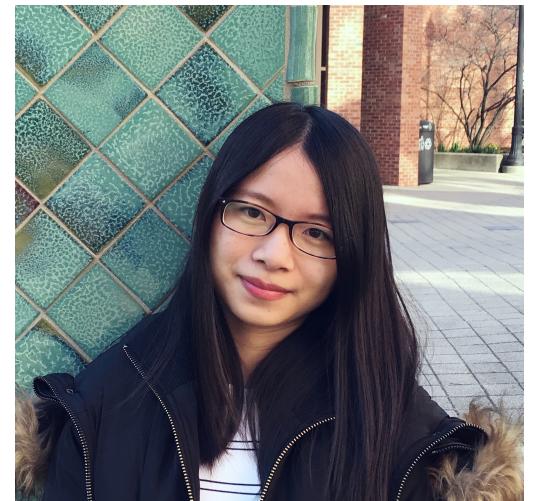


# Tutorial on Mechanistic Interpretability of Language Models

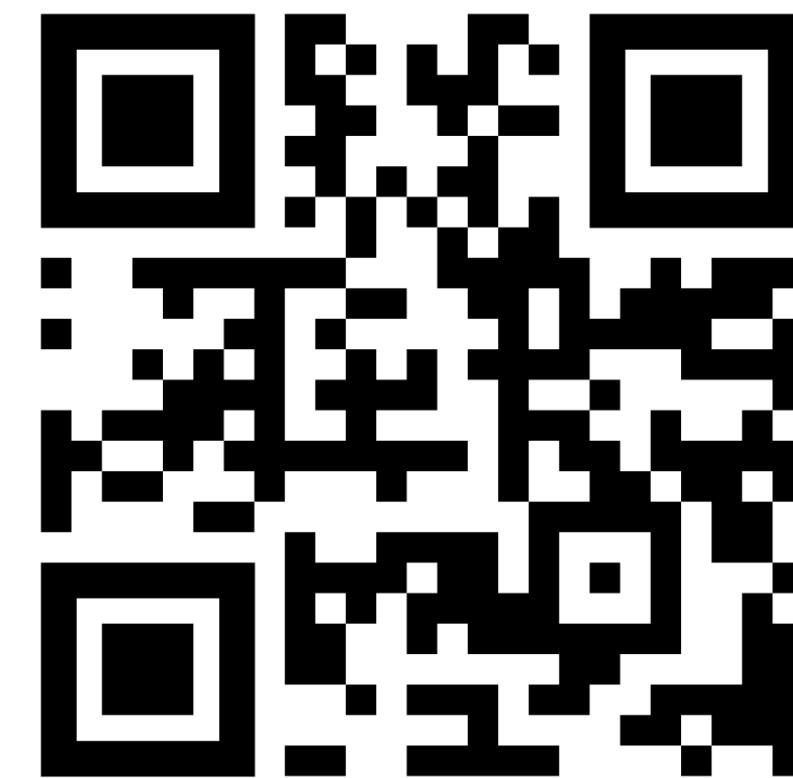


**Ziyu Yao**  
Asst. Prof.

Department of Computer Science  
George Mason University, USA



**Daking Rai**  
PhD Student



Post-event Q&A to:  
[{ziyuyao, drai2}@gmu.edu](mailto:{ziyuyao, drai2}@gmu.edu)

<https://ziyu-yao-nlp-lab.github.io/ICML25-MI-Tutorial.github.io/>

Tutorial Website with All the Materials and Recordings

# What comes to your mind when you think about **language models**?



...

