
CS 333:
Natural Language
Processing

Fall 2025

Prof. Carolyn Anderson
Wellesley College

Reminders

- ♦ Final project presentations are next Tuesday (2-3 minutes each!)
- ♦ **All work due by 12/18 at 4pm**
- ♦ Guest speaker this Friday: Jin Zhao (Brandeis)
- ♦ My next help hours: Thursday 4-5

Final Project Presentations

- ◆ Final project presentations are next Tuesday
- ◆ No more than 3 minutes per presentation (2 minutes is fine)
- ◆ You should have 1 slide
- ◆ Your slide should show an example of your task

Recap

Aligning Models with Human Preferences

Users tend to have preferences about generated text that go beyond its statistical frequency.

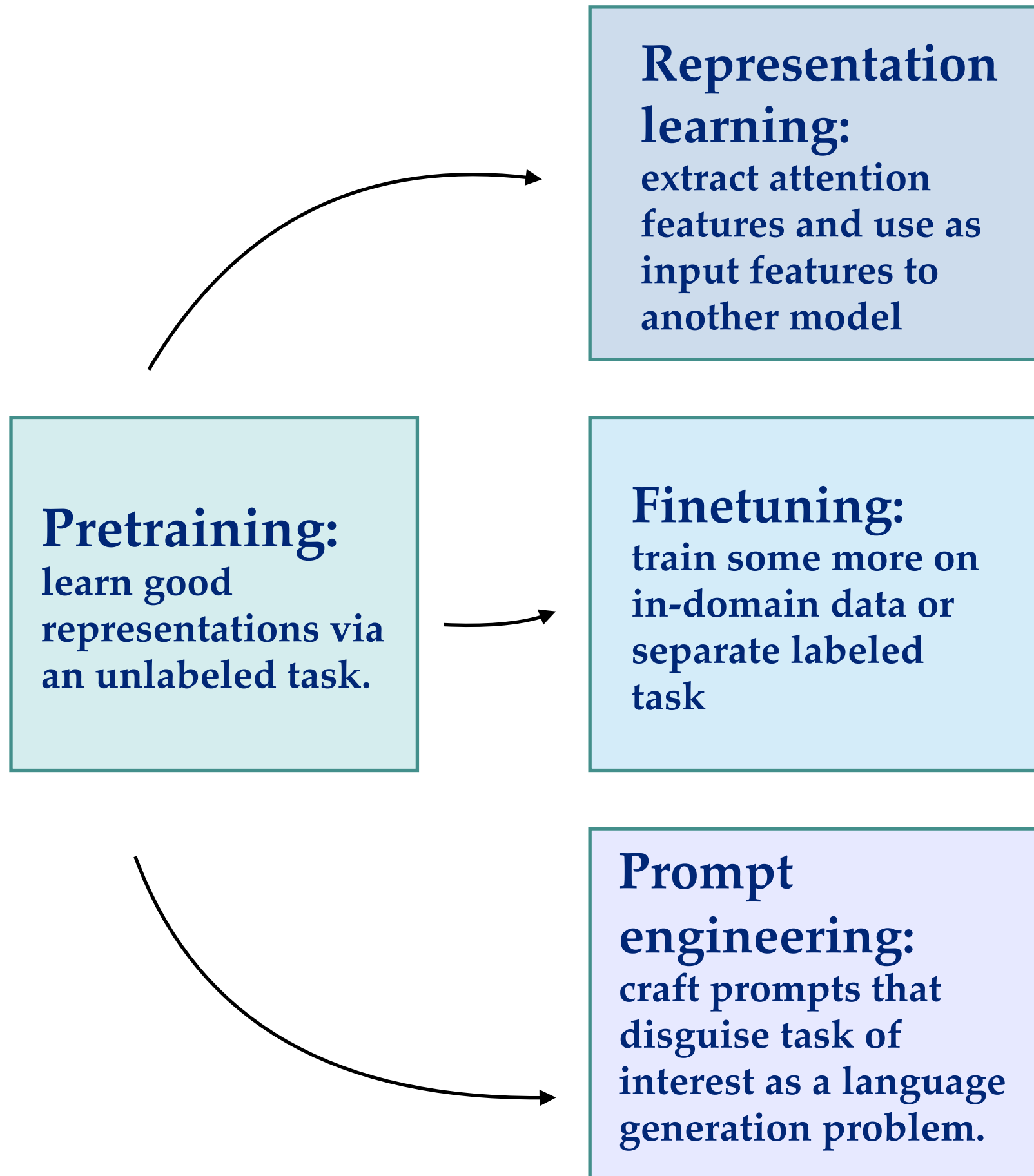
Companies also have preferences for how their models respond.



