Motivation
○○○○○○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

# Optimal ablation for interpretability
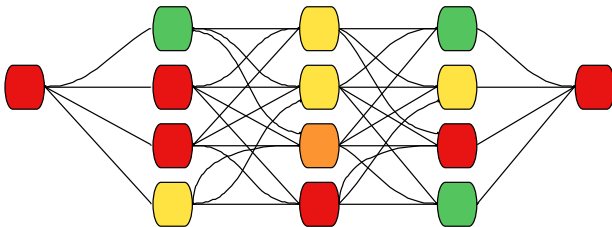
## Maximilian Li and Lucas Janson

NeurIPS 2024 (Spotlight)

November 13, 2024

Motivation
●000000000

Optimal ablation
○○○

Applications
○○○○○○○○○○○

# Interpreting neural networks

How important is a model component?

## Motivating question

Define the *ablation loss gap* $\Delta(\mathcal{M}, \mathcal{A}) := \mathcal{P}(\mathcal{M}^{\backslash \mathcal{A}}) - \mathcal{P}(\mathcal{M})$.

**What is the best performance on subtask $\mathcal{D}$ the model $M$ could have achieved without component $\mathcal{A}$?**

Motivation
○○●○○○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

## Motivating question

Define the *ablation loss gap* $\Delta(\mathcal{M}, \mathcal{A}) := \mathcal{P}(\mathcal{M}^{\backslash \mathcal{A}}) - \mathcal{P}(\mathcal{M})$.

**What is the best performance on subtask $\mathcal{D}$ the model $M$ could have achieved without component $\mathcal{A}$?**

I. <u>Performance on subtask $\mathcal{D}$</u> is measured via expected loss on the subtask, i.e. $\mathcal{P}(\tilde{\mathcal{M}}) = \mathbb{E}_{X \sim \mathcal{D}} \ \mathcal{L}(\tilde{\mathcal{M}}(X), \mathcal{M}(X))$.

Motivation
○○○○●○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

## Motivating question

Define the *ablation loss gap* $\Delta(\mathcal{M}, \mathcal{A}) := \mathcal{P}(\mathcal{M}^{\backslash \mathcal{A}}) - \mathcal{P}(\mathcal{M})$.

**What is the best performance on subtask $\mathcal{D}$ the model $M$ could have achieved without component $\mathcal{A}$?**

I. <u>Performance on subtask $\mathcal{D}$</u> is measured via expected loss on the subtask, i.e. $\mathcal{P}(\tilde{\mathcal{M}}) = \mathbb{E}_{X \sim \mathcal{D}} \mathcal{L}(\tilde{\mathcal{M}}(X), \mathcal{M}(X))$.

II. <u>Model $\mathcal{M}$ could have achieved</u>: $\mathcal{M}^{\backslash \mathcal{A}}$ is constructed solely by changing the value of $\mathcal{A}(X)$.

Motivation
○○○○●○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

## Motivating question

Define the *ablation loss gap* $\Delta(\mathcal{M}, \mathcal{A}) := \mathcal{P}(\mathcal{M}^{\backslash \mathcal{A}}) - \mathcal{P}(\mathcal{M})$.
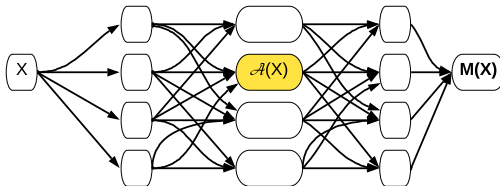
**What is the best performance on subtask $\mathcal{D}$ the model $M$ could have achieved without component $\mathcal{A}$?**

I. <u>Performance on subtask $\mathcal{D}$</u> is measured via expected loss on the subtask, i.e. $\mathcal{P}(\tilde{\mathcal{M}}) = \mathbb{E}_{X \sim \mathcal{D}} \mathcal{L}(\tilde{\mathcal{M}}(X), \mathcal{M}(X))$.
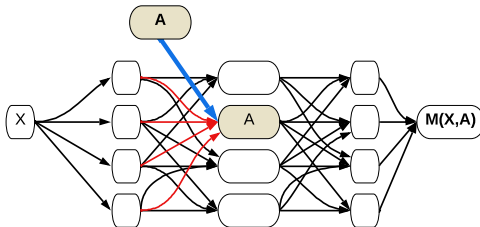
II. <u>Model $\mathcal{M}$ could have achieved</u>: $\mathcal{M}^{\backslash \mathcal{A}}$ is constructed solely by changing the value of $\mathcal{A}(X)$.

