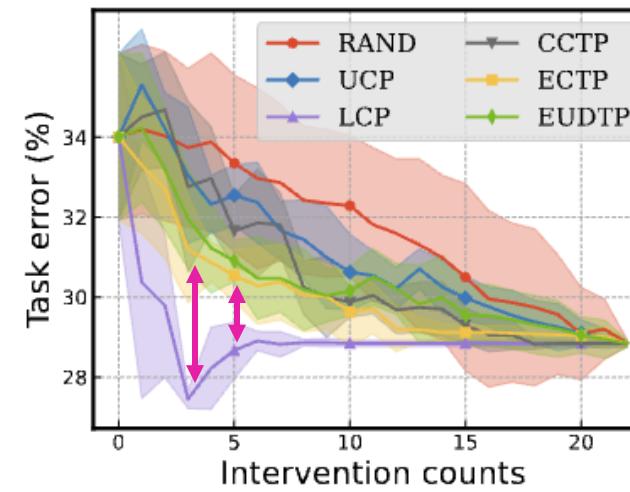


INTERVENTION POLICIES RESULTS

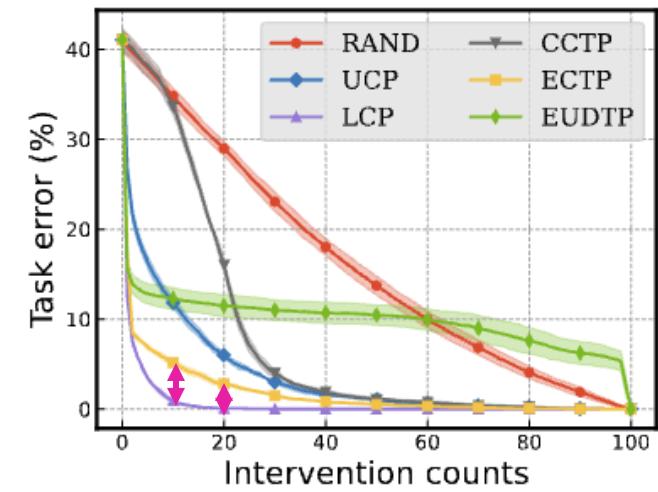
This leads to significantly different intervention curves:



(a) CUB



(b) SkinCon



(c) Synthetic

We still observe a **significant gap** between the best policy and an **optimal greedy policy (LCP)**



COMBINING POLICIES

It may be possible to shorten this gap by learning a weighting between the **concept uncertainty** and the **expected change in current prediction policies**

Cooperative Prediction Intervention Policy (CooP)

$$s_i = \alpha \mathcal{H}(\hat{c}_i) + \beta |E_{v \sim p(c_i | x)} [\hat{y}_{\hat{c}_i=v} - \hat{y}]| + \gamma q_i$$

[1] Chauhan et al. "Interactive concept bottleneck models." AAAI (2023).



COMBINING POLICIES

It may be possible to shorten this gap by learning a weighting between the **concept uncertainty** and the **expected change in current prediction policies**

Cooperative Prediction Intervention Policy (CooP)

$$s_i = \alpha \underbrace{\mathcal{H}(\hat{c}_i)}_{\text{Uncertainty of concept prediction}} + \beta \left| E_{v \sim p(c_i | x)} [\hat{y}_{\hat{c}_i=v} - \hat{y}] \right| + \gamma \underbrace{q_i}_{\text{Cost of the intervention}}$$

Uncertainty of concept prediction

Cost of the intervention

Expected change to the current predicted label if we were to intervene on c_i based on c_i 's current prediction



COMBINING POLICIES

It may be possible to shorten this gap by selecting concepts based on both

Can we do any better than this? Can we avoid the need for calculating computationally expensive scores all concepts?

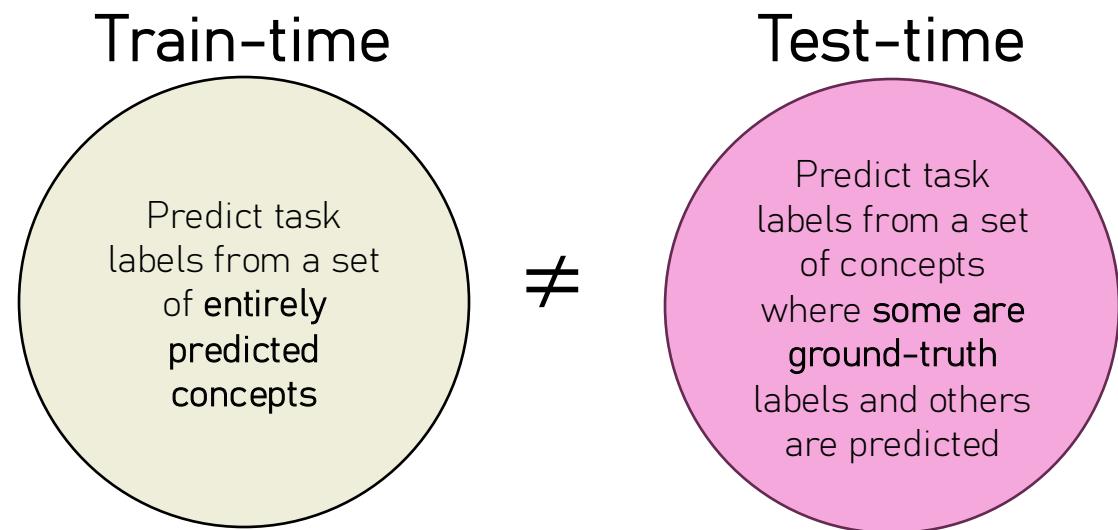
We can find the hyperparameters α , β , and γ using a concept-annotated validation set

[1] Chauhan et al. "Interactive concept bottleneck models." AAAI (2023).



INSIGHT #1: TRAINING-TIME INCENTIVES

There is a disconnect between how CBMs are trained and how they are used at test-time when they are intervened on



During training, concept-based models are **not even aware** they may be intervened on!

During testing, concept interventions may lead to **out-of-distribution bottlenecks** for a CBM!

[1] Espinosa Zarlenga et al. "Learning to receive help: Intervention-aware concept embedding models." NeurIPS (2023)



INSIGHT #2: USEFUL TRAINING FEEDBACK

If we know all task and concept labels, we can compute the optimal greedy concept intervention:

$$c_*(\mathbf{x}, \mu, \mathbf{c}, y) := \arg \max_{1 \leq i \leq k} f(\tilde{g}(\mathbf{x}, \mu \vee \mathbb{1}_i, \mathbf{c}))_y$$

(Translation: Attempt every intervention and select that one maximises the ground truth label's confidence)

This is feedback we have at training time and can use to learn an intervention policy!



INSIGHT #3: INTERVENTIONS CAN BE DIFFERENTIABLE

When modelling concepts as being a **mixture of two learnable embeddings** $\{\hat{c}_i^+, \hat{c}_i^-\}$ as in CEMs, **interventions are differentiable**:

Original Embedding Construction

$$\hat{c}_i := \hat{p}_i \hat{c}_i^+ + (1 - \hat{p}_i) \hat{c}_i^-$$

Intervened Embedding Construction

$$\hat{c}_i := (\mu_i c_i + (1 - \mu_i) \hat{p}_i) \hat{c}_i^+ + (1 - (\mu_i c_i + (1 - \mu_i) \hat{p}_i)) \hat{c}_i^-$$

Whether we intervene on the i -th concept (can be relaxed to be in $[0,1]$)



INSIGHT #3: INTERVENTIONS CAN BE DIFFERENTIABLE

When modelling concepts as being a **mixture of two learnable embeddings** $\{\hat{c}_i^+, \hat{c}_i^-\}$ as in CEMs, **interventions are differentiable**:

Original Embedding Construction

$$\hat{c}_i := \hat{p}_i \hat{c}_i^+ + (1 - \hat{p}_i) \hat{c}_i^-$$

Intervened Embedding Construction

$$\hat{c}_i := (\mu_i c_i + (1 - \mu_i) \hat{p}_i) \hat{c}_i^+ + (1 - (\mu_i c_i + (1 - \mu_i) \hat{p}_i)) \hat{c}_i^-$$

Whether we intervene on the i -th concept (can be relaxed to be in $[0,1]$)

This means an intervention policy deciding μ_i can be learnt via gradient descent!

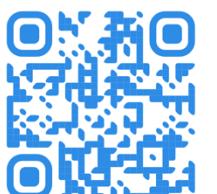


INTERVENTION-AWARE MODELS

Intervention-Aware Concept Embedding Models (IntCEMs)

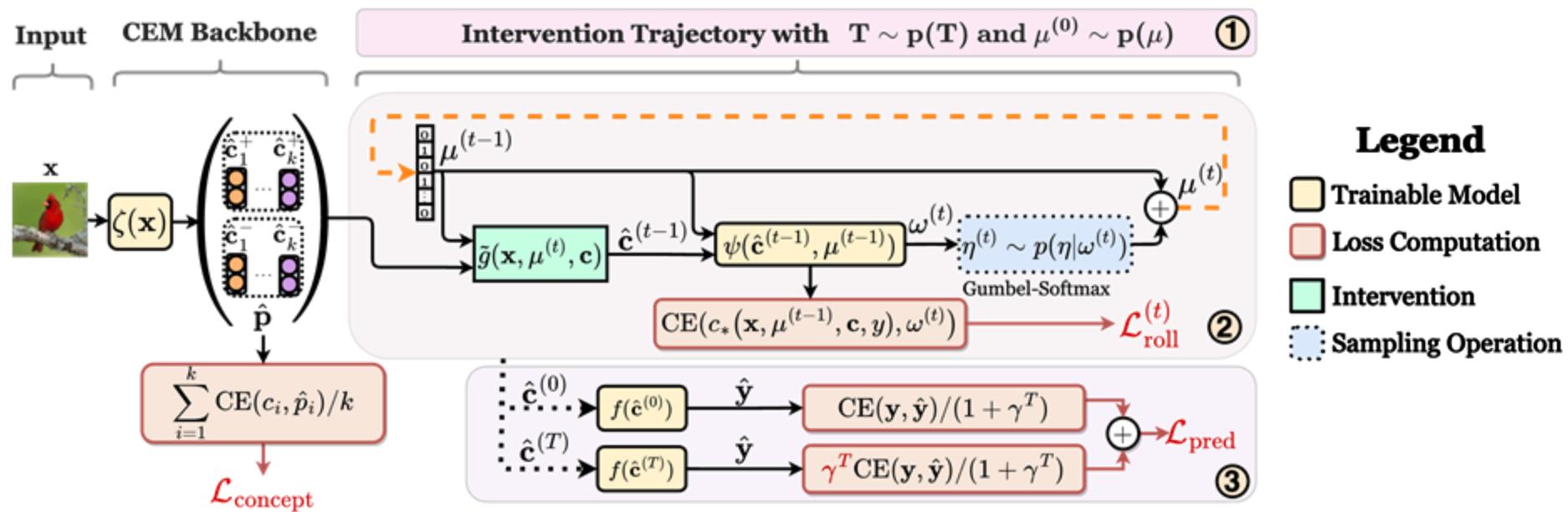
incorporate these insights into an end-to-end architecture that:

1. Introduces an **intervention-aware training loss** that encourages receptiveness to concept interventions at test-time
2. Learns **an efficient intervention policy** in an end-to-end fashion.



HOW TO TRAIN YOUR INTCEM?

This can be done using an end-to-end neural architecture:

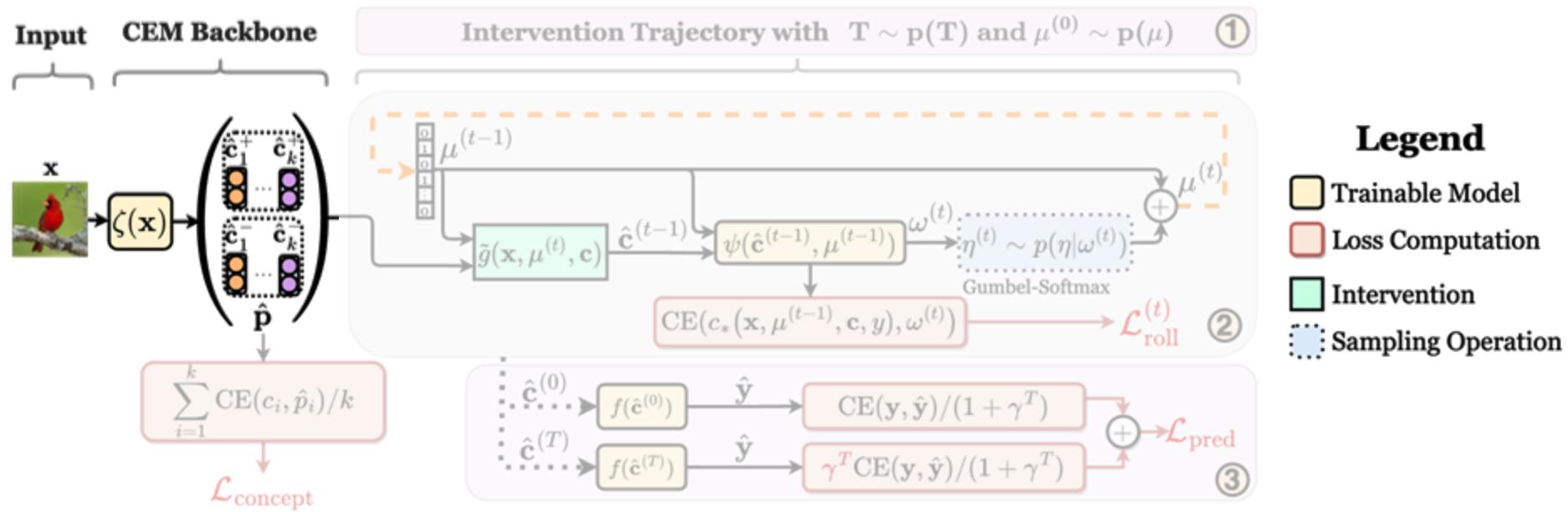


[1] Espinosa Zarlenga et al. "Learning to receive help: Intervention-aware concept embedding models." NeurIPS (2023)



HOW TO TRAIN YOUR INTCEM?

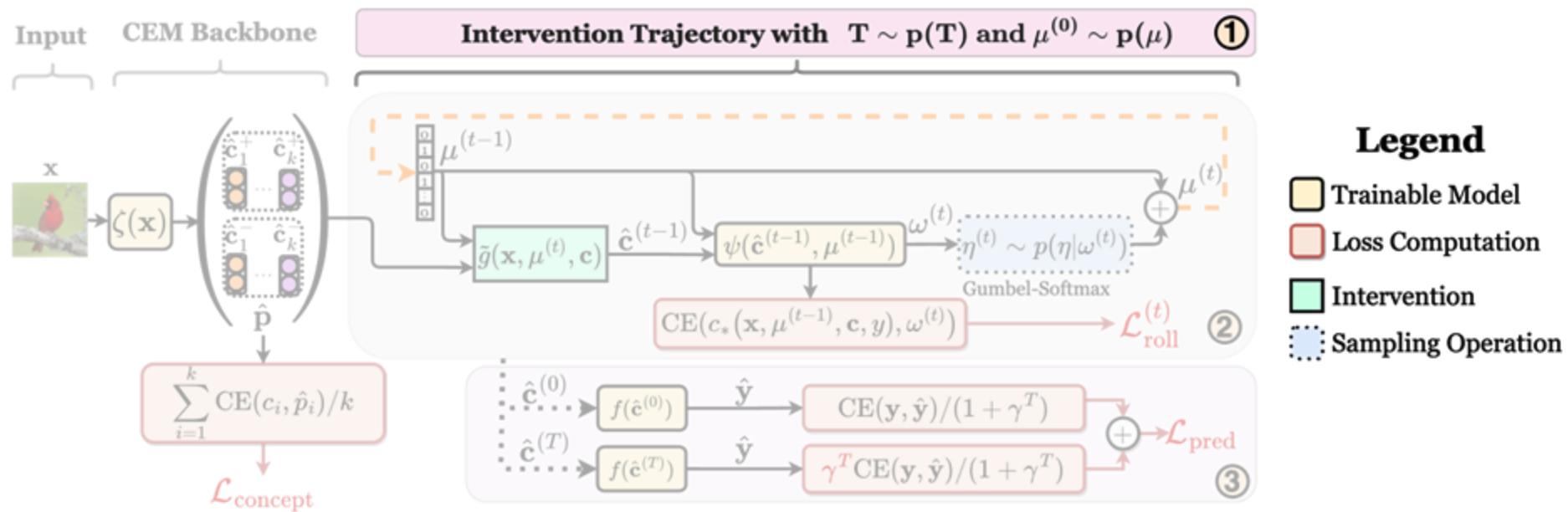
This can be done using an end-to-end neural architecture:



(1) Construct a positive and negative embedding for each training concept

HOW TO TRAIN YOUR INTCEM?

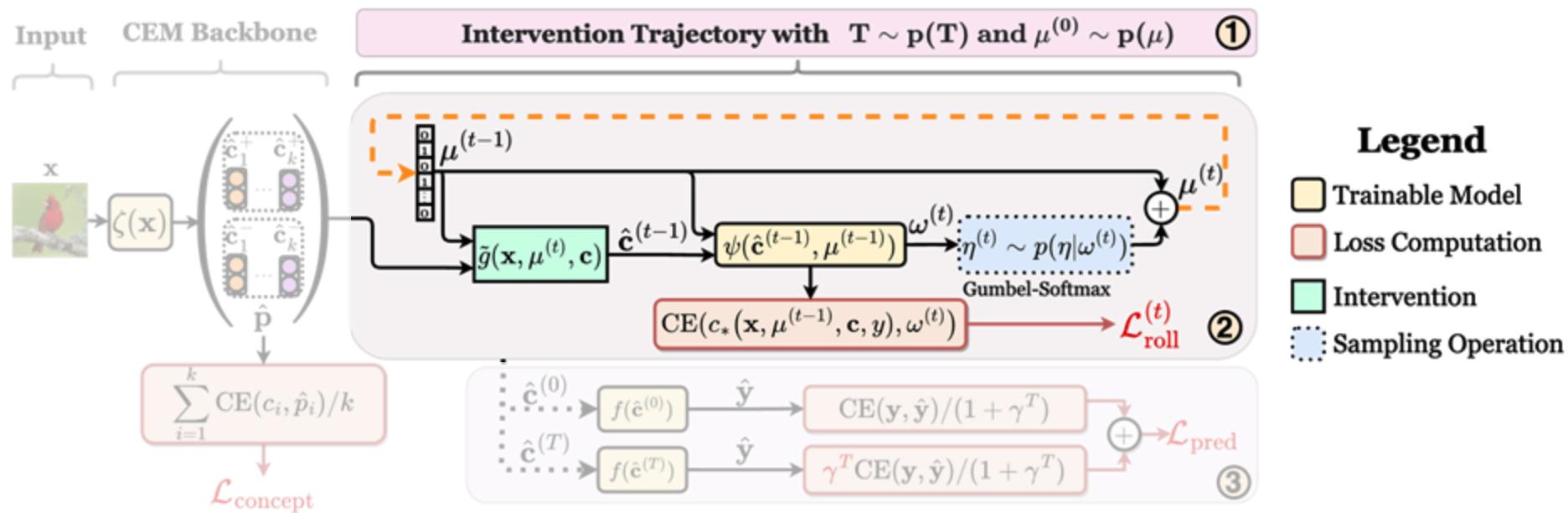
This can be done using an end-to-end neural architecture:



(2) Randomly select **a subset of concepts** which we will initially intervene on and a **number of interventions T** we will perform in this training step

HOW TO TRAIN YOUR INTCEM?

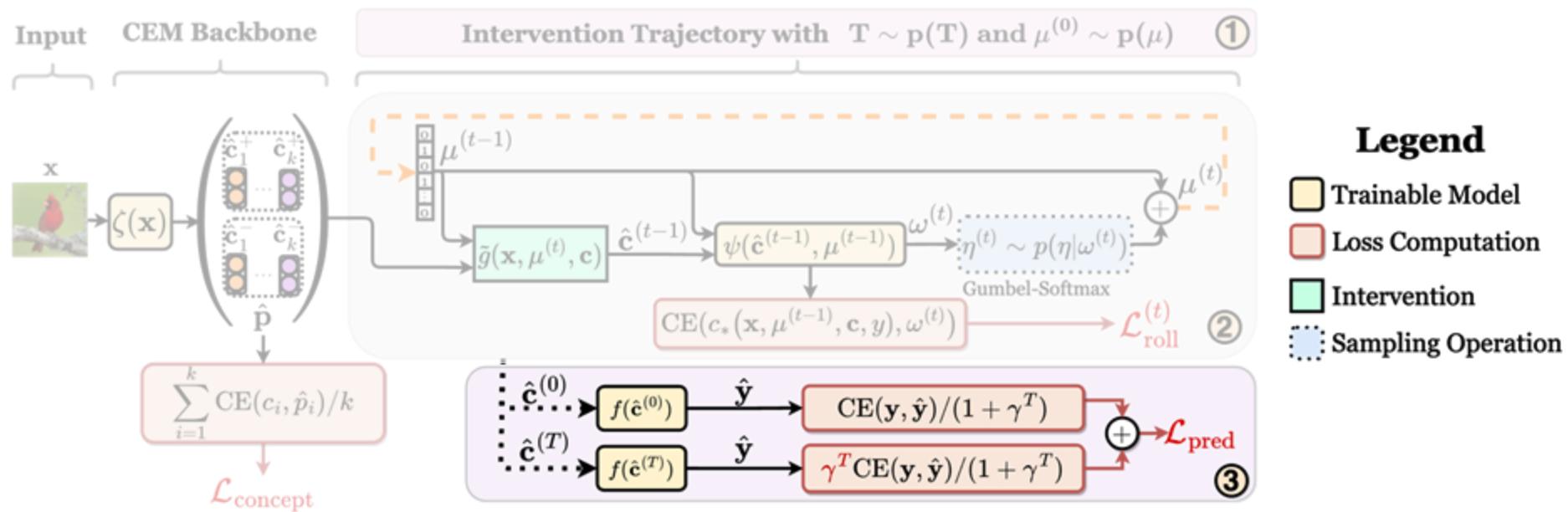
This can be done using an end-to-end neural architecture:



- (3) Recursively sample a **trajectory of T interventions** from this set using a **learnable intervention policy**. We train this policy to **align to the “oracle” optimal policy**.

HOW TO TRAIN YOUR INTCEM?

This can be done using an end-to-end neural architecture:



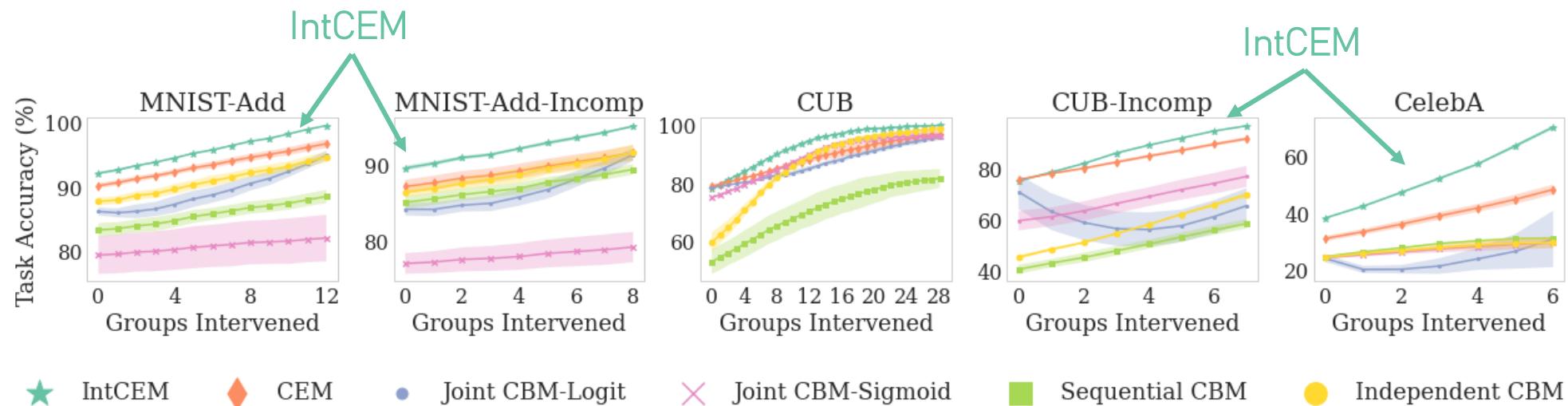
(4) Penalise the model more heavily for mispredicting the task label at the end of the intervention trajectory vs mispredicting the task label at the start of the trajectory

$$\mathbb{E}_{(x,c,y) \sim \mathcal{D}} \left[\frac{\mathcal{L}_{\text{task}}(y, f(\hat{c}^{(0)})) + \gamma^T \mathcal{L}_{\text{task}}(y, f(\hat{c}^{(T)}))}{1 + \gamma^T} \right]$$

WHAT DOES ALL OF THIS GIVE YOU?

WHAT DOES ALL OF THIS GIVE YOU?

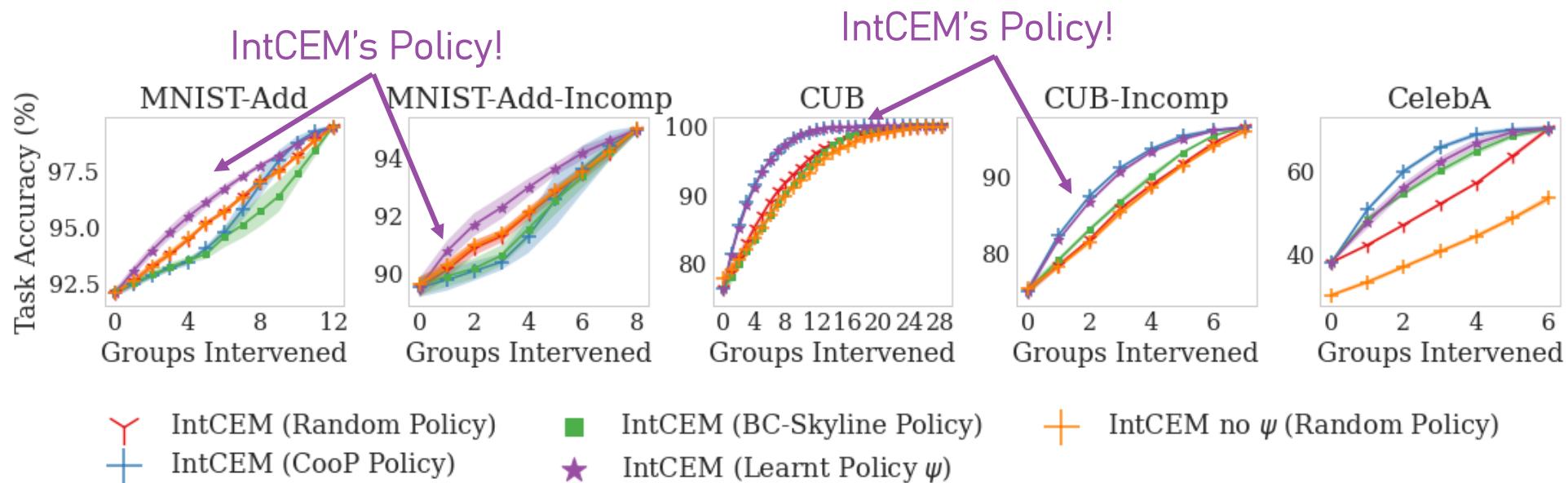
(1) A model that is much better at receiving test-time feedback
even if concepts are intervened in a random order



Up to 9% in absolute improvement when 25% of concepts are randomly selected to be intervened on!

WHAT DOES ALL OF THIS GIVE YOU?

(2) An **efficient intervention policy** that selects useful concepts to intervene on next

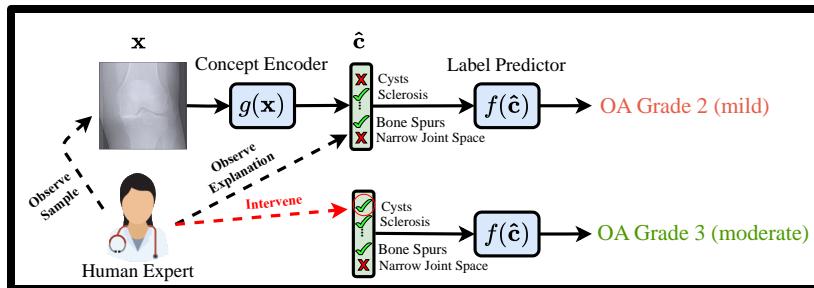


CRITICAL LIMITATIONS OF INTERVENTIONS

When intervening, we assume that concept interventions are:

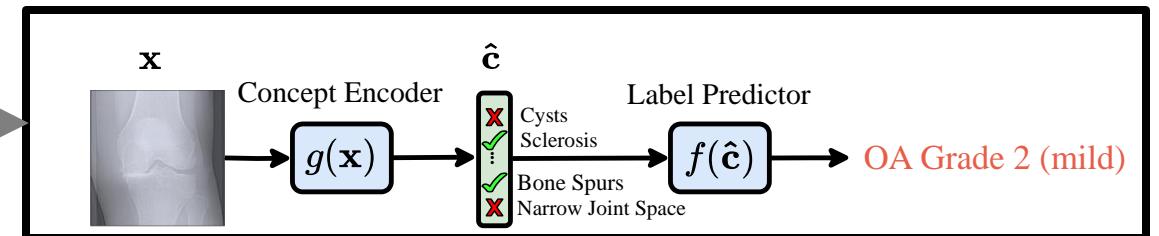
1. **Transient**: after an intervention is made, it is **forgotten**

If an intervention is made on a sample x



Some time later...

The same mistake will be made if x is seen again



CRITICAL LIMITATIONS OF INTERVENTIONS

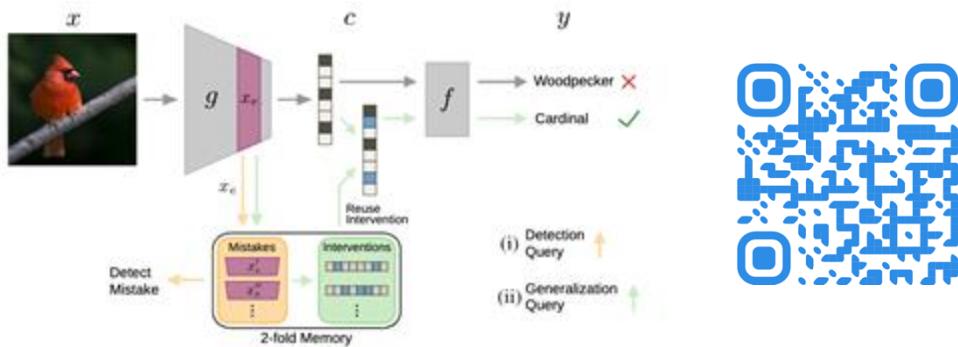
When intervening, we assume that concept interventions are:

1. **Transient**: after an intervention is made, it is **forgotten**
2. **Independent**: intervening on concept c_i will **not affect other concepts' values**

RELAXING KEY ASSUMPTIONS

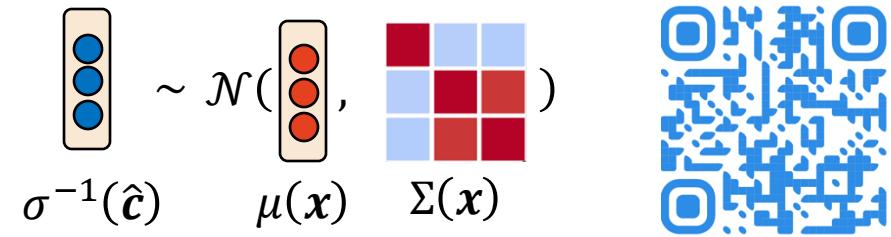
These constraints can be relaxed via clever modelling:

Concept Bottleneck Memory Models



Addresses: Transient nature of a concept intervention
Approach: Introduce a **learnable memory module** that keeps previously seen interventions and re-applies them in the future.

Stochastic Concept Bottleneck Models



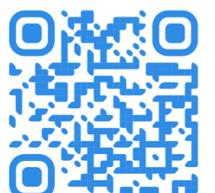
Addresses: the assumption that concepts are independent
Approach: Model the predicted concept logits as a **normal distribution** with a (learnable) non-diagonal covariance.

