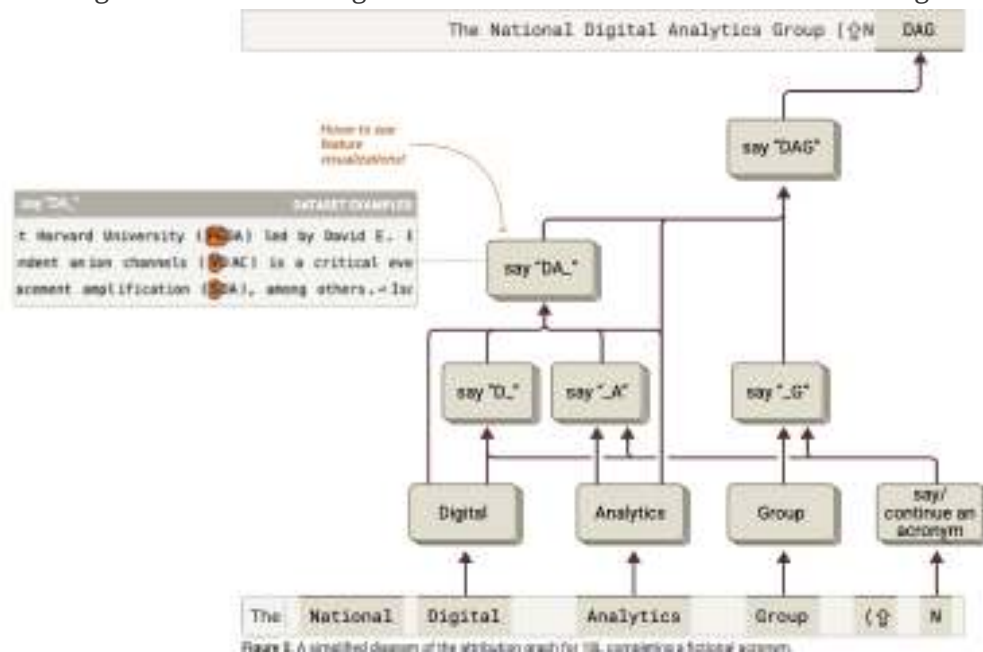# AI/ML Model Interpretability
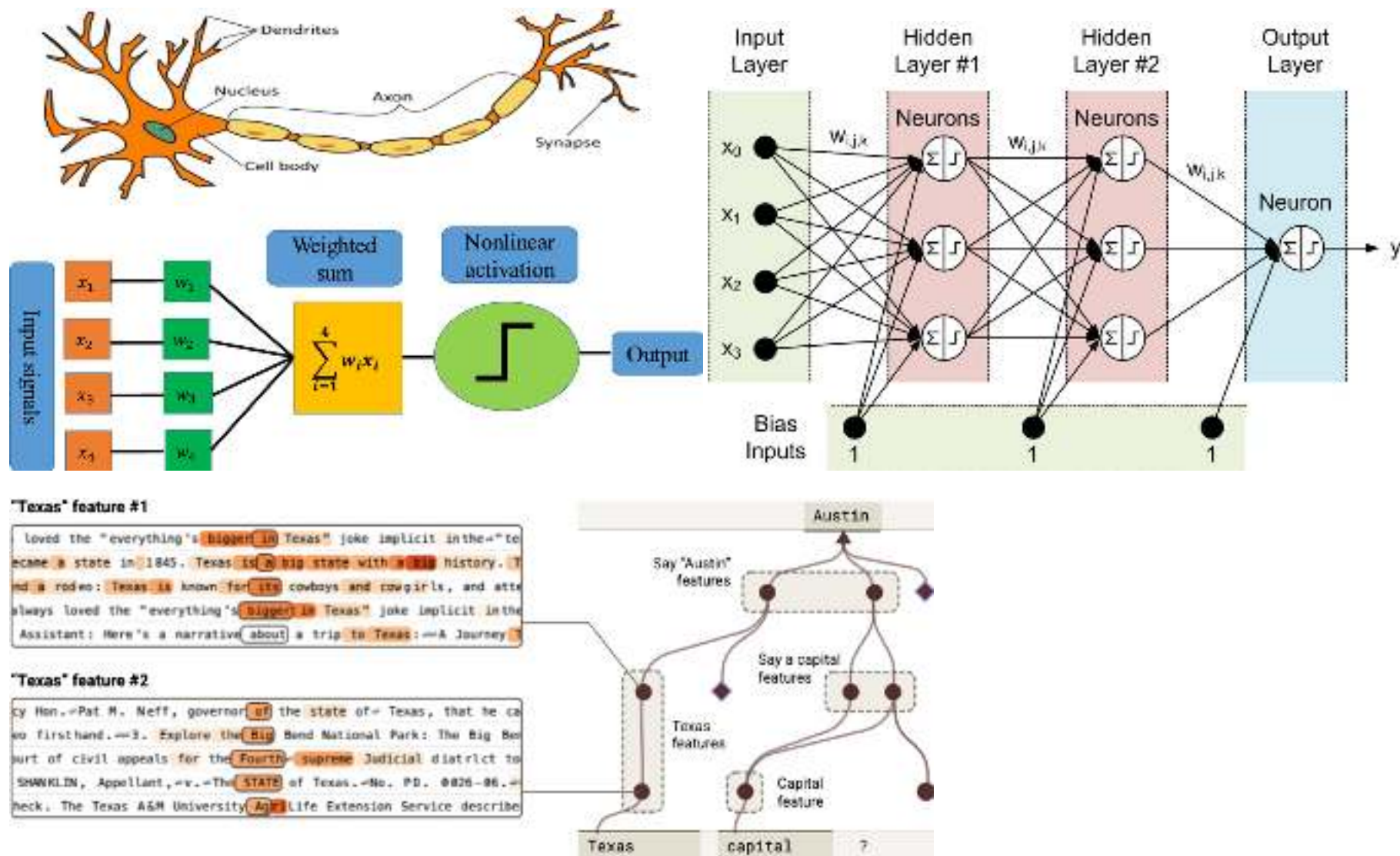
## Understanding AI Model Internals

