

L193 – Lecture 6 – Lent 2025







Are <20 interactions sufficient to fully understand this ML model?

If yes, raise your hand!

Disclaimer: Unfortunately for you, I'm always right during the lecture—but I'm happy to be proven wrong afterward! 😊





How many interactions do we need to fully understand the model?



 $c_i \in \{0,1\}$

ML Model

How many interactions do we need to fully understand the model?



How many interactions do we need to fully understand the model?



How many interactions do we need to fully understand the model?

We can recover the full Conditional Probability / Truth Table!



How many interactions do we need to fully understand the model?

We can recover the full Conditional Probability / Truth Table!

#interactions required to extract full CPT/TT is exponential in #inputs!

Can we do better?



(INFORMAL) OBJECTIVE

Build a **general-purpose** neural model that is:

Expressive (as DNNs)

and whose inference mechanism is



Functionally transparent (we fully understand CPT/TT)

- □ Tractable (CPT/TT size << exp)
- Semantically transparent (concept-based)
- **Causally transparent** (based on non-trivial cause-effect chains)



[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025)...



- General methods (foundations)
- Neural Interpretable Reasoning paradigm
- Grounding NIR (Concept Memory Reasoning example)
- Causal reasoning (metrics and models)





Functional transparency: understanding of inference mechanism is tractable

 $|\mathsf{CPT}(f)| \ll \exp(|Z|)$

Methods:

• Filter irrelevant features out





[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025).

Functional transparency: understanding of inference mechanism is tractable

$|\operatorname{CPT}(f)| \ll \exp(|Z|)$

Methods:

- Filter irrelevant features out
- **Re-parametrize** global CPT as a **mixture** of simple (e.g., linear) "local CPTs"









[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025)_.

Semantic transparency: input features are aligned with human semantics

Methods:

- Work with tabular data
- Use **concept-based** approach

... and then write the model as P(Y | C)!





[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025).

General-purpose: model design should be applicable to any data type





[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025)...

General-purpose: model design should be applicable to any data type

Interpretability is a **Markovian property**:

 $(Y \perp X) \mid N_y$

Methods:

• Re-parametrize the inputs of an ML model *f* without affecting its interpretability (idea behind CBMs)!

$$P(Y \mid X) = \sum_{C} P(Y \mid C) P(C \mid X)$$

[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025)...



Introducing x does not affect the understanding of f!



Semantic opacity: data representations are not aligned with human concepts



Semantic opacity: data representations are not aligned with human concepts



Functional opacity: CPT is unknown or intractable to reconstruct



510

Does **semantic** transparency imply **functional** transparency?

In theory it does, but...

... with *n* inputs the size of a TT is $2^n!$ \rightarrow intractable!

Does **functional** transparency imply **semantic** transparency?

Potential solution: DeepProbLog







$\mathsf{FUNCTIONAL} \rightarrow \mathsf{SEMANTIC} \mathsf{TRANSP.}?$

Can we prove that the classification head is safe?





If yes, raise your hand!

Disclaimer: Unfortunately for you, I'm always right during the lecture—but I'm happy to be proven wrong afterward! 😊



Safety depends on concepts' meaning!



Disclaimer. Unfortunately for you, I'm always right during the lecture—but I'm happy to be proven wrong afterward! 😊

NEURAL INTERPRETABLE 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)



516

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

NEURAL INTERPRETABLE REASONING

Proposed Solution





517

2] Debot. David. et al. "Interpretable concept-based memory reasoning." NeurIPS 2024.

Proposed Solution

Step 1: DNN predicts concept activations





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

Proposed Solution

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







519

Proposed Solution

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







Proposed Solution (details)

Step 4: Execute rule using rule states and concept activations

Concepts



 $\begin{array}{c} \textbf{Rule states} \\ \textbf{NOT } \textbf{C}_1 \text{ AND } \textbf{C}_2 \end{array}$



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

Proposed Solution (details)

Step 4: Execute rule using rule states and concept activations

Execution of rule states on concept activations

Concepts





 $\begin{array}{c} \textbf{Rule states} \\ \textbf{NOT } \textbf{C}_1 \text{ AND } \textbf{C}_2 \end{array}$



Proposed Solution (details)

Step 4: Execute rule using rule states and concept activations

Execution of rule states on concept activations

Concepts





 $\mathbb{I}[r_i = \text{positive}]p(c_i \mid x) + \mathbb{I}[r_i = \text{negative}]p(\neg c_i \mid x) + \mathbb{I}[r_i = \text{irrelevant}]$

$$0 \times p(c_i \mid x) + 1 \times p(\neg c_i \mid x) + 0 = p(\neg c_i \mid x) = (1 - 0.2) = 0.8$$

 $\begin{array}{c} \textbf{Rule states} \\ \textbf{NOT } \textbf{C}_1 \text{ AND } \textbf{C}_2 \end{array}$



Proposed Solution (details)

Step 4: Execute rule using rule states and concept activations

Execution of rule states on concept activations

Concepts



$$\mathbb{I}[r_{i} = \text{positive}]p(c_{i} \mid x) + \mathbb{I}[r_{i} = \text{negative}]p(\neg c_{i} \mid x) + \mathbb{I}[r_{i} = \text{irrelevant}]$$

$$\stackrel{\bullet}{\longrightarrow} 0 \times p(c_{i} \mid x) + 1 \times p(\neg c_{i} \mid x) + 0 = p(\neg c_{i} \mid x) = (1 - 0.2) = 0.8$$

$$\stackrel{\bullet}{\longrightarrow} 1 \times p(c_{i} \mid x) + 0 \times p(\neg c_{i} \mid x) + 0 = p(c_{i} \mid x) = 0.9$$