CAN INTERVENTIONS EXTEND BEYOND CBMS?

So far, the **concept intervention** strategies we have considered require one to operate on a CBM-like model

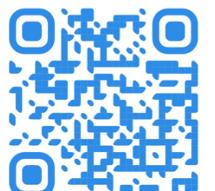
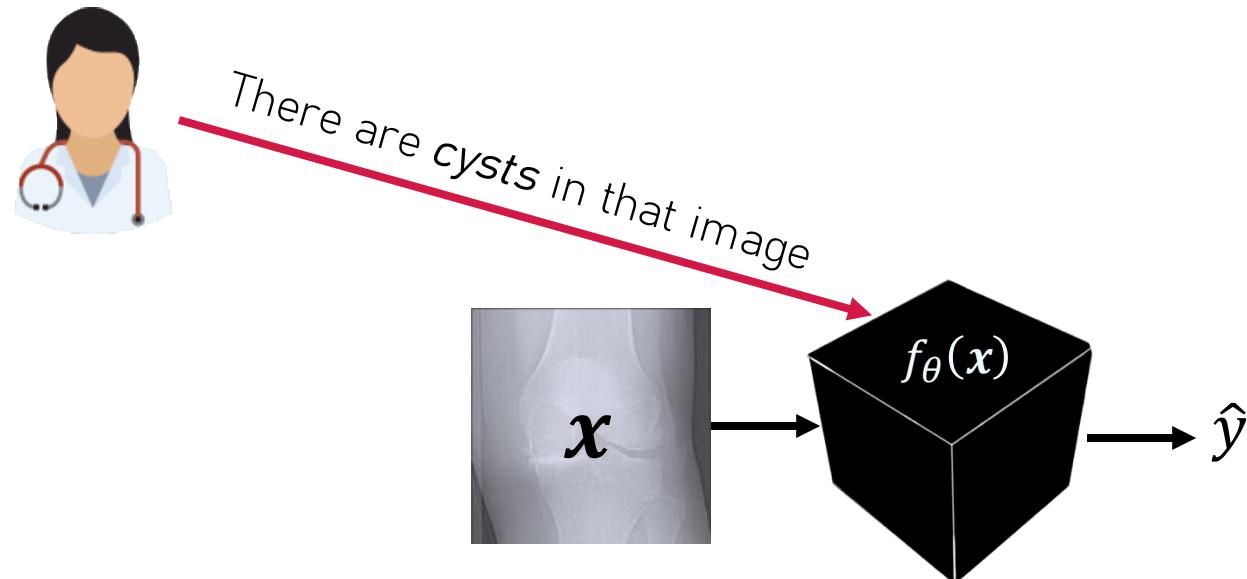
Could we potentially extend these ideas to models beyond CBMs?

[1] Laguna, Marcinkevičs, et al. "Beyond Concept Bottleneck Models: How to Make Black Boxes Intervenable?" NeurIPS (2024).



INJECTING KNOWLEDGE TO BLACK BOXES

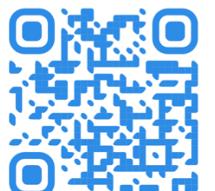
Given a black-box model $f_{\theta}(x) = g_{\psi}(h_{\phi}(x))$ and a test sample x , we may want to inject knowledge about the presence or absence of a concept in x at test time



INJECTING KNOWLEDGE TO BLACK BOXES

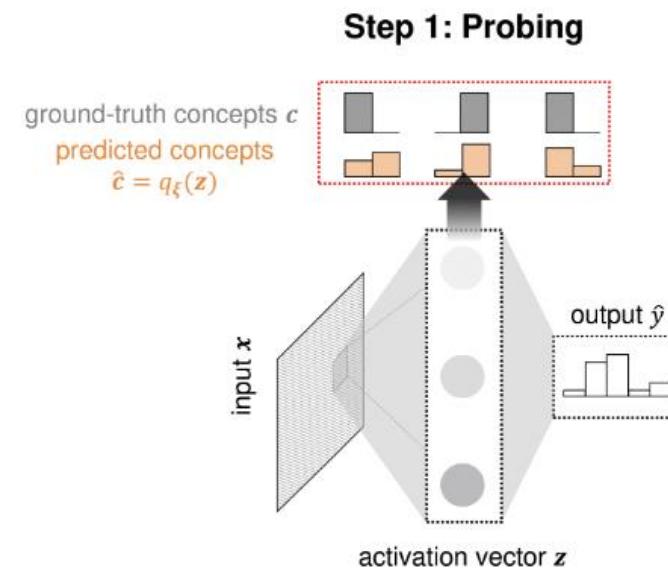
Given a black-box model $f_{\theta}(\mathbf{x}) = g_{\psi}(h_{\phi}(\mathbf{x}))$ and a test sample \mathbf{x} , we may want to inject knowledge about the presence or absence of a concept in \mathbf{x} at test time

If we have a **concept-annotated validation set** $\{(\mathbf{x}^{(i)}, \mathbf{c}^{(i)}, y^i)\}_i$, we can do this!



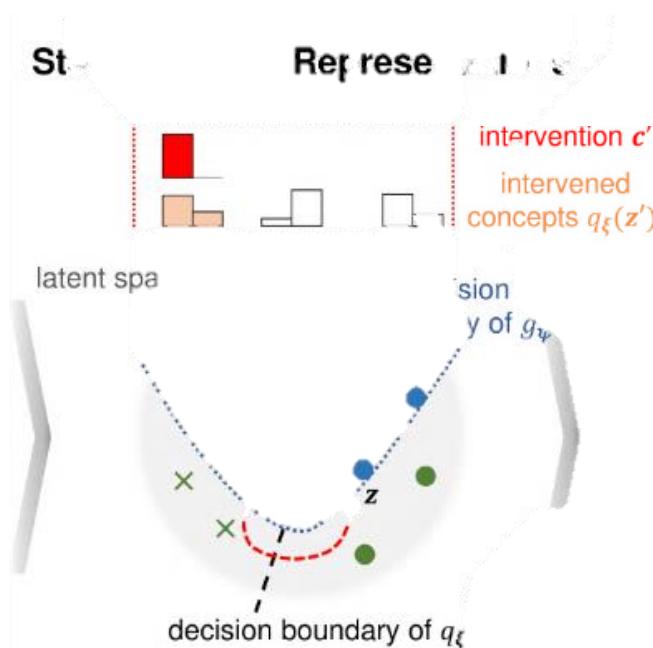
BLACK-BOX INTERVENABILITY: PROBING

We first **learn a multivariate probe** $\hat{\mathbf{c}} = \xi(\mathbf{z})$ that predicts **all** concepts given the latent space $\mathbf{z} = h_\phi(x)$ using the **annotated validation set**



BLACK-BOX INTERVENABILITY: EDITING

Given **user-provided concept labels** c' for sample \mathbf{x} , we edit the representation $\mathbf{z} = h_\phi(\mathbf{x})$ so that it **maps to** c' as predicted by the probe $\xi(\mathbf{z})$

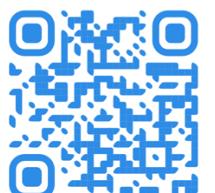


We find a new latent representation \mathbf{z}' by solving

$$\arg \min_{\mathbf{z}'} \lambda \mathcal{L}^c(q_\xi(\mathbf{z}'), \mathbf{c}') + d(\mathbf{z}, \mathbf{z}')$$

Concept alignment
Make the latent representation map to the set of user-provided concepts c'

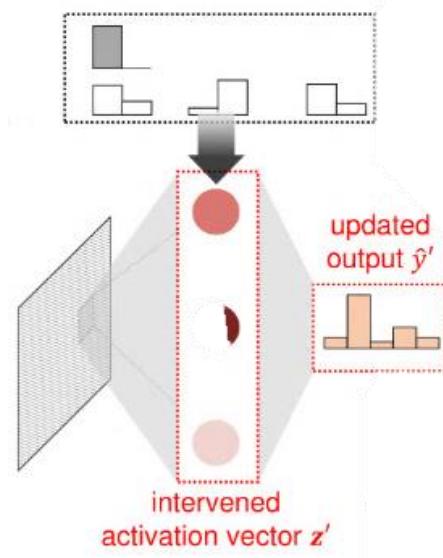
Distance Penalty
Keep the new latent representation as close as possible to the original



BLACK-BOX INTERVENABILITY: OUTPUT

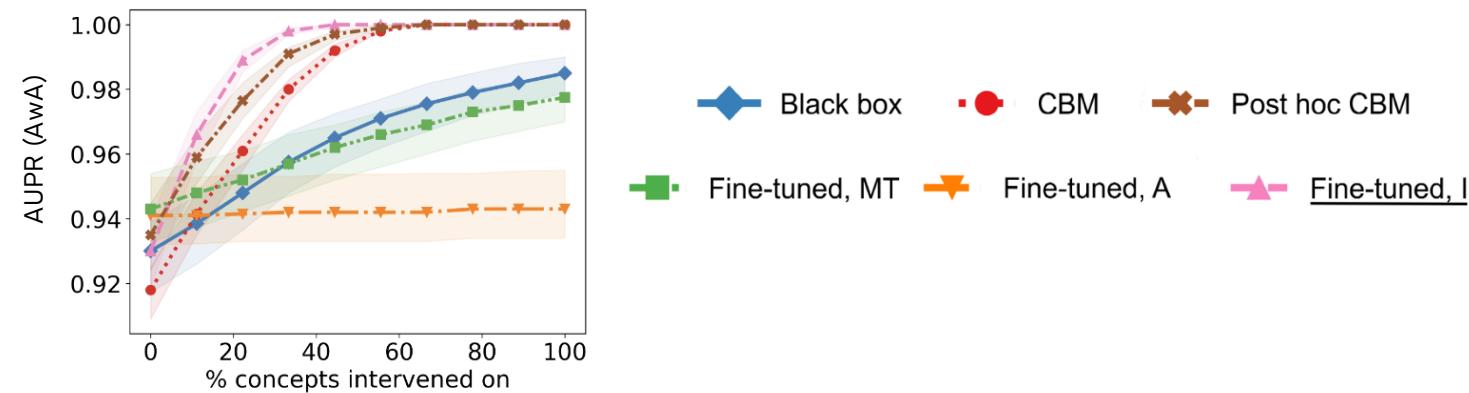
Finally, fed the edited representation z' to the second part of the DNN to obtain an updated prediction $\hat{y}' = g_\psi(z')$

Step 3: Updating Output



WHAT THIS GIVES YOU

This process allows you to **improve the task accuracy of a black-box model when you have extra test-time knowledge** in the form of concepts labels



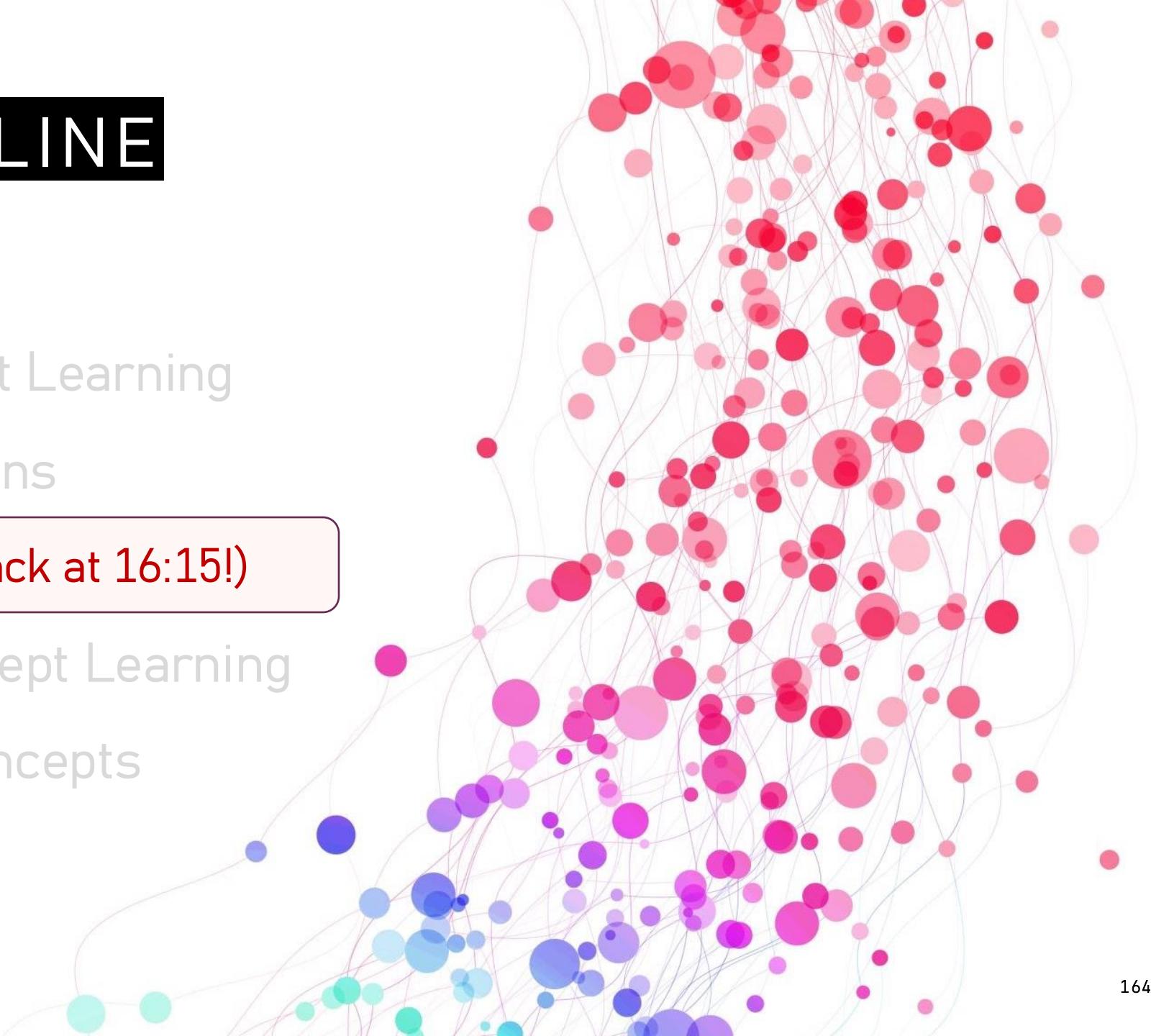
More importantly, you **can fine-tune a model** to be more receptive to this type of interventions by directly optimizing for an edit's positive effect

[1] Laguna, Marcinkevičs, et al. "Beyond Concept Bottleneck Models: How to Make Black Boxes Intervenable?" NeurIPS (2024).



TUTORIAL OUTLINE

1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
- 4. Q&A + Break (Back at 16:15!)**
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



Q&A + BREAK

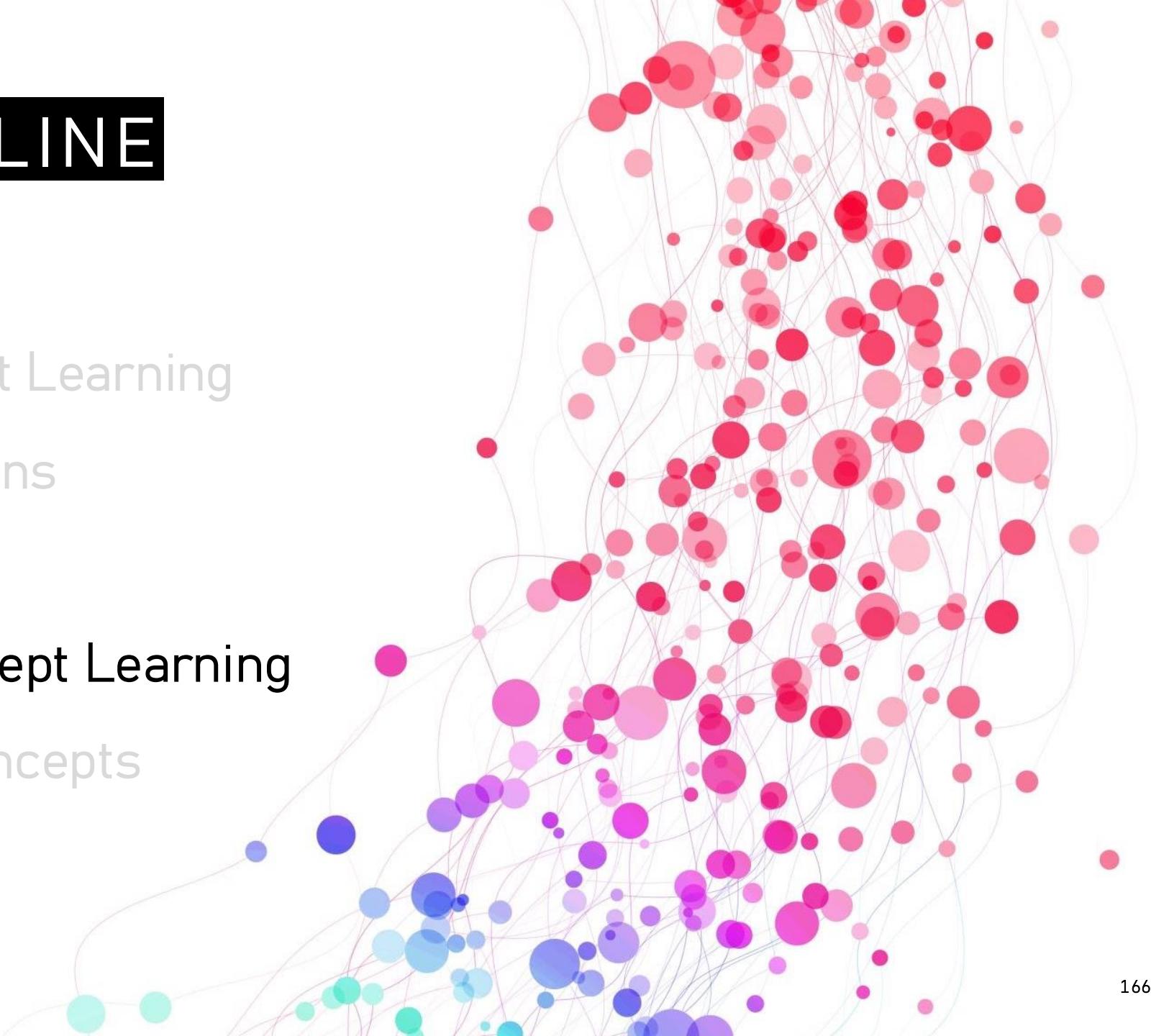


conceptlearning.github.io/

Back at 16:15 for part III!

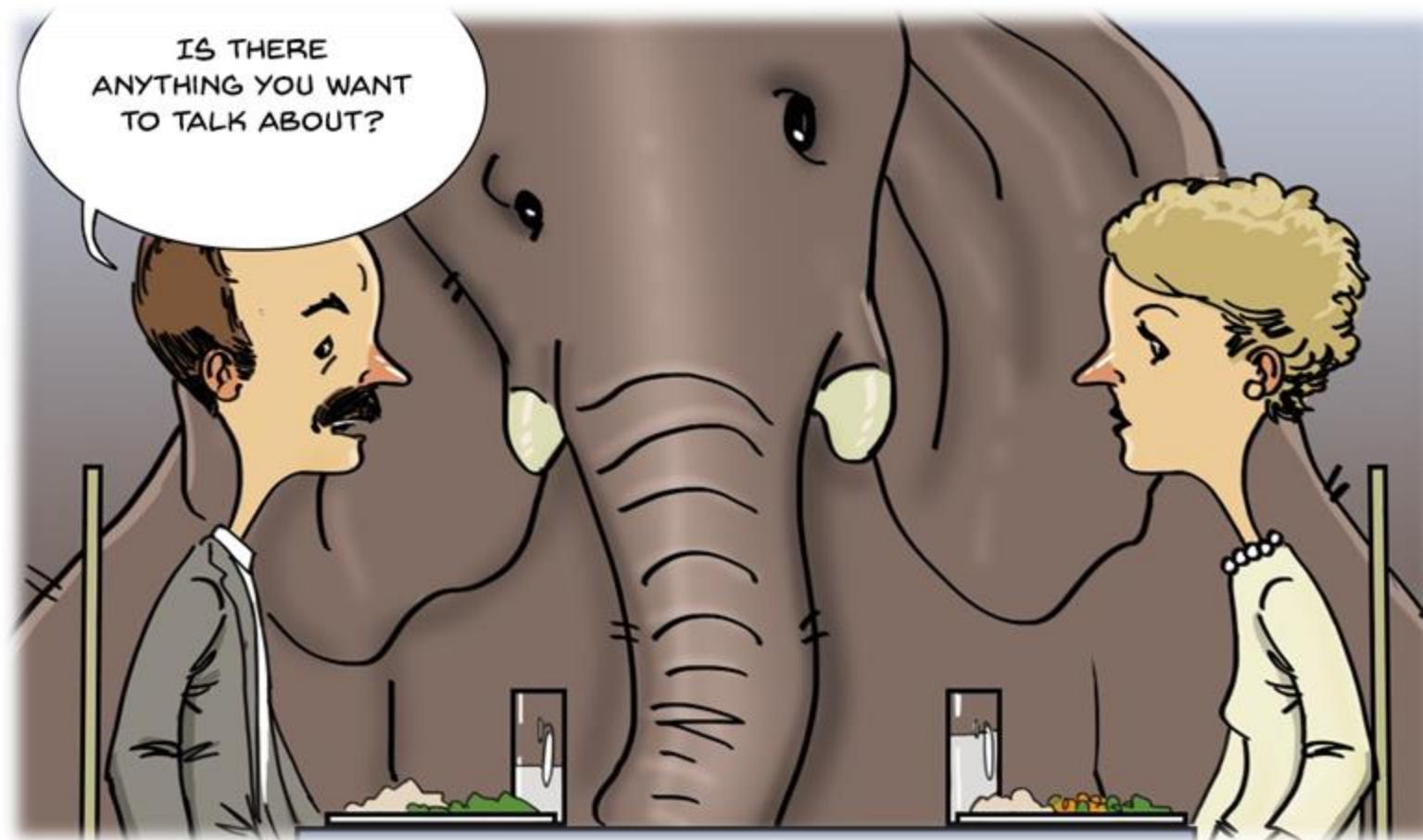
TUTORIAL OUTLINE

1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. **Unsupervised Concept Learning**
6. Reasoning With Concepts
7. Future Directions
8. Q&A



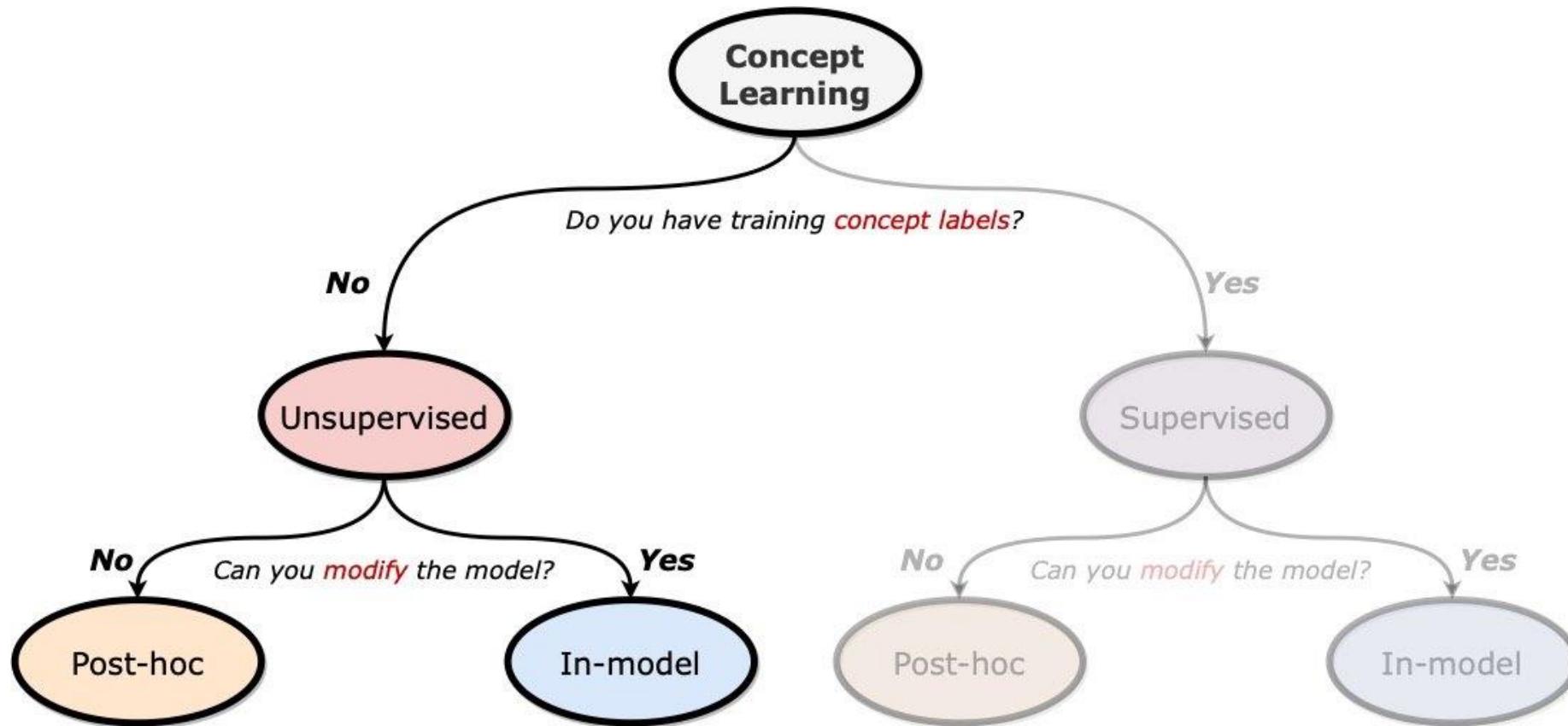
THE COST OF BEING GREAT

What if you **don't have access to concept supervisions?**



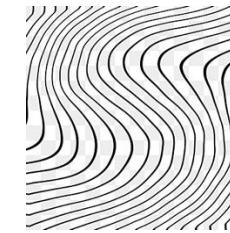
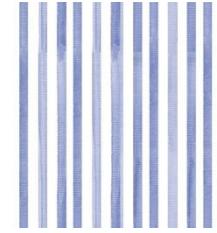
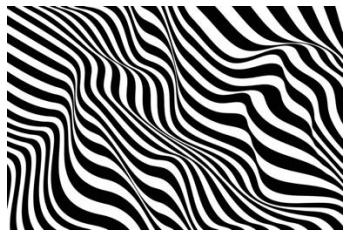
THE COST OF BEING GREAT

What if you **don't have access to concept supervisions?**



THE COST OF BEING GREAT

T-CAV requires large sets of examples of each concept of interest:



For example, when finding the influence of the concept “stripes” for a DNN, T-CAV requires a set of samples that all have the concept “stripes”

But, obtaining concept labels can be **expensive** and **intractable**

THE COST OF BEING GREAT

T-CAV requires large sets of examples of each concept of interest:

Can we extract patches automatically?

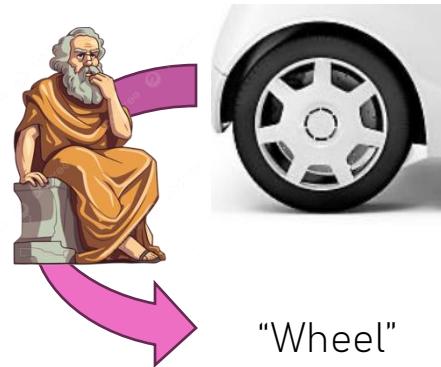
For example, when finding the instances of the concept "stripes" for a DNN, T-CAV requires a set of samples that all have the concept "stripes"

But, obtaining concept labels can be **expensive** and **intractable**

AUTOMATIC CONCEPT EXTRACTION (ACE)

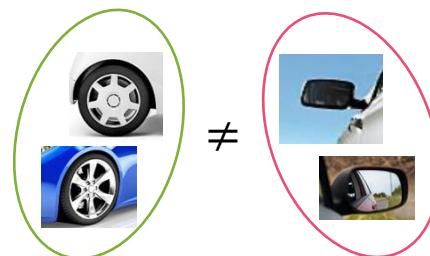
Desiderata: We would like to discover concepts / patches that are:

Meaningful



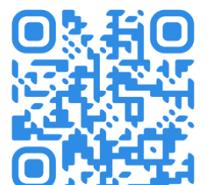
"Wheel"

Coherent



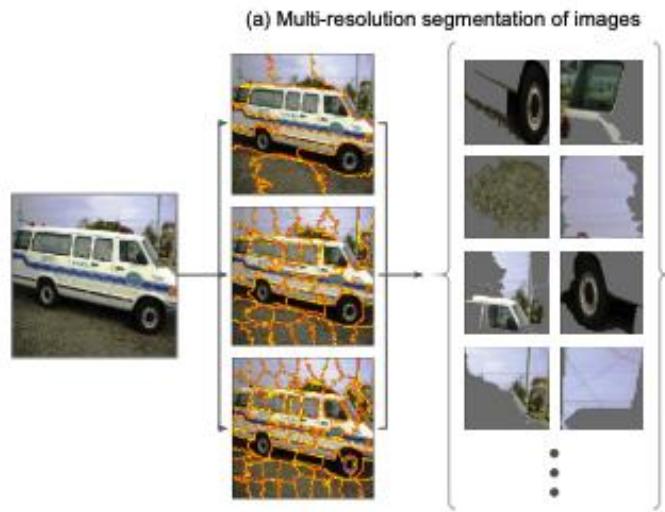
Important

$$\text{"Car"} \sim f(\text{Wheel})$$



AUTOMATIC CONCEPT EXTRACTION (ACE)

Proposed Solution



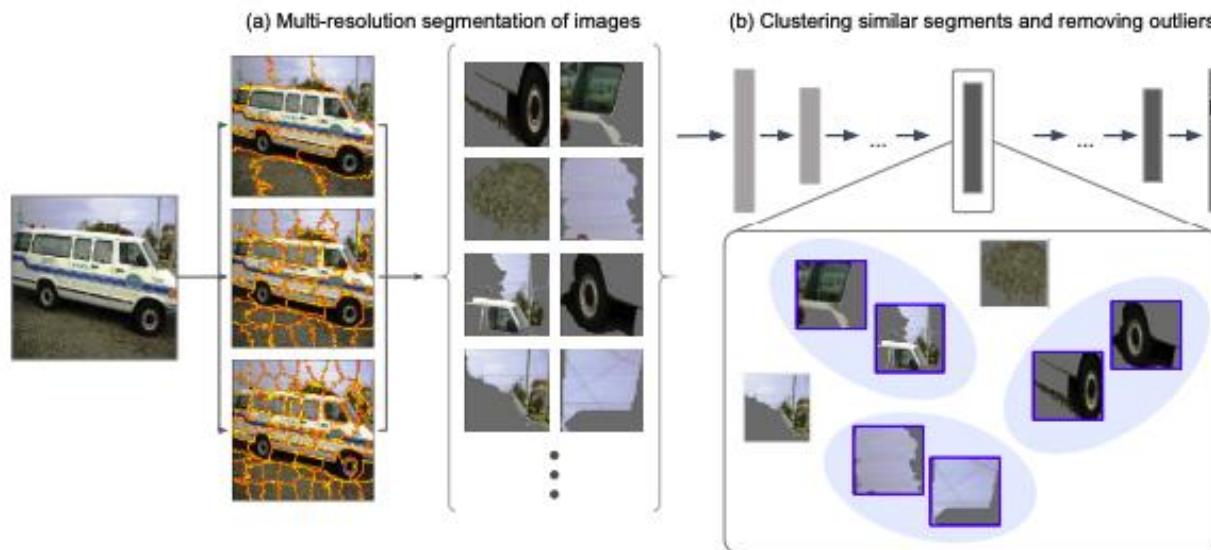
Step 1: Multi-resolution segmentation (**why?** concepts have different granularities)

Desiderata enforced: meaningfulness



AUTOMATIC CONCEPT EXTRACTION (ACE)

Proposed Solution



Step 2: cluster extracted segments using a hidden layer (**which one?**) of a CNN as a feature extractor (**why?** ensure invariances). Then get rid of outliers (**why?** noisy!).

