



# The Search for Foundations in AI Interpretability

Maxime Peyrard



**Why Explain?**

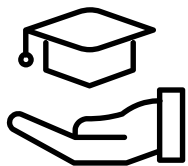
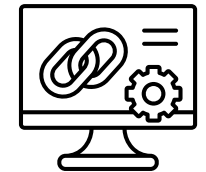
# Practical Reasons

Transparency, Trust



Accountability, Legal Compliance

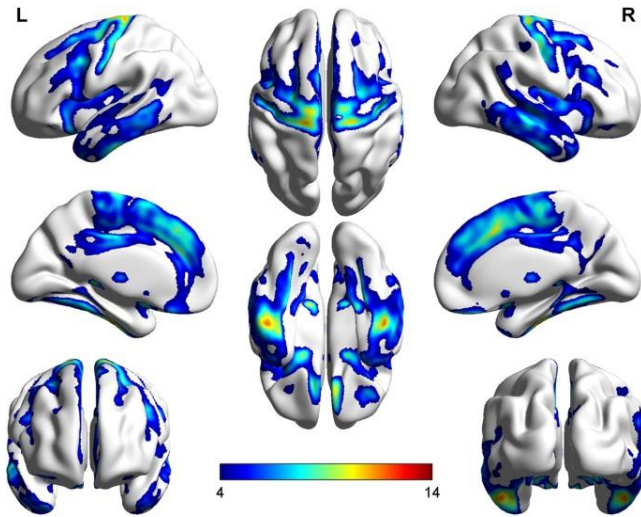
Error Diagnostic, Continuous Improvement, Robustness



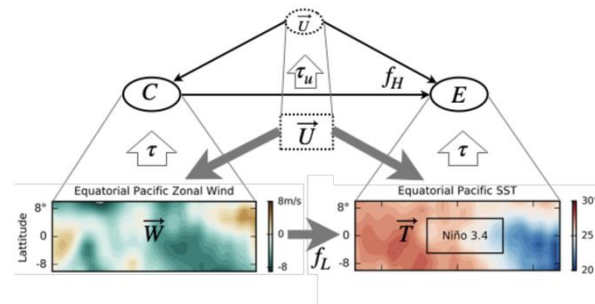
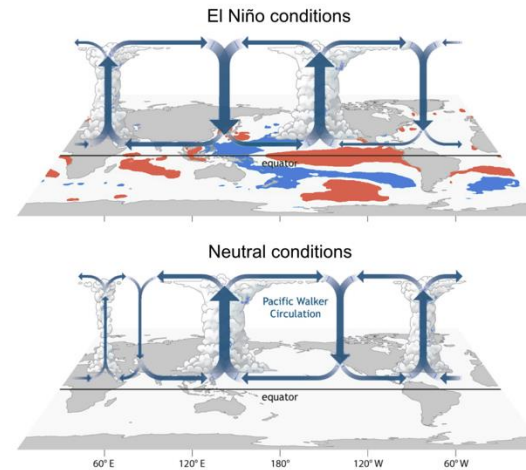
Communication, Education, Human-AI collaboration

...

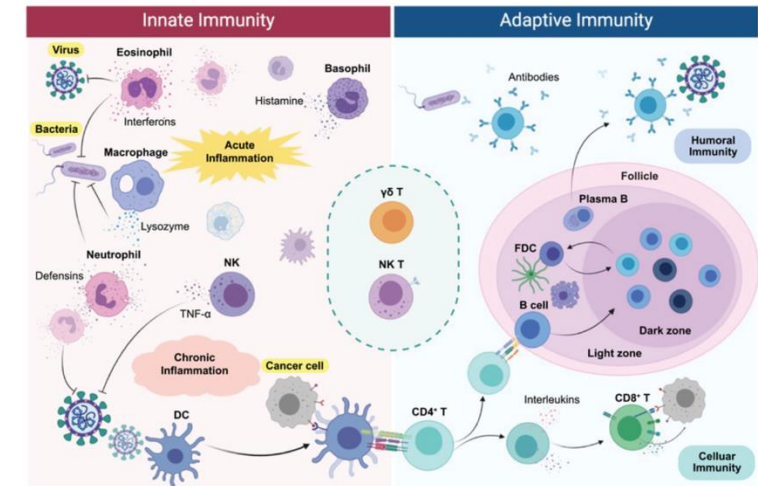
# Scientific Question: Explaining Complex Systems



Wei, et al 2018  
Nature

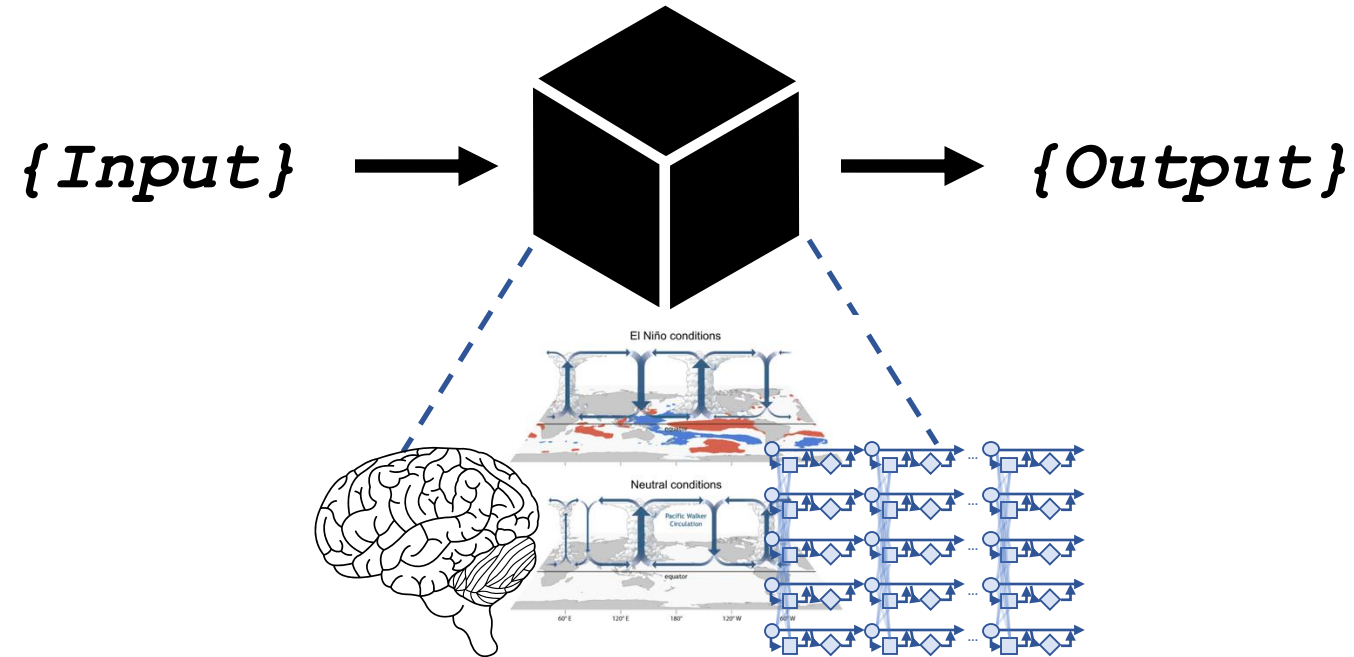


Chalupka, et al  
2018 UAI  
Beckers, et al,  
2019 UAI



Chen, et al  
2020 AM

# Interpretability research



What ***should be*** explained? What ***is*** an explanation? What is a ***valid*** explanation?

*How to go from observations to **valid explanations**?*

# Als are the *Simplest* Complex Systems to Study

Testing the Tools of Systems Neuroscience on Artificial Neural Networks

Grace W. Lindsay

AI are fully observable and manipulable

What does it mean to understand a neural network?

Timothy P. Lillicrap & Konrad P. Kording

Forms of explanation and understanding for neuroscience and artificial intelligence

Jessica A. F. Thompson

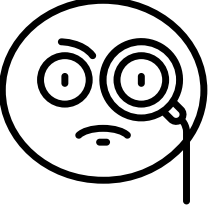
**Contributions and challenges for network models in cognitive neuroscience**

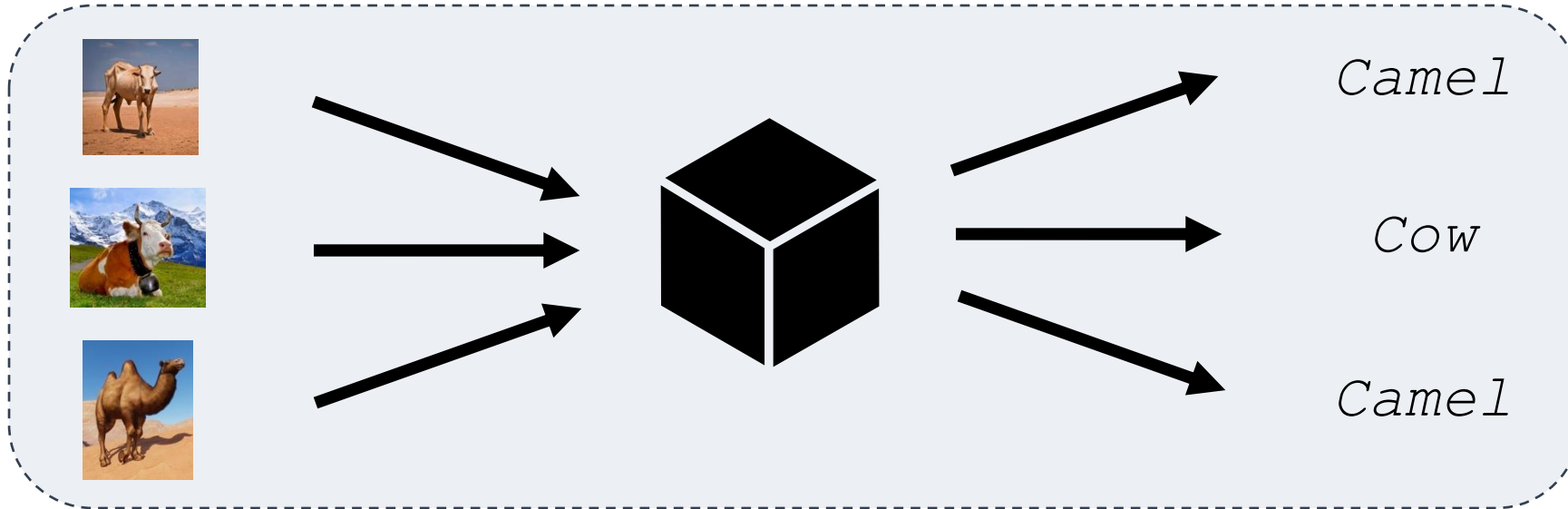
[Olaf Sporns](#) 

# **How to Explain?**

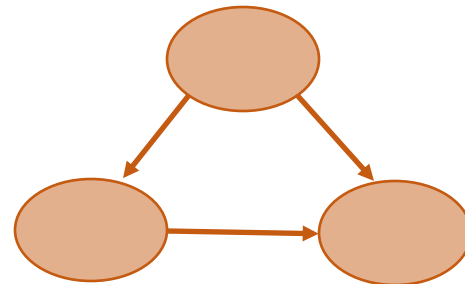
## **“Behavioral”**

# Behavioral Testing

  
Inputs  
outputs  
inspection



Carefully craft inputs, measure effects on outputs, come-up with hypothesis



High-level explanation



# Benchmarking

GAIA Leaderboard						
Model name	Average score (%)	Level 1 score (%)	Level 2 score (%)	Level 3 score (%)	organisation	Model family
<a href="#">Multi-Agent Experiment v0.1</a>	32.33	47.31	28.93	14.58	MSR AI Frontiers	GPT-4-turbo
<a href="#">MAAC_V1</a>	25.0					
<a href="#">FRIDAY</a>	24.0					
<a href="#">FRIDAY_without_learning</a>	21.0					
<a href="#">DIP</a>	15.0					
<a href="#">Chamomile</a>	14.0					
<a href="#">GPT4 + manually selected plu</a>	14.0					
<a href="#">Clarity_v1</a>	14.0					
<a href="#">Warm-up Act</a>	12.0					
<a href="#">stealth3</a>	9.3					
<a href="#">stealth2</a>	8.9					
<a href="#">stealth</a>	8.6					

Open LLM Leaderboard

LLM Benchmark

Metrics through time

About

FAQ

Submit

Search for your model (separate multiple queries with ,)

Select columns to show

☒
Average

☒
ARC

☒
HellaSwag

☒
Winogrande

☒
GSM8K

☐
Type

☐
Merged

☐
Hub License

☐
#Params

Hide models

☒
Private or deleted

☒
Contains a merge/mod

T

Model

davidkim205/Rhea-72b-v0.5

Contamination/contaminated\_proof

MTSAIR/MultiVerse\_70B

MTSAIR/MultiVerse\_70B

SF-Foundation/Ein-72B-v0.11

LMSYS Chatbot Arena Leaderboard

[Vote](#)
[Blog](#)
[GitHub](#)
[Paper](#)
[Dataset](#)
[Twitter](#)
[Discord](#)

LMSYS [Chatbot Arena](#) is a crowdsourced open platform for LLM evals. We've collected over 500,000 human preference votes to ra

Arena Elo

Full Leaderboard

Total #models: 76. Total #votes: 511252. Last updated: March 29, 2024.

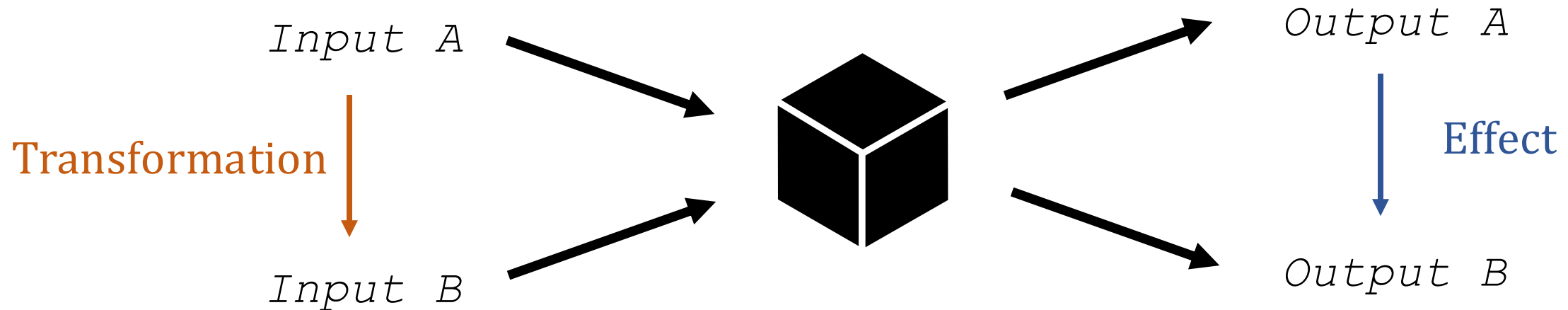
Contribute your vote at [chat.lmsys.org](#)! Find more analysis in the [notebook](#).