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

Proposed Solution (details)

Step 4: Execute rule using rule states and concept activations

Execution of rule states on concept activations

Concepts



 $\begin{array}{c} \textbf{Rule states} \\ \textbf{NOT } \textbf{C}_1 \text{ AND } \textbf{C}_2 \end{array}$

$$\mathbb{I}[r_{i} = \text{positive}]p(c_{i} \mid x) + \mathbb{I}[r_{i} = \text{negative}]p(\neg c_{i} \mid x) + \mathbb{I}[r_{i} = \text{irrelevant}]$$

$$0 \times p(c_{i} \mid x) + 1 \times p(\neg c_{i} \mid x) + 0 = p(\neg c_{i} \mid x) = (1 - 0.2) = 0.8$$

$$1 \times p(c_{i} \mid x) + 0 \times p(\neg c_{i} \mid x) + 0 = p(c_{i} \mid x) = 0.9$$

$$0 \times p(c_{i} \mid x) + 0 \times p(\neg c_{i} \mid x) + 1 = 1$$



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

Proposed Solution

Step 5: Execute the rule combining concept states and activations to predict the output label



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

CMR has 3 key features:

• Universal approximator akin to opaque DNNs (Theorem 4.1)



CMR has 3 key features:

• Universal approximator akin to opaque DNNs (Theorem 4.1)

[°]‱∧ ¬C

• Minimize #relevant concepts \rightarrow understanding is **tractable**



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



CMR has 3 key features:

- Universal approximator akin to opaque DNNs (Theorem 4.1)
- Minimize #relevant concepts \rightarrow understanding is tractable
- The concept memory allows **formal verification** of properties





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

(INFORMAL) OBJECTIVE

Build a **general-purpose** neural model that is:

Expressive (as DNNs)

and whose inference mechanism is



Functionally transparent (we fully understand CPT/TT)
 Tractable (CPT/TT size << exp)
 Semantically transparent (concept-based)
 Causally transparent (based on non-trivial cause-effect chains)



[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025)...

COMPARING BOTTLENECKS

Is CBM #1 more causally transparent than CBM #2?







If yes, raise your hand!

Disclaimer: Unfortunately for you, I'm always right during the lecture—but I'm happy to be proven wrong afterward! 😊

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



CBM #1 is potentially equally transparent w.r.t. CBM #2!

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



Association

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

534

Sometimes intervening on wrongly predicted concepts helps...



Sometimes intervening on wrongly predicted concepts helps...

and sometimes it doesn't! 😢

Causal analysis can provide us with insights!



Sometimes intervening on wrongly predicted concepts helps... and sometimes it doesn't! 😢

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





Proposed Solution

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





Proposed Solution

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



[1] Goyal. Yash. et al. "Explaining classifiers with causal concept effect (cace)." arXiv preprint 2019.

Proposed Solution

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

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



[1] Goval, Yash, et al. "Explaining classifiers with causal concept effect (cace)," arXiv preprint 2019.



Proposed Solution

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



[1] Goval. Yash. et al. "Explaining classifiers with causal concept effect (cace)," arXiv preprint 2019.



Proposed Solution

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

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



[1] Goyal. Yash. et al. "Explaining classifiers with causal concept effect (cace)." arXiv preprint 2019.

COUNTERFACTUAL CBMS

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(brake | light, crash)

 Intervention

 What if I set the light color to red?

 P(brake | do(light))

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

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

[2] Abid. Abubakar. et al. "Meaningfully debugging model mistakes using conceptual counterfactual explanations." International Conference on Machine Learning. PMLR. 2022

COUNTERFACTUAL CBMS

Proposed Solution

Step 1: Generate counterfactual concept activations





ICLR25

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

[1] Dominici. Gabriele. et al. "Counterfactual Concept Bottleneck Models." ICLR 2025.
[2] Abid. Abubakar. et al. "Meaningfully debugging model mistakes using conceptual counterfactual explanations." International Conference on Machine Learning. PMLR. 2022.

ICMI 22

COUNTERFACTUAL CBMS



Step 2: Compute causal effect on the task!





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

Association

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

[1] Dominici, Gabriele, et al. "Counterfactual Concept Bottleneck Models." ICLR 2025.
[2] Abid. Abubakar. et al. "Meaningfully debugging model mistakes using conceptual counterfactual explanations." International Conference on Machine Learning. PMLR, 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



Proposed Solution

Enforce inference through a **concept graph**!



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

549

CLeaR24

Proposed Solution

The concept graph can be:

- Given as a prior





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

550

Proposed Solution

The concept graph can be:

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



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

Proposed Solution

The concept graph can be:

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



552

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

NEURAL INTERPRETABLE REASONING

Build a **general-purpose** neural model that is:

Expressive (as DNNs)

and whose inference mechanism is



Semantically transparent (concept-based)
 Functionally transparent (we fully understand CPT/TT)
 Tractable (CPT/TT size << exp)
 Causally transparent (based on non-trivial cause-effect chains)



[1] Barbiero, Pietro, et al. "Neural Interpretable Reasoning." arXiv preprint (2025)...