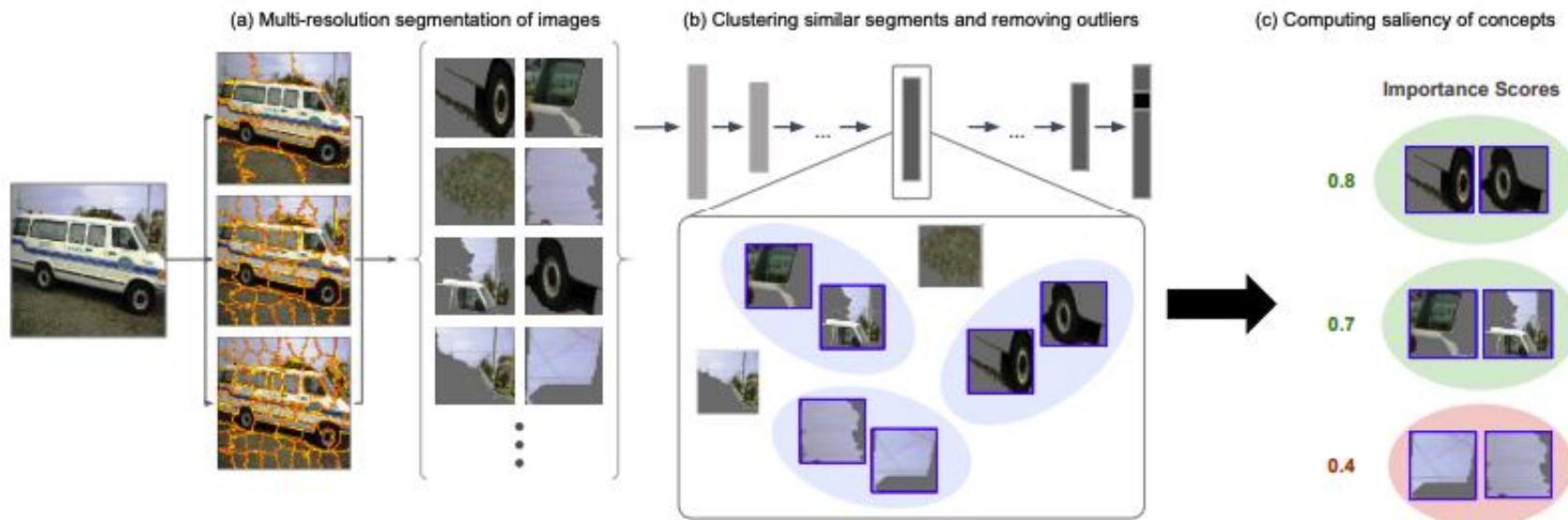
Desiderata enforced: coherence

[1] Ghorbani et al. "Towards automatic concept-based explanations." *Advances in Neural Information Processing Systems 32* (2019).



AUTOMATIC CONCEPT EXTRACTION (ACE)

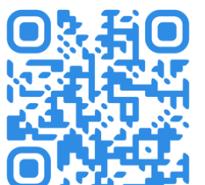
Proposed Solution



Step 3: use T-CAV with the newly discovered concepts to explain the prediction of the sample of interest!

Desiderata enforced: importance

[1] Ghorbani et al. "Towards automatic concept-based explanations." *Advances in Neural Information Processing Systems 32* (2019).



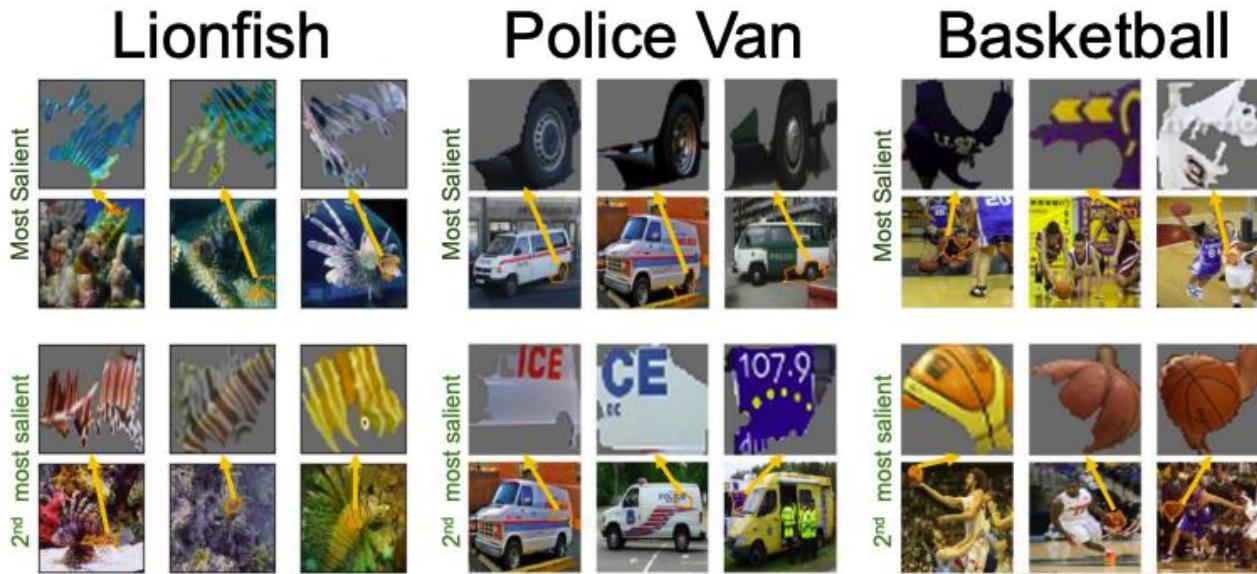
AUTOMATIC CONCEPT EXTRACTION (ACE)



What are the most **salient discovered concepts** for some of the ImageNet classes?



AUTOMATIC CONCEPT EXTRACTION (ACE)



What are the most **salient discovered concepts** for some of the ImageNet classes?

ACE has also been generalised to learn concepts in Graph Neural Networks in GCExplainer (Magister et al. 2021) [2]



[1] Ghorbani et al. "Towards automatic concept-based explanations" *Advances in Neural Information Processing Systems* 32 (2019)

[2] Magister et al. "GCExplainer: Human-in-the-loop Concept-based Explanations for Graph Neural Networks" *arXiv preprint arXiv:2107.11889* (2021)

AUTOMATIC CONCEPT EXTRACTION (ACE)

ACE's hyperparameters and processing steps have **several limitations**:

1. We can never be certain that we properly **cover all useful concepts**



Important concepts for **underrepresented populations** could be removed as outliers!

AUTOMATIC CONCEPT EXTRACTION (ACE)

ACE's hyperparameters and processing steps have **several limitations**:

1. We can never be certain that we properly **cover all useful concepts**
2. We won't detect concepts that **interact non-linearly** with the output labels

Looking at the gradients provides understanding of **local (linear) sensitivity**

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

AUTOMATIC CONCEPT EXTRACTION (ACE)

ACE's hyperparameters and processing steps have **several limitations**:

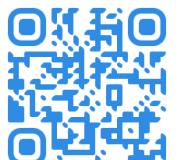
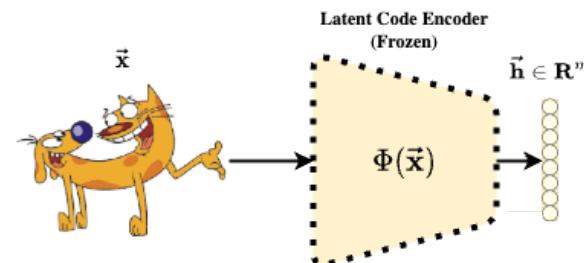
1. We can never be certain that we properly **cover all useful concepts**
2. We won't detect concepts that **interact non-linearly** with the output labels

Can we optimize accounting for concept usefulness and non-linear interactions?

COMPLETENESS-AWARE CONCEPT EXTRACTION

Proposed Solution

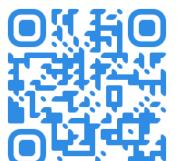
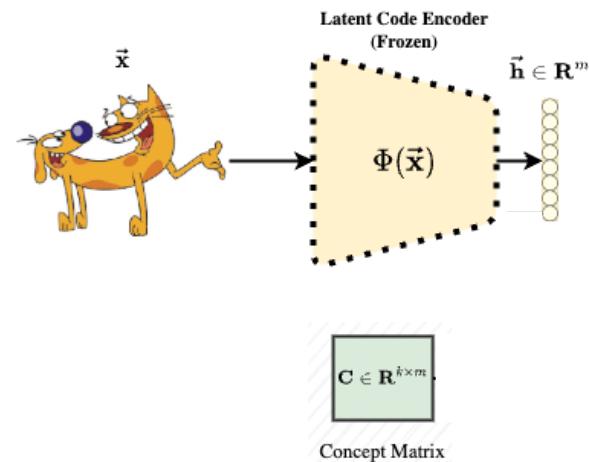
Step 1: project the input sample to DNN's intermediate hidden layer $\Phi(\vec{x})$



COMPLETENESS-AWARE CONCEPT EXTRACTION

Proposed Solution

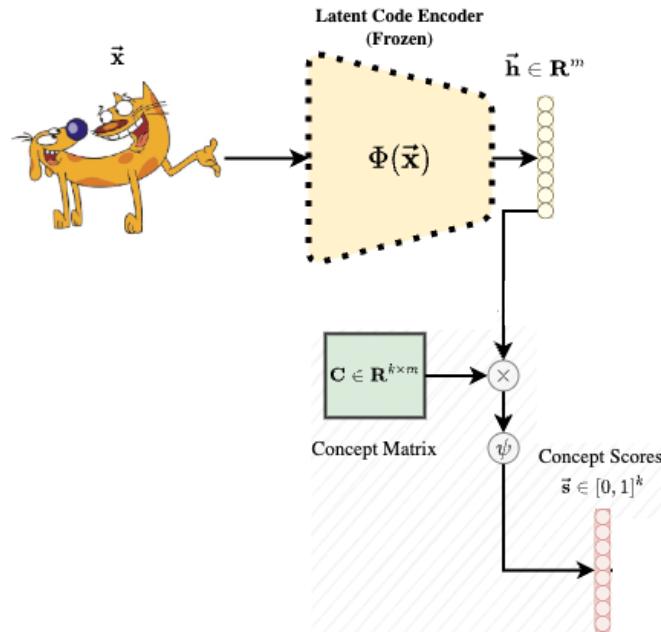
Step 2: randomly initialize a latent, learnable concept bank of k concepts $\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k]^T$



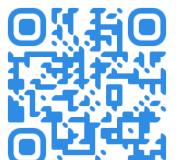
COMPLETENESS-AWARE CONCEPT EXTRACTION

Proposed Solution

Step 3: compute a set of concept scores by projecting the input embedding into the concept space



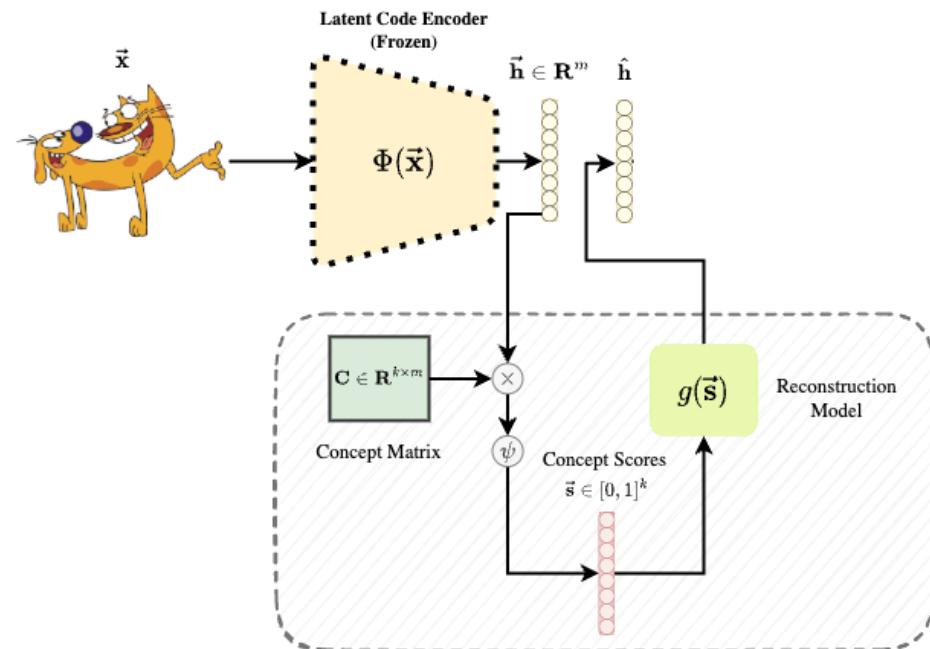
[1] Yeh et al. "On completeness-aware concept-based explanations in deep neural networks" NeurIPS (2020)



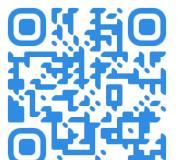
COMPLETENESS-AWARE CONCEPT EXTRACTION

Proposed Solution

Step 4: pass the concepts scores to a learnable model $g(\vec{s}) = \hat{h}$ that aims to reconstruct \vec{h} from \vec{s}



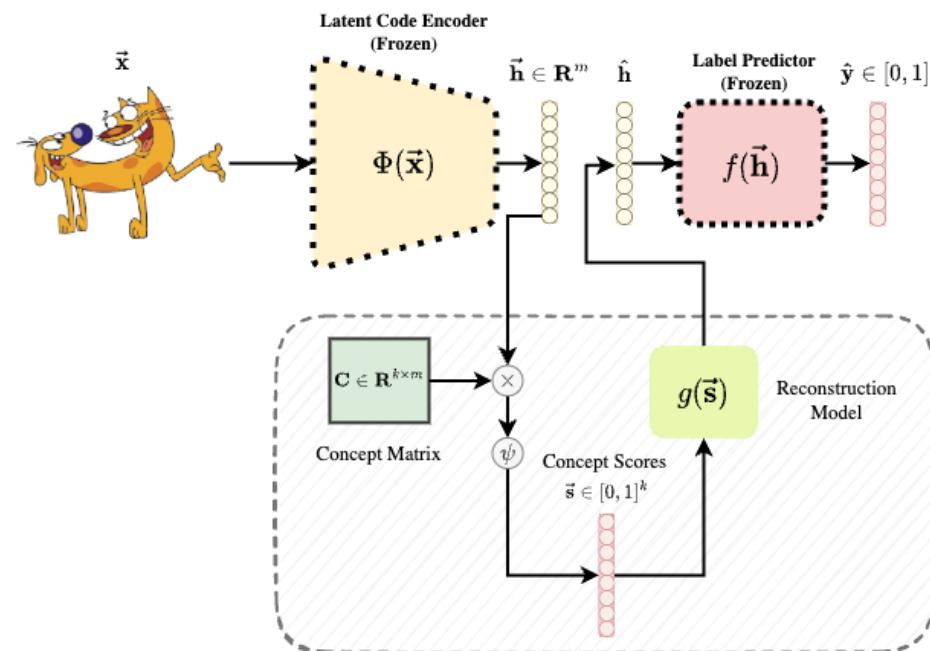
[1] Yeh et al. "On completeness-aware concept-based explanations in deep neural networks" NeurIPS (2020)



COMPLETENESS-AWARE CONCEPT EXTRACTION

Proposed Solution

Step 5: use $\hat{\mathbf{h}}$ as the reconstructed hidden layer and predict an output class using f



[1] Yeh et al. "On completeness-aware concept-based explanations in deep neural networks" NeurIPS (2020).



COMPLETENESS-AWARE CONCEPT EXTRACTION

Proposed Solution

Step 6: maximise a “concept completeness score”

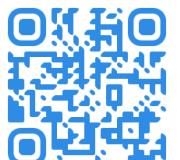
DNN's accuracy via concept projection

$$n_f(\mathbf{c}_1, \dots, \mathbf{c}_m) = \frac{\sup_g \mathbb{P}_{x,y \sim V} \left[y = \operatorname{argmax}_{y'} f_{y'} \left(g(\mathbf{C} \phi(x)) \right) \right] - a_r}{\mathbb{P}_{x,y \sim V} \left[y = \operatorname{argmax}_{y'} f_{y'}(x) \right] - a_r}$$

Original DNN's accuracy

Score is ~ 1 if and only if the projection in the concept space preserves all the information needed to predict y !

[1] Yeh et al. "On completeness-aware concept-based explanations in deep neural networks." NeurIPS (2020).



COMPLETENESS-AWARE CONCEPT EXTRACTION

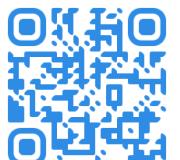
CCE further encourages discovered concepts to be:

1. **Coherent**: similar samples should remain close in concept-space
2. **Diverse**: concept vectors should be as distinct from each other as possible

Coherency Diversity

$$R(\mathbf{c}) = \lambda_1 \frac{\sum_{k=1}^m \sum_{\mathbf{x}_a^b \in T_{\mathbf{c}_k}} \Phi(\mathbf{x}_a^b) \cdot \mathbf{c}_k}{mK} - \lambda_2 \frac{\sum_{j \neq k} \mathbf{c}_j \cdot \mathbf{c}_k}{m(m-1)}$$

[1] Yeh et al. "On completeness-aware concept-based explanations in deep neural networks." NeurIPS (2020).



COMPLETENESS-AWARE CONCEPT EXTRACTION

And it can be applied to different data modalities!

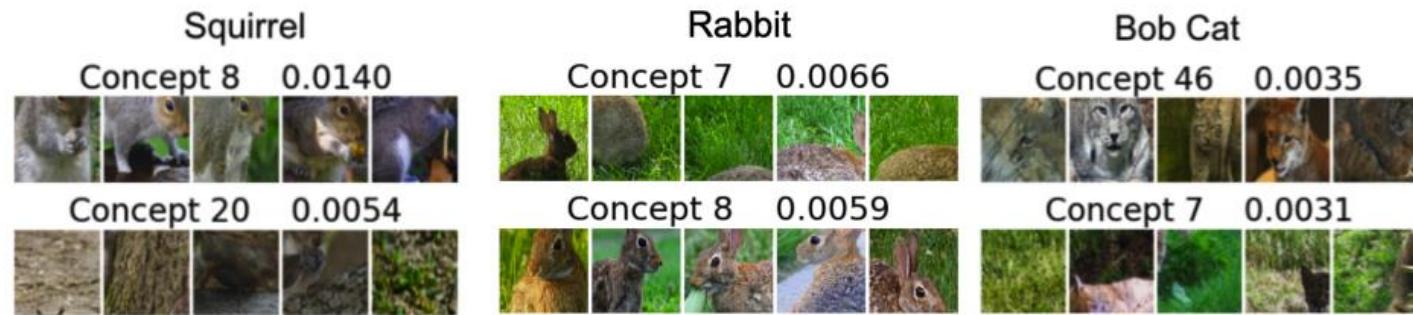


Table 2: The 4 discovered concepts and some nearest neighbors along with the most frequent words that appear in top-500 nearest neighbors.

Concept	Nearest Neighbors	Frequent words	ConceptSHAP
1	poorly constructed what comes across as interesting is the wasting my time with a comment but this movie awful in my opinion there were <UNK> and the	worst (168) ever (69) movie (61) seen (55) film (50) awful (42) time(40) waste (34) poorly (26) movies (24) films (18) long (17)	0.280
2	normally it would earn at least 2 or 3 <UNK> <UNK> is just too dumb to be called i feel like i was ripped off and hollywood	not (58) movie (39) make (25) too (23) film (22) even (19) like (18) 2 (16) never (14) minutes (13) 1 (12) doesn't (11)	0.306
3	remember awaiting return of the jedi with almost <UNK> better than most sequels for tv movies i hate male because marie has a crush on her attractive	movies (19) like (18) see (16) movie (15) love (15) good (12) character (11) life (11) little (10) ever (9) watch (9) first (9)	0.174
4	new <UNK> <UNK> via <UNK> <UNK> with absolutely hilarious homosexual and an italian clown <UNK> is an entertaining stephen <UNK> on the vampire <UNK> as a masterpiece	excellent (50) film (25) perfectly (19) wonderful (19) perfect (16) hilarious (15) best (13) fun (12) highly (11) movie (11) brilliant (9) old (9)	0.141



REMEMBER CONCEPT BOTTLENECKS?

We took care of T-CAV & friends...



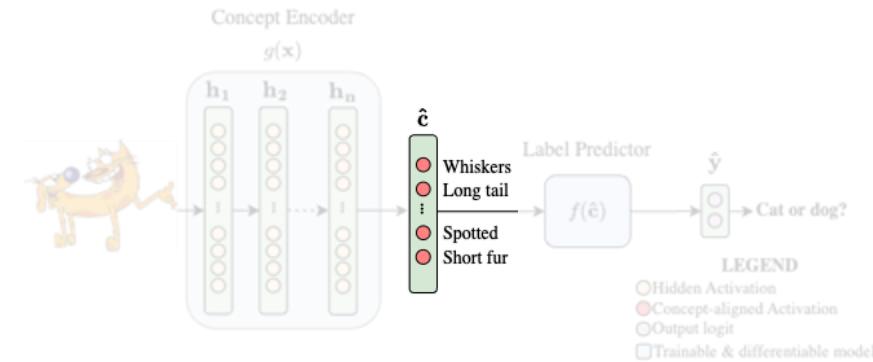
What about CBMs family?



LABEL-FREE CONCEPT BOTTLENECKS

Limitation Being Addressed

CBMs & co **require some known concepts**, or we have no bottleneck at all!



And post-hoc CBMs still **require one to know which concepts are potentially useful for a downstream task!**



ICLR23



CVPR23

[1] Oikarinen et al. "Label-Free Concept Bottleneck Models." ICLR (2023)

[2] Yang et al. "Language in a bottle: Language model guided concept bottlenecks for interpretable image classification." CVPR (2023)

LABEL-FREE CONCEPT BOTTLENECKS

Proposed Solution

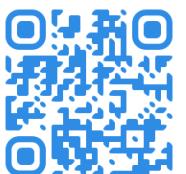
Why not **simply ask GPT** for a set of useful concepts for a specific class?



"List the most important features for recognizing something as a {class}:"



ICLR23



CVPR23

[1] Oikarinen et al. "Label-Free Concept Bottleneck Models." ICLR (2023)

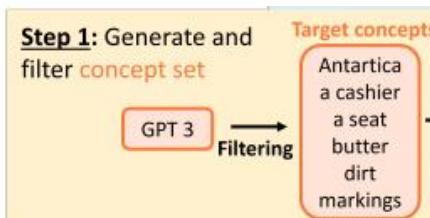
[2] Yang et al. "Language in a bottle: Language model guided concept bottlenecks for interpretable image classification." CVPR (2023)

LABEL-FREE CONCEPT BOTTLENECKS

Proposed Solution

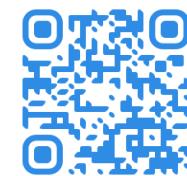
Step 1: Generate a concept set by "asking" an LLM

Label-free CBM



[1] Oikarinen et al. "Label-Free Concept Bottleneck Models." ICLR (2023)

[2] Yang et al. "Language in a bottle: Language model guided concept bottlenecks for interpretable image classification." CVPR (2023)



ICLR23

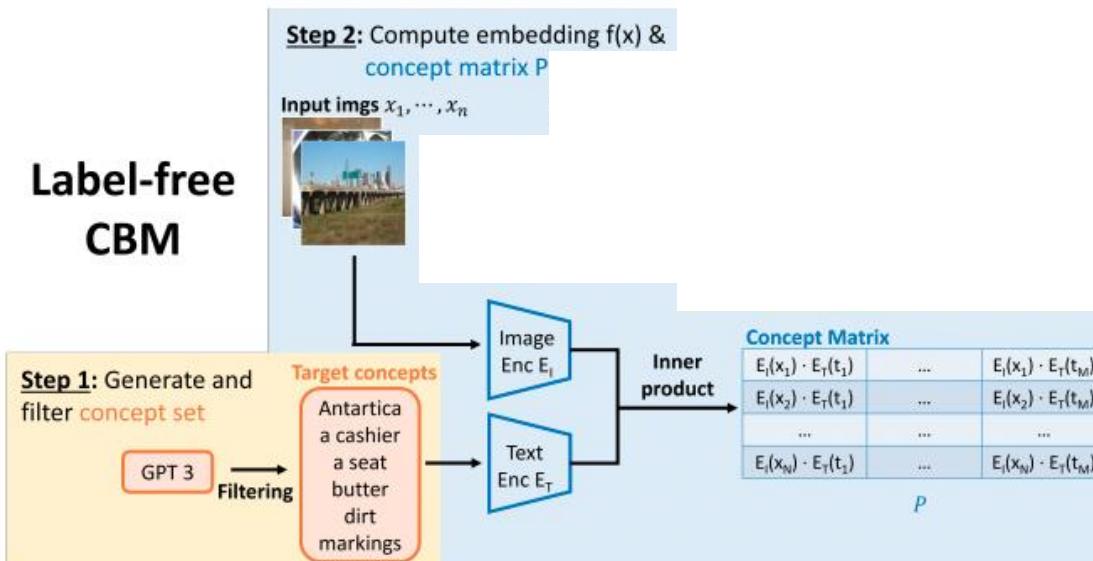


CVPR23

LABEL-FREE CONCEPT BOTTLENECKS

Proposed Solution

Step 2: Use multi-modal contrastive language model (e.g., CLIP) to compute similarity of image-text embeddings



[1] Oikarinen et al. "Label-Free Concept Bottleneck Models." ICLR (2023)

[2] Yang et al. "Language in a bottle: Language model guided concept bottlenecks for interpretable image classification." CVPR (2023)



ICLR23

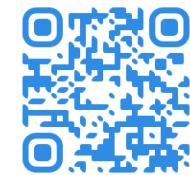
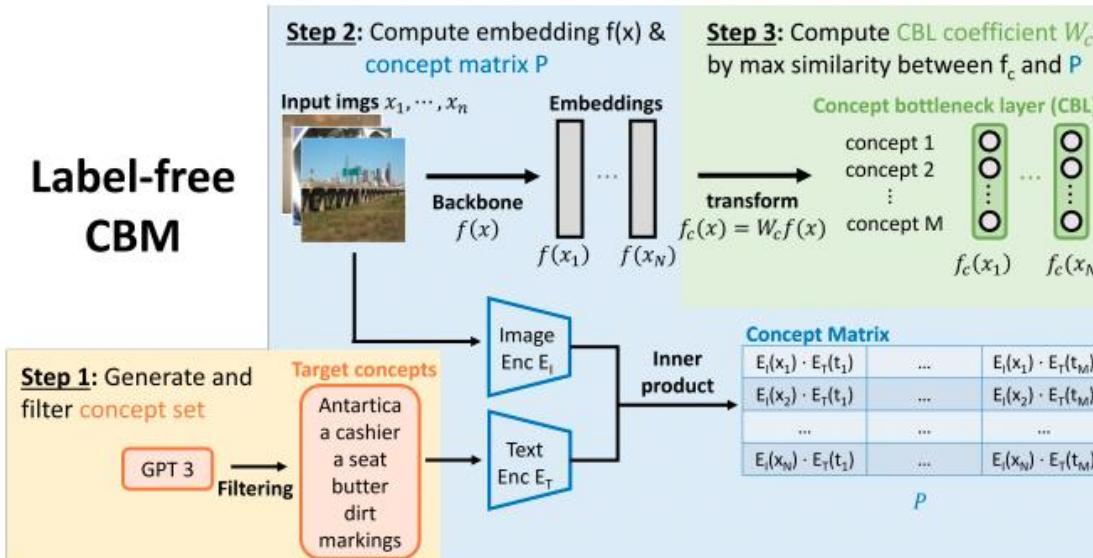


CVPR23

LABEL-FREE CONCEPT BOTTLENECKS

Proposed Solution

Step 3: Train DNN activations to **align with similarity scores** predicted by the contrastive LM



ICLR23

CVPR23

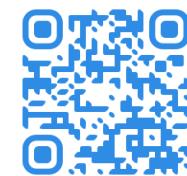
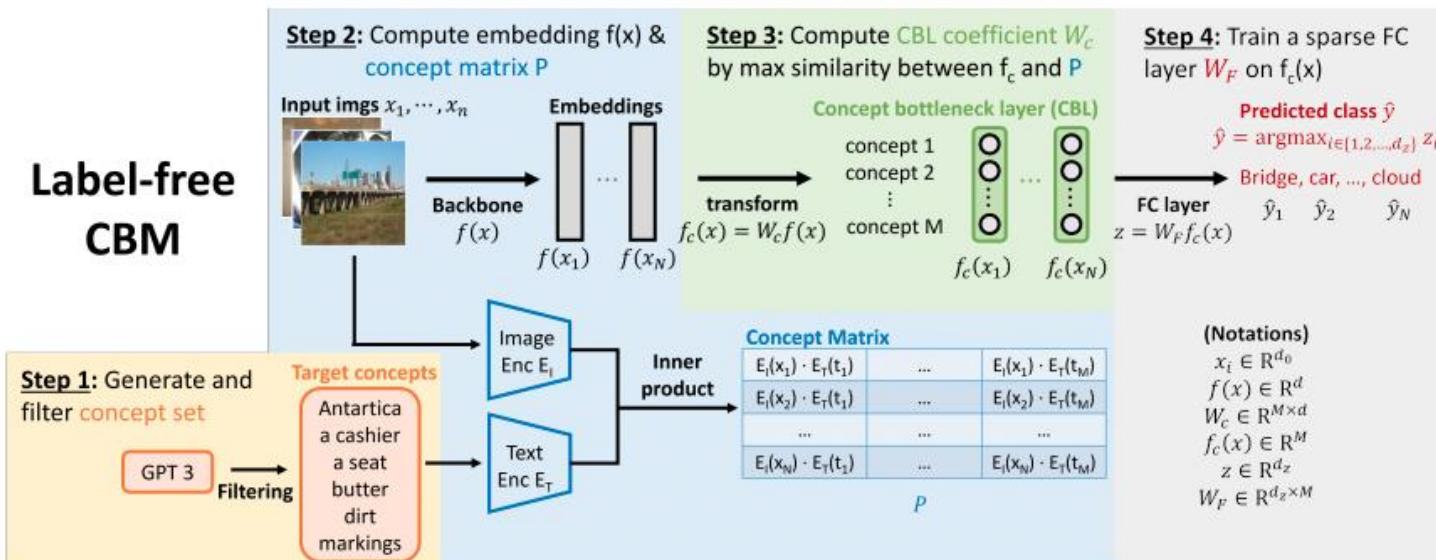
[1] Oikarinen et al. "Label-Free Concept Bottleneck Models." ICLR (2023)

[2] Yang et al. "Language in a bottle: Language model guided concept bottlenecks for interpretable image classification." CVPR (2023)

LABEL-FREE CONCEPT BOTTLENECKS

Proposed Solution

Step 4: Train a simple (linear) model to map predicted concept scores to tasks



ICLR23

CVPR23

[1] Oikarinen et al. "Label-Free Concept Bottleneck Models." ICLR (2023)

[2] Yang et al. "Language in a bottle: Language model guided concept bottlenecks for interpretable image classification." CVPR (2023)

STRIPPING CBMS TO THEIR BONES

What if we want a CBM, but... we **don't have**:

- Concept supervisions
- Pre-trained contrastive LMs

What's left??