- Goal: Make AI models transparent and interpretable through systematic analysis

- Challenge: Understanding AI models is similar to biological research - complex systems requiring sophisticated tools



Figure 1: A simplified diagram of the attribution graph for 10L completing a fictional acronym.

- Predictions about unexpected AI outputs

- "Microscopes" for AI model internals

*Lindsey, et al., "On the Biology of a Large Language Model", Transformer Circuits, 2025.*
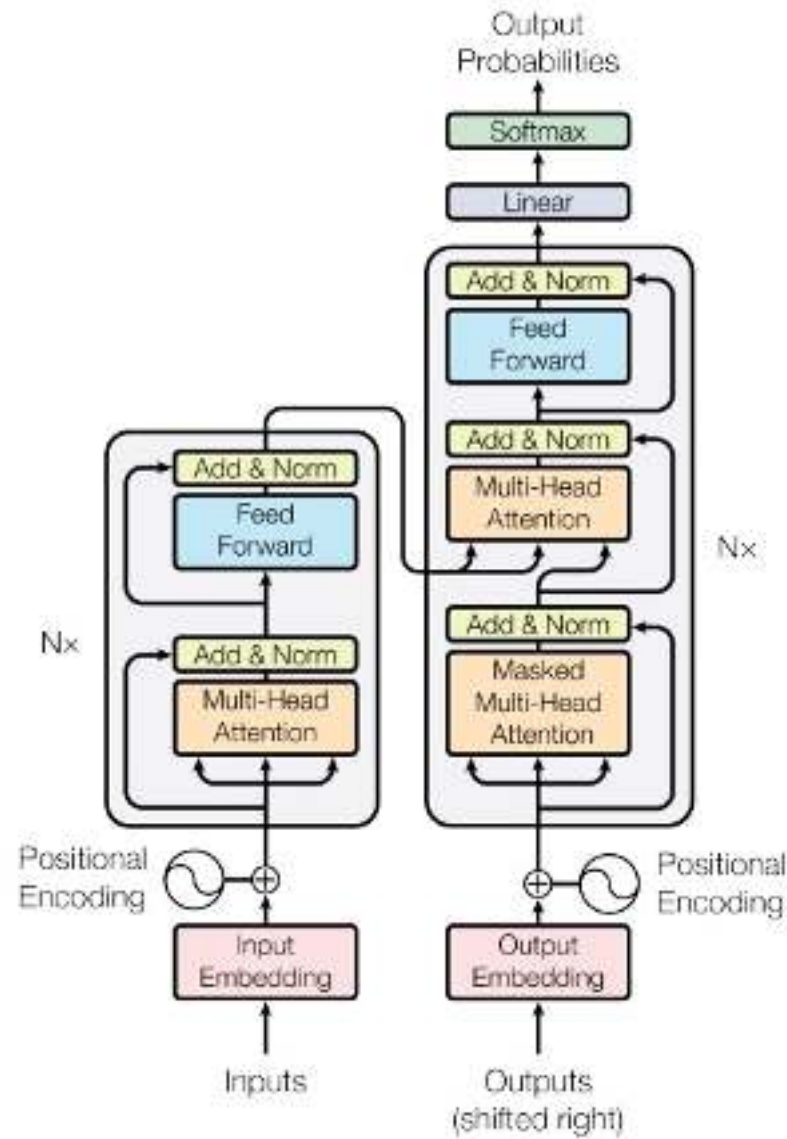
# One slide overview (MLP=ANN=FF)