III. <u>Without component $\mathcal{A}$</u>: $\mathcal{M}^{\backslash \mathcal{A}}(x)$ uses a value for $\mathcal{A}$ that conveys no information about $x$.
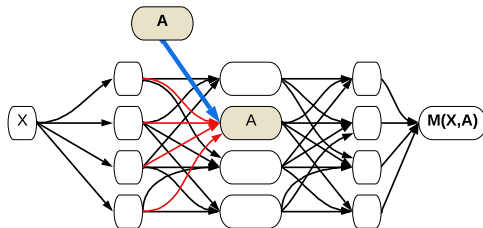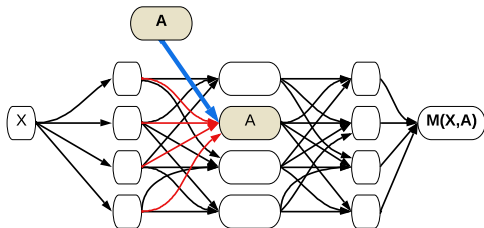
◀ □ ▶ ◀ 🗗 ▶ ◀ 🗉 ▶ ◀ 🗉 ▶ 🗉 🔗 ᘓ ℚ ℚ

Motivation
○○○○○●○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

# Example

Motivation
○○○○○○○●○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

# Example

Motivation
○○○○○○○●○○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

## Example



Definition: A **total ablation method** satisfies
$\mathcal{M}^{\setminus \mathcal{A}}(X) = \mathcal{M}^{\setminus \mathcal{A}}(X, A)$ for $A \perp\!\!\!\perp X$.

Motivation
○○○○○○○○○●○

Optimal ablation
○○○

Applications
○○○○○○○○○○○

# Example



Zero ablation: $A = 0$.

Mean ablation: $A = \mathbb{E}_{X' \sim \mathcal{D}}[\mathcal{A}(X')]$

Resample ablation: $A = \mathcal{A}(X'), X' \perp\!\!\!\perp X$.

Motivation
○○○○○○○○○●

Optimal ablation
○○○

Applications
○○○○○○○○○○○

## Motivating question

Define the *ablation loss gap* $\Delta(\mathcal{M}, \mathcal{A}) := \mathcal{P}(\mathcal{M}^{\setminus \mathcal{A}}) - \mathcal{P}(\mathcal{M})$.

**What is the best performance on subtask $\mathcal{D}$ the model $M$ could have achieved without component $\mathcal{A}$?**

IV. "Best" performance: we want to understand how much performance degrades *because* we had to ablate $\mathcal{A}$.

Seeking best performance avoids interventions that "spoof" the model by causing it to confuse $x$ for a different input, or treat $x$ in a way that it never treated any training input.

Motivation
○○○○○○○○○○

Optimal ablation
●○○

Applications
○○○○○○○○○○○

## Optimal ablation

<u>Definition</u>:

$$\mathcal{M}_{(\text{opt})}^{\backslash \mathcal{A}}(x) := \mathcal{M}_{\mathcal{A}}(x, a^*),$$
$$a^* := \arg \min_{a} \mathbb{E}_{X \sim \mathcal{D}} \ \mathcal{L}(\mathcal{M}_{\mathcal{A}}(X, a), \mathcal{M}(X))$$

### Proposition

*Let $\Delta(\mathcal{M}, \mathcal{A})$ be the ablation loss gap for some component $\mathcal{A}$ measured with any total ablation method. Then*

$$\Delta_{\text{opt}}(\mathcal{M}, \mathcal{A}) \le \Delta(\mathcal{M}, \mathcal{A})$$

# Comparison to counterfactual ablation

**Counterfactual ablation** (CF) considers pairs of parallel inputs.

- CF requires manual effort for each subtask and may not be possible for complex subtasks. OA is more versatile than CF.

