

# MEASURING THE RELIABILITY OF CAUSAL PROBING METHODS: TRADEOFFS, LIMITATIONS, AND THE PLIGHT OF NULLIFYING INTERVENTIONS

Marc E. Canby\*, Adam Davies\*, Chirag Rastogi, Julia Hockenmaier



\*Equal Contribution

University of Illinois Urbana-Champaign



# SUMMARY

Which features do LLMs use to perform a given task?

- **Causal interventions** on feature representations

# SUMMARY

Which features do LLMs use to perform a given task?

- **Causal interventions** on feature representations

How ***reliable*** are interventions?

- **Completeness:** Is the intended intervention carried out?
- **Selectivity:** Are we damaging non-targeted features?

# SUMMARY

Which features do LLMs use to perform a given task?


- **Causal interventions** on feature representations

How **reliable** are interventions?


- **Completeness:** Is the intended intervention carried out?
- **Selectivity:** Are we damaging non-targeted features?

We define an **evaluation framework** to compare *different classes* of interventions


# CAUSAL PROBING




# CAUSAL PROBING



# CAUSAL PROBING




# CAUSAL PROBING




# CAUSAL PROBING

The boy with the  
keys [MASK] the  
door




# CAUSAL PROBING

The boy with the  
keys [MASK] the  
door




# CAUSAL PROBING

The boy with the  
keys [MASK] the  
door



# CAUSAL PROBING


The boy with the  
keys [MASK] the  
door




# OUR EVALUATION FRAMEWORK

oracle  
probes


# OUR EVALUATION FRAMEWORK



# OUR EVALUATION FRAMEWORK



# OUR EVALUATION FRAMEWORK




# OUR EVALUATION FRAMEWORK


completeness




# OUR EVALUATION FRAMEWORK



# OUR EVALUATION FRAMEWORK



# OUR EVALUATION FRAMEWORK

$$\text{completeness} \left\{ 1 - \text{dist}\left(\text{counterfactual}, \right. \right.$$

$$\left. \left. \right) \right)$$

# OUR EVALUATION FRAMEWORK


$$\text{completeness} \left\{ \begin{array}{l} \text{counterfactual} \\ 1 - \text{dist} \left( \left[ \begin{array}{c} \text{red bar} \\ \text{blue bar} \end{array} \right], \left[ \begin{array}{c} \text{yellow bar} \\ \text{green bar} \end{array} \right] \right) \end{array} \right. \in [0, 1]$$

# OUR EVALUATION FRAMEWORK

completeness

$$\left\{ \begin{array}{l} \text{counterfactual} \\ 1 - \text{dist}\left(\left[ \begin{array}{c} \text{red bar} \\ \text{red bar} \end{array} \right], \left[ \begin{array}{c} \text{red bar} \\ \text{yellow bar} \end{array} \right] \right) \in [0, 1] \\ \text{removal} \\ 1 - \text{dist}\left(\left[ \begin{array}{c} \text{red bar} \\ \text{red bar} \end{array} \right], \left[ \begin{array}{c} \text{yellow bar} \\ \text{yellow bar} \end{array} \right] \right) \in [0, 1] \end{array} \right.$$

# OUR EVALUATION FRAMEWORK



# CAUSAL PROBING INTERVENTIONS

**Concept Removal:** Remove representation of target feature  
– Linear: INLP, RLACE

# CAUSAL PROBING INTERVENTIONS

**Concept Removal:** Remove representation of target feature

- **Linear:** INLP, RLACE

**Counterfactual:** Swap representation from one value to another

- **Linear:** AlterRep
- **Nonlinear (GBIs):** FGSM, PGD, AutoAttack

# CAUSAL PROBING INTERVENTIONS

**Concept Removal:** Remove representation of target feature

- **Linear:** INLP, RLACE

**Counterfactual:** Swap representation from one value to another

- **Linear:** AlterRep
- **Nonlinear (GBIs):** FGSM, PGD, AutoAttack


Intervene on [MASK] token in final layer of BERT

# OUR EVALUATION FRAMEWORK


Evaluate interventions according to

- **Completeness:** Is the intended intervention carried out?
- **Selectivity:** Are we damaging non-targeted features?
- **Reliability:** Harmonic mean of completeness and selectivity


# RELIABILITY: ALTERREP IS MOST RELIABLE!




# COMPLETENESS: REMOVAL METHODS ARE NOT COMPLETE!




# SELECTIVITY: GBIs ARE NOT SELECTIVE; LINEAR METHODS ARE




# COMPLETENESS AND SELECTIVITY ARE A TRADEOFF!




# COMPLETENESS AND SELECTIVITY ARE A TRADEOFF!



# COMPLETENESS AND SELECTIVITY ARE A TRADEOFF!



# FGSM (NONLINEAR) MORE RELIABLE IN EARLIER LAYERS ALTERREP (LINEAR) MORE RELIABLE IN LATER LAYERS



# TAKEAWAYS

We introduce an **evaluation framework** to compare **different classes** of causal probing interventions

- **Tradeoff** between *completeness* and *selectivity*
- Concept removal is **not reliable** (for causal probing)
- Linear interventions better in later layers (less collateral damage)



# Questions?

# SUPPLEMENTARY MATERIALS FOR QA

# MORE RELIABLE METHODS → GREATER $\Delta$ IN TASK ACCURACY

