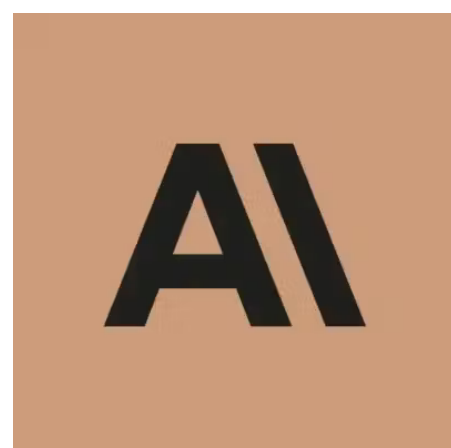# Sparse Feature Circuits

## Discovering and Editing Interpretable Causal Graphs in Language Models

Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, **Aaron Mueller**
2025 International Conference on Learning Representations (ICLR)
26 April 2025

# Interpretability

For a model to generalize, it must achieve right answers *for the right reasons.*

# Interpretability

For a model to generalize, it must achieve right answers *for the right reasons*.

**How** do neural networks (NNs) perform particular behaviors?

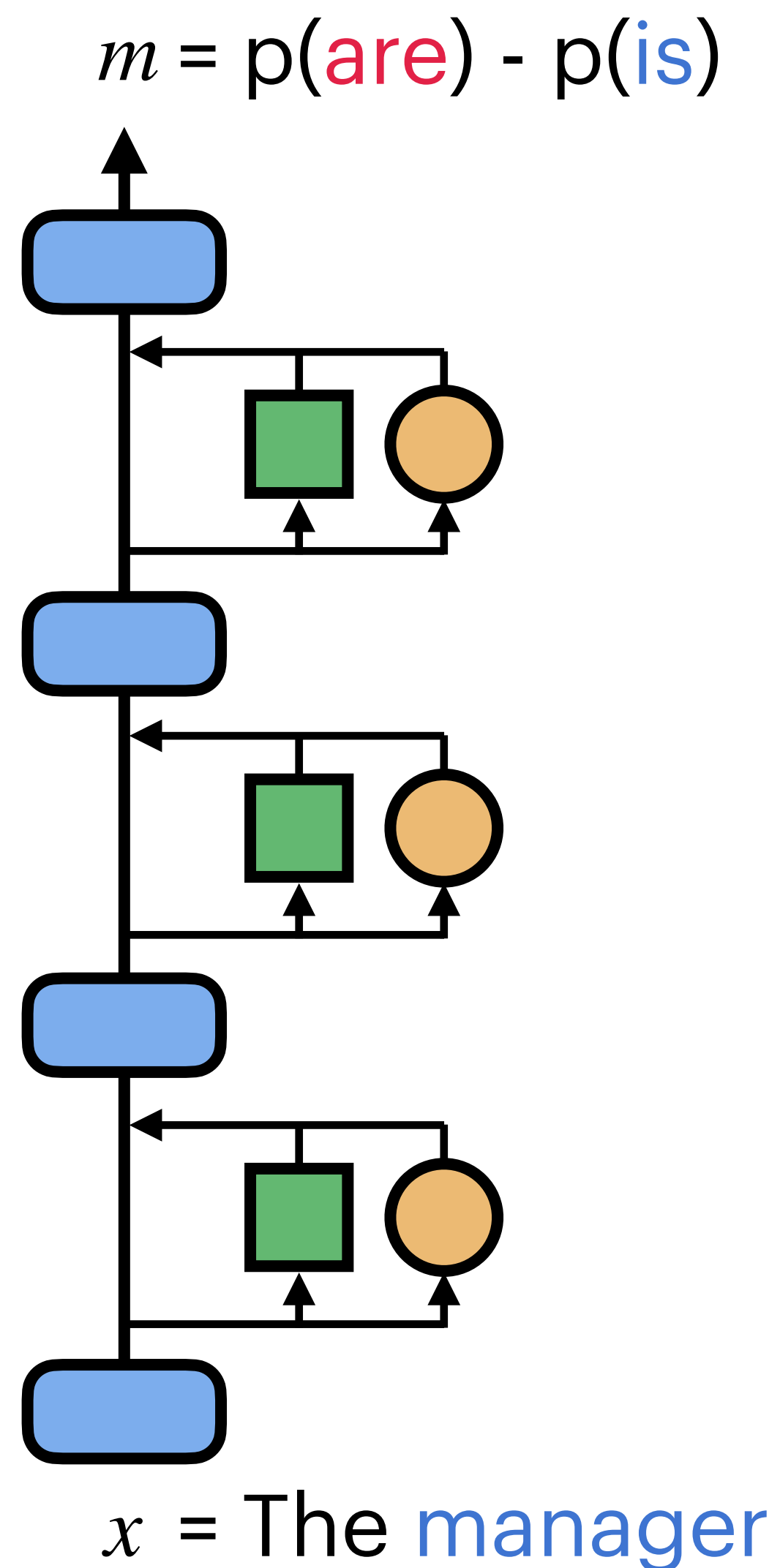**Why** do they behave in certain ways on certain inputs?