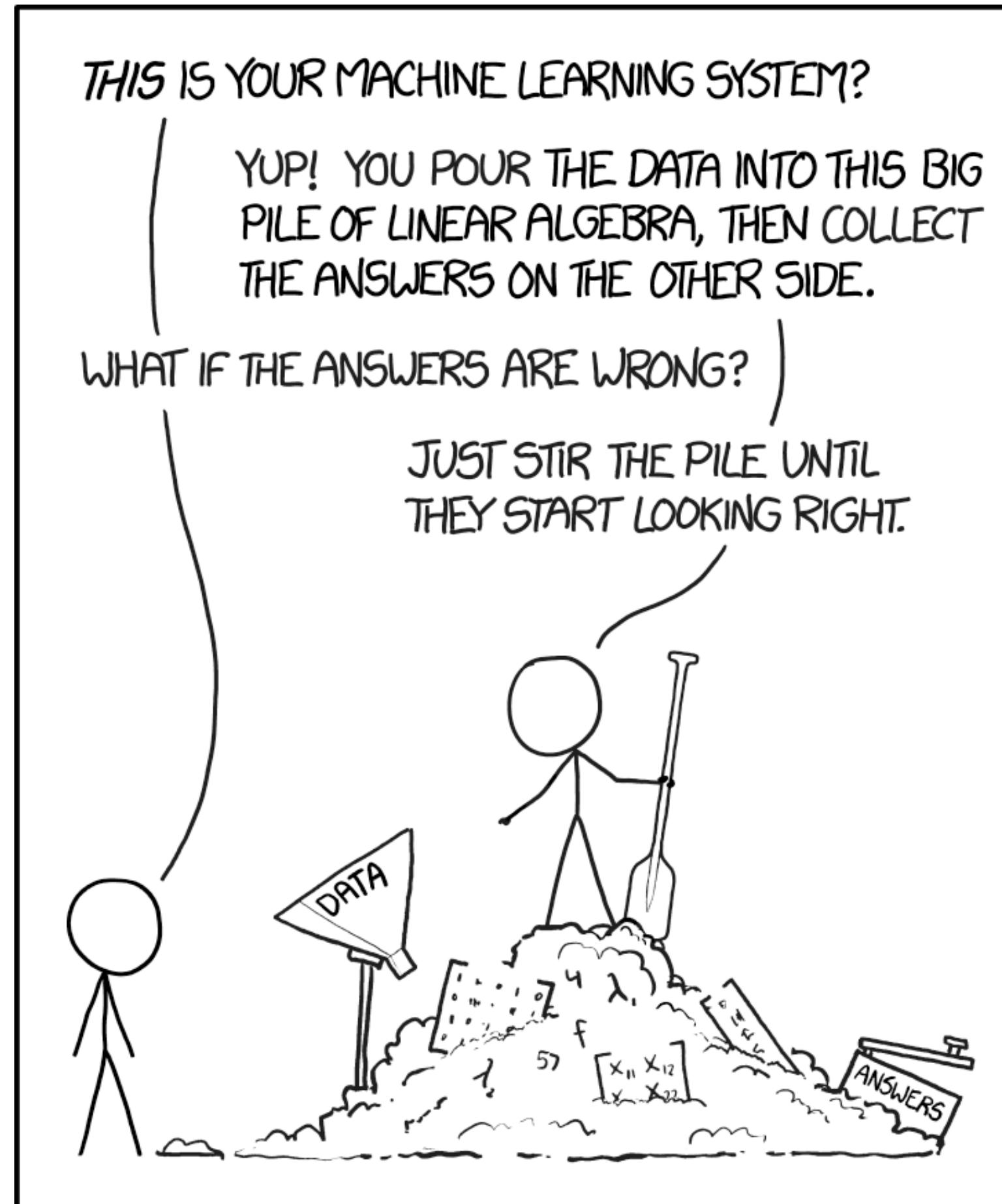
These?

- A **language model (LM)** is a probabilistic model that learns the likelihood of next-token prediction in a language sequence.