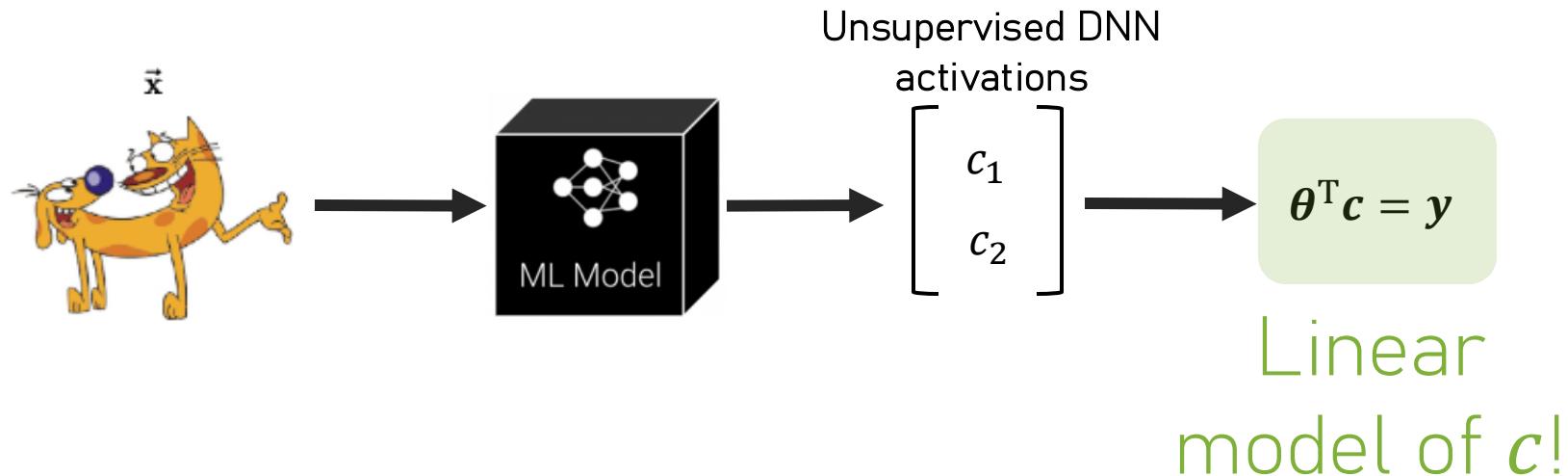


STRIPPING CBMS TO THEIR BONES

What if we want a CBM, but... we **don't have**:

- Concept supervisions
- Pre-trained contrastive LMs

What's left??



STRIPPING CBMS TO THEIR BONES

What if we want a CBM, but... we **don't have**:

- Concept supervisions
- Pre-trained contrastive LMs

What's left??

*Can we make a DNN behave
like a proper linear model?*

LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

1. **[Expressiveness]** The relevance weights used to make the output prediction must be able to **dynamically adapt depending on the input**:

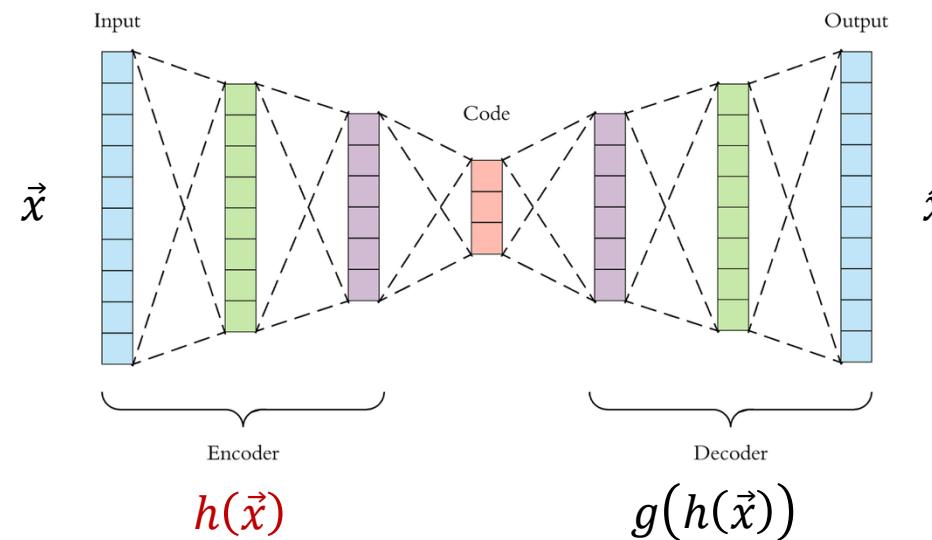
$$\begin{array}{ll} \text{Linear Model Output: } & f(\vec{x}) = \theta^T \vec{x} \\ \text{"Linear-ish DNN" Model output: } & f(\vec{x}) = \theta(\vec{x})^T \vec{x} \end{array}$$

where $\theta: \mathcal{X} \rightarrow \mathcal{W}$ is parameterised as a learnable DNN!

LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

2. [*Interpretability*] If the features are not interpretable (e.g., individual pixels), then **we should learn a high-level “concept” representation $h(\vec{x})$** :



LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

2. [*Interpretability*] If the features are not interpretable (e.g., individual pixels), then **we should learn a high-level “concept” representation $h(\vec{x})$** :

$$\begin{array}{ll} \text{Linear Model Output: } & f(\vec{x}) = \theta^T \vec{x} \\ \text{“Linear-ish DNN” Model output: } & f(\vec{x}) = \theta(\vec{x})^T h(\vec{x}) \end{array}$$

where $\theta: \mathcal{X} \rightarrow \mathcal{W}$ and $h: \mathcal{X} \rightarrow \mathcal{Z}$ are parameterised as a learnable DNNs!

LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

3. **[*Local Linearity*]** The model should behave, at least in the neighborhood of a sample, as a linear classifier.

What does this imply?

LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

3. **[Local Linearity]** The model should behave, at least in the neighborhood of a sample, as a linear classifier.

$$\nabla_{h(\vec{x})} f(\vec{x}) \approx \theta(\vec{x})$$

Relevance coefficients **adapt with the inputs** but they do so in a **stable/slow** manner

LINEARIZING A DNN

If we want to make a DNN **act as a linear model while maintaining its expressive power**, we need a few things:

3. **[Local Linearity]** The model should behave, at least in the neighborhood of a sample, as a linear classifier.

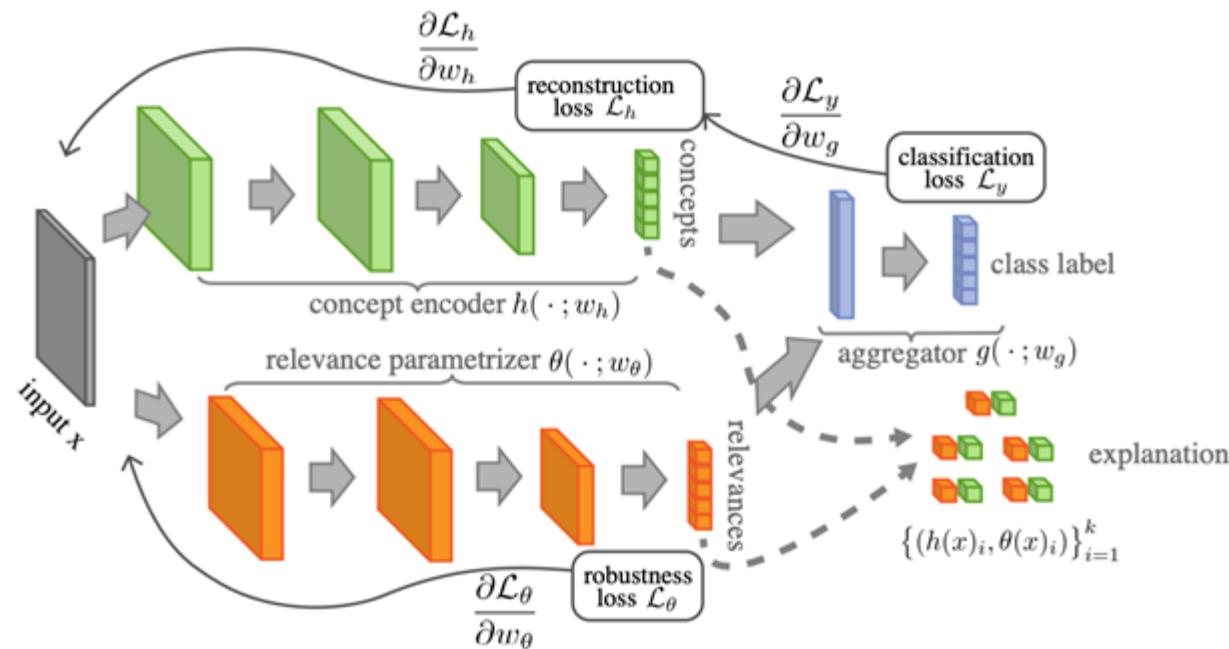
$$\nabla_{h(\vec{x})} f(\vec{x}) \approx \theta(\vec{x}) \quad \begin{array}{l} \text{Relevance coefficients adapt with the inputs} \\ \text{but they do so in a stable/slow manner} \end{array}$$

We can encourage this local linearity by **including the following training regulariser**:

$$\mathcal{L}_{reg}(\vec{x}) := \left\| \nabla_{\vec{x}} f(\vec{x}) - J_{\vec{x}}^{h(\vec{x})}(\vec{x}) \right\|$$

SELF-EXPLAINING NEURAL NETS

This is the idea behind **Self Explaining Neural Networks (SENNs)**!

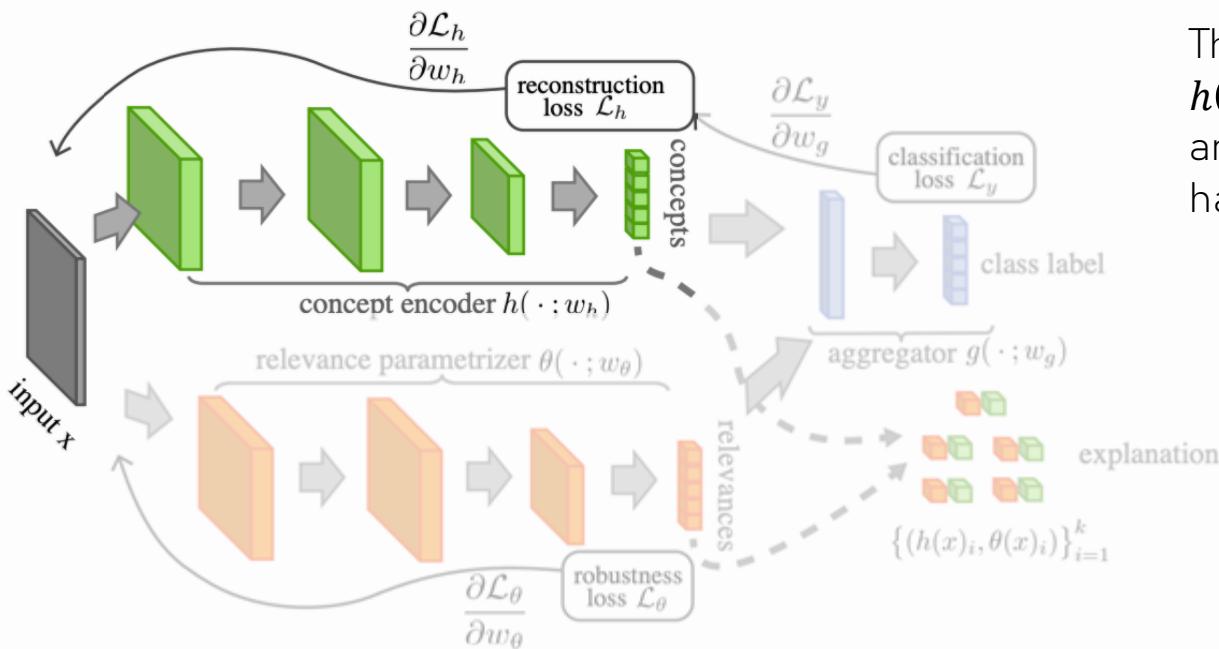


[1] Alvarez Melis et al. "Towards robust interpretability with self-explaining neural networks." *NeurIPS* (2018).



SELF-EXPLAINING NEURAL NETS

Step 1: extract **concepts** from our input distribution:

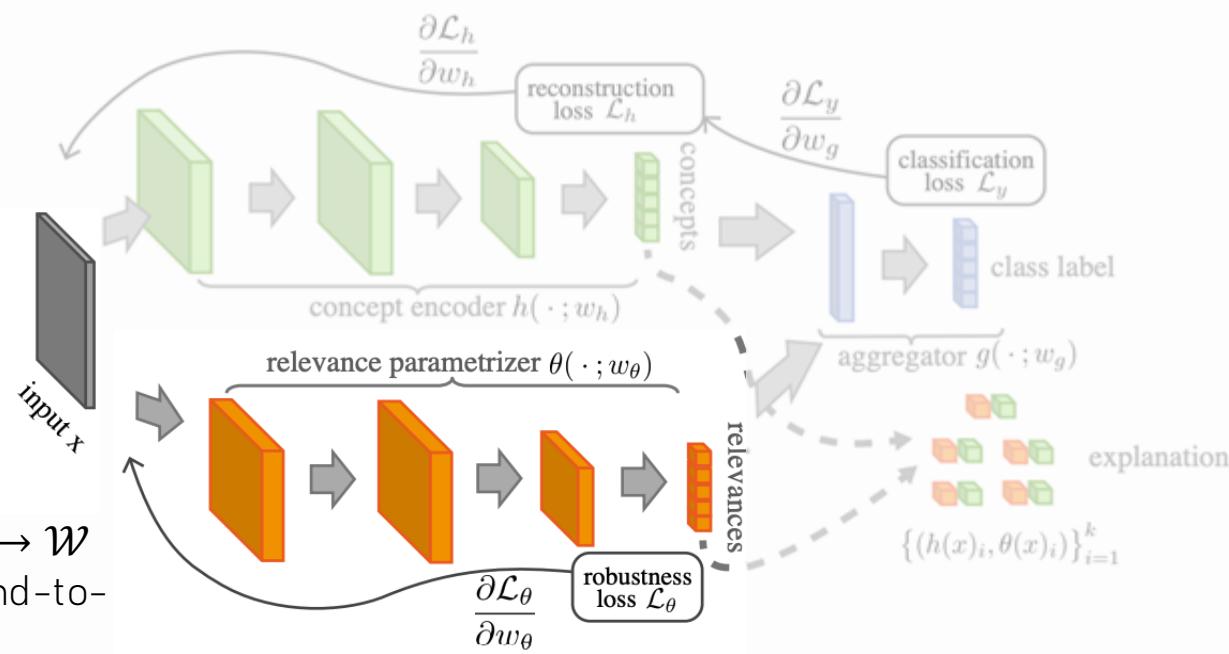


The concept extractor $h(x): \mathcal{X} \rightarrow \mathcal{Z}$ can be learnt via an **autoencoder model** or via handcrafted feature extractors



SELF-EXPLAINING NEURAL NETS

Step 2: use DNN to dynamically predict the set of linear weights for each sample:



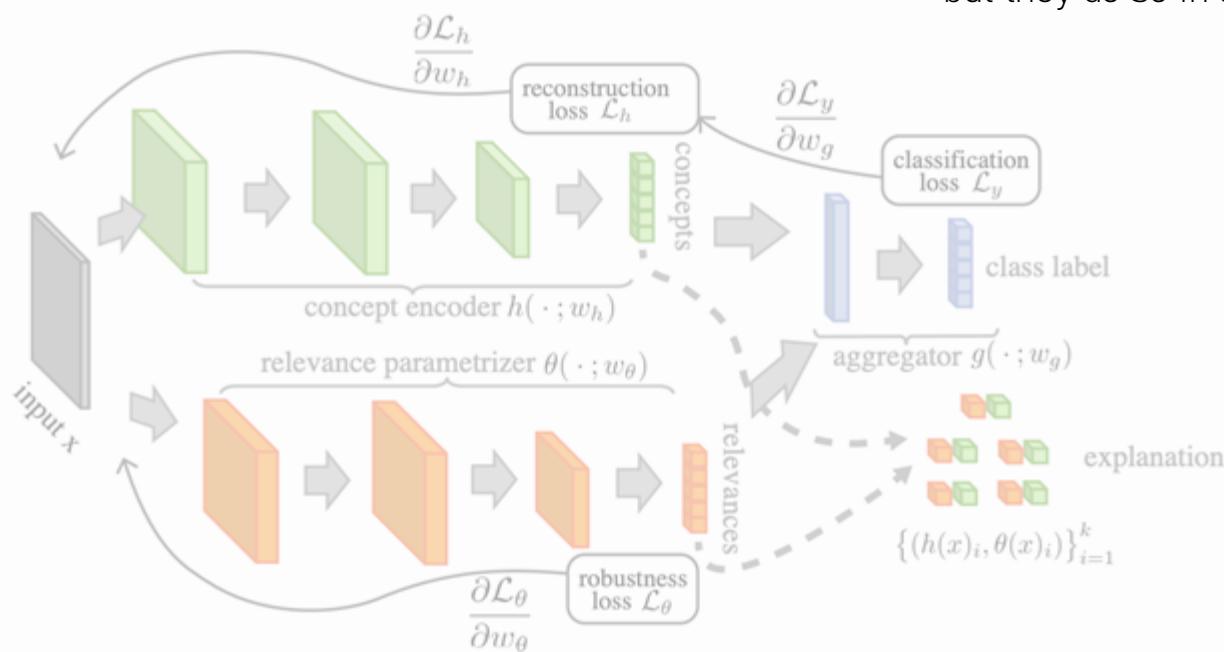
This is done via a weight relevance model $\theta(x): \mathcal{X} \rightarrow \mathcal{W}$ that can be learnt in an end-to-end fashion



SELF-EXPLAINING NEURAL NETS

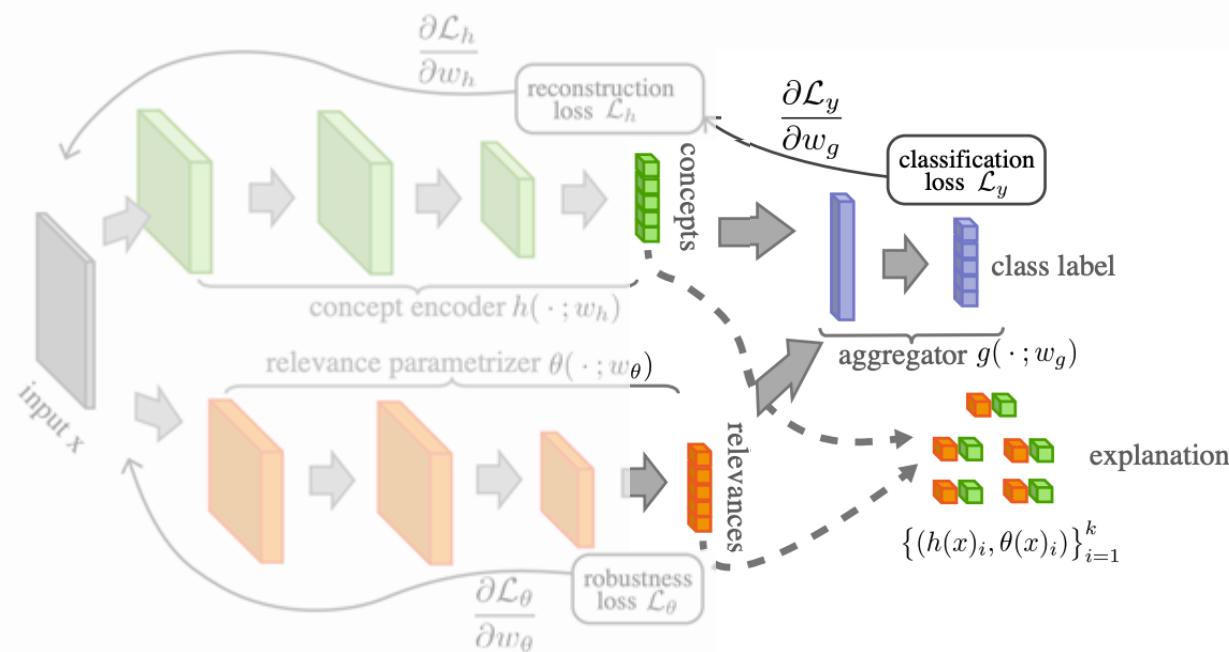
Step 3: Add regulariser that will encourage local linearity: $\mathcal{L}_\theta(f(x)) := \|\nabla_x f(x) - \theta(x)^\top J_x^h(x)\|$

Relevance coefficients **adapt with the inputs**
but they do so in a **stable/slow** manner



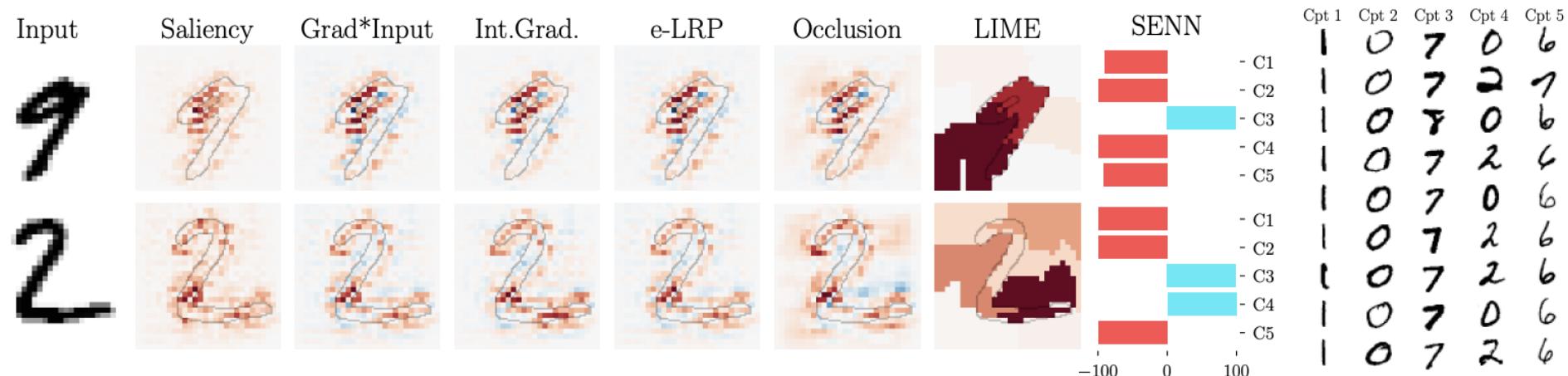
SELF-EXPLAINING NEURAL NETS

Step 4: Generate prediction with the **linear form** $\theta(x)^T h(x)$. The **explanation** is the tuple (concept, relevance weight)



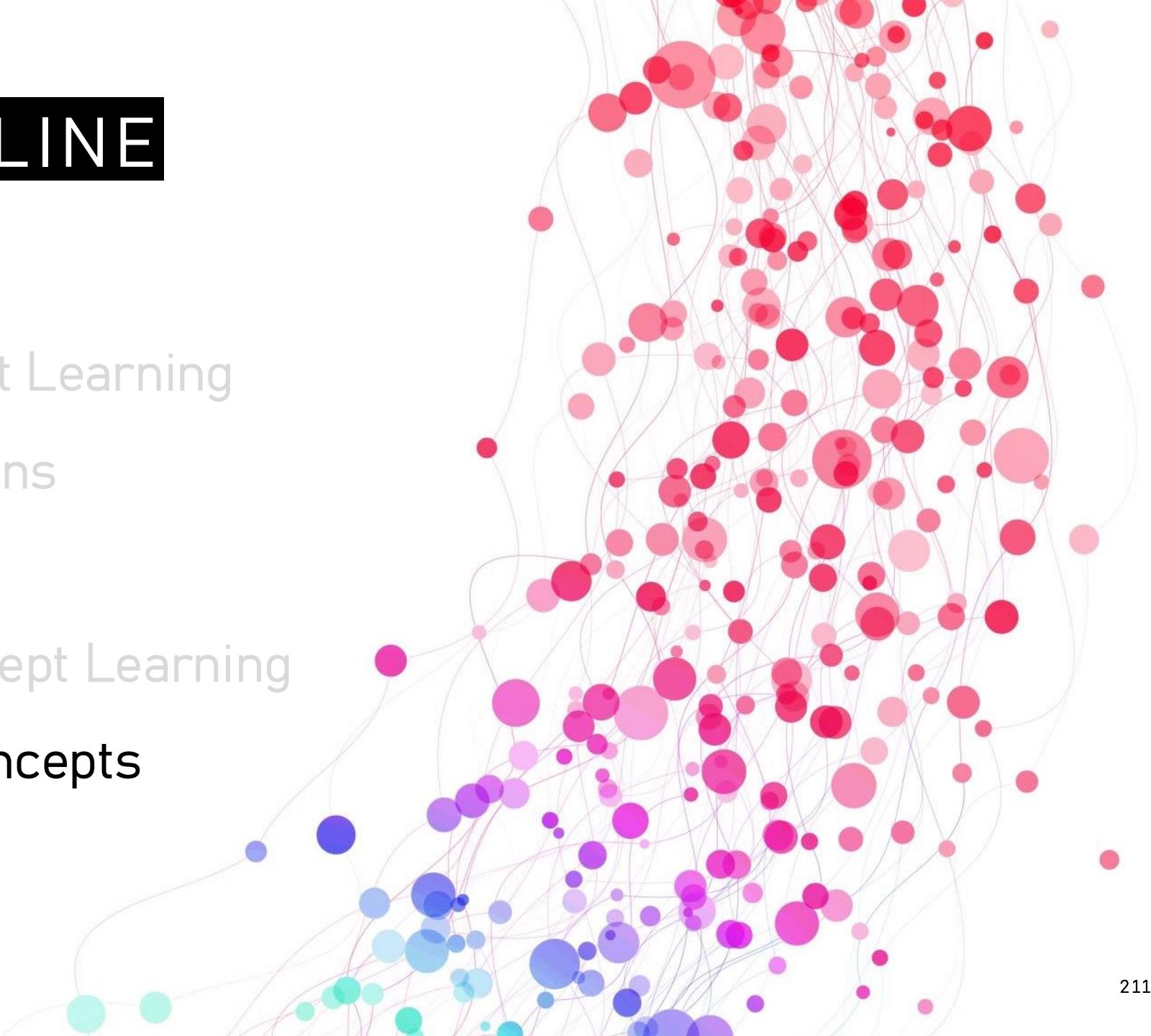
SELF-EXPLAINING NEURAL NETS

When features lack useful semantics, learnt concepts can be understood via **prototypical examples**:



TUTORIAL OUTLINE

1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
- 6. Reasoning With Concepts**
7. Future Directions
8. Q&A



CONCEPT-BASED REASONING

We'll focus on two main branches of concept-based reasoning:

Neural symbolic concept reasoning

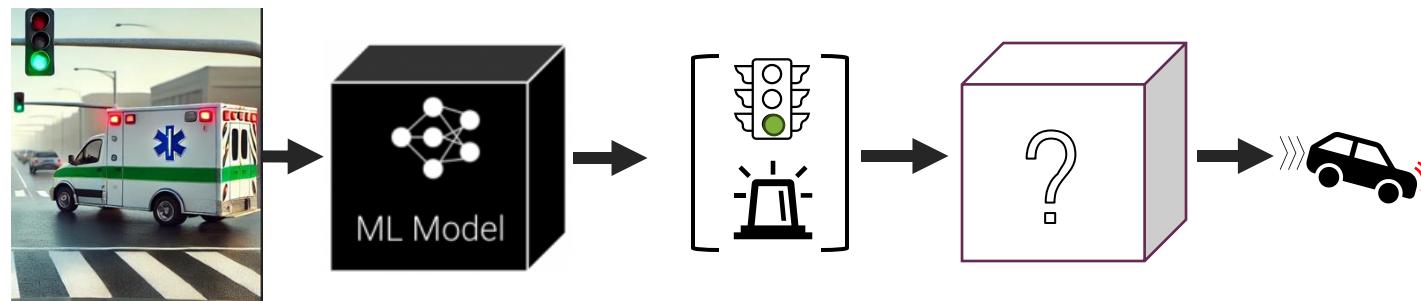


Causal concept reasoning



TIME TO GET YOUR C*EPTS TOGETHER

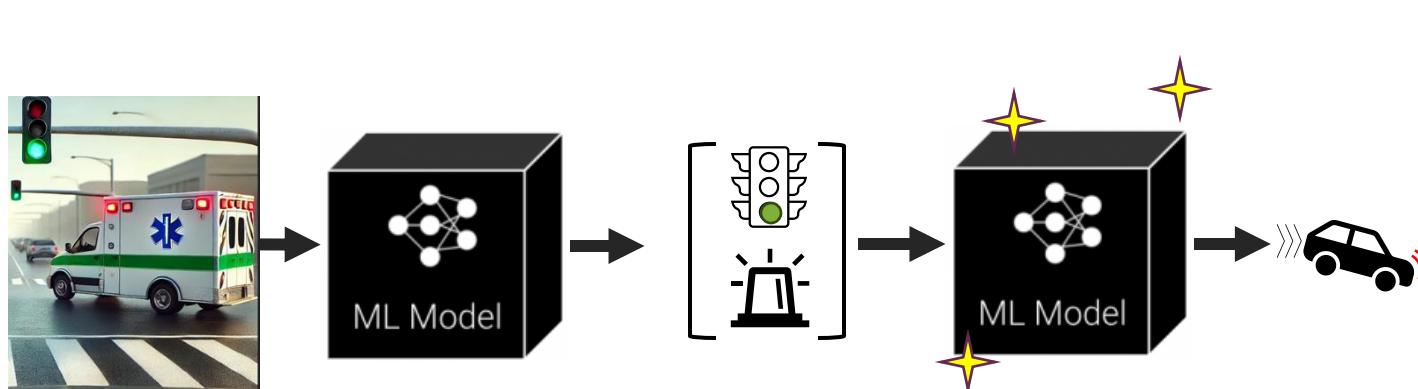
Let's say we have a nice set of concepts, what should we use as **classification head?**



TIME TO GET YOUR C*EPTS TOGETHER

Let's say we have a nice set of concepts, what should we use as **classification head**?

... what about an **opaque DNN**?



Back to square 1!

TIME TO GET YOUR C*EPTS TOGETHER

Let's say we have a nice set of concepts, what should we use as **classification head**?

... what about an **opaque DNN**?

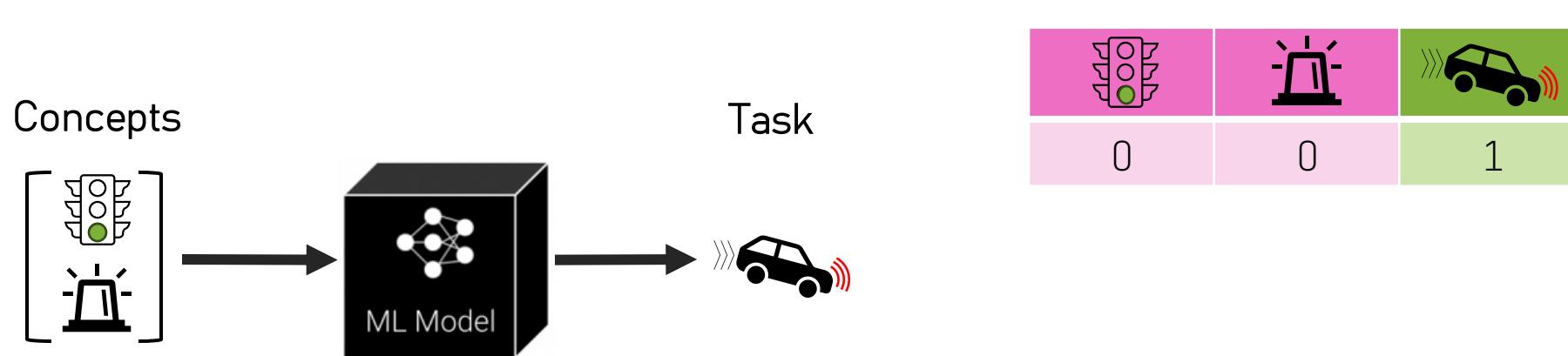
Can we do better?

Back to square 1!

FROM INTERVENTIONS TO LOGIC REASONING

Let's say we have a nice set of concepts, what should we use as **classification head**?

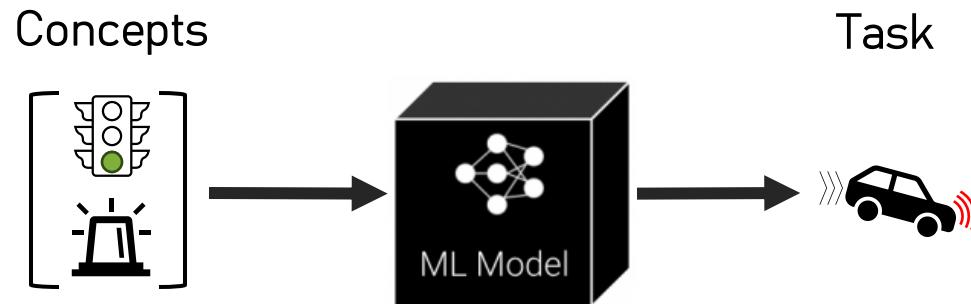
... what about an **opaque DNN**?



FROM INTERVENTIONS TO LOGIC REASONING

Let's say we have a nice set of concepts, what should we use as **classification head**?

... what about an **opaque DNN**?