Rank	Model	Arena Elo	95% CI	Votes	Organization
1	<a href="#">Claude 3 Opus</a>	1255	+3/-4	37663	Anthropic
1	<a href="#">GPT-4-1106-preview</a>	1252	+3/-3	56936	OpenAI
1	<a href="#">GPT-4-0125-preview</a>	1249	+3/-4	38105	OpenAI
4	<a href="#">Bard (Gemini Pro)</a>	1204	+5/-5	12468	Google
4	<a href="#">Claude 3 Sonnet</a>	1200	+3/-4	40389	Anthropic

LMSYS Chatbot Arena Leaderboard							
<a href="#">Vote</a>   <a href="#">Blog</a>   <a href="#">GitHub</a>   <a href="#">Paper</a>   <a href="#">Dataset</a>   <a href="#">Twitter</a>   <a href="#">Discord</a>							
LMSYS <a href="#">Chatbot Arena</a> is a crowdsourced open platform for LLM evals. We've collected over 500,000 human preference votes to rank LLMs with the Elo ranking system.							
Arena Elo   Full Leaderboard							
Total #models: 76. Total #votes: 511252. Last updated: March 29, 2024.							
Contribute your vote 🗳️ at <a href="#">chat.lmsys.org</a> ! Find more analysis in the <a href="#">notebook</a> .							
Rank	Model	Arena Elo	95% CI	Votes	Organization	License	Knowledge Cutoff
1	<a href="#">Claude 3 Opus</a>	1255	+3/-4	37663	Anthropic	Proprietary	2023/8
1	<a href="#">GPT-4-1106-preview</a>	1252	+3/-3	56936	OpenAI	Proprietary	2023/4
1	<a href="#">GPT-4-0125-preview</a>	1249	+3/-4	38105	OpenAI	Proprietary	2023/12
4	<a href="#">Bard (Gemini Pro)</a>	1204	+5/-5	12468	Google	Proprietary	Online
4	<a href="#">Claude 3 Sonnet</a>	1200	+3/-4	40389	Anthropic	Proprietary	2023/8
6	<a href="#">GPT-4-0314</a>	1185	+4/-4	35803	OpenAI	Proprietary	2021/9
7	<a href="#">Claude 3 Haiku</a>	1177	+3/-4	26773	Anthropic	Proprietary	2023/8
8	<a href="#">GPT-4-0613</a>	1160	+3/-5	54509	OpenAI	Proprietary	2021/9

# Controlled Setups

Idea: controlled **transformation** of the inputs, and measure **effects** on outputs



Evaluating Models' Local Decision  
Boundaries via Contrast Sets  
EMNLP 2020

Learning What Makes a Difference from  
Counterfactual Examples  
ECCV 2021

# Controlled Setups - Examples

## Noisy Exemplars Make Large Language Models More Robust: A Domain-Agnostic Behavioral Analysis

**Hongyi Zheng**  
New York University  
hz2212@nyu.edu

**Abulhair Saparov**  
New York University  
as17582@nyu.edu

### Original

Natalia sold 48 clips in April. Then she sold half as many clips in May. How many clips did Natalia sell altogether?



She sold  $48 \div 2 = 24$  clips in May, so she sell  $48 + 24 = 72$  clips altogether. ✓

### Perturbed

Natalia sold 48 clips in April. Then she sold half as many clips in May. **Then she sold half as many clips in May.** How many clips did Natalia sell altogether?

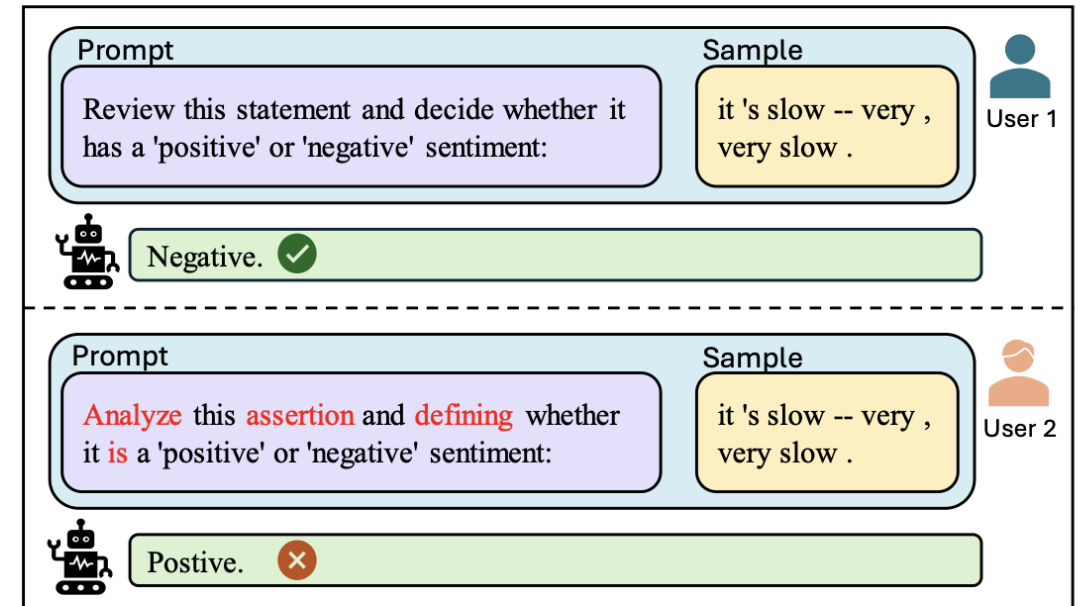


She sold  $48 \div 2 \div 2 = 12$  clips in May, so she sell  $48 + 12 = 60$  clips altogether. ✗

## PromptBench: Towards Evaluating the Robustness of Large Language Models on Adversarial Prompts

Kaijie Zhu<sup>1,2\*</sup>, Jindong Wang<sup>1†</sup>, Jiaheng Zhou<sup>2</sup>, Zeek Wang<sup>1</sup>, Hao Chen<sup>3</sup>, Yidong Wang<sup>4</sup>, Linyi Yang<sup>5</sup>, Wei Ye<sup>4</sup>, Yue Zhang<sup>5</sup>, Neil Zhenqiang Gong<sup>6</sup>, Xing Xie<sup>1</sup>

<sup>1</sup>Microsoft Research <sup>2</sup>Institute of Automation, CAS <sup>3</sup>Carnegie Mellon University  
<sup>4</sup>Peking University <sup>5</sup>Westlake University <sup>6</sup>Duke University



(b) Synonyms lead to errors in sentiment analysis problems.

# Input Feature Attributions



controlled changes on the input features  $\rightarrow$  effects on output

**LIME:** local approximation of the boundary around an input

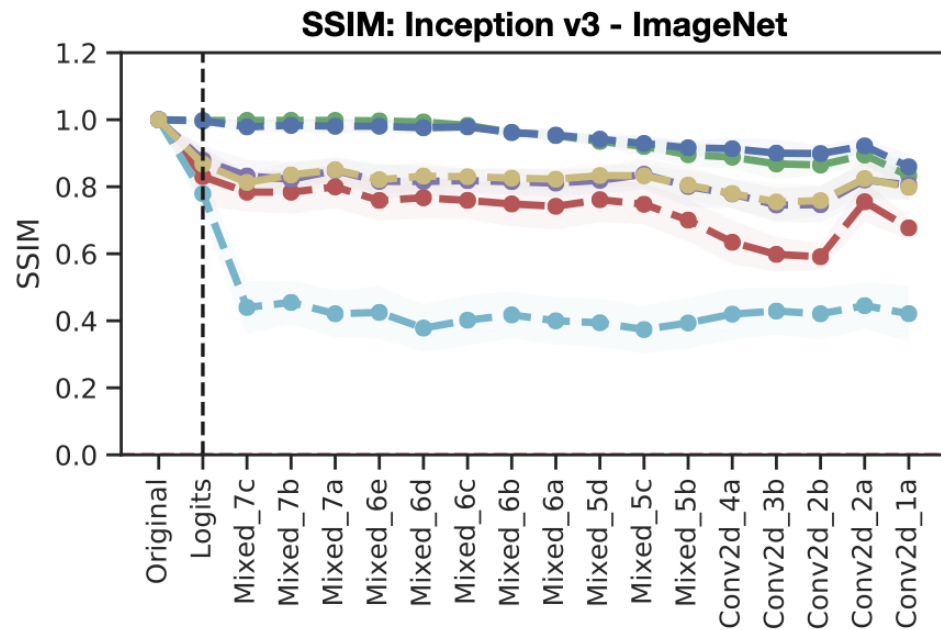
**SHAP:** measure each feature contribution relative to others

**Integrated Gradient:** use the gradient information backpropagated in the input features





# Problem with Feature Attributions



## Sanity Checks for Saliency Maps

**Julius Adebayo<sup>\*</sup>, Justin Gilmer<sup>#</sup>, Michael Muelly<sup>#</sup>, Ian Goodfellow<sup>#</sup>, Moritz Hardt<sup>#†</sup>, Been Kim<sup>#</sup>**  
juliusad@mit.edu, {gilmer,muelly,goodfellow,mrtz,beenkim}@google.com

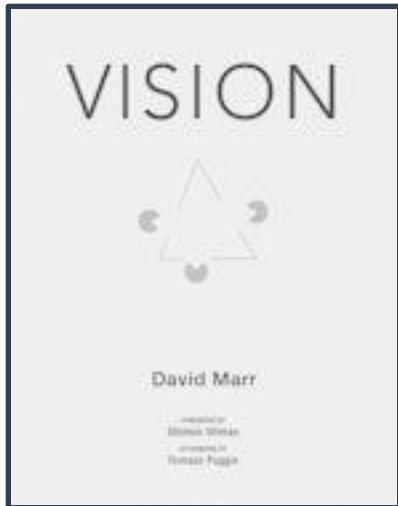
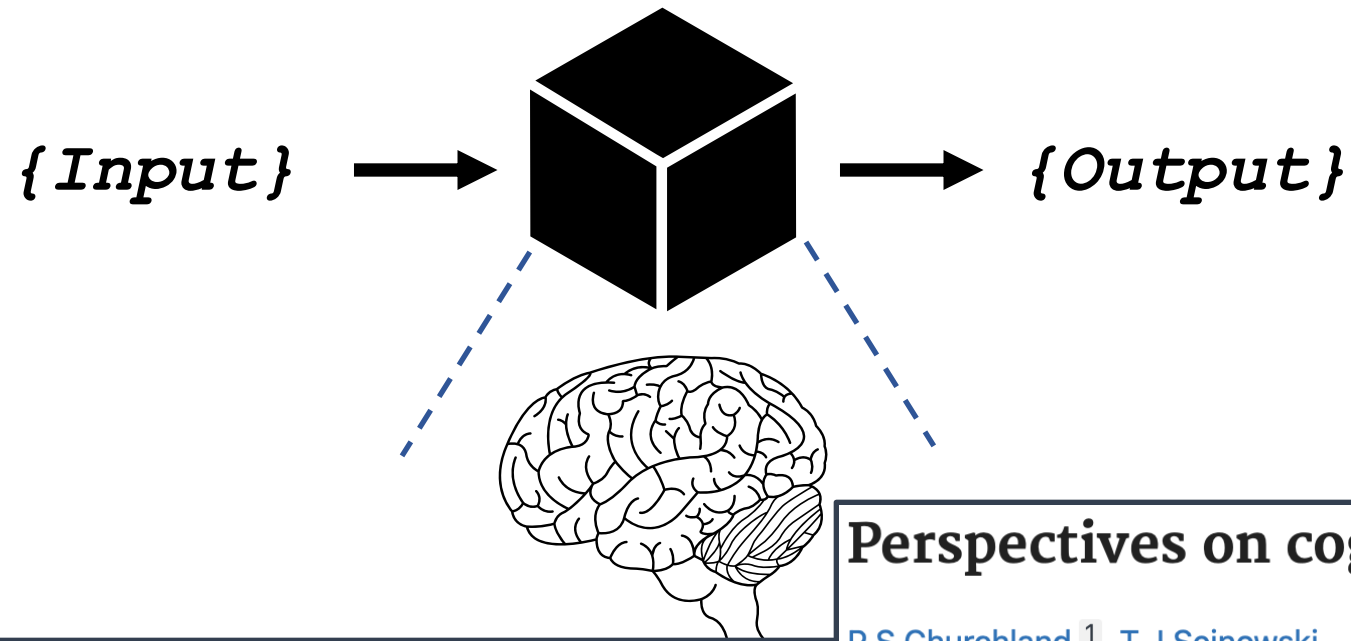
<sup>#</sup>Google Brain

<sup>†</sup>University of California Berkeley

Similar feature attributions for  
randomly initialized networks  
compared to trained ones

(

# Neuroscience detour I



Review Article | Published: 12 June 2008

## What we can do and what we cannot do with fMRI

[Nikos K. Logothetis](#)

[Nature](#) **453**, 869–878 (2008) | [Cite](#)

## Perspectives on cognitive neuroscience

[P S Churchland](#) <sup>1</sup>, [T J Sejnowski](#)

## Neural representation and the cortical code

[R C deCharms](#) <sup>1</sup>, [A Zador](#)

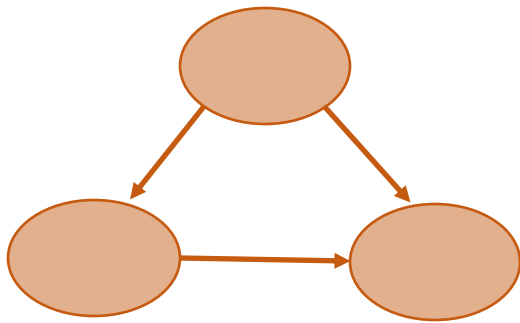
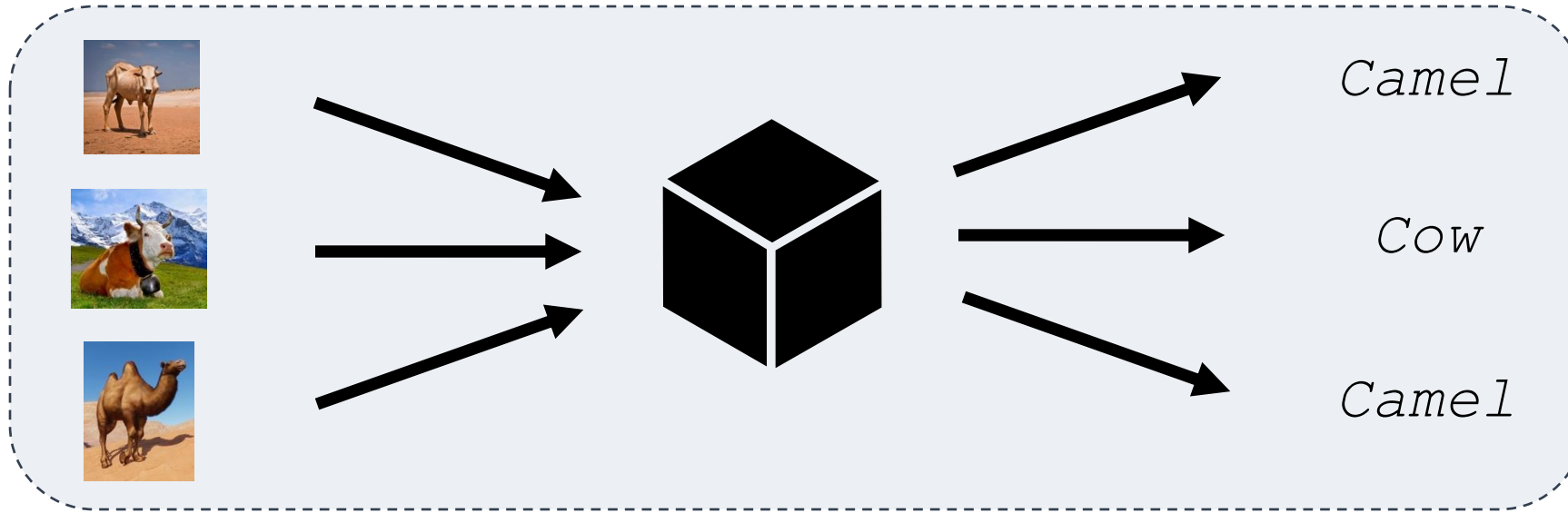
Behavior is not enough; we have to look at the computation to find objective, measurable, and generalizable predictors