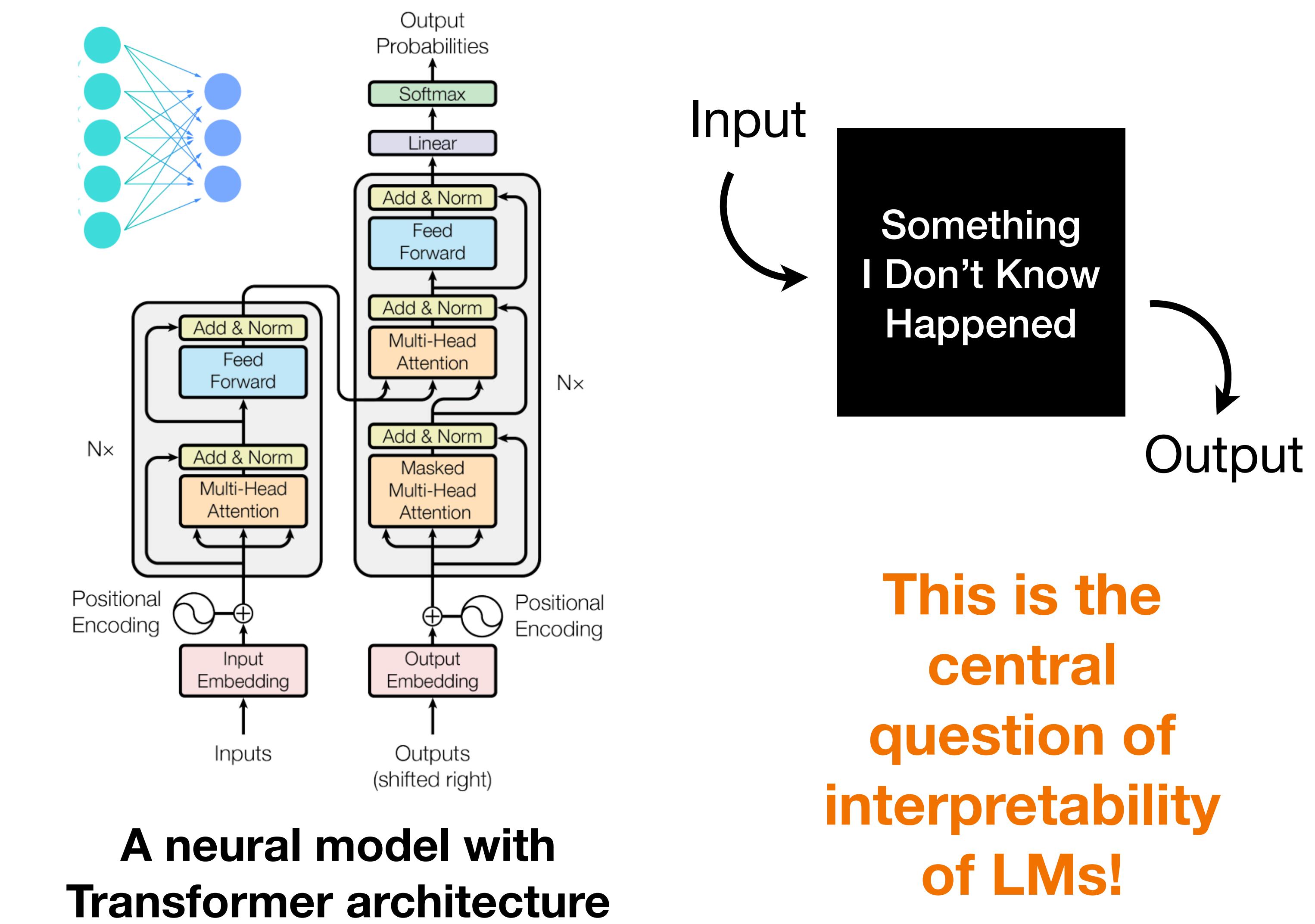
$$P(x_n | x_1, x_2, \dots, x_{n-1}), x_n \in \mathcal{V} \quad (\mathcal{V}: \text{the vocabulary})$$

- This tutorial: LMs implemented by **the transformer architecture** (Vaswani+17)

# What comes to your mind when you think about how LM works?



As a Machine Learning model...



# Why do we want to understand how LM works?

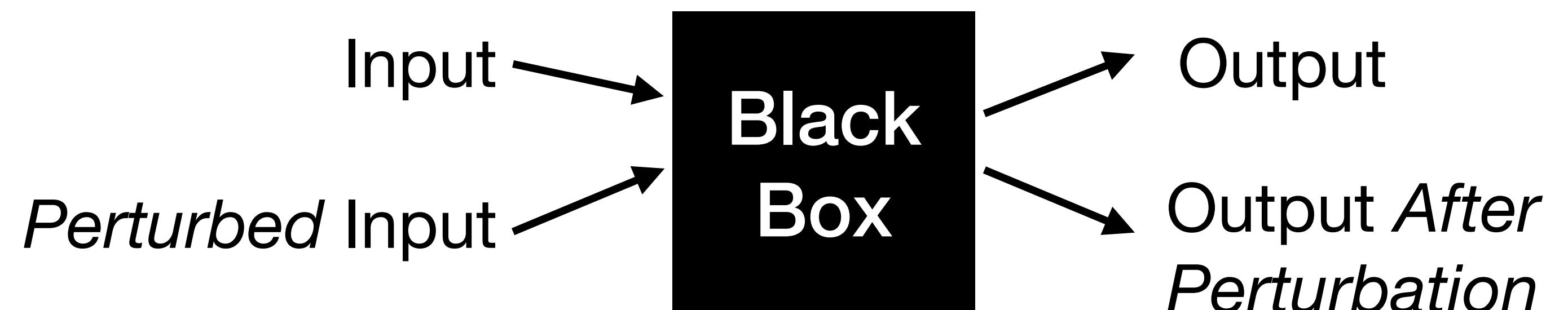
- Scientific understanding
- Practical reasons
  - AI safety and alignment
  - Human trust in AI
  - Enhancing model capability
  - Guiding model design, training, inference



I'm just curious!