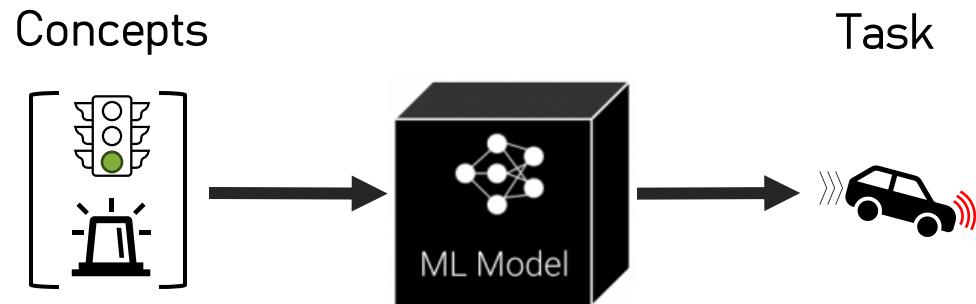
	Traffic Light	Megaphone	Car
Traffic Light	0	0	1
Megaphone	1	0	0
Car			

green light=1

FROM INTERVENTIONS TO LOGIC REASONING

Let's say we have a nice set of concepts, what should we use as **classification head**?

... what about an **opaque DNN**?



	Traffic Light	Siren	Car
Traffic Light	0	0	1
Siren	1	0	0
Car	0	1	1

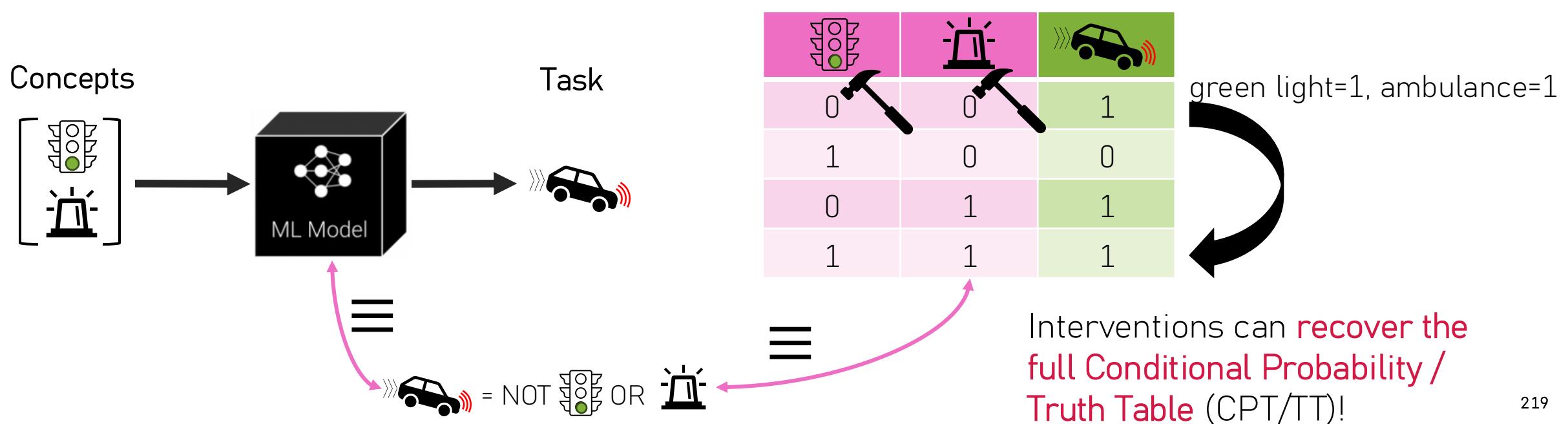
ambulance=1

ambulance=1

FROM INTERVENTIONS TO LOGIC REASONING

Let's say we have a nice set of concepts, what should we use as **classification head**?

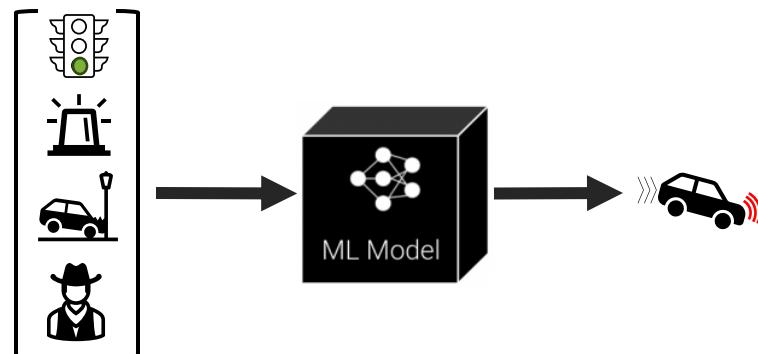
... what about an **opaque DNN**?



LOGIC-EXPLAINED NETWORKS (LENS)

Limitation Being Addressed

#interventions required to extract full CPT/TT is **exponential** in #concepts!



Can we extract a CPT/TT
more efficiently?



LOGIC-EXPLAINED NETWORKS (LENS)

Proposed Solution

Step 1: Filter concept activations using learnable **attention weights** α

Concept activations	Learnable "attention" weights	Filtered concept activations
	1	
	0.9	
	0.1	
	0.2	



LOGIC-EXPLAINED NETWORKS (LENS)

Proposed Solution

Step 2: Minimize the entropy of the **attention weights** α . Why? Concept set should be **small!**

$$\min H(\alpha)$$

Concept activations	Learnable "attention" weights	Filtered concept activations

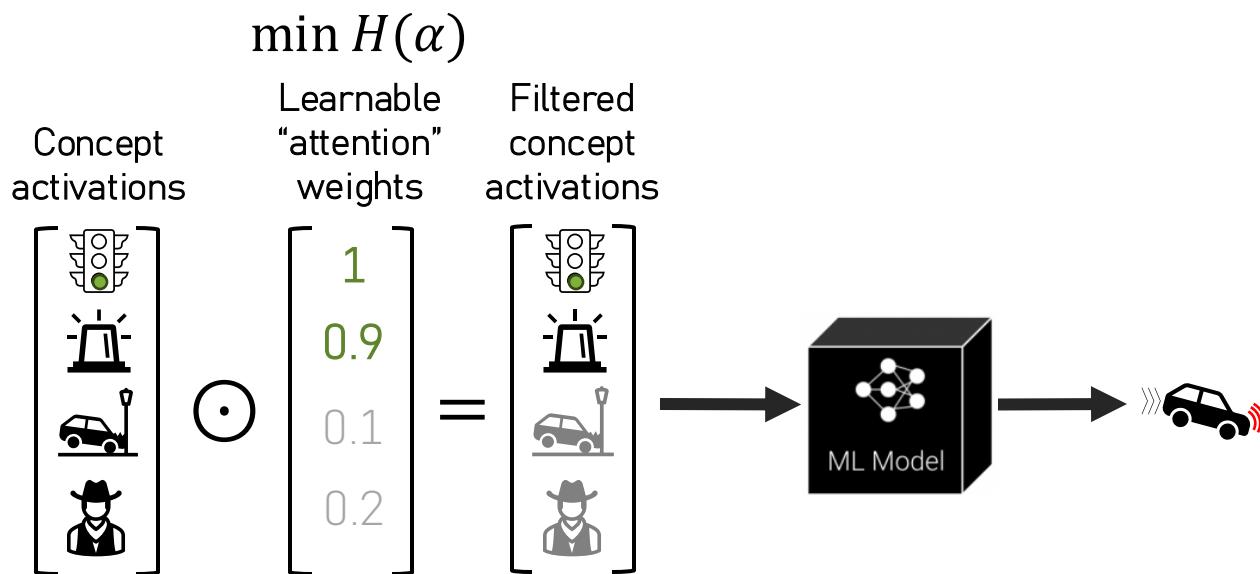
The diagram illustrates a linear transformation where a vector of concept activations is multiplied by a matrix of learnable "attention" weights to produce a filtered vector of concept activations. The matrix has four columns corresponding to the concepts in the vector. The first column has a value of 1, indicating it is the primary focus (the most "attentive"). The other three columns have values of 0.9, 0.1, and 0.2 respectively, indicating they are secondary or tertiary concepts.



LOGIC-EXPLAINED NETWORKS (LENS)

Proposed Solution

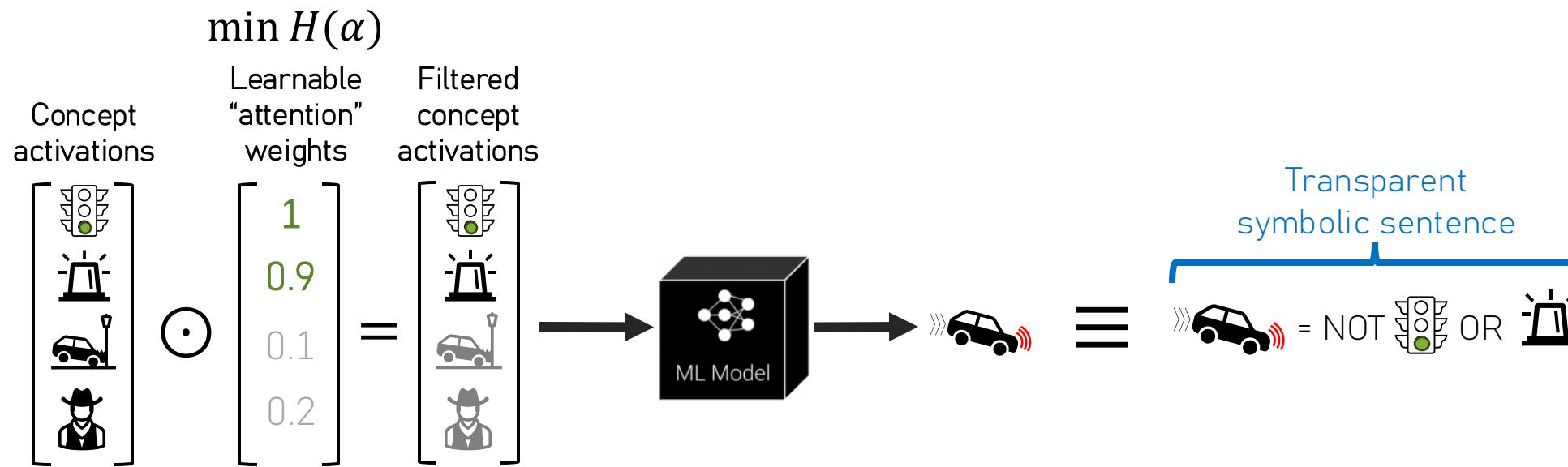
Step 3: Solve task with the selected concepts. Why? Concept set should be **relevant!**



LOGIC-EXPLAINED NETWORKS (LENS)

Proposed Solution

Step 4: Derive explanation in DNF from the (empirical) truth table



LOGIC-EXPLAINED NETWORKS (LENS)

Proposed Solution

Step 4: Derive explanation in DNF from the (empirical) truth table

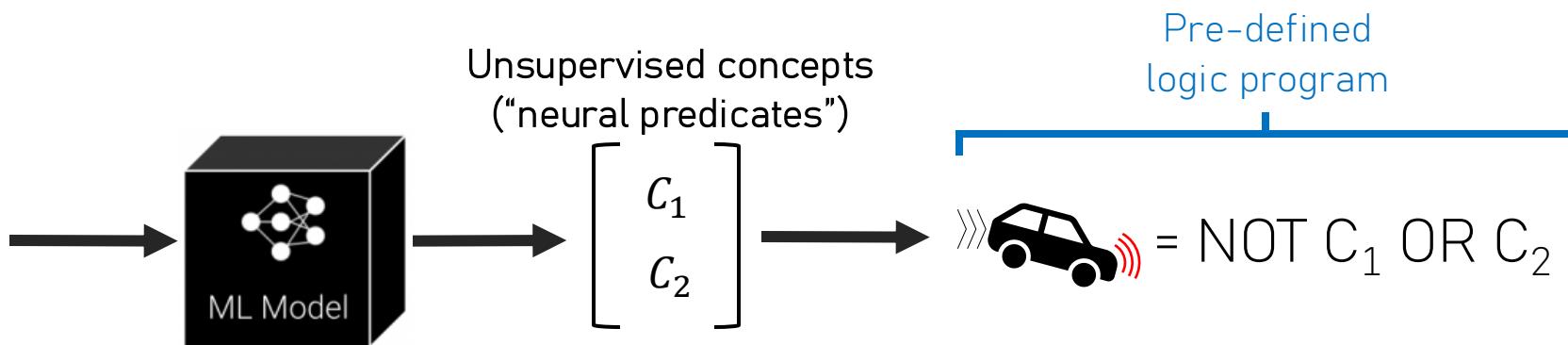
*What if we know the logic program,
but we don't have concept supervisions?*



NEURAL PROBABILISTIC LOGIC PROGRAMMING

Proposed Solution

Replace task predictor with a **pre-defined logic program!**



NEURAL PROBABILISTIC LOGIC PROGRAMMING

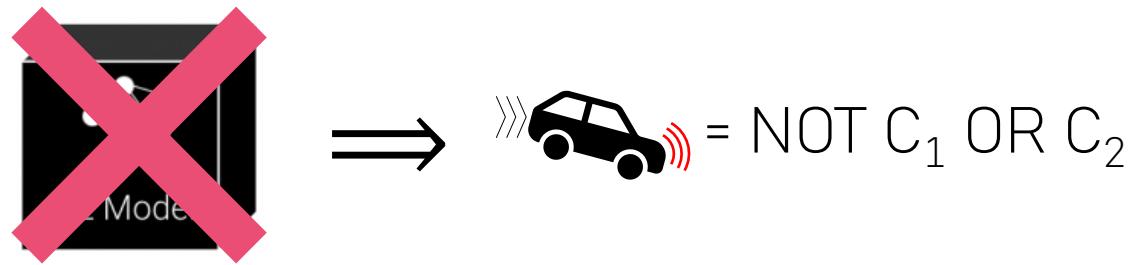
Proposed Solution

Replace task predictor with a **pre-defined logic program!**

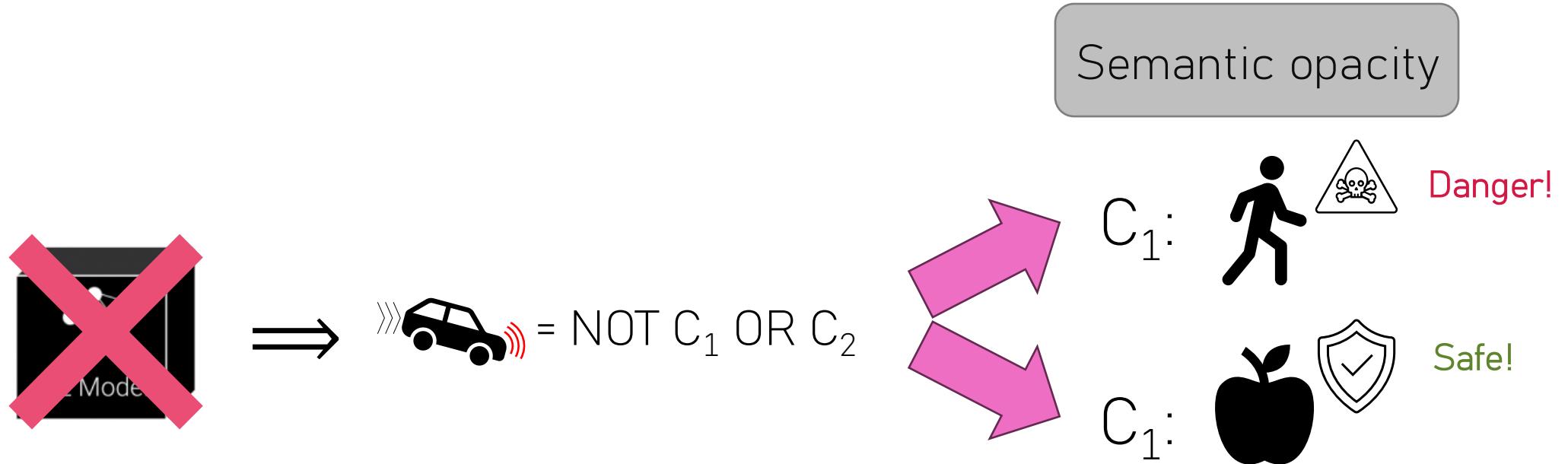
*Are symbolic classification heads sufficient
for a model to be interpretable?*



SEMANTIC & FUNCTIONAL OPACITY

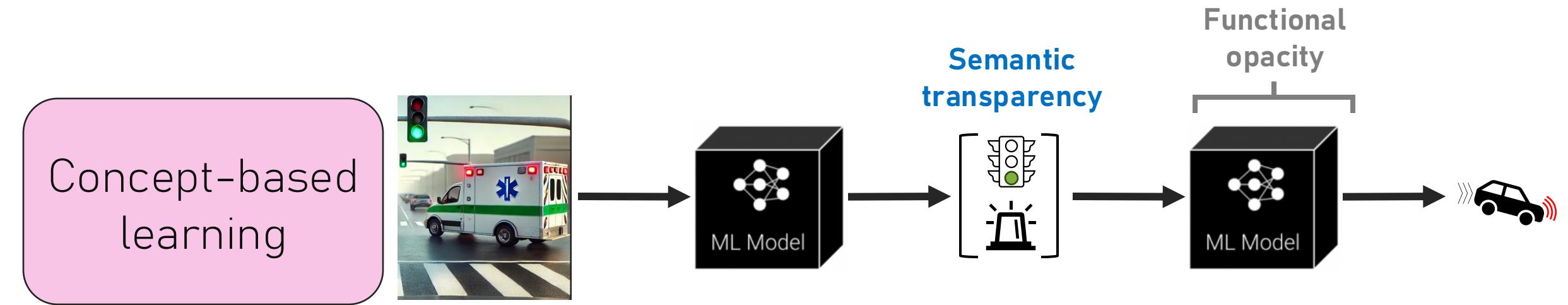


SEMANTIC & FUNCTIONAL OPACITY



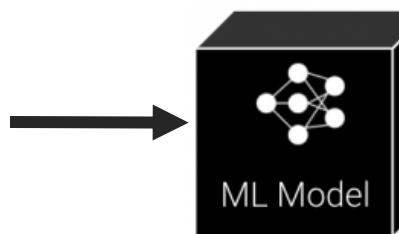
A symbolic classification head alone **does not guarantee semantic transparency**
(... as well as Logistic Regression, Additive Models, Decision Trees, etc...)!

SEMANTIC & FUNCTIONAL OPACITY

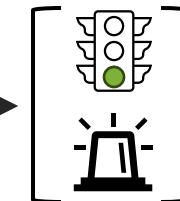


SEMANTIC & FUNCTIONAL OPACITY

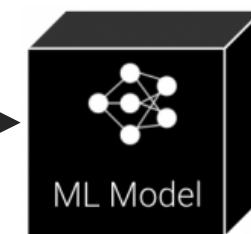
Concept-based learning



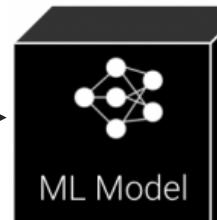
Semantic transparency



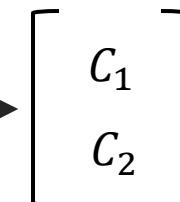
Functional opacity



Symbolic reasoning



Semantic opacity



Functional transparency



$= \text{NOT } C_1 \text{ OR } C_2$

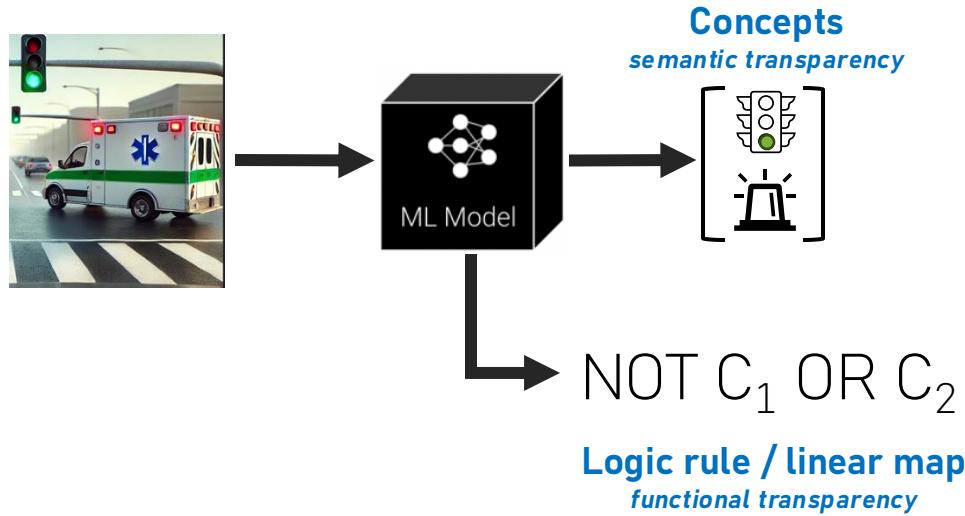
SEMANTIC & FUNCTIONAL OPACITY

*Can we combine concept-based learning
with symbolic reasoning?*

NEURAL-SYMBOLIC CONCEPT REASONING

Proposed Solution

Step 1: DNN generates both concept activations & rule parameters (*neural generation*)



[1] Barbiero et al. "Interpretable neural-symbolic concept reasoning" International Conference on Machine Learning, PMLR 2023.
[2] Debot et al. "Interpretable concept-based memory reasoning" NeurIPS 2024.



ICML23



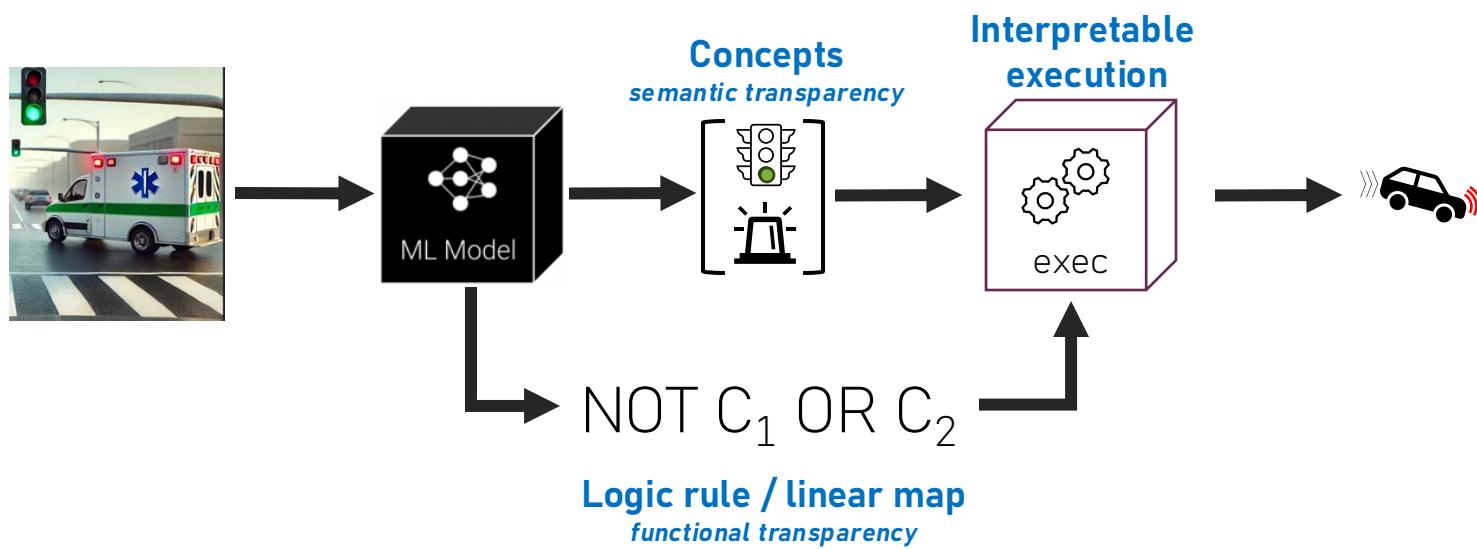
NeurIPS24

NEURAL-SYMBOLIC CONCEPT REASONING

Proposed Solution

Step 1: DNN generates both concept activations & rule parameters (*neural generation*)

Step 2: Symbolic engine executes the rule using concept activations (*interpretable execution*)



[1] Barbiero et al. "Interpretable neural-symbolic concept reasoning" International Conference on Machine Learning, PMLR 2023
[2] Debot et al. "Interpretable concept-based memory reasoning" NeurIPS 2024



ICML23

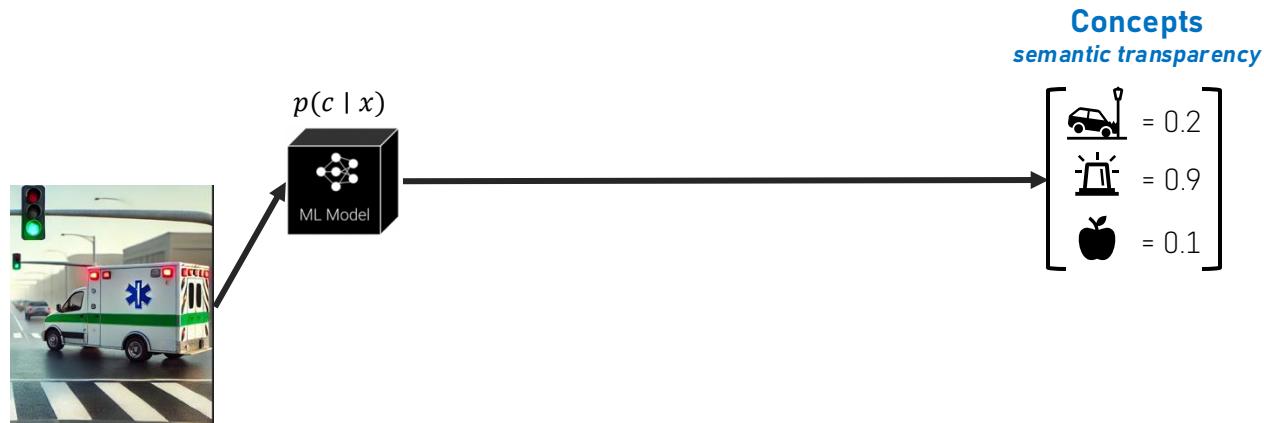


NeurIPS24

CONCEPT-BASED MEMORY REASONING

Proposed Solution

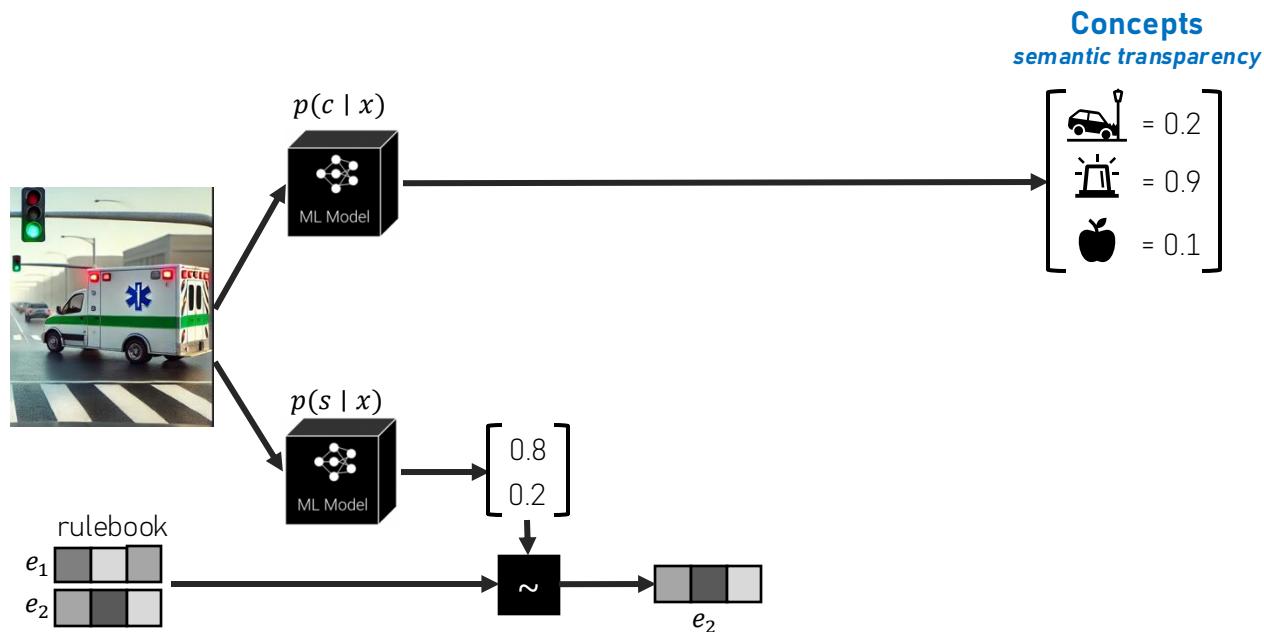
Step 1: DNN predicts concept activations



CONCEPT-BASED MEMORY REASONING

Proposed Solution

Step 2: DNN predicts embedding to be selected from the latent rulebook



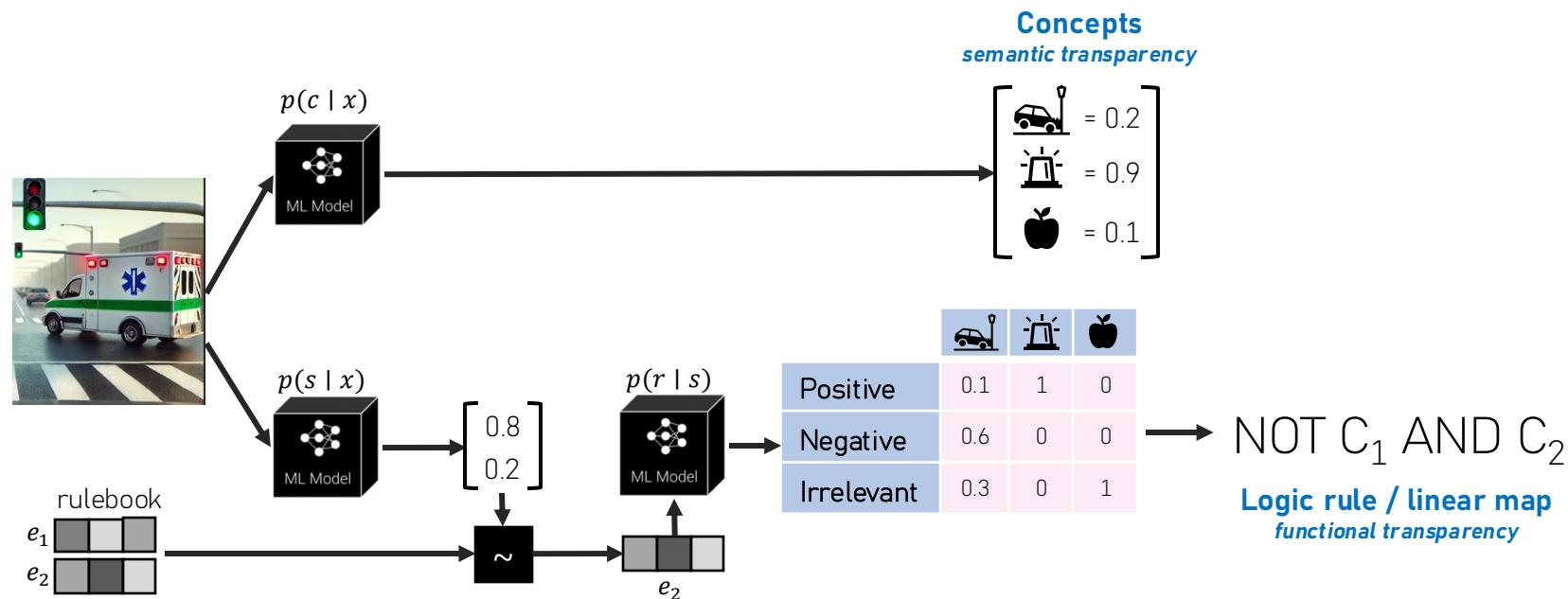
[1] Dehghani et al. "Interpretable concept-based memory reasoning." NeurIPS 2024



CONCEPT-BASED MEMORY REASONING

Proposed Solution

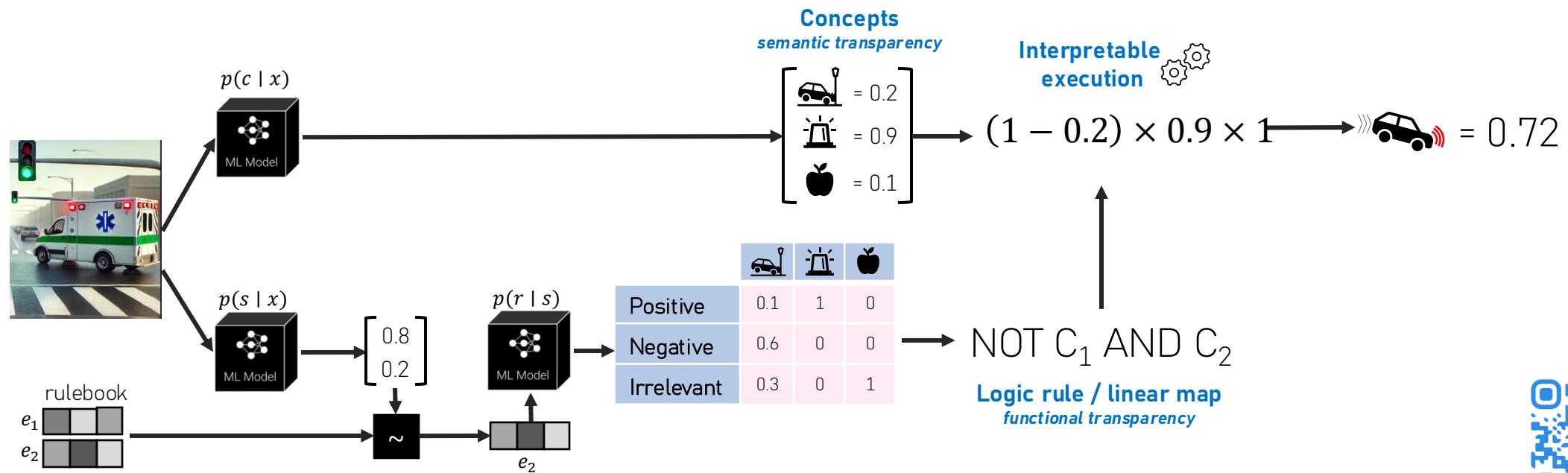
Step 3: DNN decodes selected embedding into 3 states: positive, negative, irrelevant



CONCEPT-BASED MEMORY REASONING

Proposed Solution

Step 4: Execute the rule combining concept states and activations to predict the output label



CONCEPT-BASED MEMORY REASONING

CMR has 3 key features:

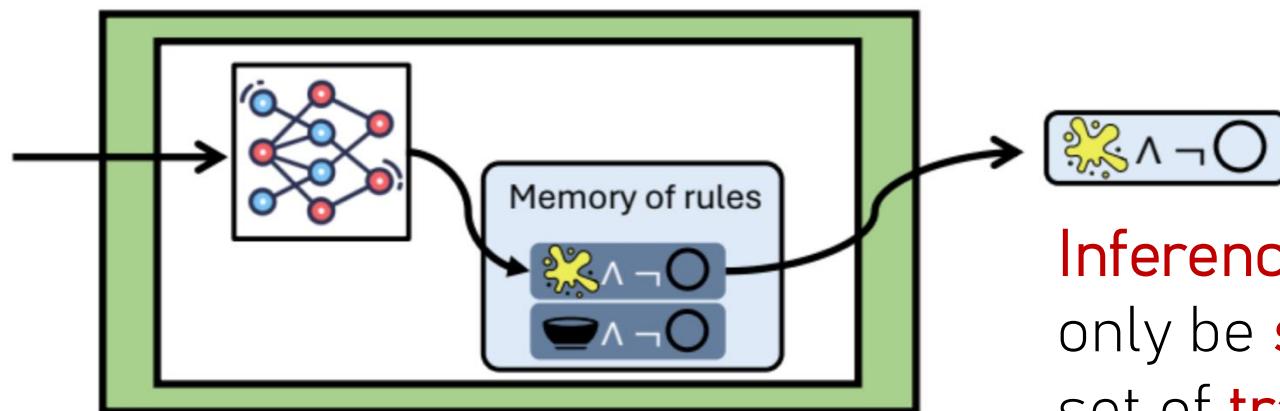
- **Universal approximator** akin to opaque DNNs (Theorem 4.1)



CONCEPT-BASED MEMORY REASONING

CMR has 3 key features:

- **Universal approximator** akin to opaque DNNs (Theorem 4.1)
- Provides both **local and global interpretability** by design



Inference mechanisms can
only be selected from a finite
set of transparent rules!