)

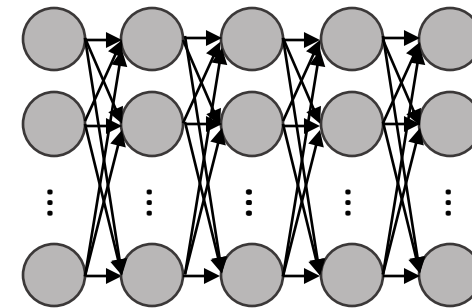


# Behavior vs Computation

Inputs  
outputs  
inspection



High-level explanation

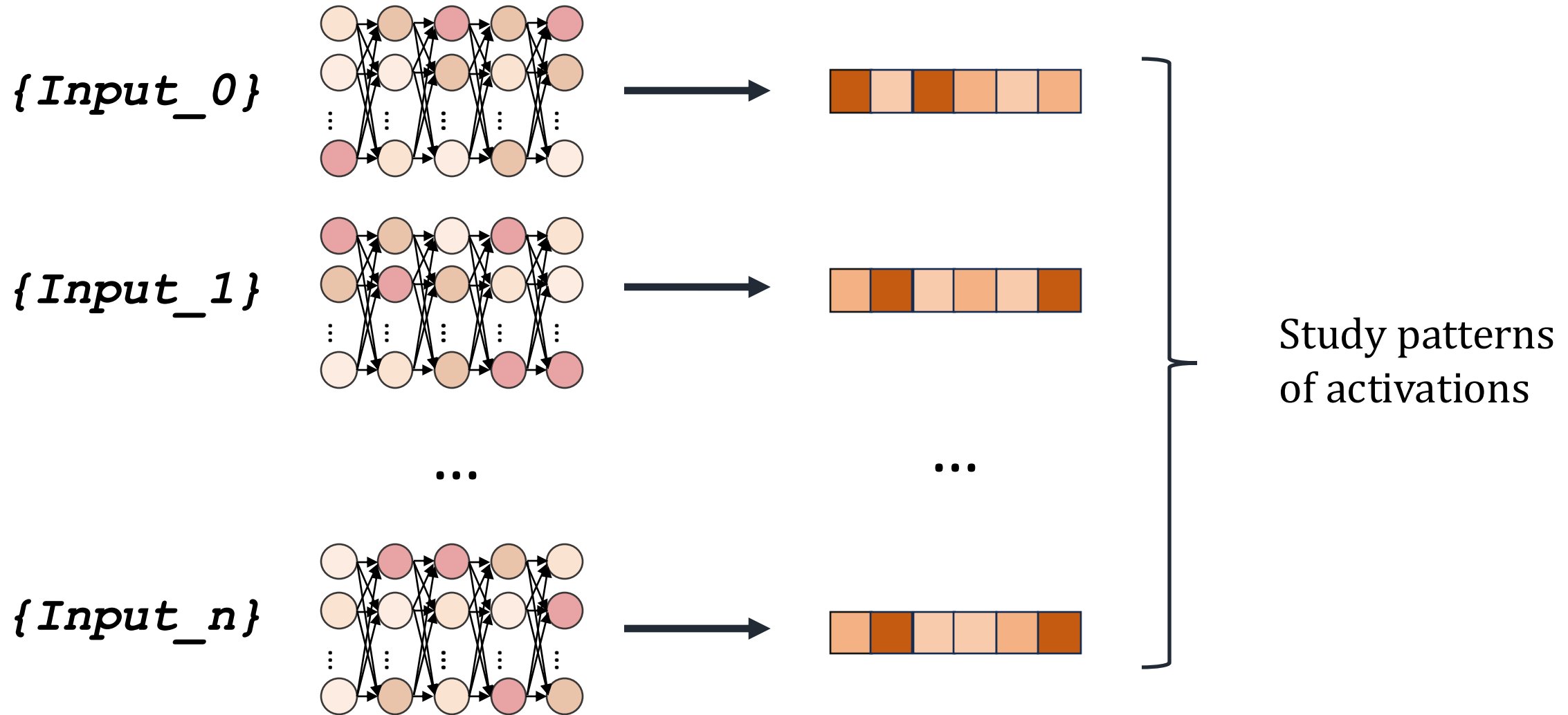


But is it **consistent with** the low-level implementation?

# **How to Explain?**

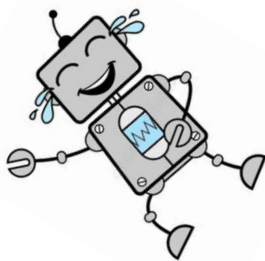
## **“Neural Correlates”**

# “Neural Correlates”



# Laughing Heads: Can Transformers Detect What Makes a Sentence Funny?

Maxime Peyrard, Beatriz Borges, Kristina Gligorić and Robert West  
EPFL

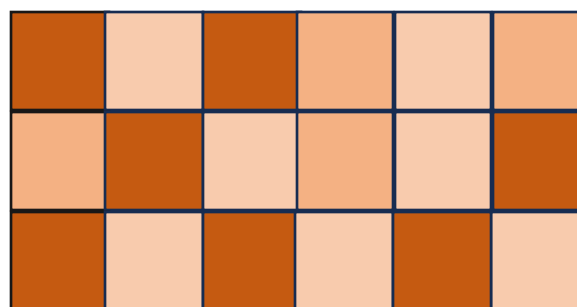


Satirical

City opens new art {jail}

Non-modified  
chunk in funny

Modified chunk  
in funny

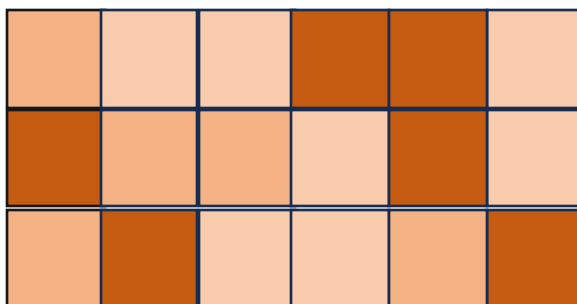


Serious

City opens new art {museum}

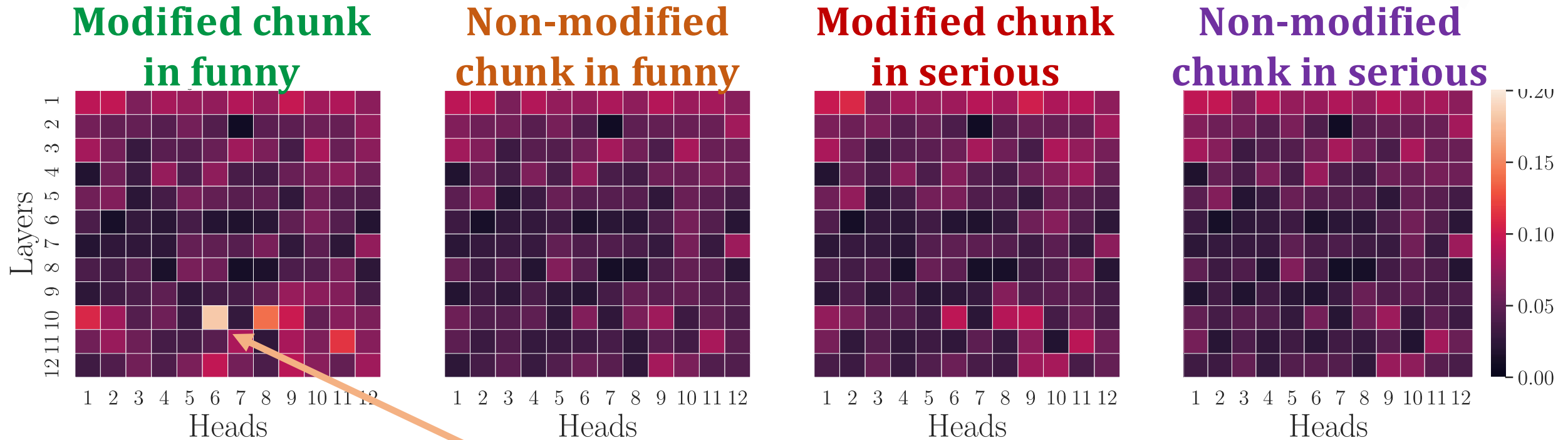
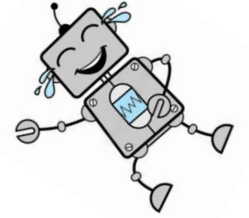
Non-modified chunk  
in serious

Modified chunk  
in serious



Matched  
comparison

# Laughing Heads



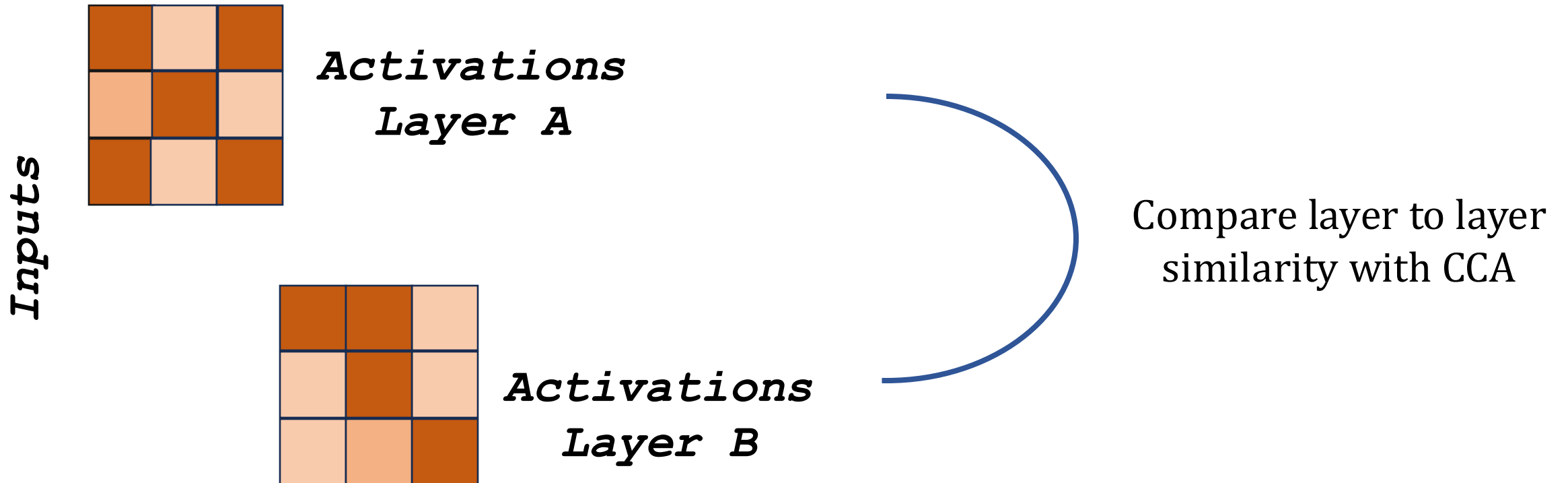
Surprisingly, one head attends a lot to modified chunk in funny sentence and only in this case

---

# SVCCA: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability

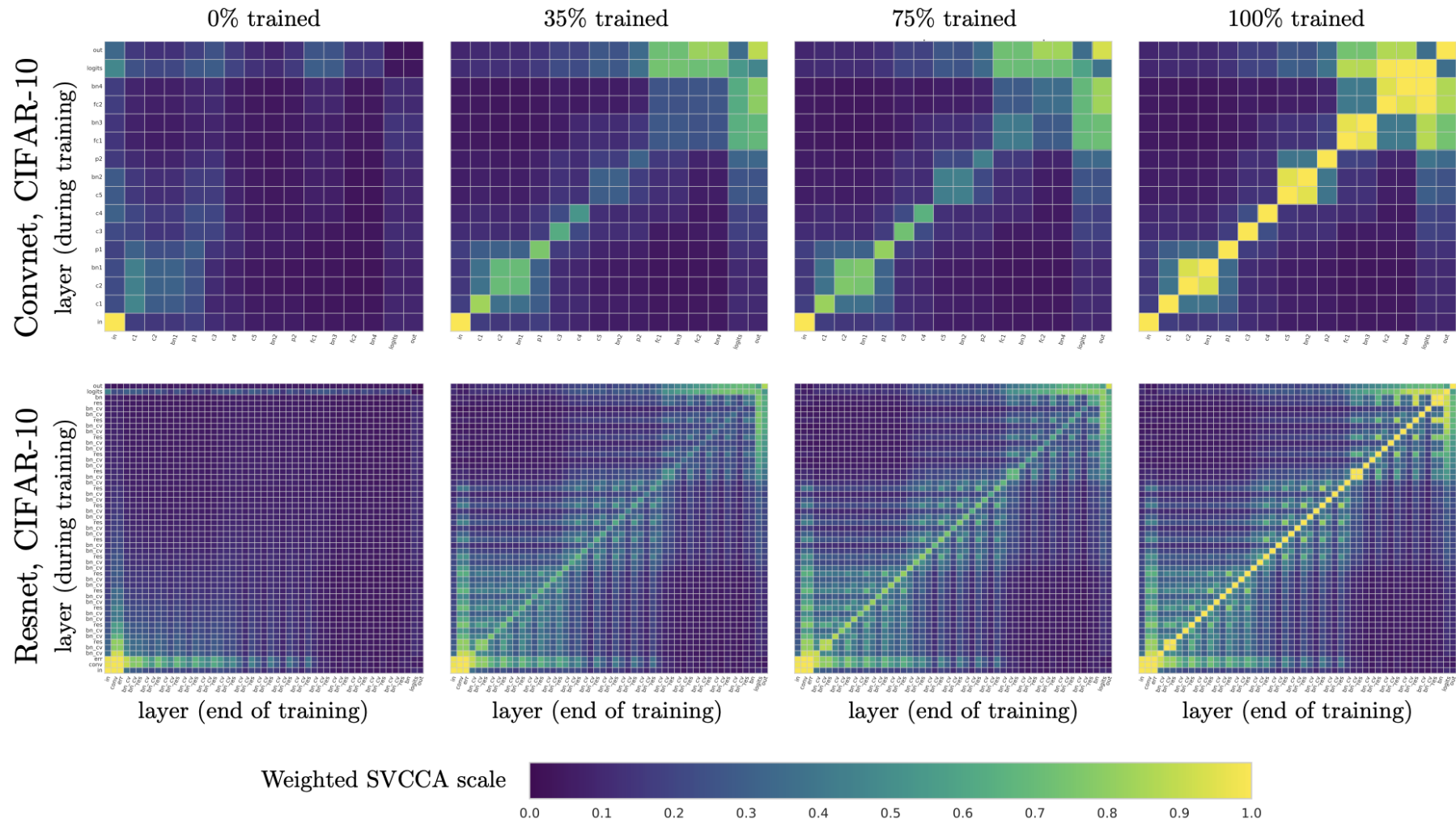
---

Maithra Raghu,<sup>1,2</sup> Justin Gilmer,<sup>1</sup> Jason Yosinski,<sup>3</sup> & Jascha Sohl-Dickstein<sup>1</sup>  
<sup>1</sup>Google Brain <sup>2</sup>Cornell University <sup>3</sup>Uber AI Labs