**How can we align
model behavior
with these
preferences?**



Mark Dredze
@mdredze



Fine-tuning = Further training

Fine-tuning:
train some more on
in-domain data or
separate labeled
task

**Simple
Fine-tuning**

*Just train on more
(labeled) data*

**Proxy
Tuning**

*Fine-tune a smaller
model and use it
steer a larger model*

**Reinforcement
Learning**

*Tune the model with
a reward function*

Interpretability

Interpretability

Mechanistic interpretability is the study of how models make their decisions.

How can we figure out what is happening inside the model, given that contemporary models have *billions* of parameters?

Common Tools for Interpretability

Logit Lens: **apply the output layer** to an intermediate layer's output in order to understand what prediction the model would make at that point in its computation.

Common Tools for Interpretability

Logit Lens: apply the output layer to an intermediate layer's output in order to understand what prediction the model would make at that point in its computation.

Activation Patching: **replace** a specific set of the model's activations with some other vector to see if it changes the model's prediction.

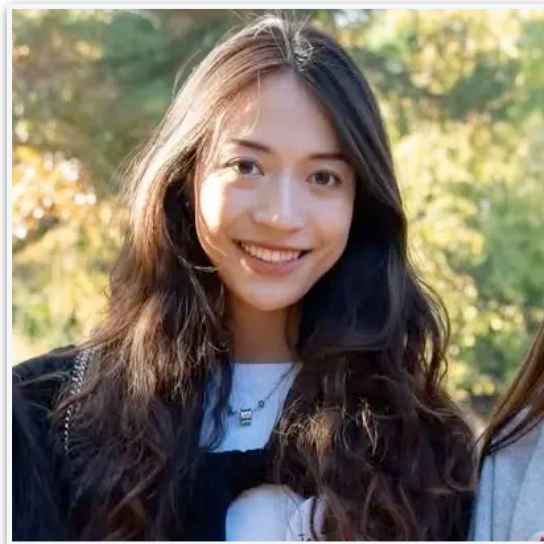
Common Tools for Interpretability

Logit Lens: apply the output layer to an intermediate layer's output in order to understand what prediction the model would make at that point in its computation.

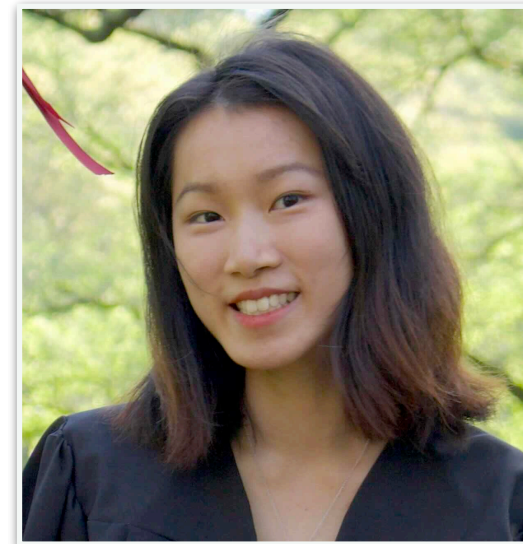
Activation Patching: replace a specific set of the model's activations with some other vector to see if it changes the model's prediction.

Activation Steering: add a vector to the model's activations at a specific point to see if it changes its behavior.