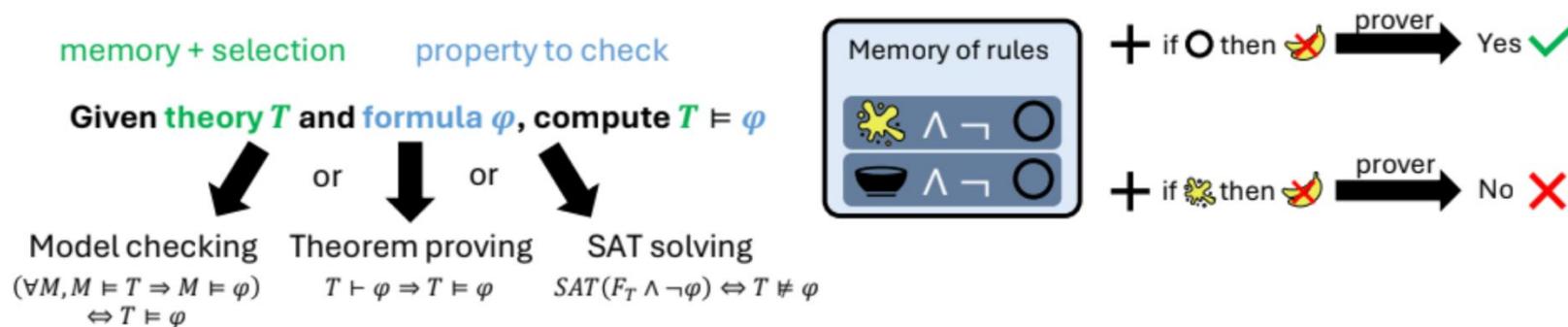


CONCEPT-BASED MEMORY REASONING

CMR has 3 key features:

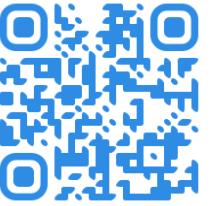
- **Universal approximator** akin to opaque DNNs (Theorem 4.1)
- Provides both **local and global interpretability** by design
- The concept memory allows **formal verification** of properties

"Does a property hold no matter which rule is selected?"



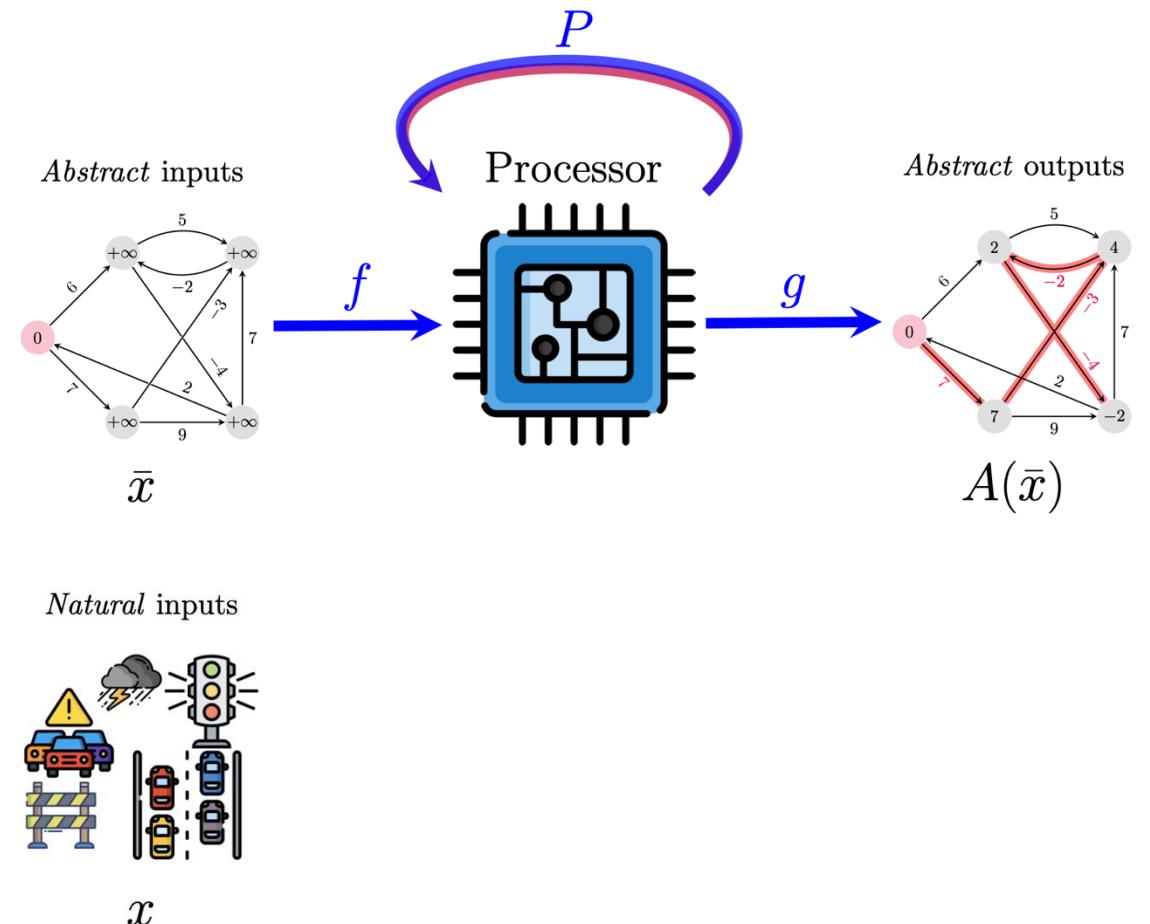
ARE WE JUST TALKING HOT AIR?





NEURAL ALGORITHMIC REASONING

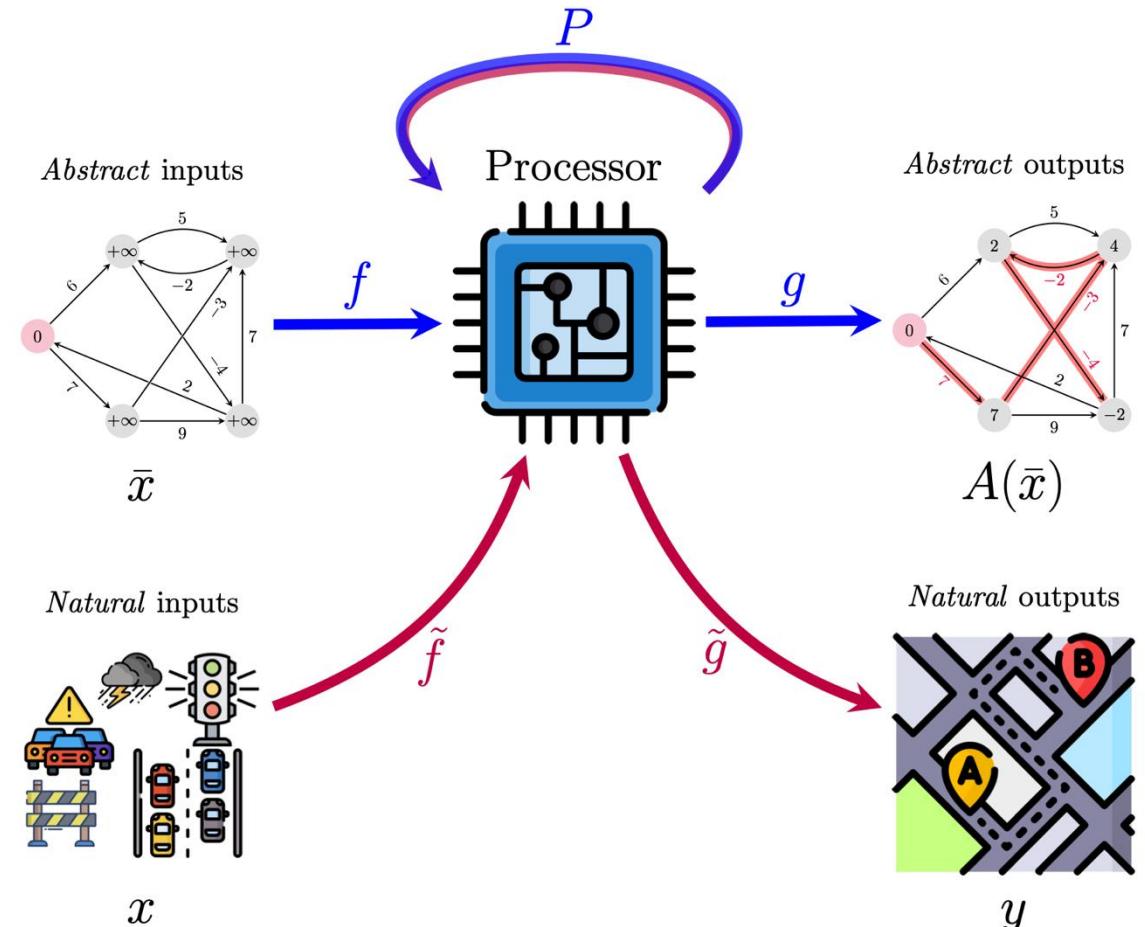
- Algorithmic reasoning
 - + OOD generalization
 - - discrete representations
- Neural nets
 - - OOD generalization
 - + continuous representations



NEURAL ALGORITHMIC REASONING



- Execute algorithms with DNNs
 - + OOD generalization (from algorithm exec)
 - + adapt to real-world inputs (e.g., images)

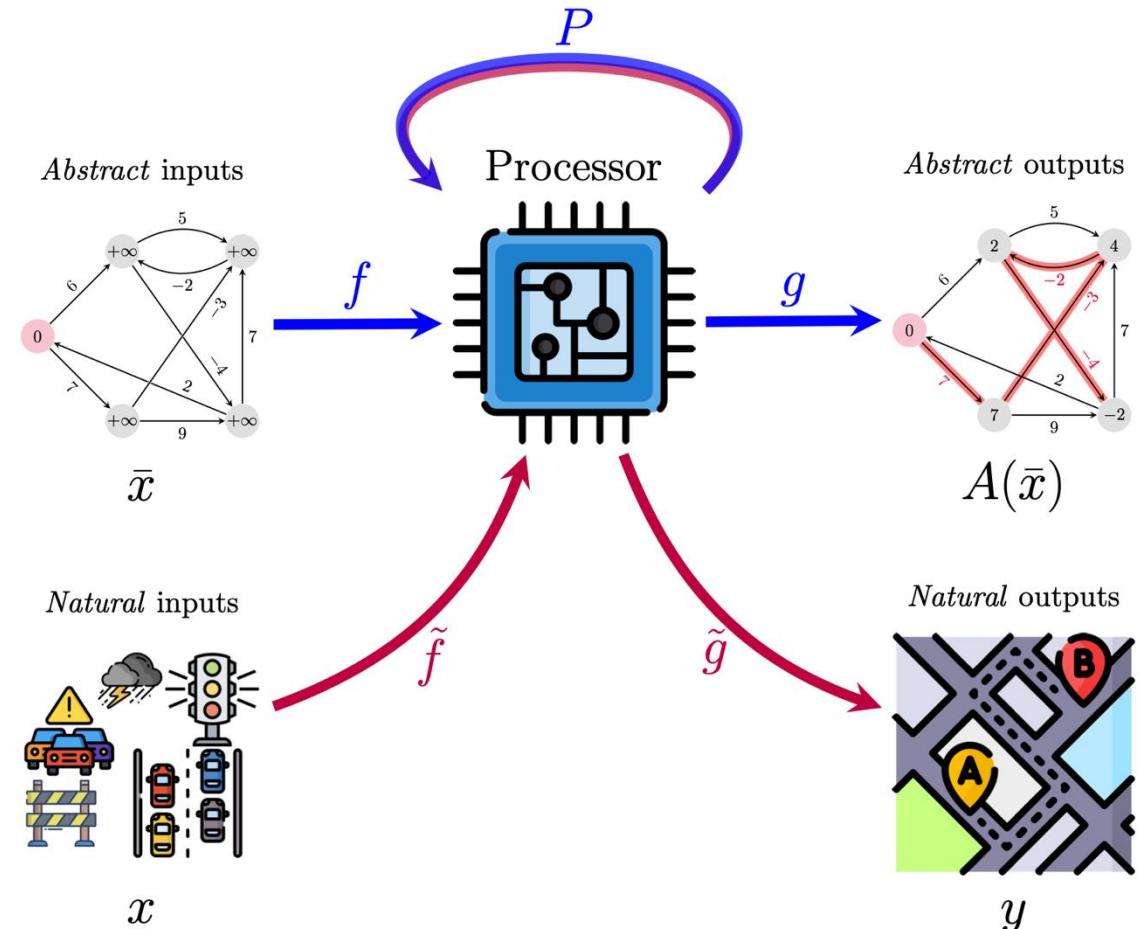


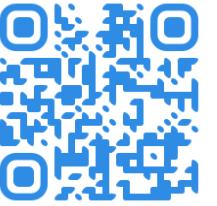
NEURAL ALGORITHMIC REASONING



- Execute algorithms with DNNs
 - + OOD generalization (from algorithm exec)
 - + adapt to real-world inputs (e.g., images)
 - + (potentially) find new heuristics!

How?

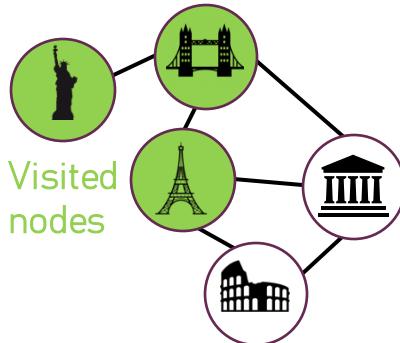


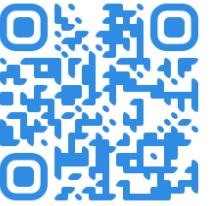


NEURAL ALGORITHMIC REASONING

- Execute algorithms with DNNs
 - + OOD generalization (from algorithm exec)
 - + adapt to real-world inputs (e.g., images)
 - + (potentially) find new heuristics!

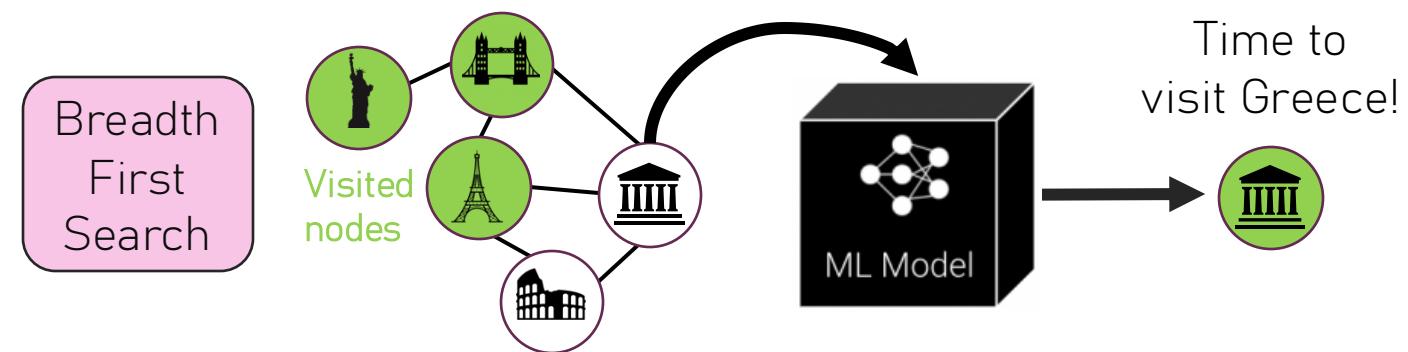
Breadth
First
Search





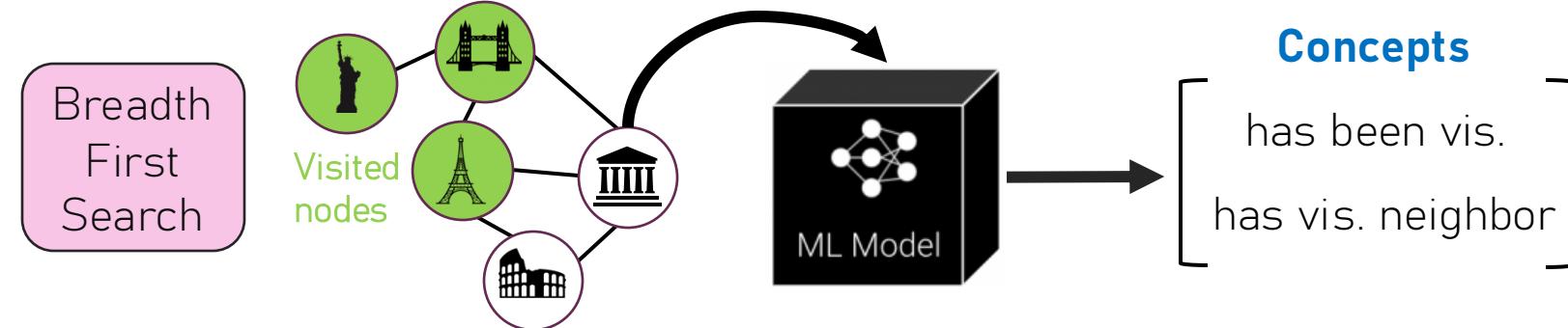
NEURAL ALGORITHMIC REASONING

- Execute algorithms with DNNs
 - + OOD generalization (from algorithm exec)
 - + adapt to real-world inputs (e.g., images)
 - + (potentially) find new heuristics!



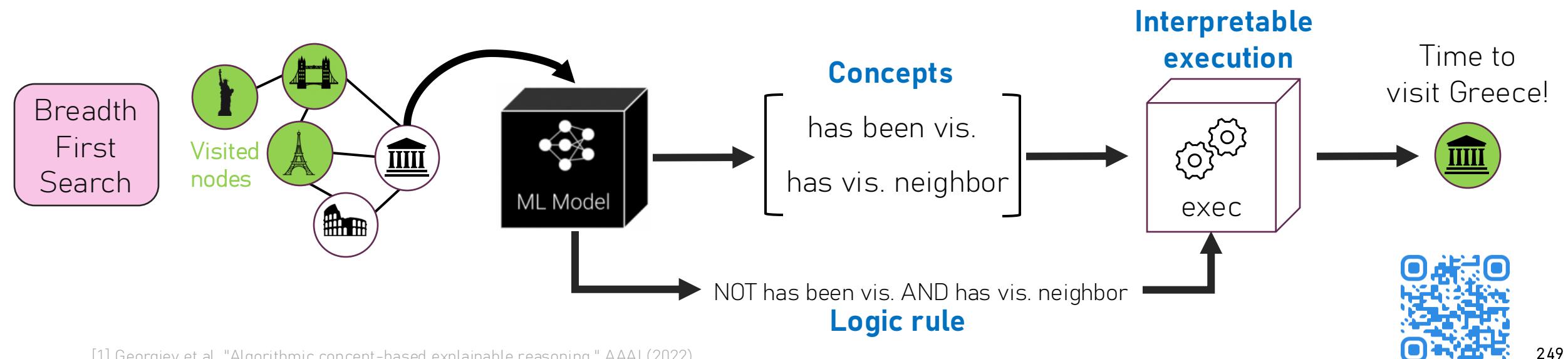
CONCEPT-BASED NEURAL ALGORITHMIC REASONING

- Execute algorithms with DNNs
 - + OOD generalization (from algorithm exec)
 - + adapt to real-world inputs (e.g., images)
 - + (potentially) find new heuristics!



CONCEPT-BASED NEURAL ALGORITHMIC REASONING

- Execute algorithms with DNNs
 - + OOD generalization (from algorithm exec)
 - + adapt to real-world inputs (e.g., images)
 - + (potentially) find new heuristics!



SHOULD INTERPRETABILITY BOTHER ABOUT CAUSALITY?

We'll focus on two main branches of concept-based reasoning:

Neural symbolic
concept reasoning

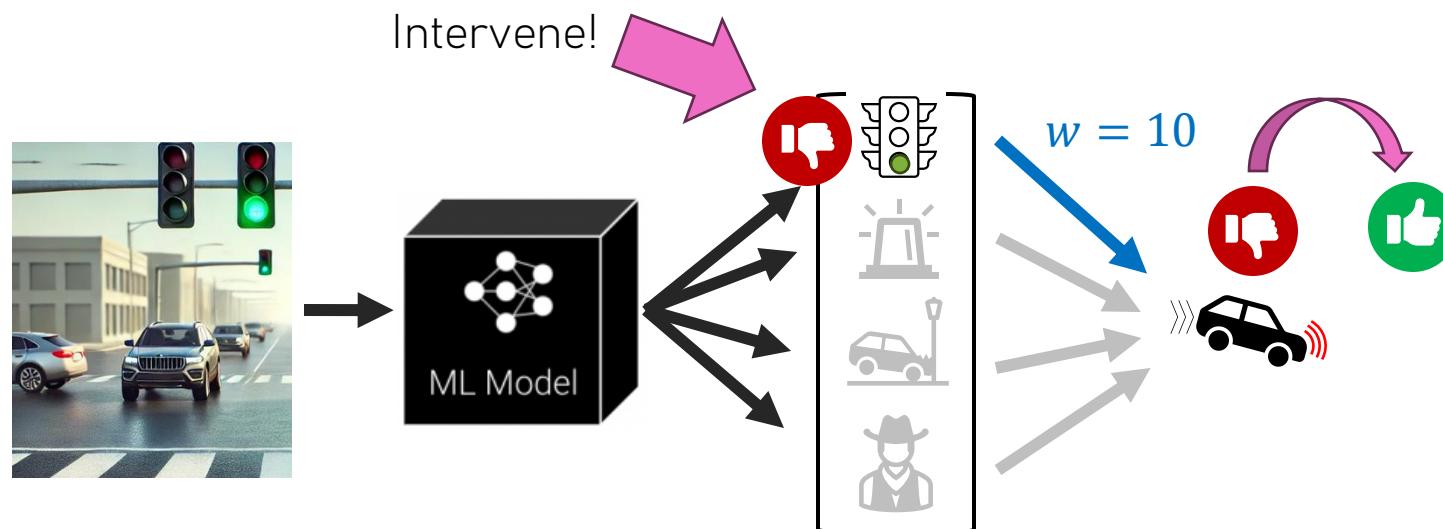


Causal
concept reasoning



SHOULD INTERPRETABILITY BOTHER ABOUT CAUSALITY?

Sometimes intervening on wrongly predicted concepts helps...

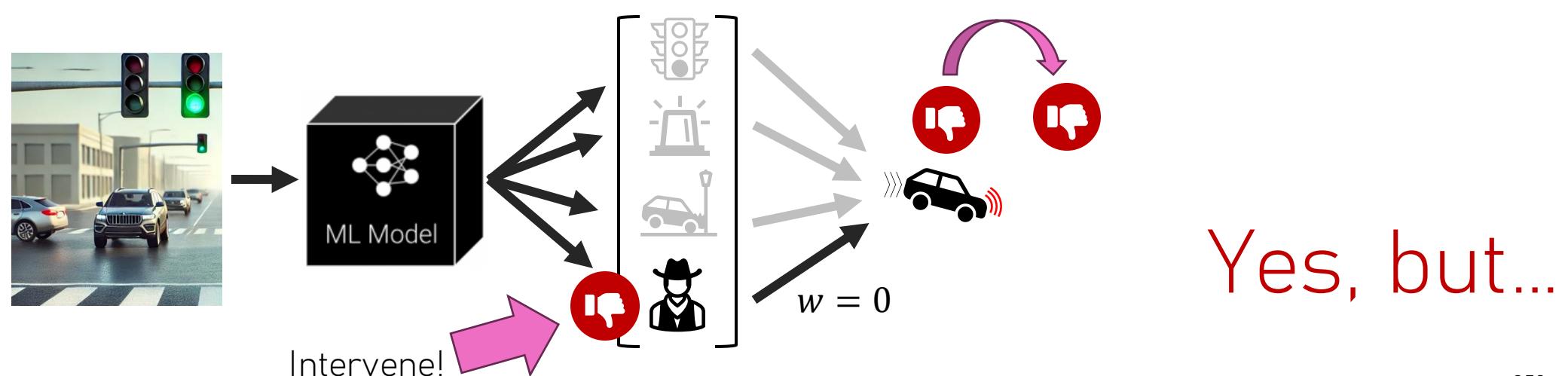


SHOULD INTERPRETABILITY BOTHER ABOUT CAUSALITY?

Sometimes intervening on wrongly predicted concepts helps...

and sometimes it doesn't! 😢

Causal analysis can provide us with insights!



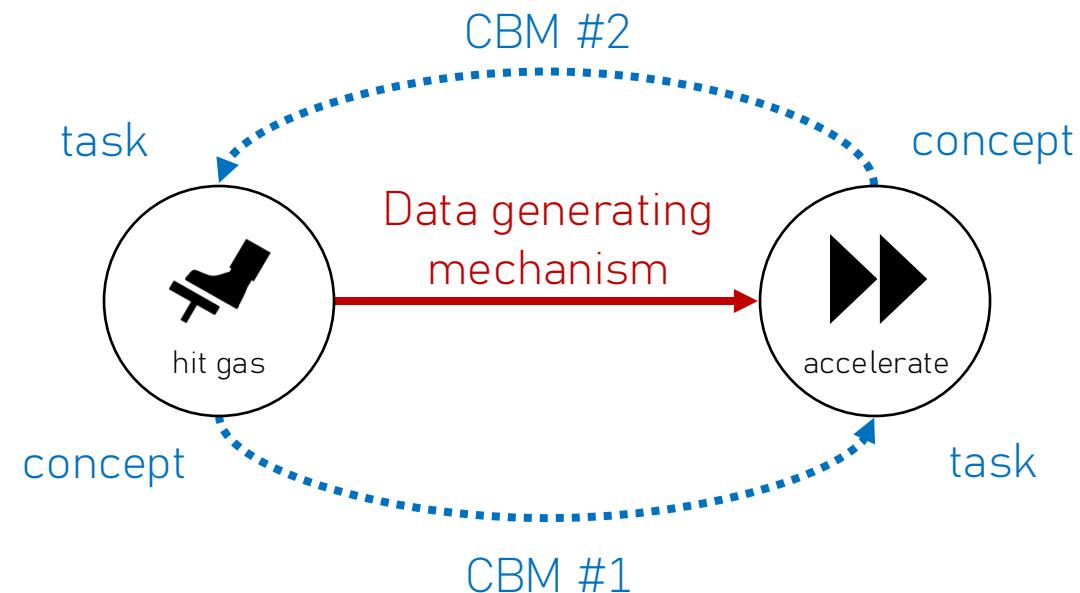
CAUSAL OPACITY

- **Causal reliability:** discover causal mechanisms of the **data generating process**



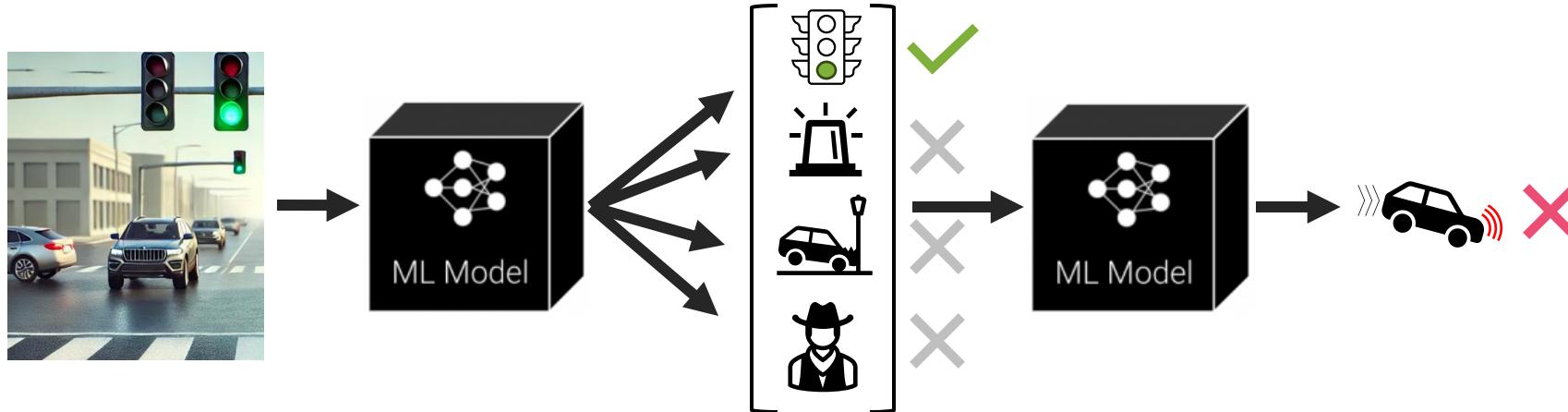
CAUSAL OPACITY

- **Causal reliability**: discover causal mechanisms of the **data generating process**
- **Causal opacity**: discover causal mechanism of a **model's inference process**



CONCEPT-BASED CAUSAL REASONING

CBMs can **answer association** queries (duh...)