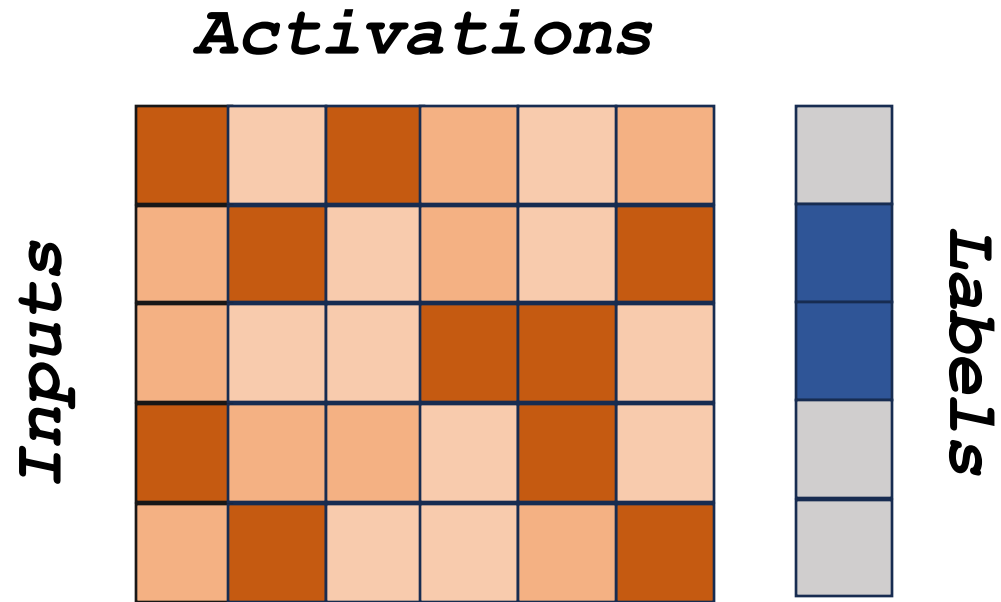


# SV-CCA

Redundant layers appearing during training → possibilities for pruning



# Probes



A model  $f$  that **reliably** predict some behavior labels from **activations**

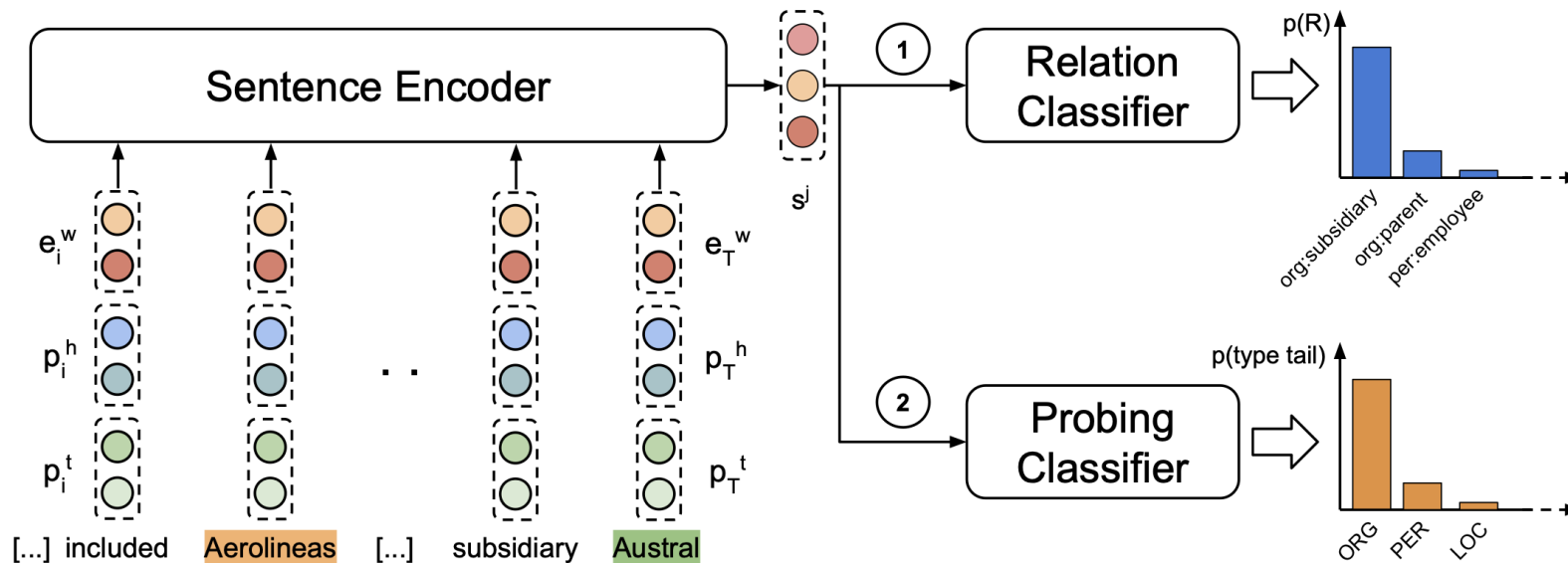
$$f\left(\begin{array}{|c|c|c|c|c|c|}\hline \text{orange} & \text{light orange} & \text{brown} & \text{orange} & \text{light orange} & \text{orange} \\ \hline\end{array}\right) = \begin{array}{|c|}\hline \text{dark blue} \\ \hline\end{array}$$



# Example of Probing: Linguistic Features

**Labels: linguistic features**

*Probing Linguistic Features of Sentence-Level Representations in Neural Relation Extraction*  
ACL 2020



Use the activations to predict linguistic features



# Problems with Probing



## **Probing the Probing Paradigm: Does Probing Accuracy Entail Task Relevance?**

**Abhilasha Ravichander<sup>1</sup>   Yonatan Belinkov<sup>2\*</sup>   Eduard Hovy<sup>1</sup>**

<sup>1</sup>Language Technologies Institute, Carnegie Mellon University

<sup>2</sup>Technion – Israel Institute of Technology

## **Probing Classifiers: Promises, Shortcomings, and Advances**

**Yonatan Belinkov\***

Technion – Israel Institute of Technology

## **An information theoretic view on selecting linguistic probes**

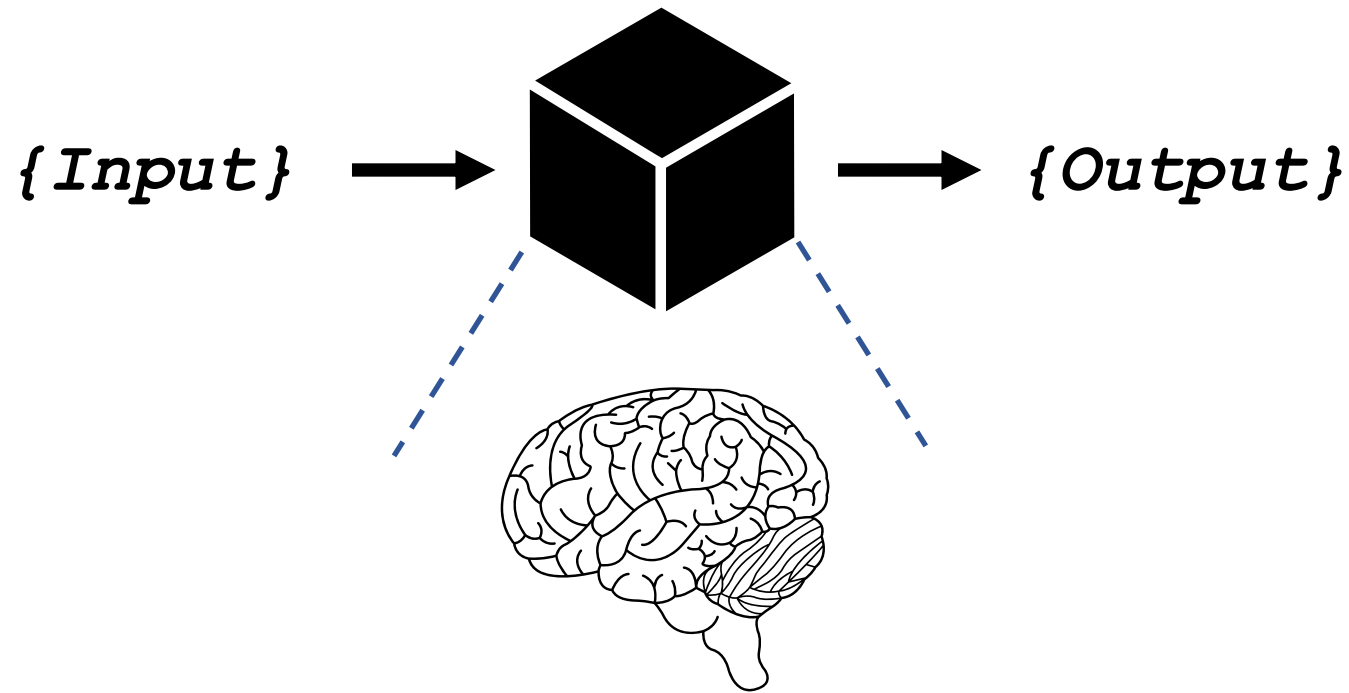
**Zining Zhu<sup>1,2</sup>, Frank Rudzicz<sup>3,4,1,2</sup>**

<sup>1</sup> University of Toronto, <sup>2</sup> Vector Institute, <sup>3</sup> Surgical Safety Technologies

<sup>4</sup> Li Ka Shing Knowledge Institute, St Michael's Hospital

(

# Neuroscience detour II

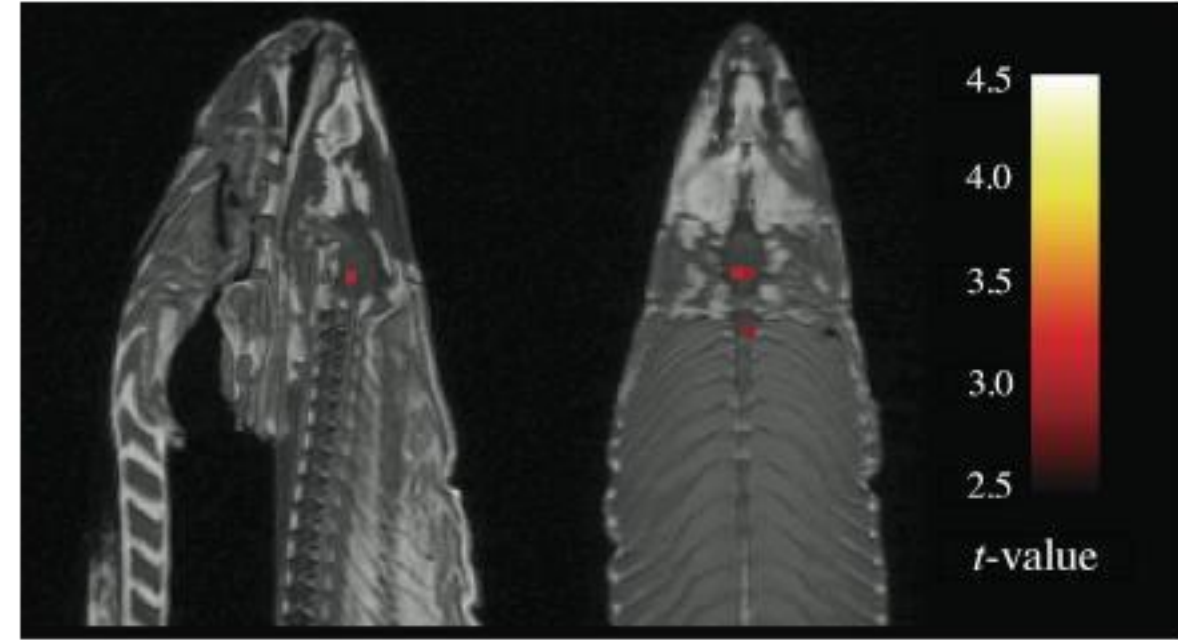
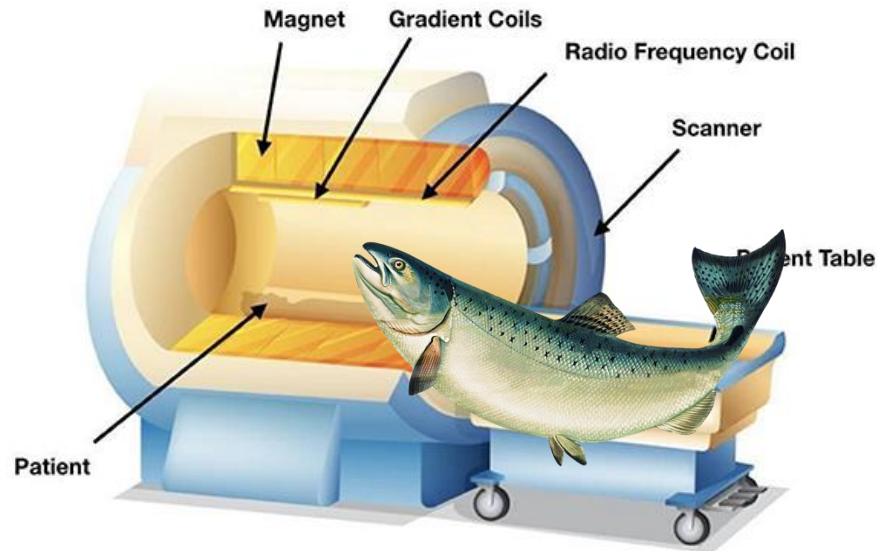


Representation, Pattern Information, and Brain Signatures:  
From Neurons to Neuroimaging