*Ameisen, et al., "Circuit Tracing: Revealing Computational Graphs in Language Models", Transformer Circuits, 2025.*

# Transformer Architecture of LLMs

•



Attention layer followed by **FF Layer**

- Attention Heads: long-range connections within text

- Feed Forward: same as in previous slide

- N x times : multiple heads run in parallel

- If you trained a neural network, similar plus attention

*Vaswani A et al. Attention is all you need. arXiv:1706.03762. 2017;30*

# The Landscape of Interpretability

- Neurons are "Poly-semantic" (superposition phenomenon)

- 

**Cross-Layer Transcoder**
Features read from one layer and write to all following ones
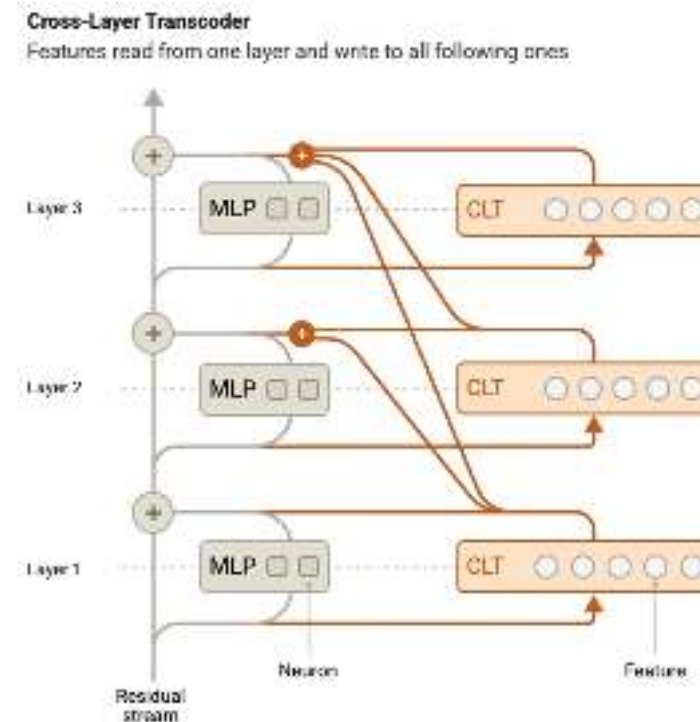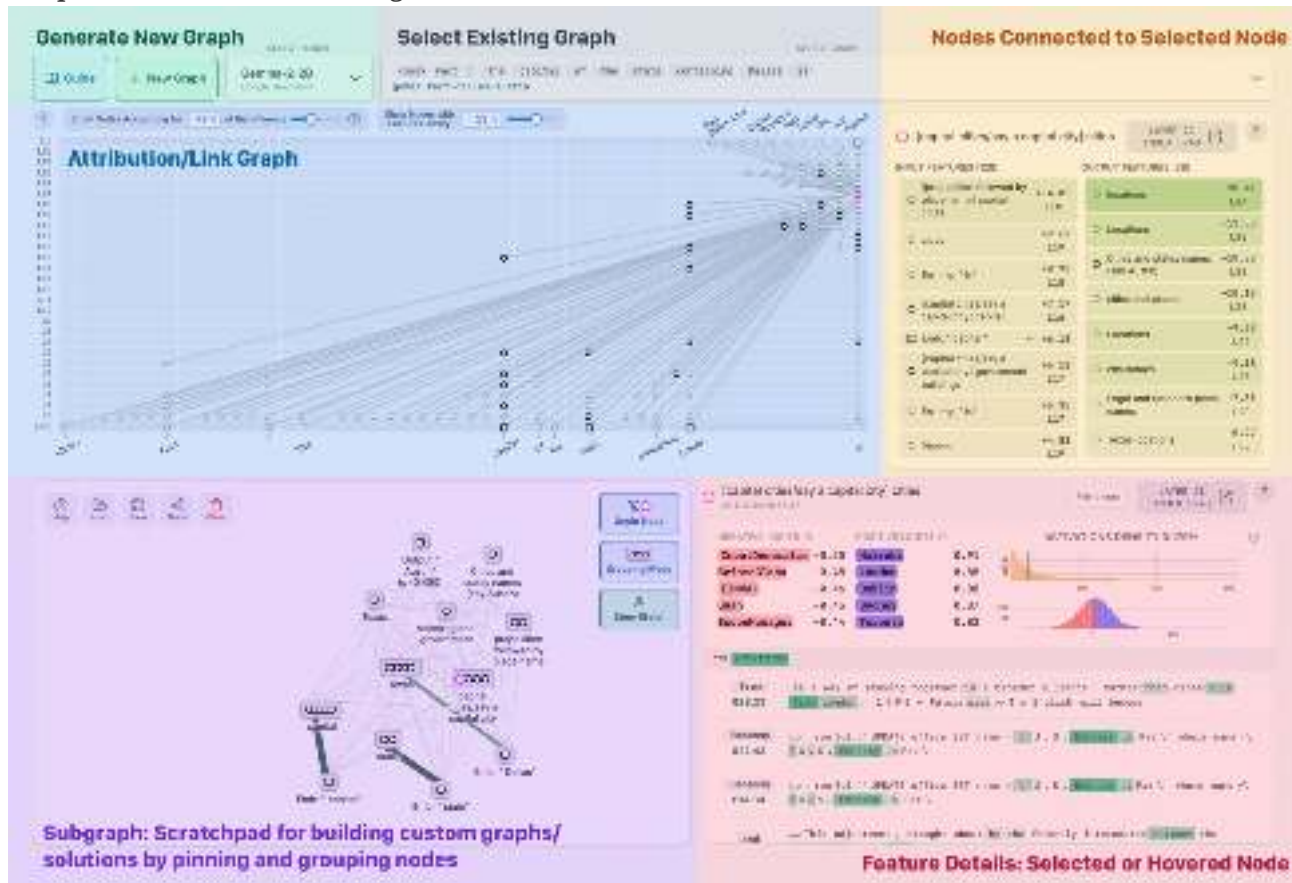


Figure 1: The cross-layer transcoder (CLT) forms the core architecture of our replacement model.

Sparse Autoencoders: Identifying knowledge (features) stored in the neurons

- Linear Probes: Internal linear representations of specific concepts

- Intervention Experiments: Steering, neural activation patching, and ablations

# Attribution Graphs

- Interactions between features (text from input) activating neurons

- Graphs showing feature-feature interactions on specific prompts



- Interaction chains influencing model output

- Prompt analysis revealing circuits

https://www.neuronpedia.org/

# Thanks !

- https://github.com/KKrampis/presentations

- Transformers the tech behind LLMs https://tinyurl.com/3b1b-Transformers

- https://tinyurl.com/alignment-nanda-papers

- https://tinyurl.com/nanda-become-interp-researcher

- https://www.neelnanda.io/mechanistic-interpretability/quickstart