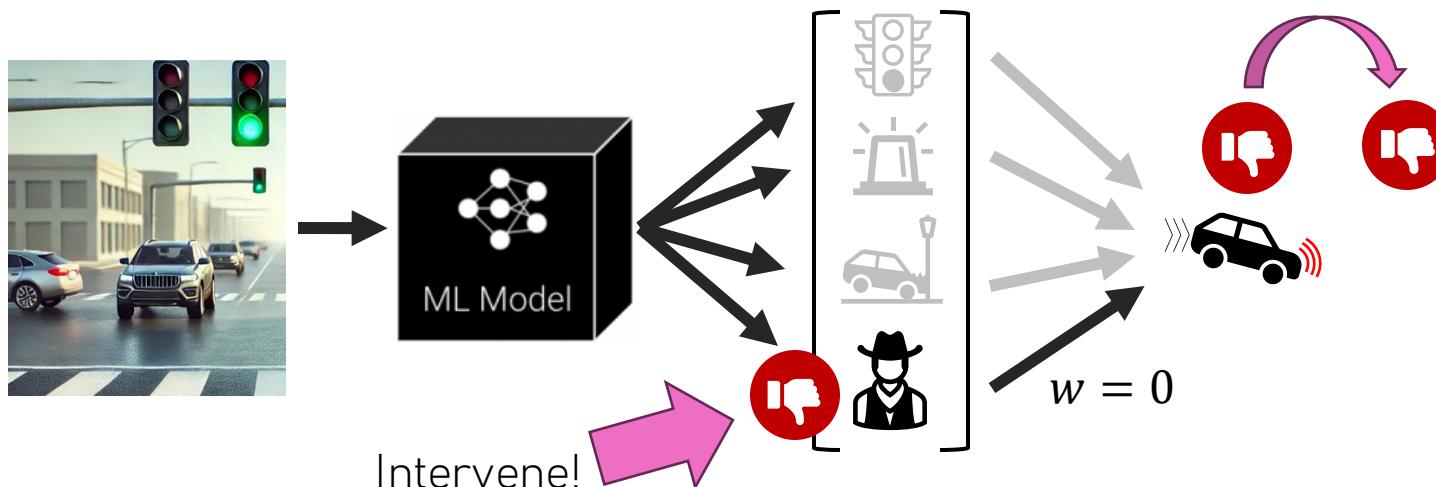
Association

What if the model sees a green light?
 $P(brake \mid light)$

CONCEPT-BASED CAUSAL REASONING

CBMs can **answer association** queries (duh...)

However, intervening on  influences the task, while intervening on  does not!



Intervention
What if I set the light color to red? $P(brake do(light))$
Association
What if the model sees a green light? $P(brake light)$

CONCEPT-BASED CAUSAL REASONING

CBMs can **answer association** queries (duh...)

However, intervening on  influences the task, while intervening on  does not!

*Can we measure the causal influence of a concept
on the task?*



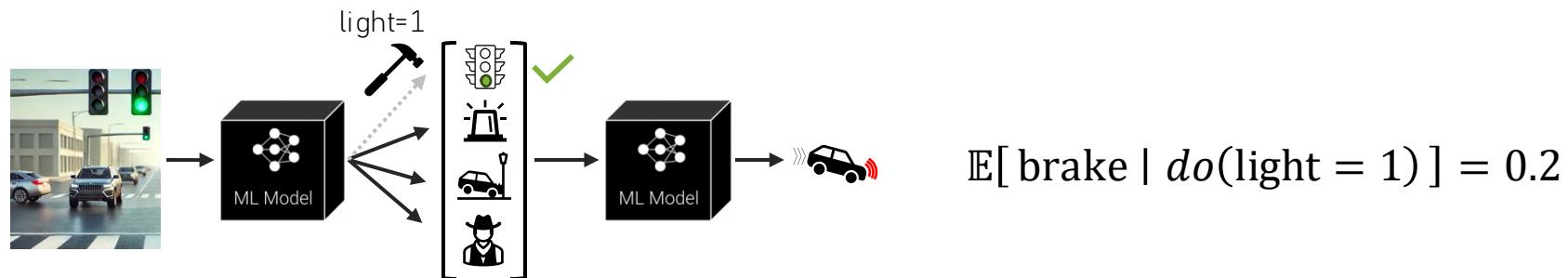
What if the model sees a green light?
 $P(brake \mid light)$

CAUSAL CONCEPT EFFECT



Proposed Solution

Step 1: Compute **expected value** of the task with $do(c_i = 1)$



Intervention
What if I set the light color to red?
 $P(\text{brake} \mid do(\text{light}))$

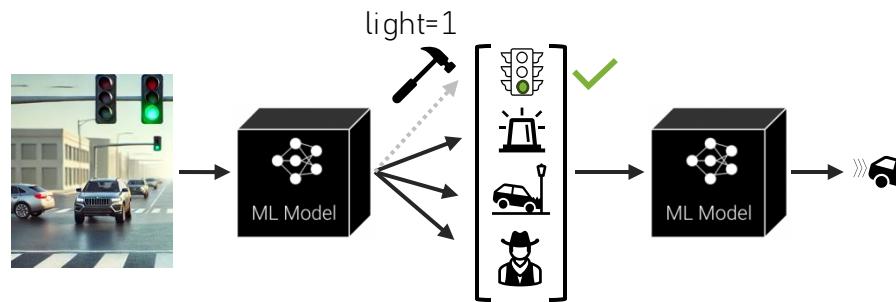
Association
What if the model sees a green light?
 $P(\text{brake} \mid \text{light})$

CAUSAL CONCEPT EFFECT

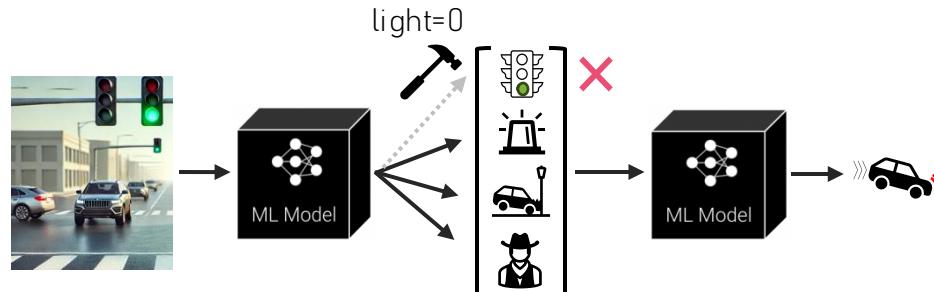


Proposed Solution

Step 2: Compute **expected value** of the task with $do(c_i = 0)$



$$\mathbb{E}[\text{brake} | do(\text{light} = 1)] = 0.2$$



$$\mathbb{E}[\text{brake} | do(\text{light} = 0)] = 1$$

Intervention
What if I set the light color to red?
 $P(\text{brake} | do(\text{light}))$

Association
What if the model sees a green light?
 $P(\text{brake} | \text{light})$

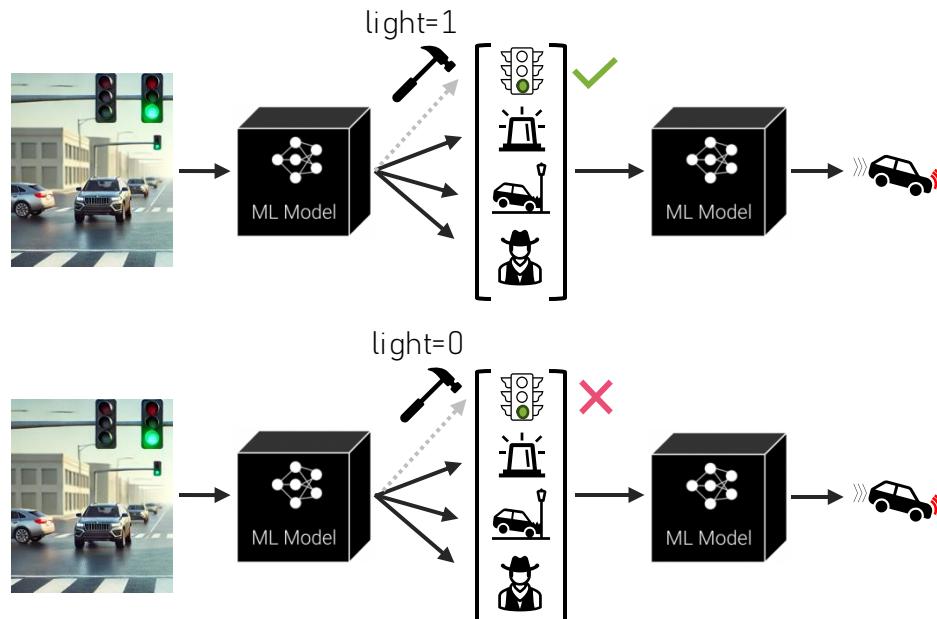
CAUSAL CONCEPT EFFECT



Proposed Solution

Step 3: Compute **difference of expected values**: absolute value is proportional to **causal effect**

$$\text{CaCE} = \mathbb{E}[\text{brake} | do(\text{light} = 1)] - \mathbb{E}[\text{brake} | do(\text{light} = 0)] = -0.8$$



$$\mathbb{E}[\text{brake} | do(\text{light} = 1)] = 0.2$$

$$\mathbb{E}[\text{brake} | do(\text{light} = 0)] = 1$$

Intervention

What if I set the light color to red?
 $P(\text{brake} | do(\text{light}))$

Association

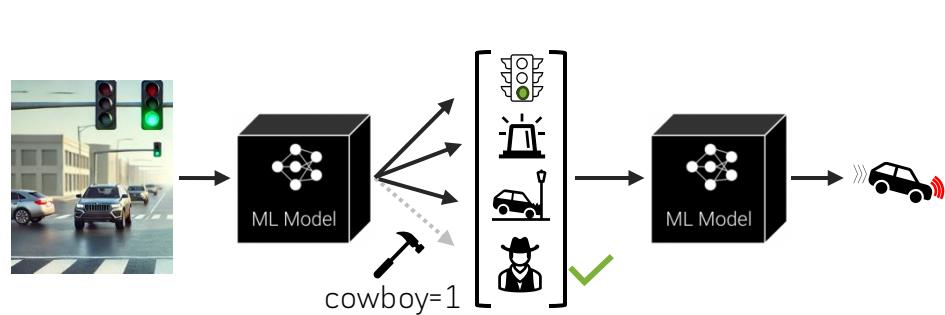
What if the model sees a green light?
 $P(\text{brake} | \text{light})$

CAUSAL CONCEPT EFFECT

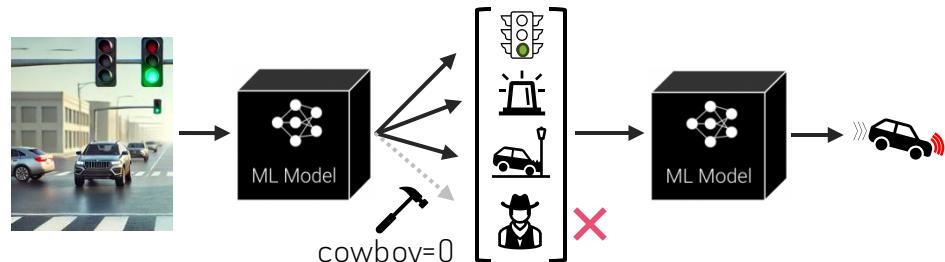


Proposed Solution

Step 3: Compute **difference of expected values**: absolute value is proportional to **causal effect**



$$\mathbb{E}[\text{brake} \mid \text{do}(\text{cowboy} = 1)] = 0.5$$



$$\mathbb{E}[\text{brake} \mid \text{do}(\text{cowboy} = 0)] = 0.5$$

Intervention

What if I set the light color to red?
 $P(\text{brake} \mid \text{do}(\text{light}))$

Association

What if the model sees a green light?
 $P(\text{brake} \mid \text{light})$

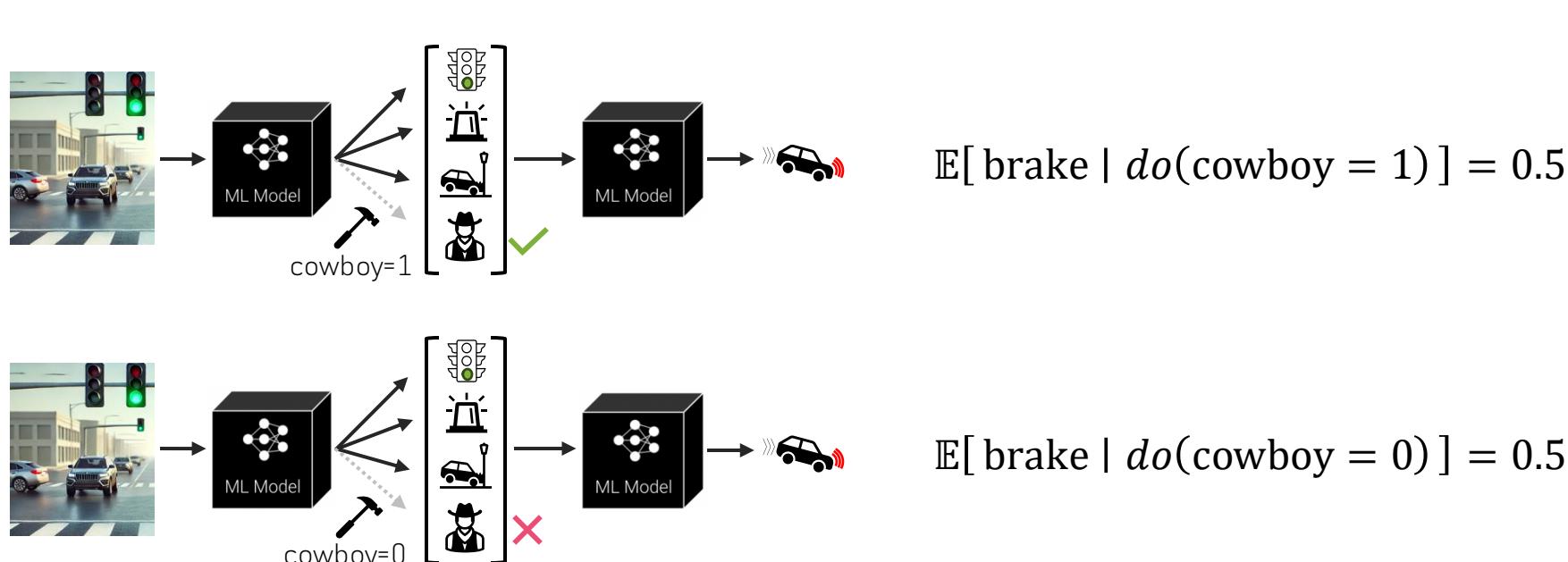
CAUSAL CONCEPT EFFECT



Proposed Solution

Step 3: Compute **difference of expected values**: absolute value is proportional to **causal effect**

$$\text{CaCE} = \mathbb{E}[\text{brake} | do(\text{cowboy} = 1)] - \mathbb{E}[\text{brake} | do(\text{cowboy} = 0)] = 0$$



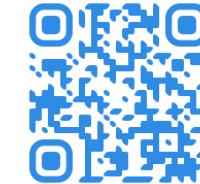
Intervention

What if I set the light color to red?
 $P(\text{brake} | do(\text{light}))$

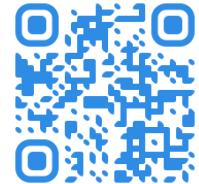
Association

What if the model sees a green light?
 $P(\text{brake} | \text{light})$

COUNTERFACTUAL CBMS



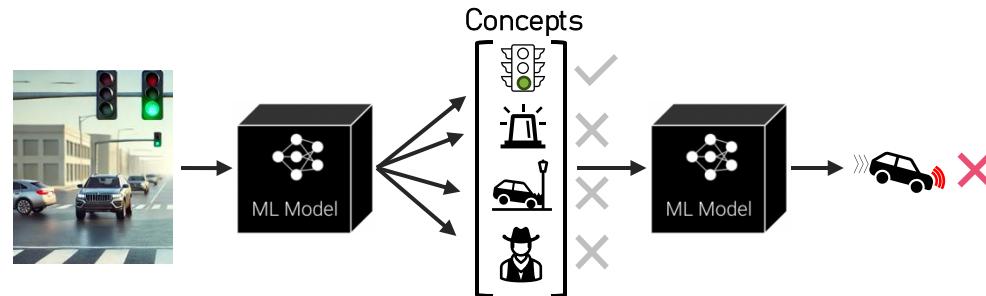
ICLR25



ICML22

Limitation Being Addressed

CBMs cannot answer **counterfactual queries**!



Counterfactual

What would have been predicted in the same circumstance had a car crash be seen?

$$P(\text{brake} \mid \text{light}, \text{crash})$$

Intervention

What if I set the light color to red?

$$P(\text{brake} \mid \text{do}(\text{light}))$$

Association

What if the model sees a green light?

$$P(\text{brake} \mid \text{light})$$

[1] Dominici et al. "Counterfactual Concept Bottleneck Models." ICLR (2025)

[2] Abid et al. "Meaningfully debugging model mistakes using conceptual counterfactual explanations." ICML (2022)

COUNTERFACTUAL CBMS

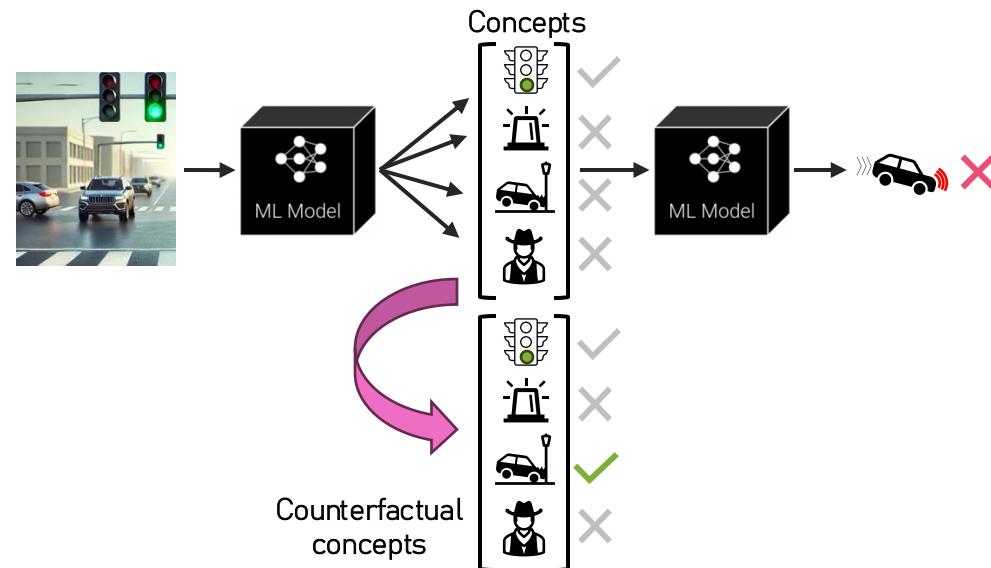


ICLR25

ICML22

Proposed Solution

Step 1: Generate **counterfactual concept activations**



Counterfactual

What would have been predicted in the same circumstance had a car crash be seen?
 $P(\text{brake} \mid \text{light}, \text{crash})$

Intervention

What if I set the light color to red?
 $P(\text{brake} \mid \text{do}(\text{light}))$

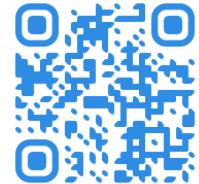
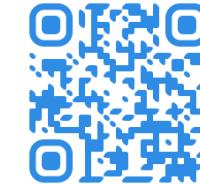
Association

What if the model sees a green light?
 $P(\text{brake} \mid \text{light})$

[1] Dominici et al. "Counterfactual Concept Bottleneck Models." ICLR (2025)

[2] Abid et al. "Meaningfully debugging model mistakes using conceptual counterfactual explanations." ICML (2022)

COUNTERFACTUAL CBMS

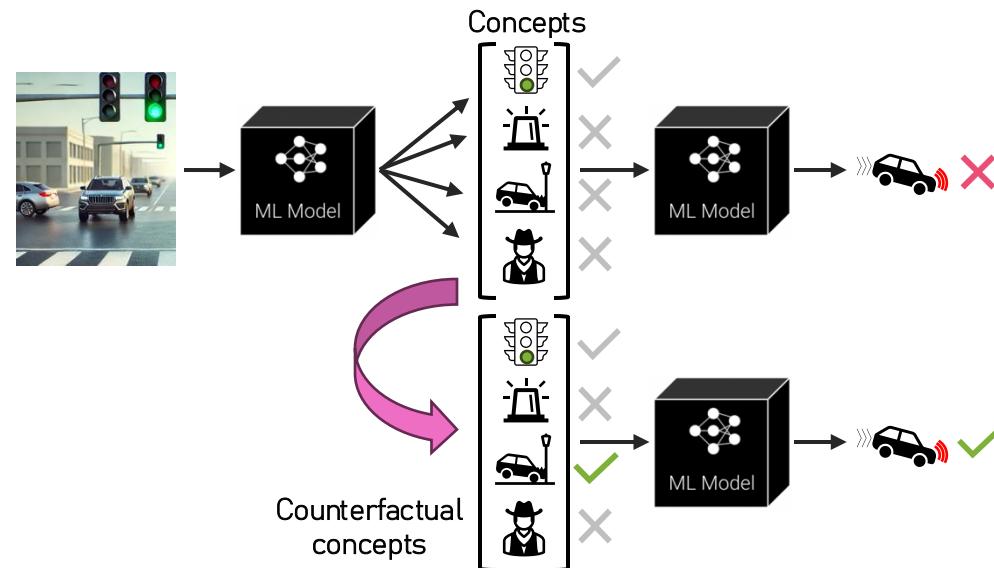


ICLR25

ICML22

Proposed Solution

Step 2: Compute **causal effect** on the task!



Counterfactual

What would have been predicted in the same circumstance had a car crash be seen?

$$P(\text{brake} \mid \text{light}, \text{crash})$$

Intervention

What if I set the light color to red?

$$P(\text{brake} \mid \text{do}(\text{light}))$$

Association

What if the model sees a green light?

$$P(\text{brake} \mid \text{light})$$

[1] Dominici et al. "Counterfactual Concept Bottleneck Models." ICLR (2025)

[2] Abid et al. "Meaningfully debugging model mistakes using conceptual counterfactual explanations." ICML (2022)

DIRECT COUNTERFACTUAL DEPENDENCE

So far, we have been making 2 **strong assumptions**...

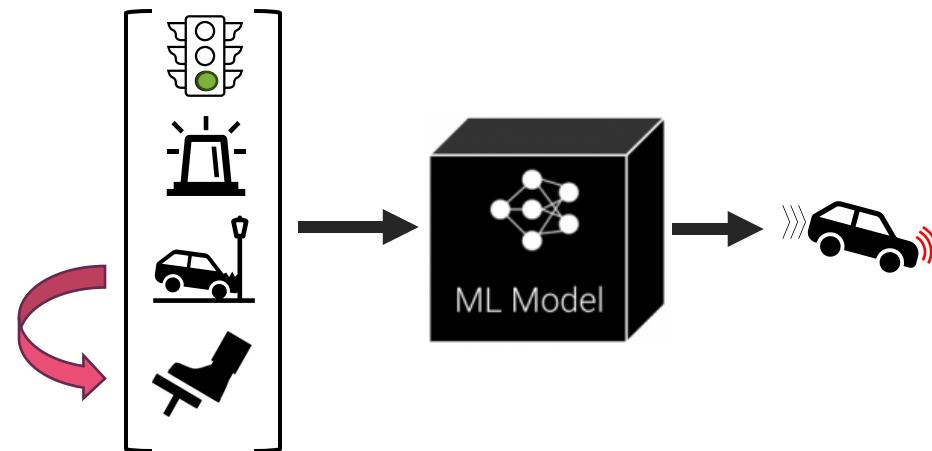


DIRECT COUNTERFACTUAL DEPENDENCE

So far, we have been making 2 **strong assumptions**:

- Concepts are **mutually independent**

Intervening on “car crash” does not increase the likelihood of hitting the brakes!



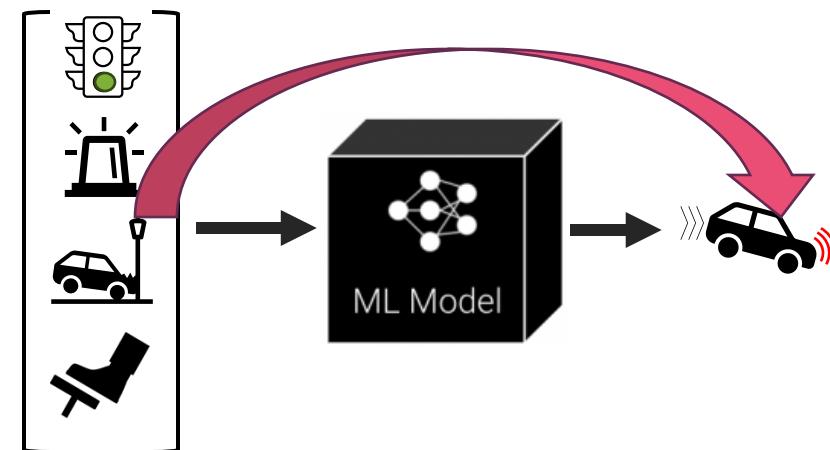
DIRECT COUNTERFACTUAL DEPENDENCE

So far, we have been making 2 **strong assumptions**:

- Concepts are **mutually independent**
- Concepts are **direct causes** of the task

Intervening on “car crash” does not increase the likelihood of hitting the brakes!

Intervening on “car crash” directly causes the car to brake!



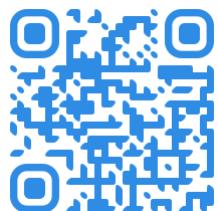
CONCEPT GRAPH MODELS

Limitation Being Addressed

CBMs (as most XAI methods) assume **direct counterfactual dependence!**



ICLR25

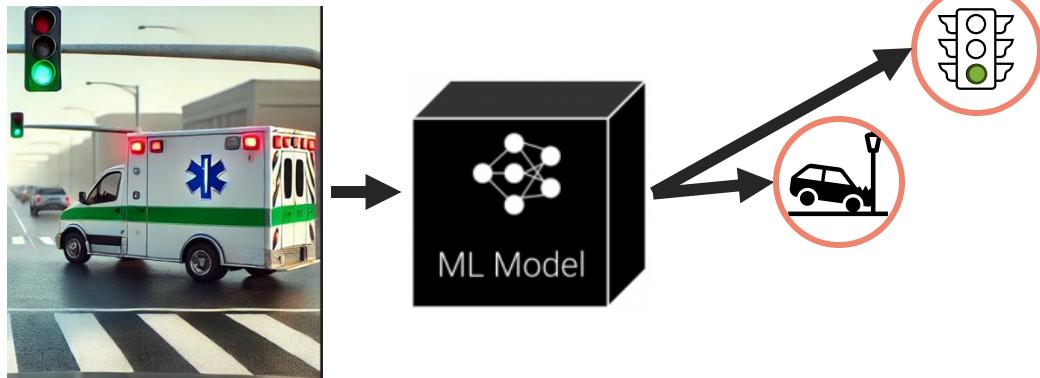


CLeaR24

CONCEPT GRAPH MODELS

Proposed Solution

Enforce inference through a **concept graph**!



- [1] Dominici et al. "Causal Concept Graph Models: Beyond Causal Opacity in Deep Learning" ICLR 2025
[2] Moreira et al. "Diconstruct: Causal concept-based explanations through black-box distillation" CLeaR 2024



ICLR25

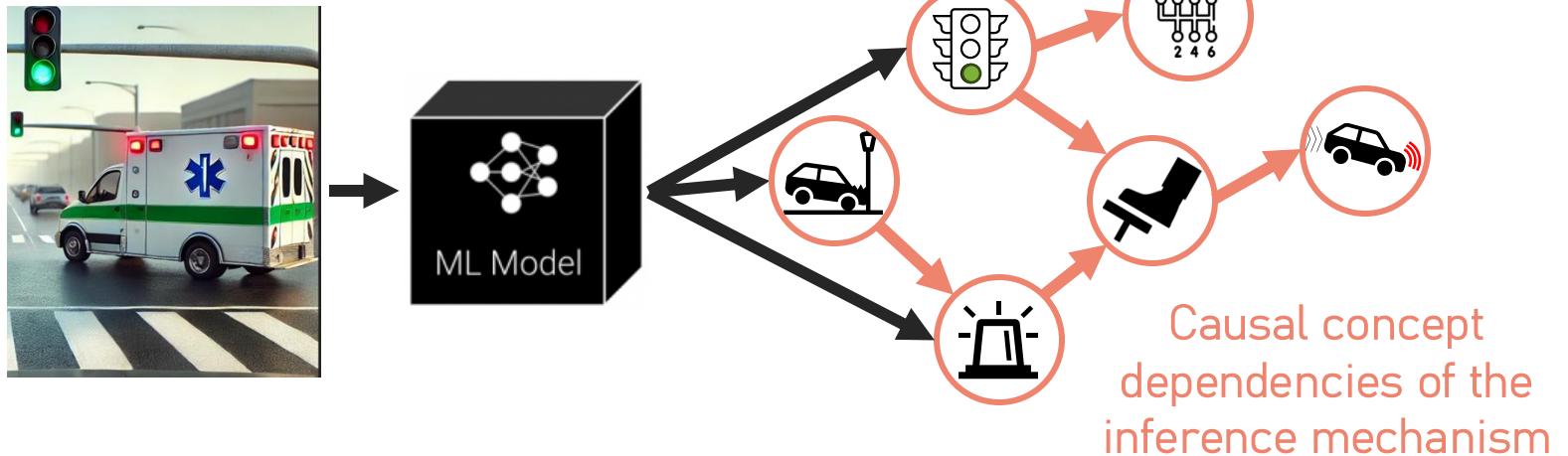


CLeaR24

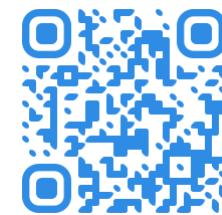
CONCEPT GRAPH MODELS

Proposed Solution

Enforce inference through a **concept graph**!