[Philip A. Kragel](#)<sup>1,2</sup> · [Leonie Koban](#)<sup>1</sup> · [Lisa Feldman Barrett](#)<sup>3,4,5</sup> · [Tor D. Wager](#)<sup>1</sup>  

Very common to do “Probing” on brain activations: *“mutivariate pattern analyses”, “brain signatures”, ...*

# The Dead Salmon



**Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon:  
An argument for multiple comparisons correction**

Craig M. Bennett<sup>1</sup>, Abigail A. Baird<sup>2</sup>, Michael B. Miller<sup>1</sup>, and George L. Wolford<sup>3</sup>

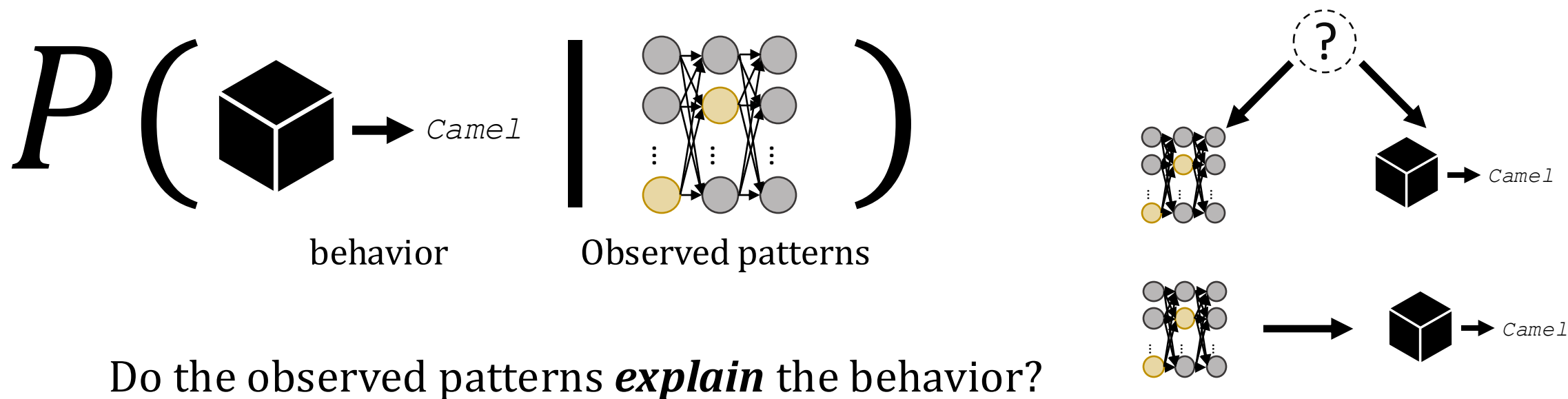
**The lure of misleading causal statements in functional connectivity research**

David Marc Anton Mehler, Konrad Paul Kording

)

# Neural correlates fall short

*Predicting is not understanding* – correlations are everywhere and do not generalize



**Predicting is not Understanding:**

Damien Teney<sup>1,3</sup>   Maxime Peyrard<sup>2</sup>   Ehsan Abbasnejad<sup>3</sup>

**Sparse Autoencoders Can Interpret Randomly Initialized Transformers**

Thomas Heap, Tim Lawson, Lucy Farnik, Laurence Aitchison

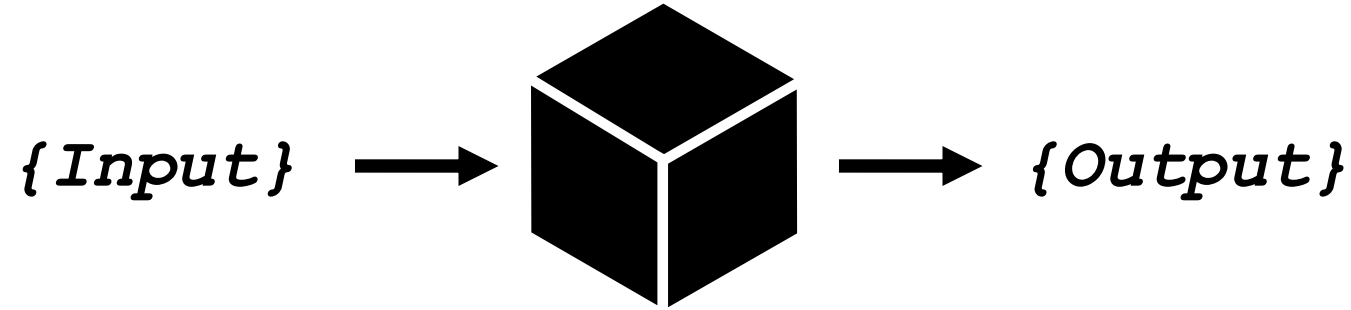
# **How to Explain?**

## **“Causal patterns”**

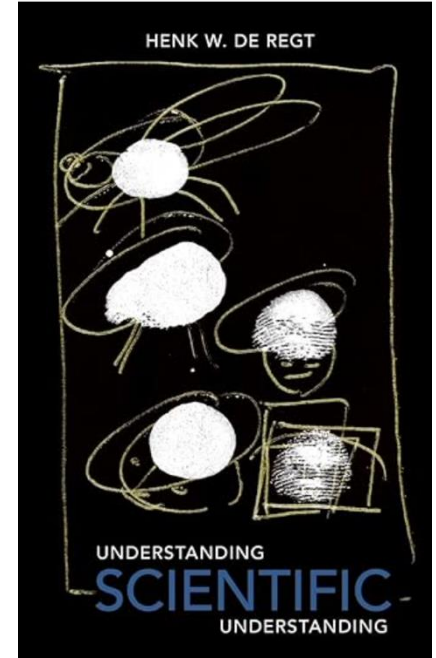




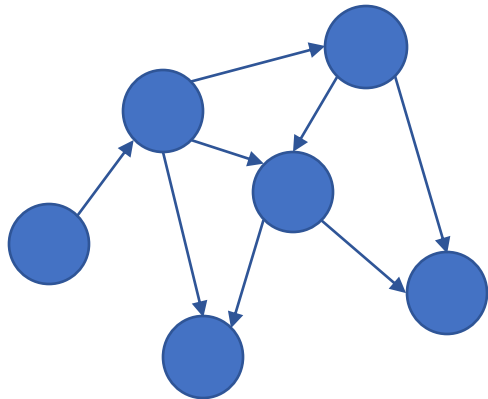
# I - Causal Models



*Why?*



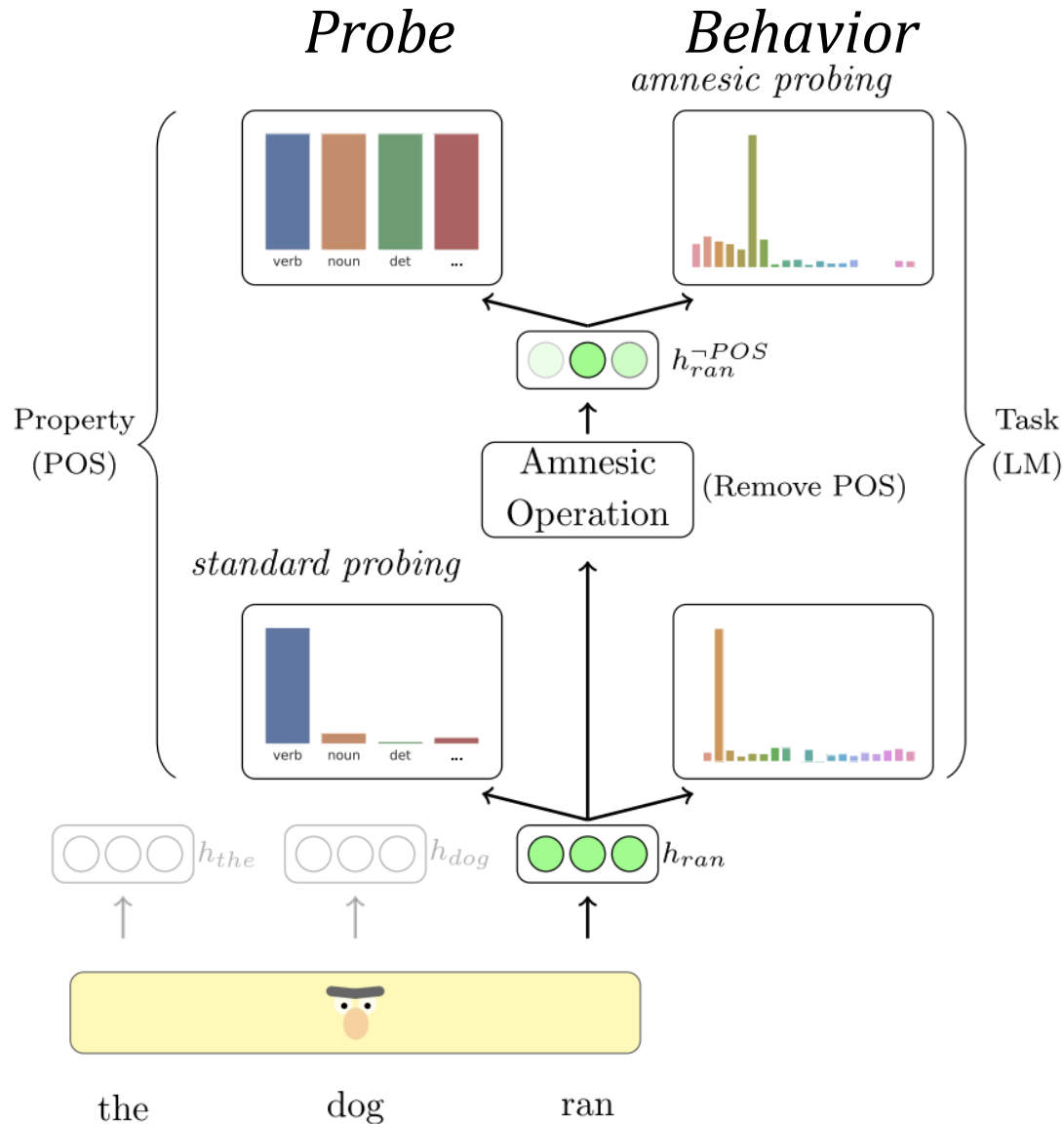
*Philosophers of Science* argue that **explanations** must be **causal analyses**



*"causes **explain** their effects »*

- Understanding: know the behavior in any scenario
- Control: know the impact of modifications on the system

# Amnesic probing: Validating with Interventions



## Amnesic Probing: Behavioral Explanation with Amnesic Counterfactuals

Yanai Elazar<sup>1,2</sup> Shauli Ravfogel<sup>1,2</sup> Alon Jacovi<sup>1</sup> Yoav Goldberg<sup>1,2</sup>

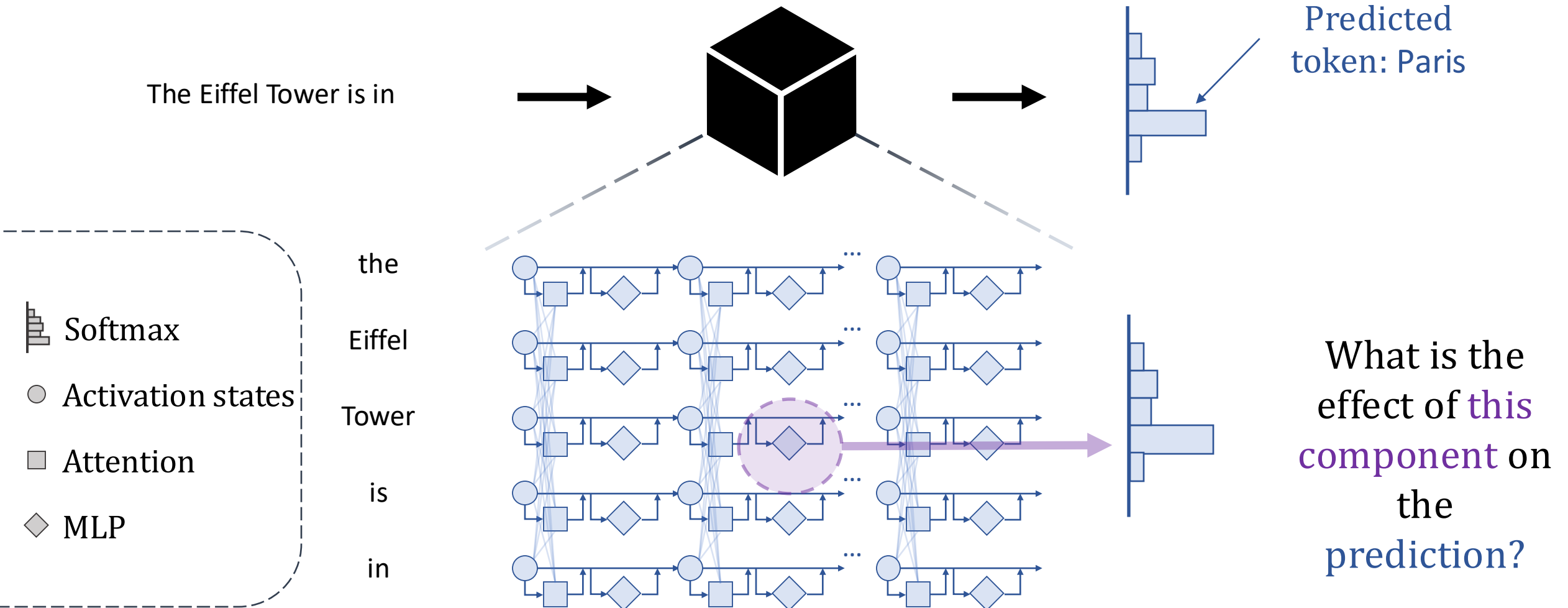
<sup>1</sup>Computer Science Department, Bar Ilan University

<sup>2</sup>Allen Institute for Artificial Intelligence

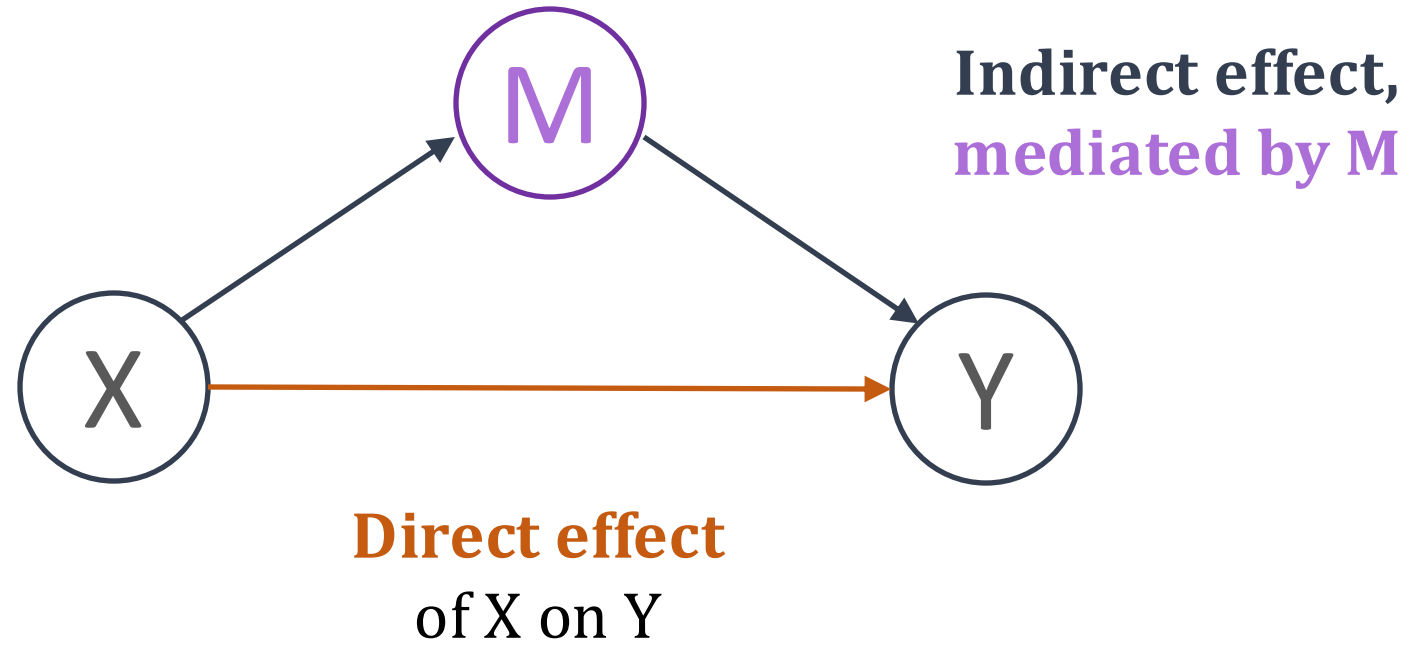
*Does modifying the activation to fool the probe, removes the behavior?*

# Causal mediation analysis

Goal: Understand the impact of *model components* on *model behavior*



# Causal Mediation Analysis



*How much of the effect of X on Y **is explained by** the path through M?*



i.e., understanding the mechanisms by which X acts on Y,  
disentangling the different paths of influences.