CS333 Alums Working on Interpretability



Jennifer Meng Lu '24



Yik Siu Chan '24

Paths Not Taken: Understanding and Mending the Multilingual Factual Recall Pipeline

Meng Lu*
Brown University
meng_lu@brown.edu

Ruochen Zhang*
Brown University
ruochen_zhang@brown.edu

Carsten Eickhoff
University of Tübingen
carsten.eickhoff@uni-tuebingen.de

Ellie Pavlick
Brown University
ellie_pavlick@brown.edu

Abstract

Multilingual large language models often exhibit factual inconsistencies

EMNLP 2025

derlying pipeline that LLMs employ involves using the English-centric call mechanism to process multilingu and then translating English answers

Pathway to Relevance: How Cross-Encoders Implement a Semantic Variant of BM25

Meng Lu*
Brown University
meng_lu@brown.edu

Catherine Chen*
Brown University
catherine_s_chen@brown.edu

Carsten Eickhoff
University of Tübingen
carsten.eickhoff@uni-tuebingen.de

Abstract

Mechanistic interpretation has greatly contributed to a more detailed understanding of generative language models, enabling significant progress in identifying structures that im

EMNLP 2025

whether a document is relevant to a query. In

principles. These models incorporate explicit components for semantic TF and IDF computations, blending neural architectures with established IR heuristics to improve relevance estimation.

However, the advent of transformer-based models revolutionized the field of IR. These models, trained end-to-end on large numbers of query-document pairs (Nguyen et al., 2016; Thakur et al., 2021), excel at extracting context-dependent semantic signals for ranking tasks. By leveraging multi-headed attention and vast parameter spaces

CAN WE PREDICT ALIGNMENT BEFORE MODELS FINISH THINKING? TOWARDS MONITORING MISALIGNED REASONING MODELS

Yik Siu Chan* Zheng-Xin Yao
Department of Computer Science
Brown University

NeurIPS 2025

Reasoning long chain adversarial for prediction of final results

SPEAK EASY: Eliciting Harmful Jailbreaks from LLMs with Simple Interactions

ICML 2025

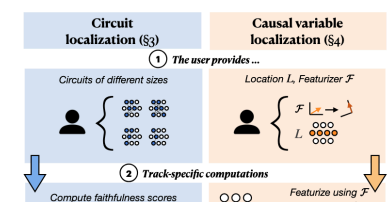
MIB: A Mechanistic Interpretability Benchmark

Aaron Mueller*¹ Atticus Geiger*² Sarah Wiegrefe³
Dana Arad⁴ Iván Arcuschin⁵ Adam Belfki⁶ Yik Siu Chan⁷ Jaden Fiotto-Kaufman⁶ Tal Haklay⁴
Michael Hanna⁸ Jing Huang⁹ Rohan Gupta¹⁰ Yaniv Nikankin⁴ Hadas Orgad⁴ Nikhil Prakash⁶
Anja Reusch⁴ Aruna Sankaranarayanan¹¹ Shun Shao¹² Alessandro Stolfo¹³ Martin Tutek⁴ Amir Zur²
David Bau⁶ Yonatan Belinkov⁴

Abstract

How can we know whether new mechanistic interpretability methods achieve real improvements? We propose a benchmark for mechanistic interpretability methods that evaluates their ability to accurately predict the causal variables in neural language models. The benchmark consists of two parts: (1) The user provides a set of circuits of different sizes, and (2) Track-specific computations. We compute faithfulness scores and featureize using \mathcal{F} .

ICML 2025



significant efforts have been made (e.g., Bai et al., 2022a,b; Gan et al., 2023). However, these “jailbreaks” (Jin et al., 2023) aim to breach safety guardrails and elicit harmful responses, often by malicious actors

LLMs by non-technical users. Existing methods for understanding LLMs (Zou et al., 2023) or for identifying harmful content (Mehrotra et al., 2023; Mehrotra et al., 2023) may not accurately capture the average user attempts to jailbreak LLMs (NPR, 2025).

The Fragility of Model Alignment

Alignment is Fragile

Today we're going to use mechanistic interpretability methods to explore a key challenge with contemporary LLMs: model alignment is very fragile.

Our goal is to *jailbreak* a model: trick the model into telling us something that it has been fine-tuned not to reveal.

Activation Patching Demo

Interpretable Neural Networks

NNsight (/ɛn.saɪt/) is a package for interpreting and manipulating the internals of models.

[Start](#)[Docs](#)[Features](#)[About](#)