# This tutorial: Mechanistic Interpretability (MI) of LMs

- Behavioral Interpretability (Ribeiro+16; Lundberg&Lee+17)



Let's see how the model reacts to changes in the input...

- Mechanistic Interpretability (Olah+20; Elhage+21)

- Understand the internals of LMs, how they function individually, and how they are connected to enable LM behaviors and capabilities
- Intuitively more insights than looking at only behaviors!



Let's open the black box and look inside!

[1] Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "" Why should i trust you?" Explaining the predictions of any classifier." SIGKDD 2016.

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

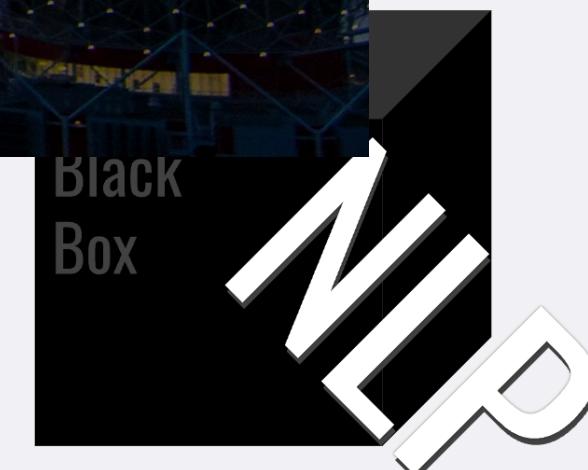
[3] Olah, Chris, et al. "Zoom in: An introduction to circuits." Distill 5.3 (2020)

[4] Elhage, Nelson, et al. "A mathematical framework for transformer circuits." Transformer Circuits Thread 1.1 (2021): 12.

# The thriving community, with doubts



**Stay Tuned!**  
<https://xllm-reasoning-planning-workshop.github.io/>



BlackboxNLP 2025

The Eight Workshop on Analyzing and Interpreting Neural Networks for NLP  
Co-located with EMNLP 2025 in Suzhou, China on November 10th, 2025

All > Technology & Research

## The Misguided Quest for Mechanistic AI Interpretability

Despite years of effort, mechanistic interpretability has failed to provide insight into AI behavior — the result of a flawed foundational assumption.

AI Frontiers

Dan Hendrycks and Laura Hiscott — May 15, 2025



**Mechanistic?**  
Naomi Saphra, Sarah Wiegreffe  
Meanwhile,  
argument about  
its scope,  
definitions, etc.

# Scope of this tutorial

Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao.  
"A practical review of mechanistic interpretability for transformer-based language models." Preprint 2025 (Version 2).

- We are not here to tell you if MI is worth pursuing or not, but we equip you with the necessary knowledge background.
  - **Part 1: Fundamental Objects of Study (7 min)**
    - What are the fundamental objects of studies of MI? What are the scopes?
  - **Part 2: MI Practices and Techniques (75 min) *10min Q&A/Break***
    - What are the common practices of applying MI to LMs? What techniques should you grasp?
  - **Part 3: Findings and Applications (25 min)**
    - What people have found out by applying MI to LMs? How has MI been helpful in applications?
  - **Part 4: Challenges and Future Work (10 min)**

# Refresher: Transformer Architecture

- An input token  $x_i$  is processed layer by layer (assuming  $L$  layers)
- $h_i^0 \in \mathbb{R}^d$ : the embedding vector looked up from  $W_E \in \mathbb{R}^{|\mathcal{V}| \times d}$
- $h_i^l \in \mathbb{R}^d$ : representation (or “activation”) at layer  $l$

$$h_i^l = h_i^{l-1} + a_i^l + f_i^l$$

## Multi-head Attention (MHA)

$$a_i^l = \text{Concat}(a_{i,1}^l, \dots, a_{i,H}^l) W_O^l$$

( $H$ : number of heads)

## Feed-Forward (FF)

$$f_i^l = f((h_i^{l-1} + a_i^l) \cdot W_k^l) \cdot W_v^l$$

- Finally, calculate the logit distribution for the next token with the unembedding matrix  $W_U \in \mathbb{R}^{|\mathcal{V}| \times d}$ :  $\text{Logit} = h_i^L \cdot W_U^\top$