# Examples – Gender Bias

**Prompt  $u$ :** The nurse said that \_\_

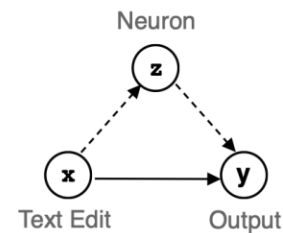
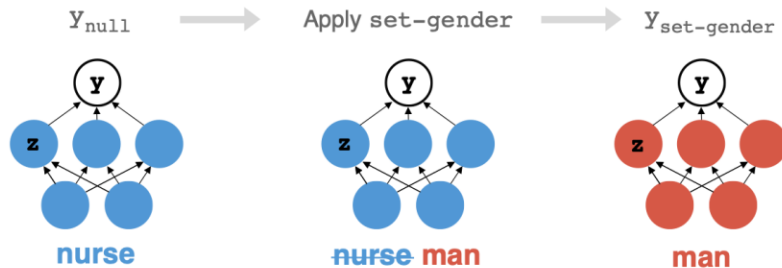
**Stereotypical candidate:** she

**Anti-stereotypical candidate:** he

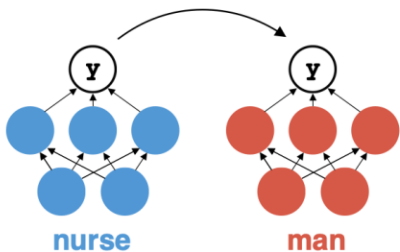
*Causal Mediation Analysis for Interpreting  
Neural NLP: The Case of Gender Bias*

NeurIPS 2020

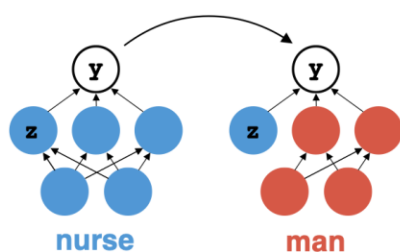
(a) Causal mechanism



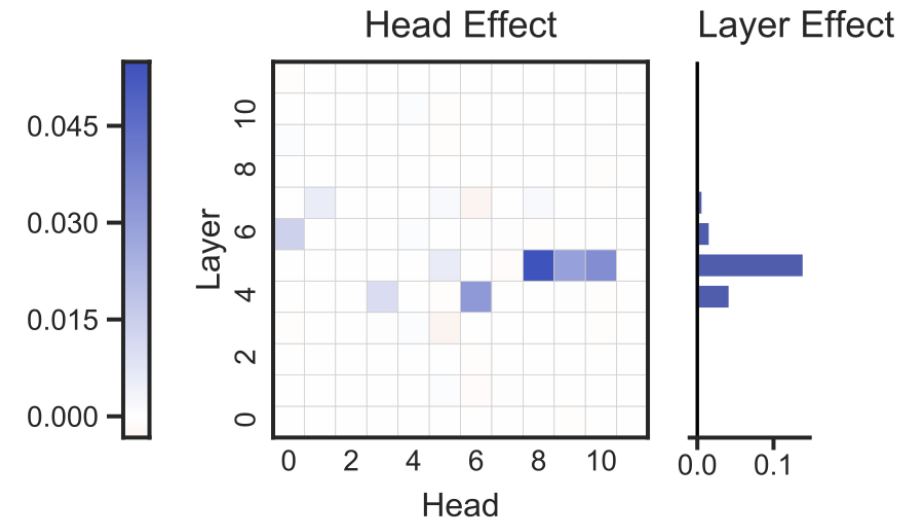
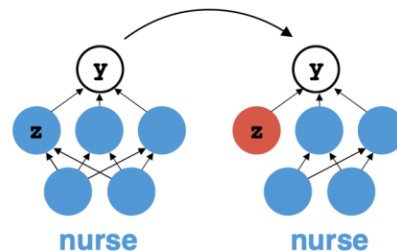
(b) Total Effect



(c) Direct Effect



(d) Indirect Effect



# Examples – factual recall

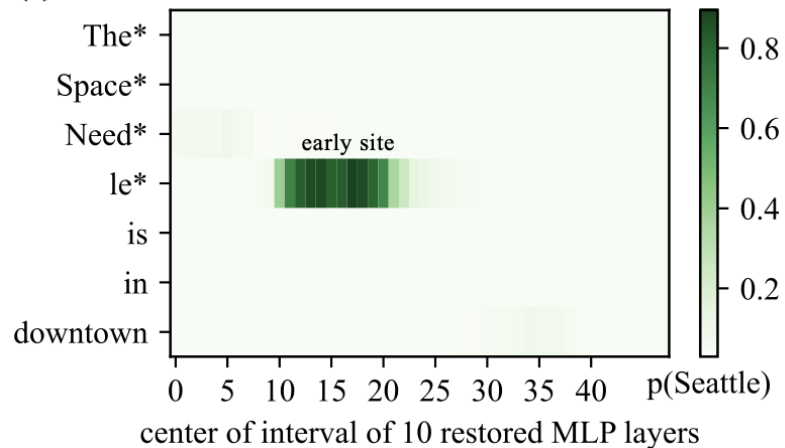
## Locating and Editing Factual Associations in GPT

Kevin Meng\* MIT CSAIL    David Bau\* Northeastern University    Alex Andonian MIT CSAIL    Yonatan Belinkov† Technion – IIT

## Transformer Feed-Forward Layers Are Key-Value Memories

Mor Geva<sup>1,2</sup>    Roei Schuster<sup>1,3</sup>    Jonathan Berant<sup>1,2</sup>    Omer Levy<sup>1</sup>  
<sup>1</sup>Blavatnik School of Computer Science, Tel-Aviv University  
<sup>2</sup>Allen Institute for Artificial Intelligence  
<sup>3</sup>Cornell Tech

(f) Impact of restoring MLP after corrupted input



## A Glitch in the Matrix? Locating and Detecting Language Model Grounding with Fakepedia

Giovanni Monea,<sup>◇</sup> Maxime Peyrard,<sup>♡</sup> Martin Josifoski,<sup>◇</sup> Vishrav Chaudhary,<sup>♣</sup>  
Jason Eisner,<sup>♣</sup> Emre Kıcıman,<sup>♣</sup> Hamid Palangi,<sup>♣</sup> Barun Patra,<sup>♣</sup> Robert West<sup>◇</sup>  
<sup>◇</sup>EPFL    <sup>♡</sup>Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG    <sup>♣</sup>Microsoft Corporation

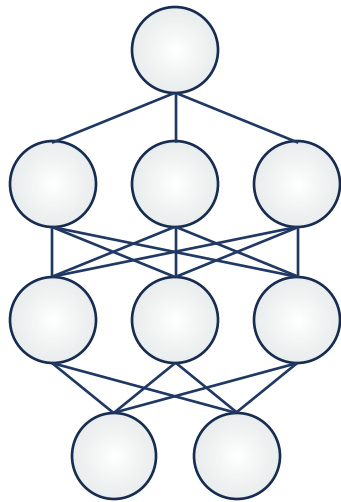
*Facts are localized in few MLPs that  
are **associative memories** for factual  
knowledge*

High causal effect on the prediction in early sites  
→ due to the activity of few MLPs

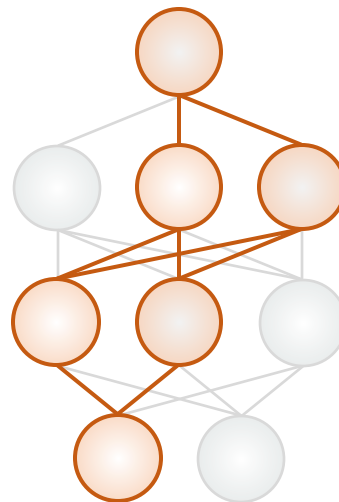
# Mechanistic Interpretability

Idea: Reverse-engineer trained neural networks to find **simple, human-interpretable, algorithms embedded in the computation**

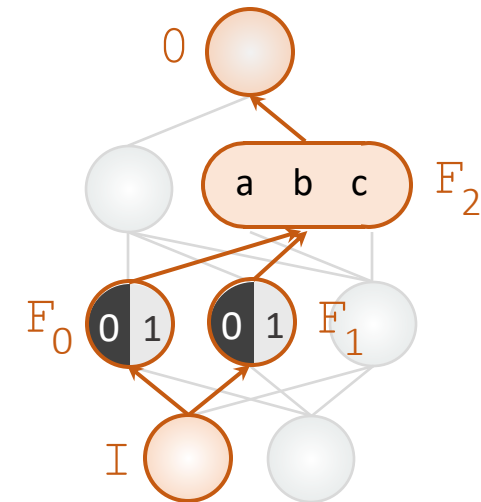
Base  
Neural Network



1. Circuit



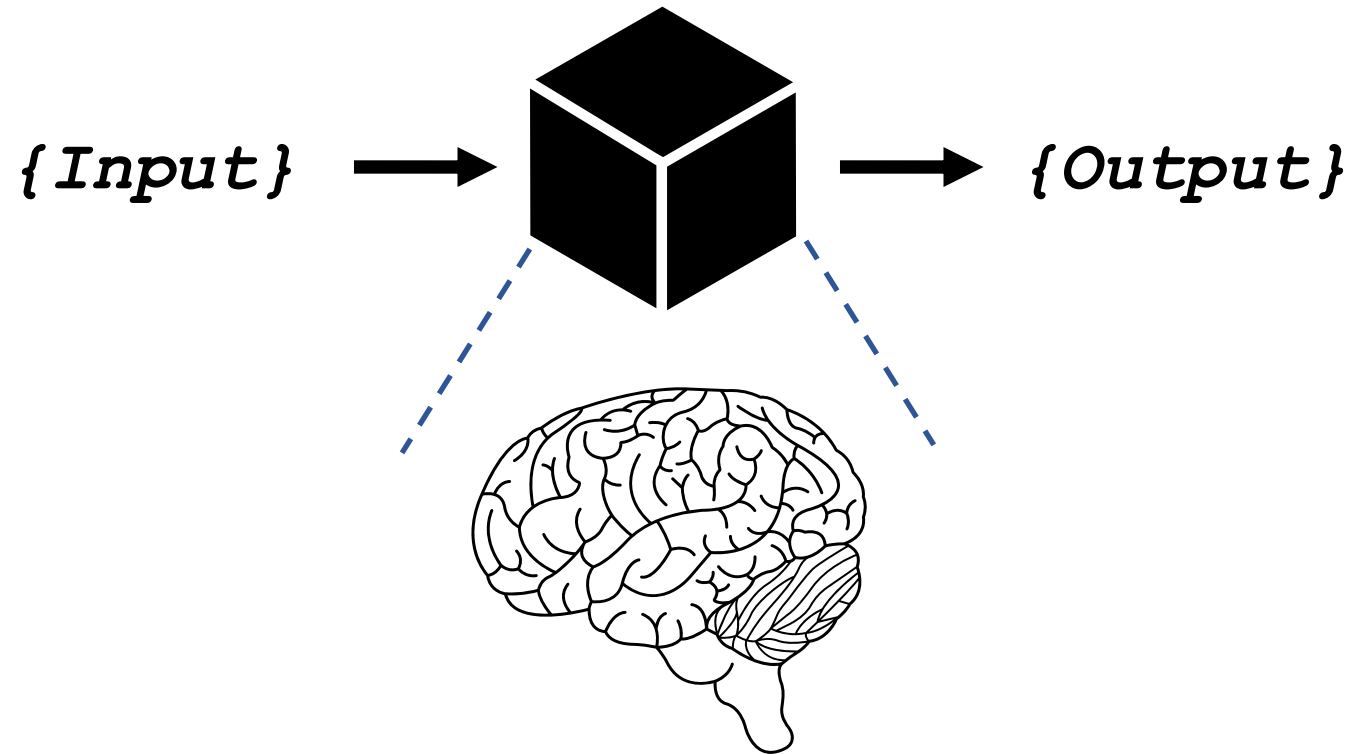
2. Interpret  
components







# Neuroscience detour III



How do we know that our explanations are correct?