Motivation
○○○○○○○○○○

Optimal ablation
○○●

Applications
○○○○○○○○○○○

## Comparison to counterfactual ablation

**Counterfactual ablation** (CF) considers pairs of parallel inputs.

- CF removes *less* information than OA, yet still achieves higher loss, which is evidence that *most* loss can be attributed to spoofing.

Table 1: Comparison of ablation loss gap $\Delta$ on IOI

|  | Zero | Mean | Resample | CF-Mean | Optimal | CF |
|---|---|---|---|---|---|---|
| Rank correlation with CF | 0.590 | 0.825 | 0.828 | 0.833 | **0.907** | 1 |
| Median ratio of $\Delta_{opt}$ to $\Delta$ | 11.1% | 33.0% | 17.7% | 31.7% | 100% | 88.9% |

Motivation
○○○○○○○○○○

Optimal ablation
○○○

Applications
●○○○○○○○○○○

# Circuit discovery

We introduce a *uniform gradient sampling* method to find circuits.

Motivation
○○○○○○○○○○

Optimal ablation
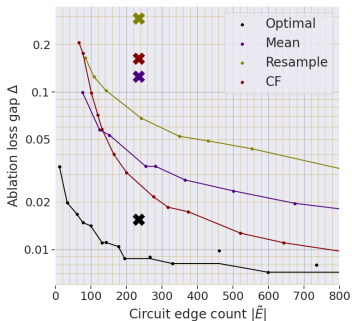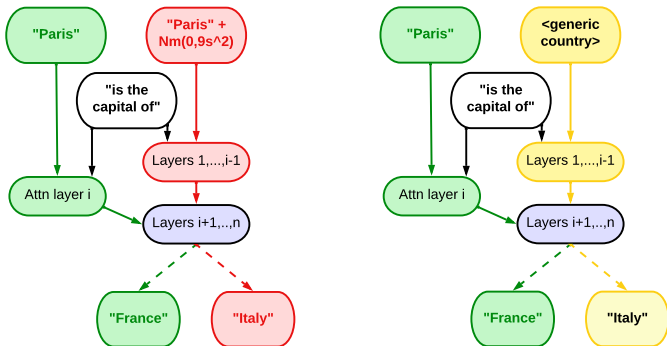○○○

Applications
○●○○○○○○○○○○

# Circuit discovery results



IOI circuits, ablation comparison

Greater-Than circuits, ablation comparison

# Causal tracing

Motivation
○○○○○○○○○○

Optimal ablation
○○○

Applications
○○○○●○○○○○○○

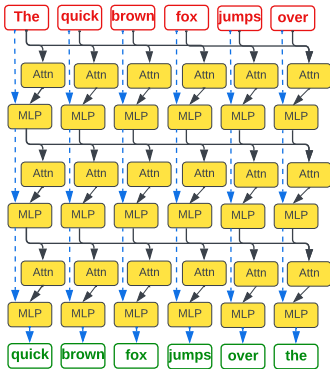# Causal tracing results

Causal tracing intervention at last token, window size 5

Motivation
○○○○○○○○○○○○

Optimal ablation
○○○

Applications
○○○○○●○○○○○○

# Latent prediction

Motivation
ooooooooooo

Optimal ablation
ooo

Applications
ooooooooooo

# Latent prediction

Motivation
oooooooooo

Optimal ablation
ooo

Applications
ooooooo●oooo

# Latent prediction: tuned lens

Motivation
○○○○○○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○●○○○

# Latent prediction: Optimal Constant Attention (OCA lens)

Motivation
○○○○○○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○●○○

# Latent prediction results



Lens loss, GPT-2-XL

Legend:
- OCA lens
- Tuned lens
- Mean
- Resample

KL-divergence vs Layer number

Motivation
○○○○○○○○○○

Optimal ablation
○○○

Applications
○○○○○○○○○○●○

# Latent prediction: causal faithfulness



Causal faithfulness: basis-aligned perturbation

Motivation
oooooooooo

Optimal ablation
ooo

Applications
ooooooooooo●

# Latent prediction: truthful elicitation



Elicitation accuracy on selected datasets with 10 demos, GPT-2-XL