Using sparse feature circuits, we can discover and edit unanticipated mechanisms in language models.

**Sparse Feature Circuits: Discovering** and Editing Interpretable Causal Graphs in Language Models Samuel Marks<sup>1</sup>, Can Rager<sup>2</sup>, Eric J. Michaud<sup>3</sup>, Yonatan Belinkov<sup>4</sup>, David Bau<sup>2</sup>, Aaron Mueller<sup>2,4</sup>



Northeastern University

TECHNION Israel Institute of Technology

# **Discovering Sparse Feature Circuits**

**Sparse** autoencoders let us decompose polysemantic



We fold SAEs into the forward pass, and attribute graphs of features in performing a behavior:



# **Case Study: Subject-Verb Agreement**



fewer nodes than neuron circuits:

Sentence boundaries, healthcare, places, i.a.?

## $\begin{array}{c} & & a_2 \\ & & b_2 \\ & & \epsilon_2 \end{array}$ Compute and filter edges. 3 Compute effects. Filter nodes. m m $a_2$ $b_2$ $(a_2)$ $(b_2)$ $(b_1) \epsilon_1$ $\widehat{\text{IE}}(a,m) = \nabla_a m \cdot (a - a)$ $\hat{\text{IE}}(a,m) > T_N$

and editor for

		Pythia-70N	M	Gemma-2-2B			
Method	<b>†Profession</b>	↓Gender	↑Worst group	<b>†Profession</b>	↓Gender	↑Worst group	
Original	61.9	87.4	24.4	67.7	81.9	18.2	
CBP	83.3	60.1	67.7	90.2	50.1	86.7	
Random	61.8	87.5	24.4	67.3	82.3	18.0	
Shift	88.5	54.0	76.0	76.0	51.5	50.0	
SHIFT + retrain	93.1	52.0	89.0	95.0	52.4	92.9	
Neuron skyline	75.5	73.2	41.5	65.1	84.3	5.6	
Feature skyline	88.5	54.3	62.9	80.8	53.7	56.7	
Oracle	93.0	49.4	91.9	95.0	50.6	93.1	

SHIFT outperforms a neuron-based approach that has an unfair advantage.



Most interpretability pipelines require us to know which behaviors we're looking for. What about unanticipated behaviors?

Given large corpus of  $\{(x_i, y_i)\}$ , encode into feature activations  $\mathbf{v}(x_i, y_i)$ , cluster  $\mathbf{v}$ , discover sparse feature circuits for each cluster.

## **Cluster 382:** Incrementing sequences

var input = [1, 2, 3, 4, 5, 6, 7, 8

Step 1. Download the latest CompsNY 3.49 Full Step 2. Double click the Setup file and follow the prompts [...] Step 3. After the main install closes, click OK [...

Example features involved:





Our judgments about feature relevance are largely informative. We get close to the performance of a classifier trained on **unbiased** data!

	12				
subervised		Interpreta		PID	eline

Narrow induction				
A3 $A \rightarrow 3$ or III or $4$				
A7 $A \rightarrow 7$ or vii or 8				

## **Cluster 475**: "to" as infinitive object

At issue, whether the defendant should be allowed to British Prime Min David Cameron says in televised remarks he would like Britain to Reader bloggers are asked to

Example features involved:

Objects which can precede object complements	Other words which precede infinitive objects				
Direct the user to It's up to you to	According to This infection leads to				