How can we trust our analysis methods if we never test on examples of behavior / true explanations

# Neuroscience detour

RESEARCH ARTICLE

## Could a Neuroscientist Understand a Microprocessor?

Eric Jonas<sup>1\*</sup>, Konrad Paul Kording<sup>2,3</sup>

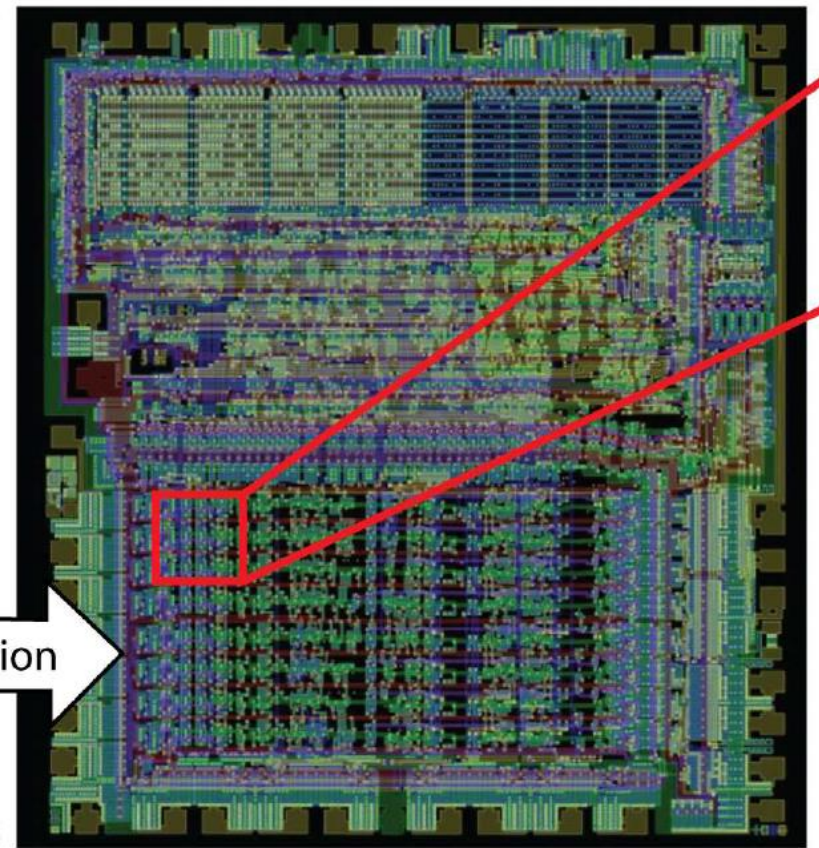


We know everything about the microprocessor,  
Let's treat it as if it was a brain (where transistor  $\approx$  neurons)

*Can analysis methods recover meaningful information about the microprocessor even with perfect observations and manipulation capabilities?*

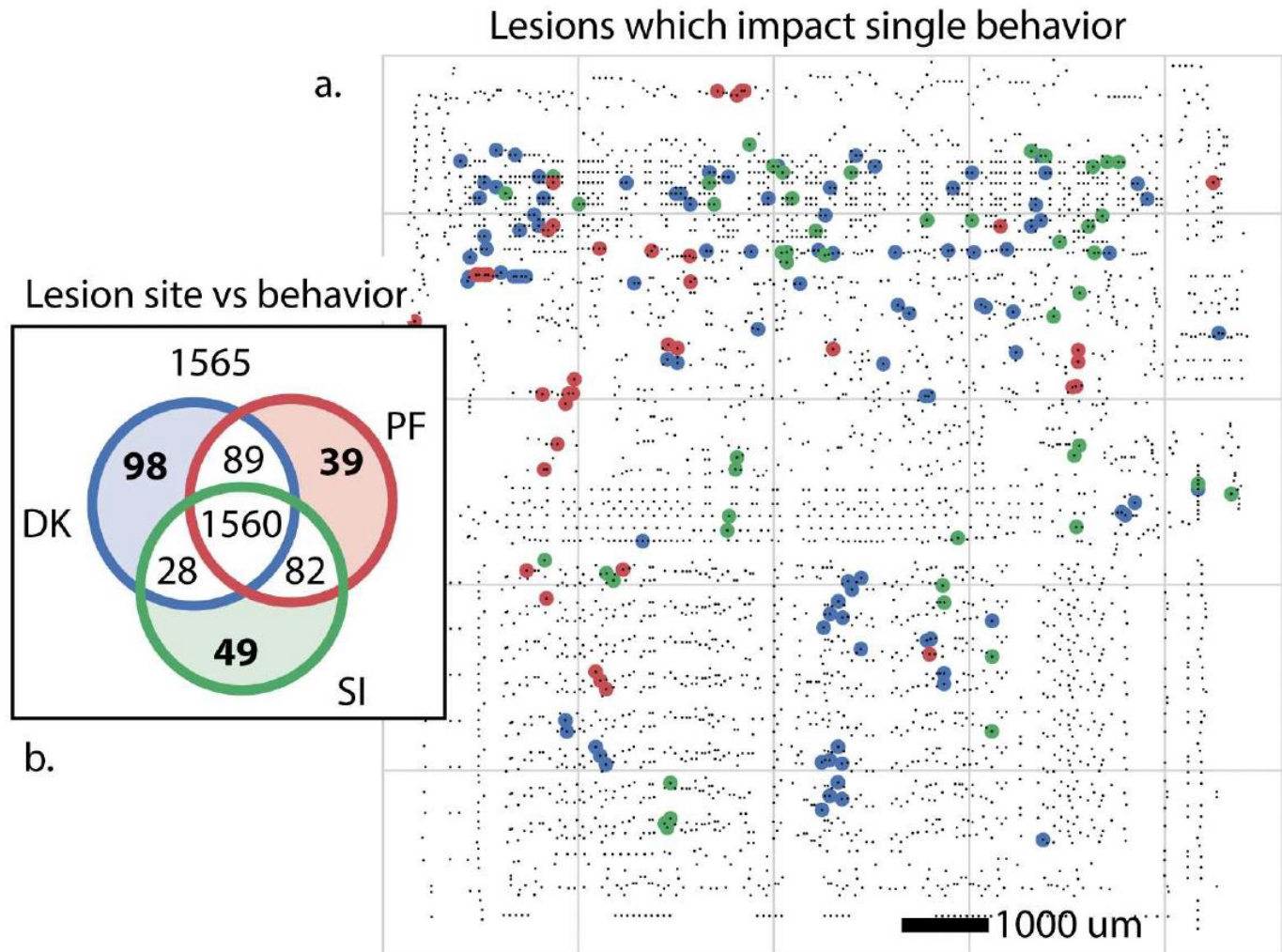
**NO**

# Neuroscience detour



ion

Even extensive intervention  
study gives no useful  
information!



**Fig 4. Lesioning every single transistor to identify function.** We identify transistors whose elimination disrupts behavior analogous to lethal alleles or lesioned brain areas. These are transistors whose elimination results in the processor failing to render the game. **(A)** Transistors which impact only one behavior, colored by behavior. **(B)** Breakdown of the impact of transistor lesion by behavioral state. The elimination of 1565 transistors have no impact, and 1560 inhibit all behaviors.

)

# ML is not the end of the story

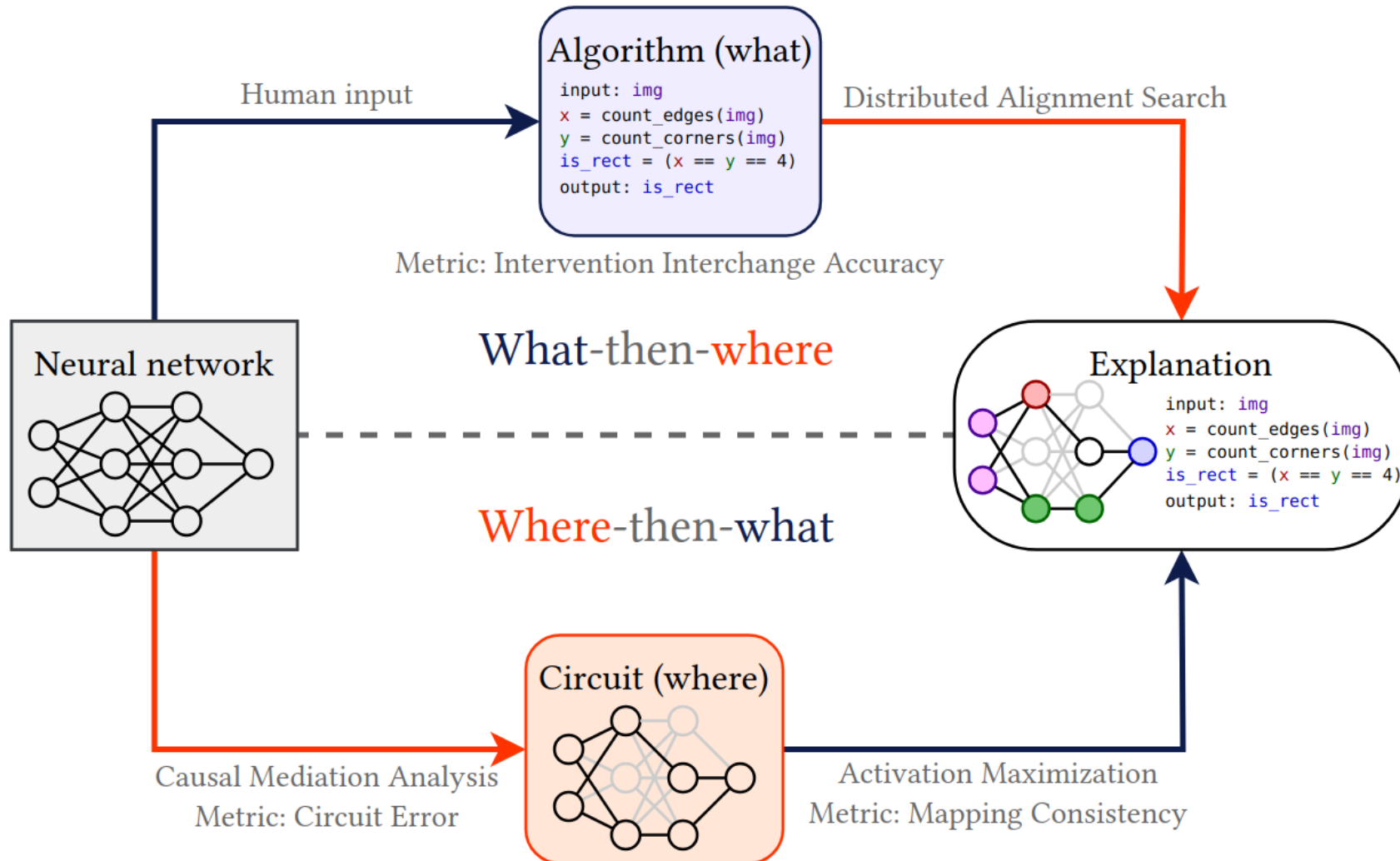
## Everything, Everywhere, All at Once: Is Mechanistic Interpretability Identifiable?

Maxime Méloux, Silviu Maniu, François Portet, Maxime Peyrard 

We can find *almost any explanation* if we look hard enough  
Even in randomly initialized neural networks

*“High-dimensional nonlinear systems may be hard to understand, but they are easy to find stories in.”* – Grace W. Lindsay → **un-identifiable**

# The Two Types of MI approaches

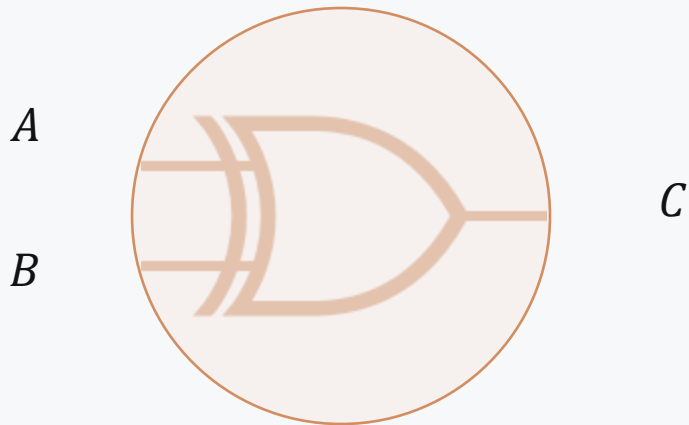




# Experimental Setup

## Base Neural Network

Train an MLP to implement XOR



$$A = 0|1 + \mathcal{N}(0, \varepsilon)$$

$$B = 0|1 + \mathcal{N}(0, \varepsilon)$$

$$C = \text{round}(A) \oplus \text{round}(B)$$

## Toy exercise in interpretability:

- **What** sequence of logic gates is implemented by the MLP?
- **Where** in the network is each gate implemented?

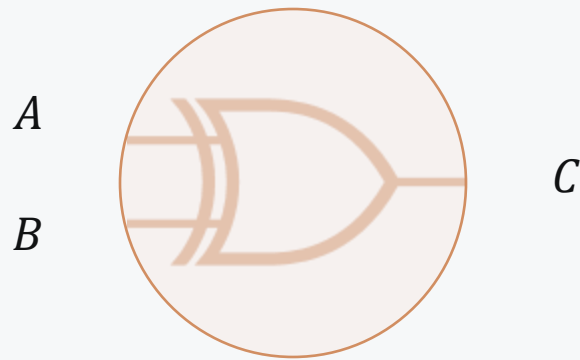
Enumerate **exhaustively** all candidate algorithms and mappings and test them with existing criteria.

# Where-then-What is not Identifiable

Do the current criteria used for selecting **circuits** and their **grounding** induce a unique solution? **NO**

## Base Neural Network

Train an MLP to implement XOR

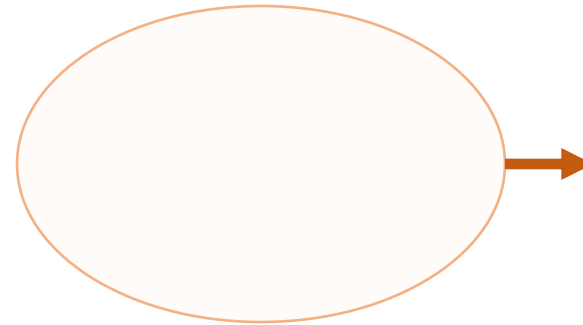


$$A = 0|1 + \mathcal{N}(0, \varepsilon)$$

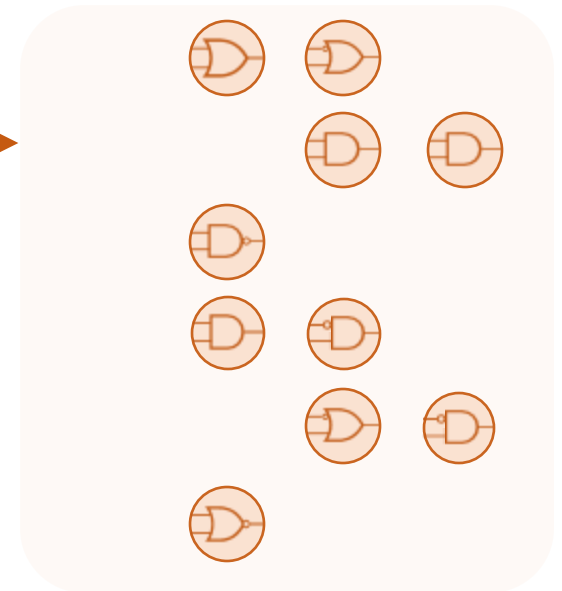
$$B = 0|1 + \mathcal{N}(0, \varepsilon)$$

$$C = \text{round}(A) \oplus \text{round}(B)$$

**85 unique circuits** with perfect accuracy



**25 unique explanations** with exact grounding for one circuit



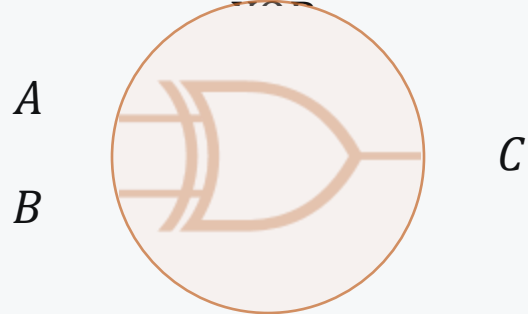


# What-then-where is not Identifiable

Do the current criteria used for assessing **causal alignment** of an **explanatory algorithm** guarantee a unique solution? **NO**