- [1] Dominici et al. "Causal Concept Graph Models: Beyond Causal Opacity in Deep Learning" ICLR 2025
[2] Moreira et al. "Diconstruct: Causal concept-based explanations through black-box distillation" CLeaR 2024



ICLR25



CLeaR24

CONCEPT GRAPH MODELS

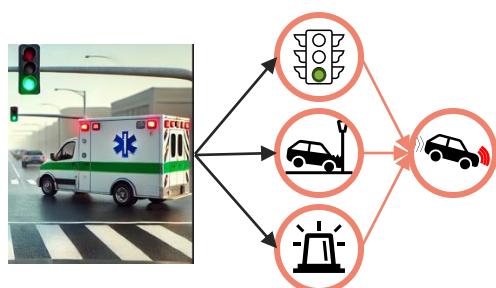
Proposed Solution

The concept graph can be:

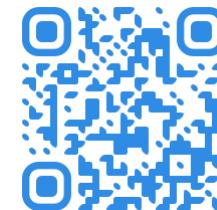
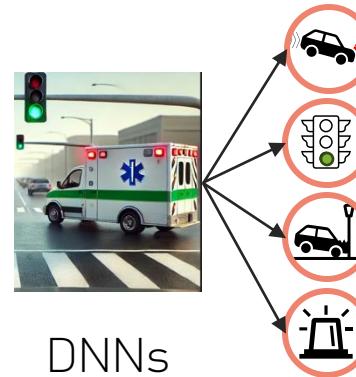
- Given as a prior



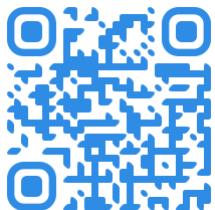
CBMs



DNNs



ICLR25



CLeaR24

[1] Dominici et al. "Causal Concept Graph Models: Beyond Causal Opacity in Deep Learning" ICLR 2025

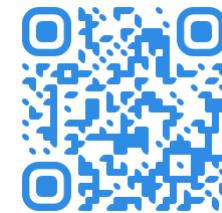
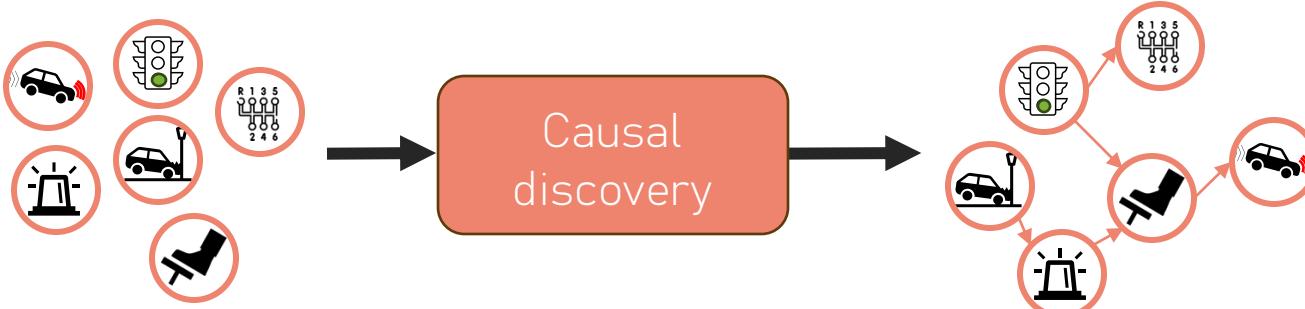
[2] Moreira et al. "Diconstruct: Causal concept-based explanations through black-box distillation" CLeaR 2024

CONCEPT GRAPH MODELS

Proposed Solution

The concept graph can be:

- Given as a prior
- Extracted from data with causal discovery techniques

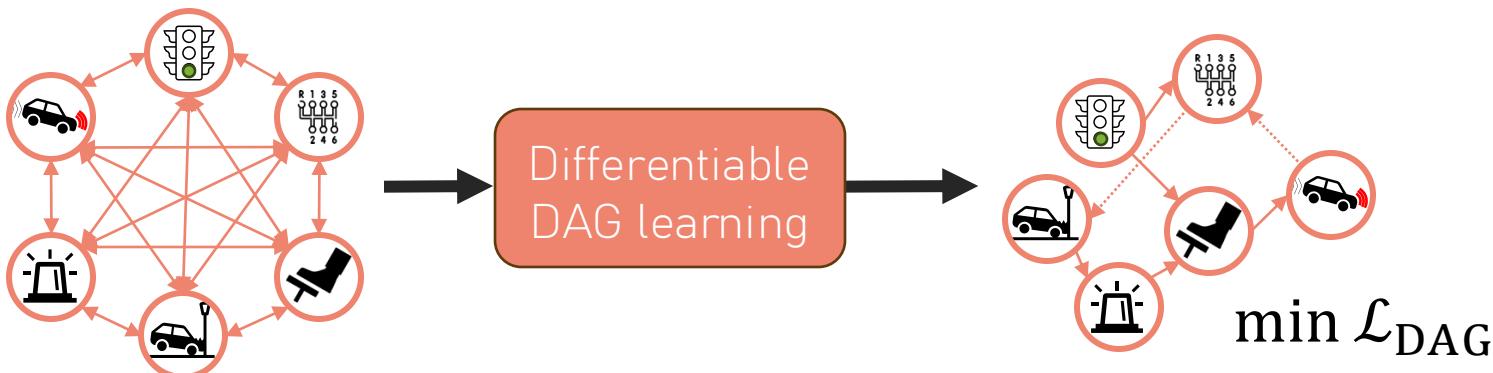


CONCEPT GRAPH MODELS

Proposed Solution

The concept graph can be:

- Given as a prior
- Extracted from data with causal discovery techniques
- Obtained with differentiable DAG learning



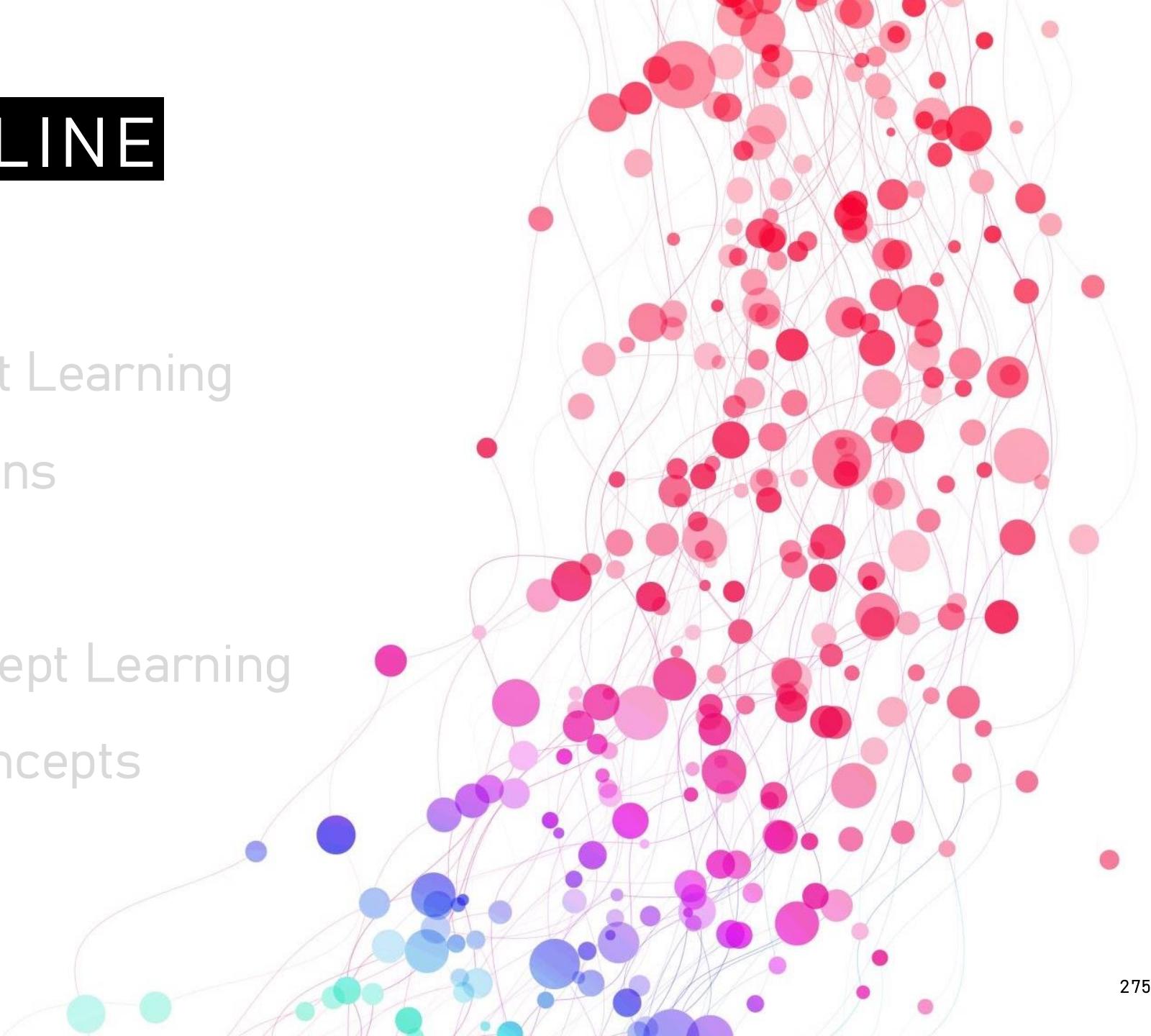
ICLR25



CLeaR24

TUTORIAL OUTLINE

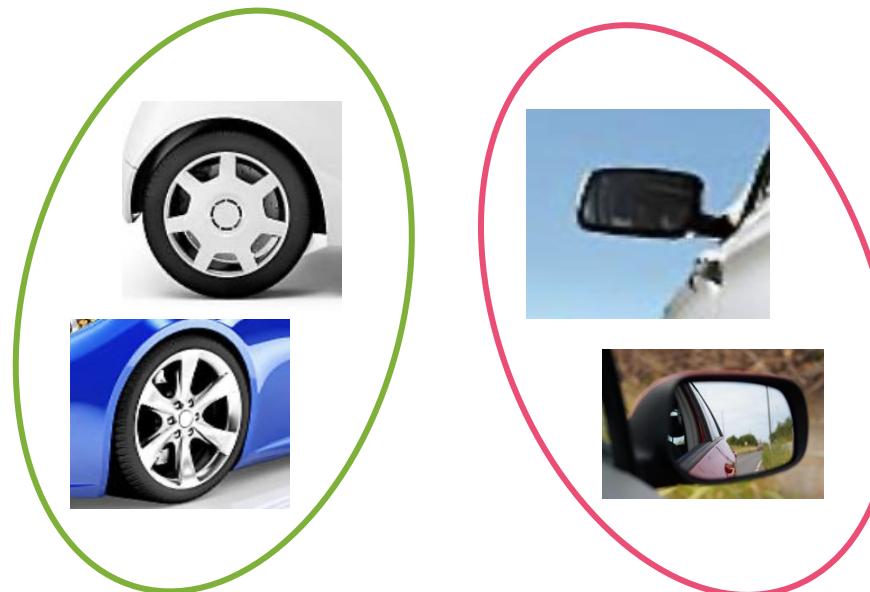
1. Introduction
2. Supervised Concept Learning
3. Concept Interventions
4. Q&A + Break
5. Unsupervised Concept Learning
6. Reasoning With Concepts
7. Future Directions
8. Q&A



AMAZING CONCEPTS & WHERE TO FIND THEM

Concept interpretability is **not the first nor the only** area focusing on concepts!

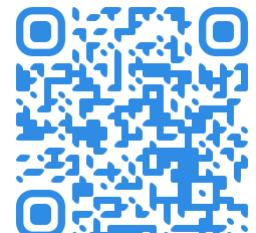
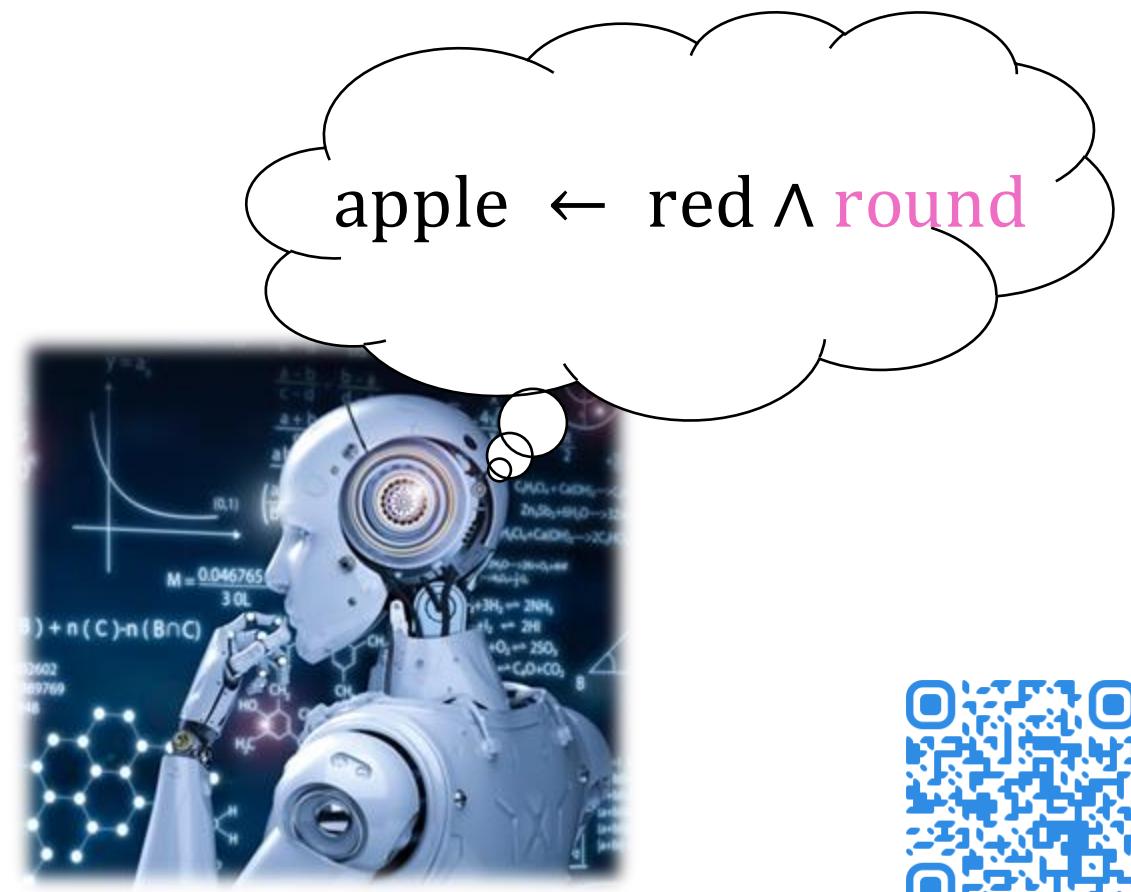
- Prototypes (clustering)



AMAZING CONCEPTS & WHERE TO FIND THEM

Concept interpretability is **not the first nor the only** area focusing on concepts!

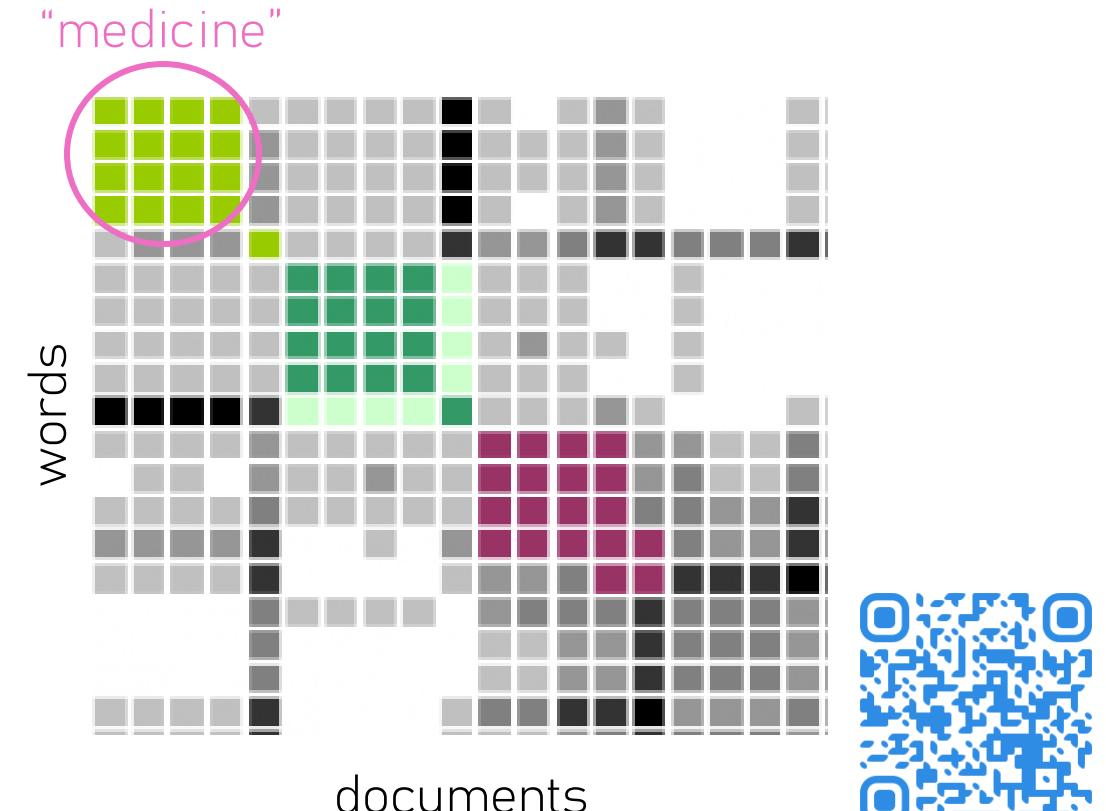
- Prototypes (clustering)
- Symbols (logic, neural-symbolic AI)



AMAZING CONCEPTS & WHERE TO FIND THEM

Concept interpretability is **not the first nor the only** area focusing on concepts!

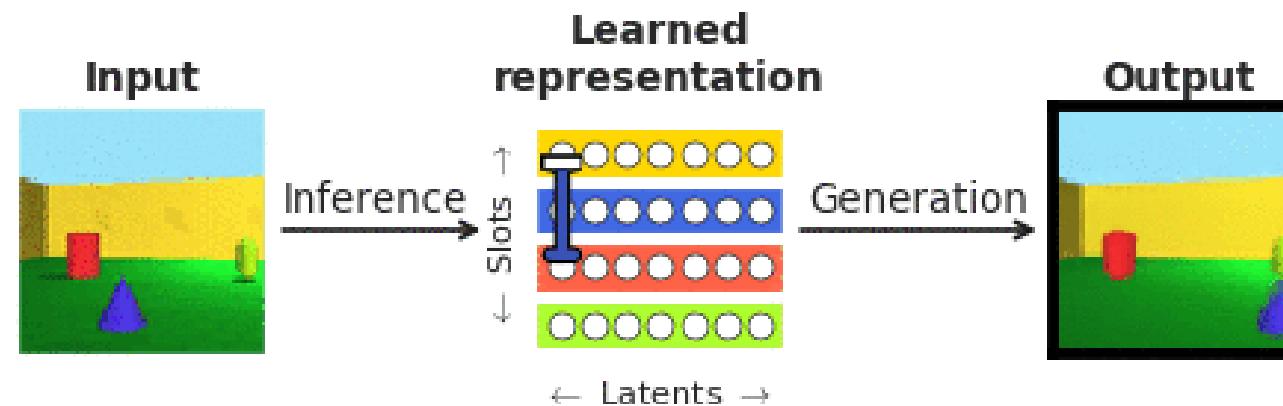
- Prototypes (clustering)
- Symbols (logic, neural-symbolic AI)
- Topic models (semantic analysis)



AMAZING CONCEPTS & WHERE TO FIND THEM

Concept interpretability is **not the first nor the only** area focusing on concepts!

- Prototypes (clustering)
- Symbols (logic, neural-symbolic AI)
- Topic models (semantic analysis)
- Factors of variation (disentanglement learning)



AMAZING CONCEPTS & WHERE TO FIND THEM

Concept interpretability is **not the first nor the only** area focusing on concepts!

... but it has a few key **differences**:

- Focus on **intervenability** & different **forms of transparency** (semantic, functional, causal)
- For this reason, often **different assumptions** hold (e.g., concepts don't have to be independent as in disentanglement learning!)

OPEN CHALLENGES

Label-free models are currently not as reliable as supervised ones

- How to effectively intervene in label-free settings?



OPEN CHALLENGES

Label-free models are currently not as reliable as supervised ones

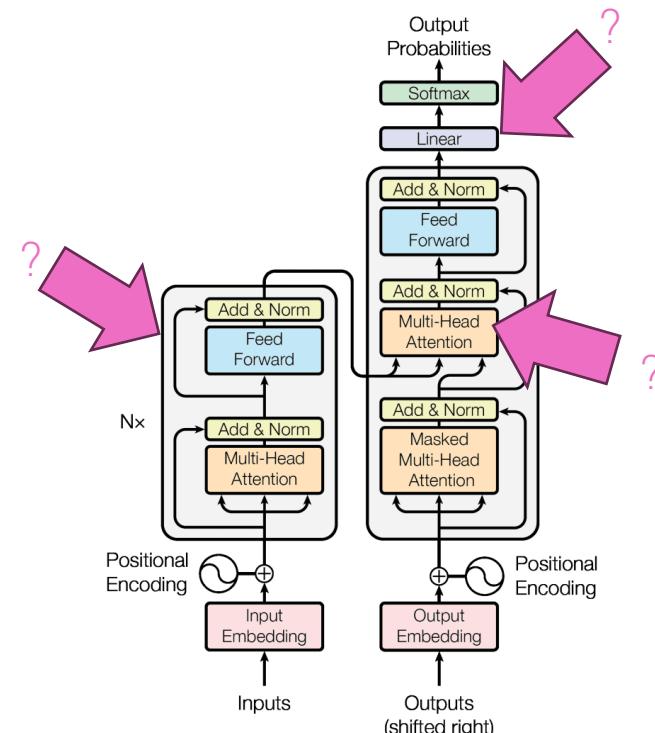
- How to effectively intervene in label-free settings?
- How to construct robust annotations without pre-trained domain-specific models?



OPEN CHALLENGES

Concept-based models are currently not designed nor integrated to **scale** to large models

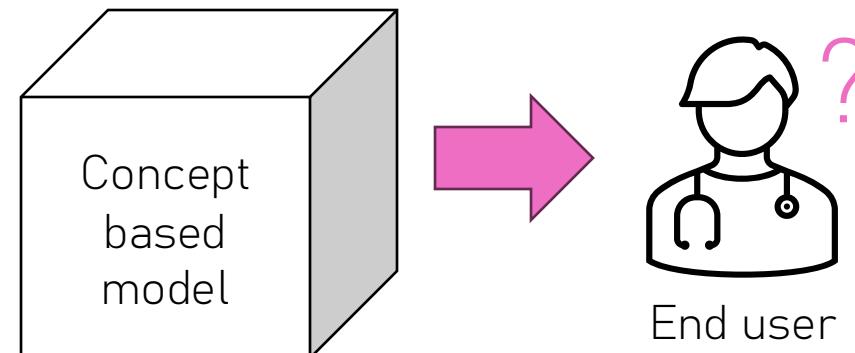
- **Where** (autoregressive, sentence, or paragraph) should we **look for/place** concepts in large models?
- **Should** large models **reason** based on concepts?



OPEN CHALLENGES

Concept-based models are currently not designed nor integrated to **scale** to large models

- **Where** (autoregressive, sentence, or paragraph) should we **look for/place** concepts in large models?
- **Should** large models **reason** based on concepts?
- **Which** guidelines should we follow to **deploy** concept-based models in the wild?



OPEN CHALLENGES

Some concepts are intrinsically hard to **represent** or **intervene on**

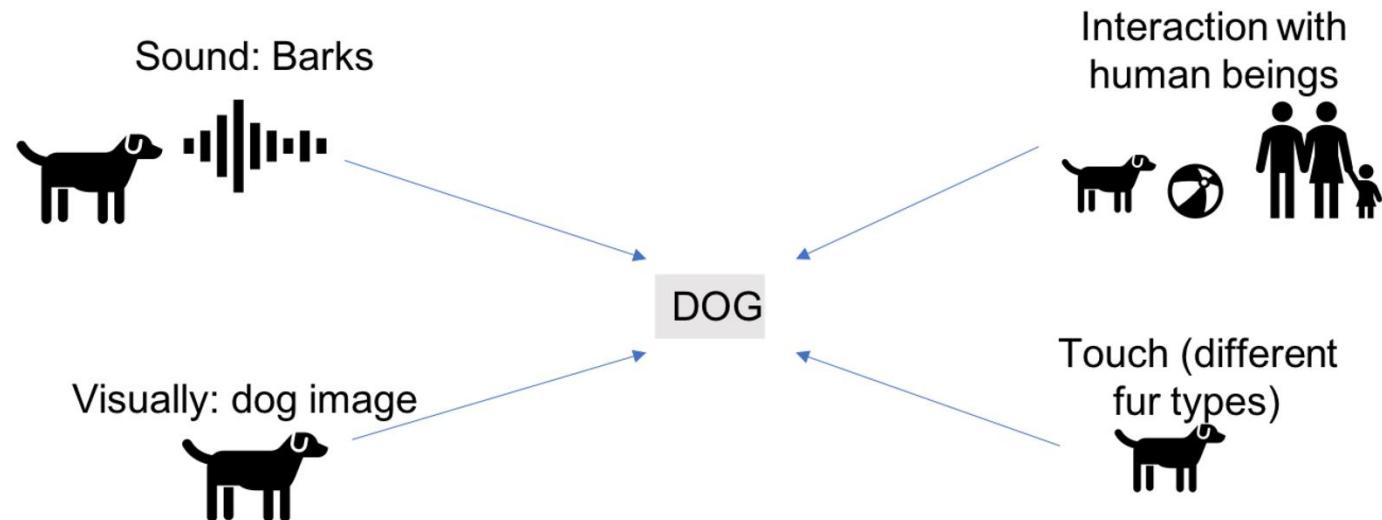
- **How** to deal with abstract (e.g., moral) or subjective concepts (e.g., aesthetics)?



OPEN CHALLENGES

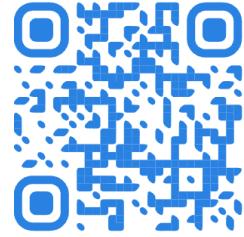
Some concepts are intrinsically hard to **represent** or **intervene on**

- **How** to deal with abstract (e.g., moral) or subjective concepts (e.g., aesthetics)?
- **How** to construct and intervene on multi-modal concepts?



RESOURCES

conceptlearning.github.io/



Cornerstone papers highlighted in this presentation

Extended bibliography on the tutorial website and in the slide deck's appendix

Explaining Classifiers with Causal Concept Effect (CaCE)

Yash Goyal^{1*}, Amir Feder², Uri Shalit², Been Kim³

model's prediction, their explanations might be confounded correlations present in the data which are not causally related to the model, as we describe now.

what drives classification decisions in a neural network, e.g. "what caused the output 'BICYCLE'?", "what caused the output 'BIRD'?"

Abstract

How can we understand classification decisions made by deep neural networks? Many explanation methods rely solely on correlations and fail to account for confounding, which result in potentially misleading explanations. To overcome this problem, we define the Causal Concept Effect (CaCE) as the causal effect of a concept on a deep neural net's measure of the presence or absence of a human concept. We show that the CaCE measure stemming from confounding is difficult in situations where we do not have access to the do-operator. To mitigate this, we propose a method to learn a causal concept effect (CaCE) that can be used to explain the predictions of a deep neural network.

Concept Bottleneck Models

Pang Wei Koh^{*1}, Thao Nguyen^{*1,2}, Yew Siang Tang^{*1}, Stephen Mussmann¹, Emma Pierson¹, Been Kim², Percy Liang¹

We seek to learn models that we can interact with using high-level concepts: if the model did not think there was a bone spur in the x-ray, would it still predict severe arthritis? State-of-the-art models today do not typically support the manipulation of concepts like "the existence of bone spurs", as they are trained end-to-end to go directly from raw input (e.g., pixels) to output (e.g.,

ON COMPLETENESS-AWARE CONCEPT-BASED EXPLANATIONS IN DEEP NEURAL NETWORKS

A PREPRINT

Chih-Kuan Yeh¹, Been Kim², Sercan Ö. Arık³, Chun-Liang Li³, Tomas Brodin¹, Michael Riegler¹, Tommi Jaakkola¹, Zoubin Ghahramani¹

¹Machine Learning Department, Carnegie Mellon University
²Google Brain
³Google Cloud AI

Interpretable Concept-Based Memory Reasoning

David Debof
KU Leuven
david.debof@kuleuven.be

Pietro Barbiero
Università della Svizzera Italiana
University of Cambridge
barbiero@tutanota.com

Francesco Giannini
Scuola Normale Superiore
francesco.giannini@sns.it

Gabriele Ciravegna
DAUNI, Politecnico di Torino
gabriele.ciravegna@polito.it

Giuseppe Marra
KU Leuven
giuseppe.marra@kuleuven.be

Michelangelo Diligenti
University of Siena
michelangelo.diligenti@unisi.it

Abstract

Concept Bottleneck Models (CBMs) tackle the opacity of neural architectures by reasoning and explaining their predictions using a set of high-level concepts. A key property of these models is that they permit *concept interventions*, wherein incorrect predicted concepts are identified and thus improved the model's performance. However, this has shown that intervention efficacy can be highly dependent on which concepts are intervened on and on the model's architecture parameters. We argue that this is rooted in a CBM's lack of transparency in the decision-making processes of deep learning systems presents a significant challenge in modern artificial intelligence (AI), as it impairs users' ability to rely on and verify these systems. To address this challenge, Concept-Based Models (CBMs) have made significant progress by incorporating human-interpretable concepts into deep learning architectures. This approach proposes the model to be appropriately receptive to concept interventions. We propose Intervention-aware Concept Embedding models (I-CE), a CBM-based architecture and training paradigm that improves

Learning to Receive Help: Intervention-Aware Concept Embedding Models

Mateo Espinosa Zarlenga
University of Cambridge
me466@cam.ac.uk

Katherine M. Collins
University of Cambridge
kmc61@cam.ac.uk

Adrian Weller
University of Cambridge
965@cam.ac.uk

Zohreh Shams
University of Cambridge
zs315@cam.ac.uk

Mateja Jamnik
University of Cambridge
mateja.jamnik@cl.cam.ac.uk

Krishnamurthy (Dj) Dvijotham
Google DeepMind
dvij@google.com

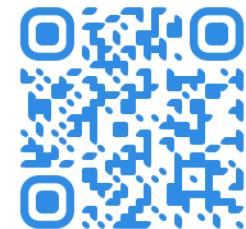
Abstract

Bottleneck Models (CBMs) tackle the opacity of neural architectures by reasoning and explaining their predictions using a set of high-level concepts. A key property of these models is that they permit *concept interventions*, wherein incorrect predicted concepts are identified and thus improved the model's performance. However, this has shown that intervention efficacy can be highly dependent on which concepts are intervened on and on the model's architecture parameters. We argue that this is rooted in a CBM's lack of transparency in the decision-making processes of deep learning systems presents a significant challenge in modern artificial intelligence (AI), as it impairs users' ability to rely on and verify these systems. To address this challenge, Concept-Based Models (CBMs) have made significant progress by incorporating human-interpretable concepts into deep learning architectures. This approach proposes the model to be appropriately receptive to concept interventions. We propose Intervention-aware Concept Embedding models (I-CE), a CBM-based architecture and training paradigm that improves

RESOURCES



PyC



@github

@medium

Some **XAI libraries** implement concept-based techniques (check out tutorial website!)

We are working on **PyTorch Concepts (PyC)**, a library dedicated to concept-based interpretability

- APIs are designed to **implement existing models**, but also to **support the development** of new ones
- Currently **supports** concept-based: data types, layers, interventions, metrics, models
- The PyC team is publishing **hands-on tutorials** on Medium!

```
encoder = torch.nn.Sequential(  
    torch.nn.Linear(n_features, latent_dims),  
    torch.nn.LeakyReLU(),  
)  
concept_bottleneck = LinearConceptLayer(latent_dims, [concept_names])  
y_predictor = torch.nn.Sequential(  
    torch.nn.Flatten(),  
    torch.nn.Linear(n_concepts, latent_dims),  
    torch.nn.LeakyReLU(),  
    LinearConceptLayer(latent_dims, [task_names]),  
)  
model = torch.nn.Sequential(encoder, concept_bottleneck, y_predictor)  
  
# generate concept and task predictions  
emb = encoder(x_train)  
c_emb = concept_emb_bottleneck(emb)  
c_pred = concept_score_bottleneck(c_emb)  
c_intervened = CF.intervene(c_pred, c_train, intervention_indexes)  
c_mix = CF.concept_embedding_mixture(c_emb, c_intervened)  
y_pred = y_predictor(c_mix)
```

[1] https://github.com/pyc-team/pytorch_concepts

[2] https://medium.com/@pyc_devteam

A FEW THINGS TO BRING BACK HOME!

Thank you for your time! Before leaving, remember that concept-based interpretability:

- is **connected** to other AI areas, but it focuses on **specific research questions** (intervenability and different forms of opacity)
- can make things **easier** (human interaction)... or **worse** (need for annotations)
- is a relatively young research field, so there's a lot of **work to do** for all of us!

Read our Medium stories to implement your first concept-based model in **<15 minutes!**



PyC



conceptlearning.github.io/