## Base Neural Network

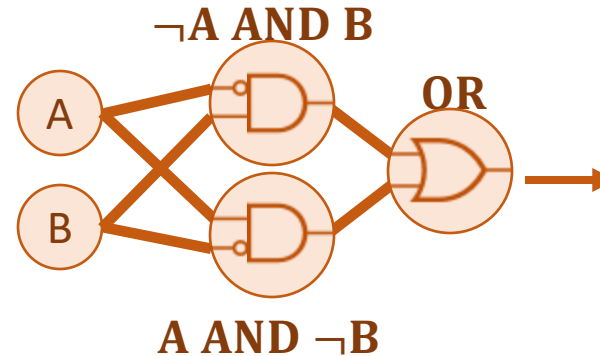
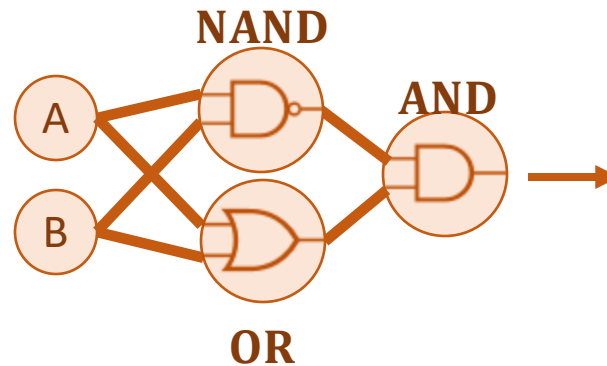
Train an MLP to implement



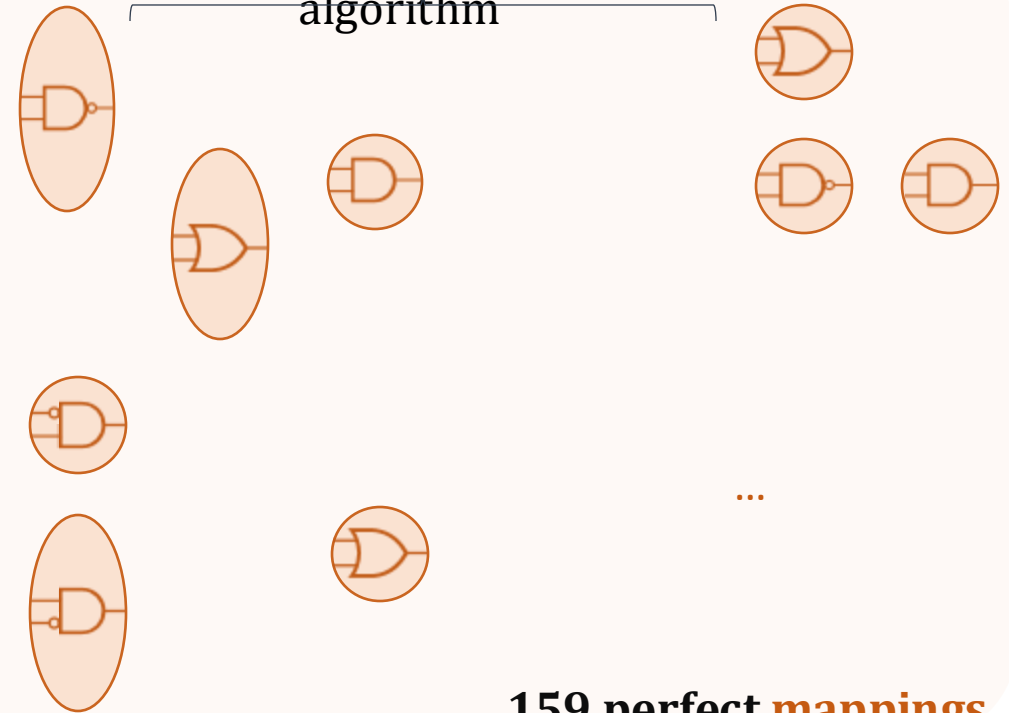
$$A = 0|1 + \mathcal{N}(0, \varepsilon)$$

$$B = 0|1 + \mathcal{N}(0, \varepsilon)$$

$$C = \text{round}(A) \oplus \text{round}(B)$$

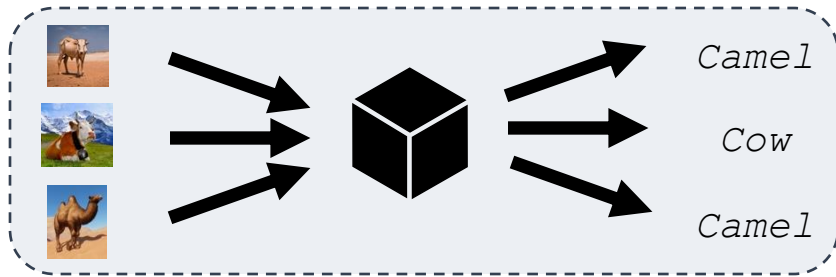


Example of 2 **perfect mappings** for one algorithm

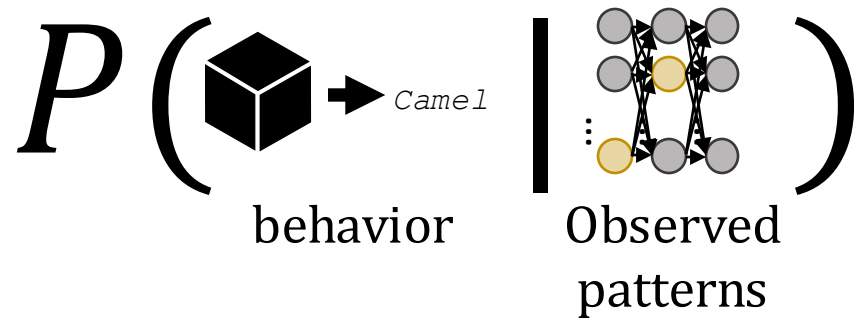


**159 perfect mappings**  
(IIA=1)

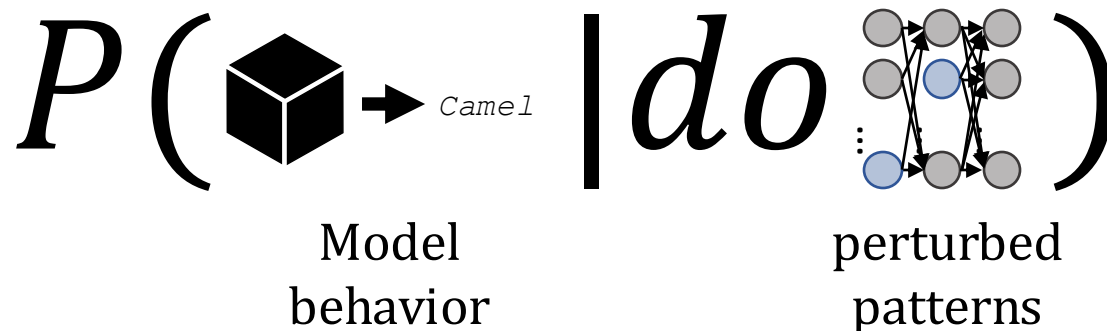
# Identifiability Issues → Generalization Issues



**Behavior:** Many explanations compatible with observed behavior. *Which one matches the computation? (Which one generalizes?)*



**Computational Correlate:** Many causal mechanisms compatible with observed correlations. *Which one generalizes?*

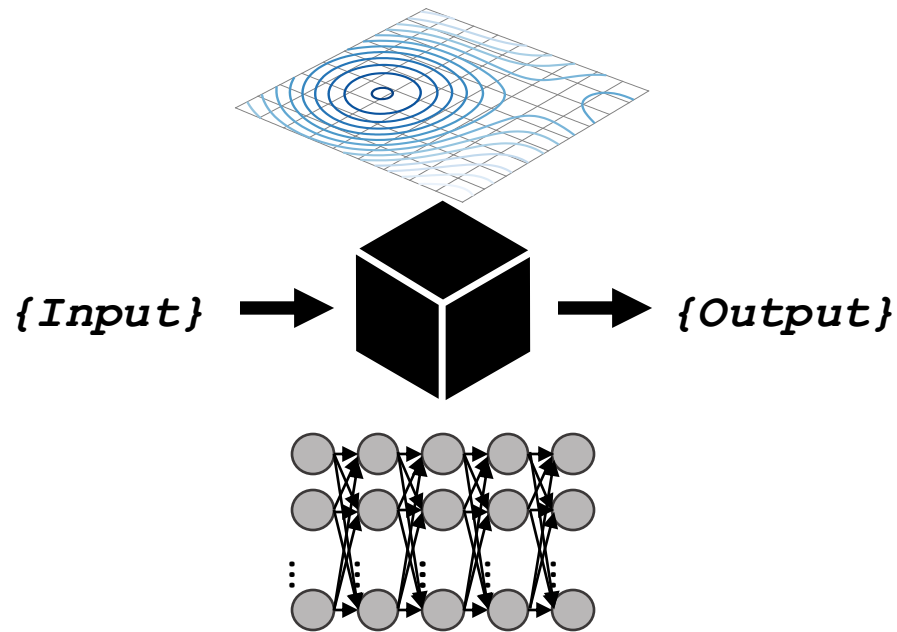


**Computational Causal Mechanisms:** Many causally aligned explanations!! *(Which one generalizes?)*

**What to do?**

# Instrumentalism: Statistical (Causal) Inference on computational data

**Distribution** over computational traces



Target property of interest  
*(estimand)*

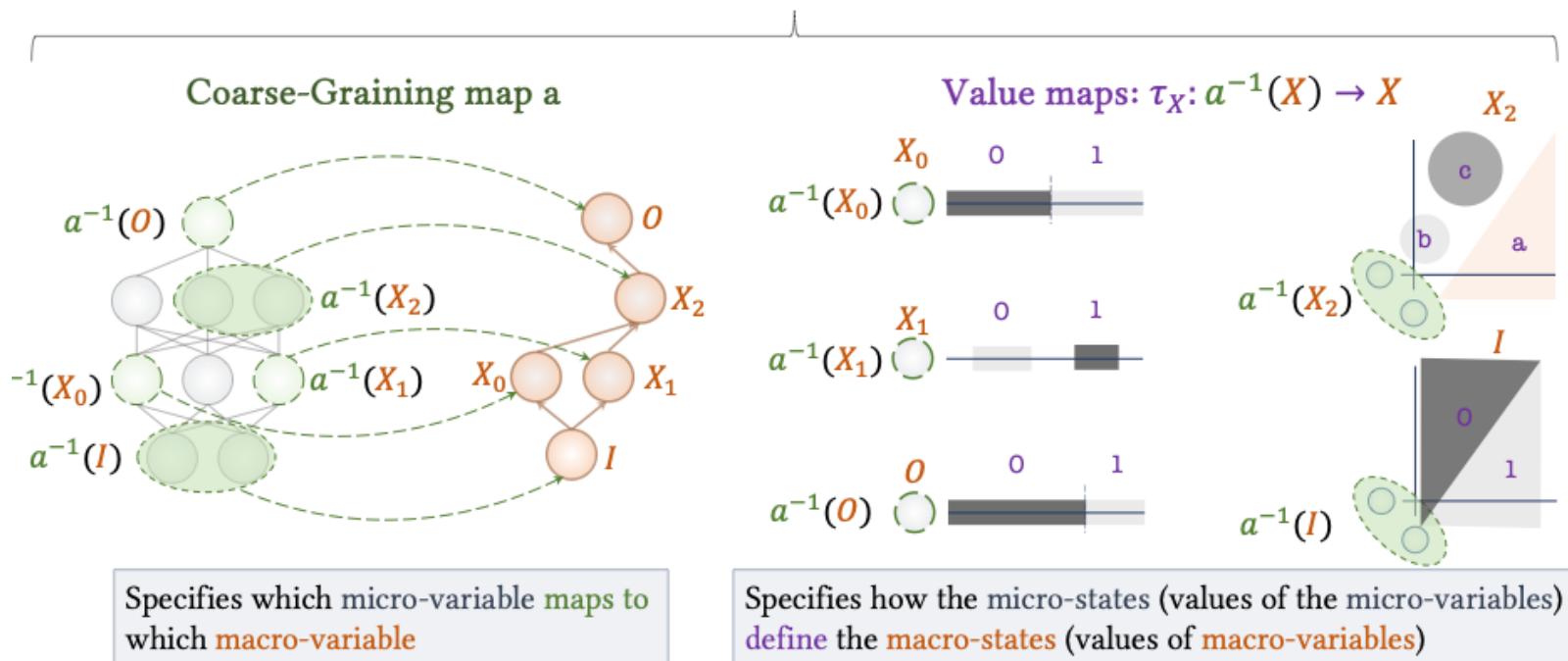
Estimator

Usual questions:

- **Estimand properties** (e.g., identifiability)
- **Estimators' properties:**  
Bias, variance, consistency, ...
- **Distributional properties:**  
Generalization, uncertainty

# Computational Summarization aka Causal Abstraction

(Constructive) Abstraction map  $\tau$



## Abstracting Causal Models

Sander Beckers

Joseph Y. Halpern

### Causal Abstraction:

A Theoretical Foundation for Mechanistic Interpretability

Atticus Geiger<sup>\*◇</sup>, Duligur Ibeling<sup>♣</sup>, Amir Zur<sup>◇</sup>, Maheep Chaudhary<sup>◇</sup>,  
Sonakshi Chauhan<sup>◇</sup>, Jing Huang<sup>♣</sup>, Aryaman Arora<sup>♣</sup>, Zhengxuan Wu<sup>♣</sup>,  
Noah Goodman<sup>♣</sup>, Christopher Potts<sup>♣</sup>, Thomas Icard<sup>\*♣</sup>

The **high-level model**  $\mathcal{A}$  is a **causal abstraction** of the low-level implementation  $\mathcal{L}$  if the variables in  $\mathcal{A}$  play the same causal role as their associated low-level variables.

# Collaborators: Thank you!



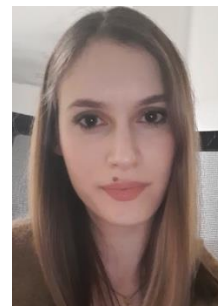
Damien Teney



Giovanni  
Monea



Robert West



Marija Sakota



Martin Josifoski



Emre Kiciman



Wei Zhao



Jason Eisner



Debjit Paul



Saibo Geng



Kristina Gligoric



Maxime  
Méloux



Francois  
Portet



Fei liu

# Thank you! Questions?

Contact: [maxime.peyrard@univ-grenoble-alpes.fr](mailto:maxime.peyrard@univ-grenoble-alpes.fr)