# Interpretability

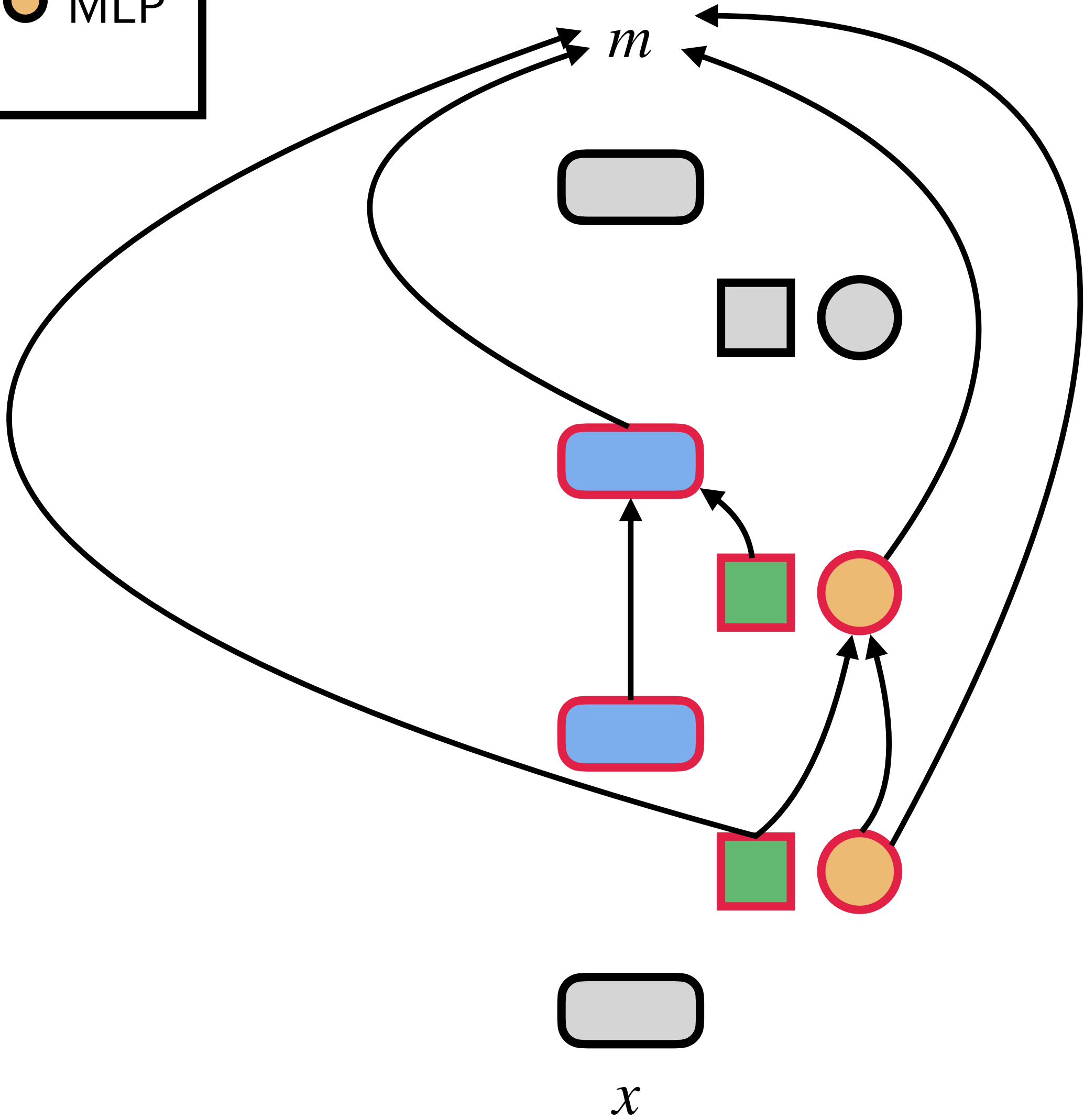For a model to generalize, it must achieve right answers *for the right reasons*.

**How** do neural networks (NNs) perform particular behaviors?
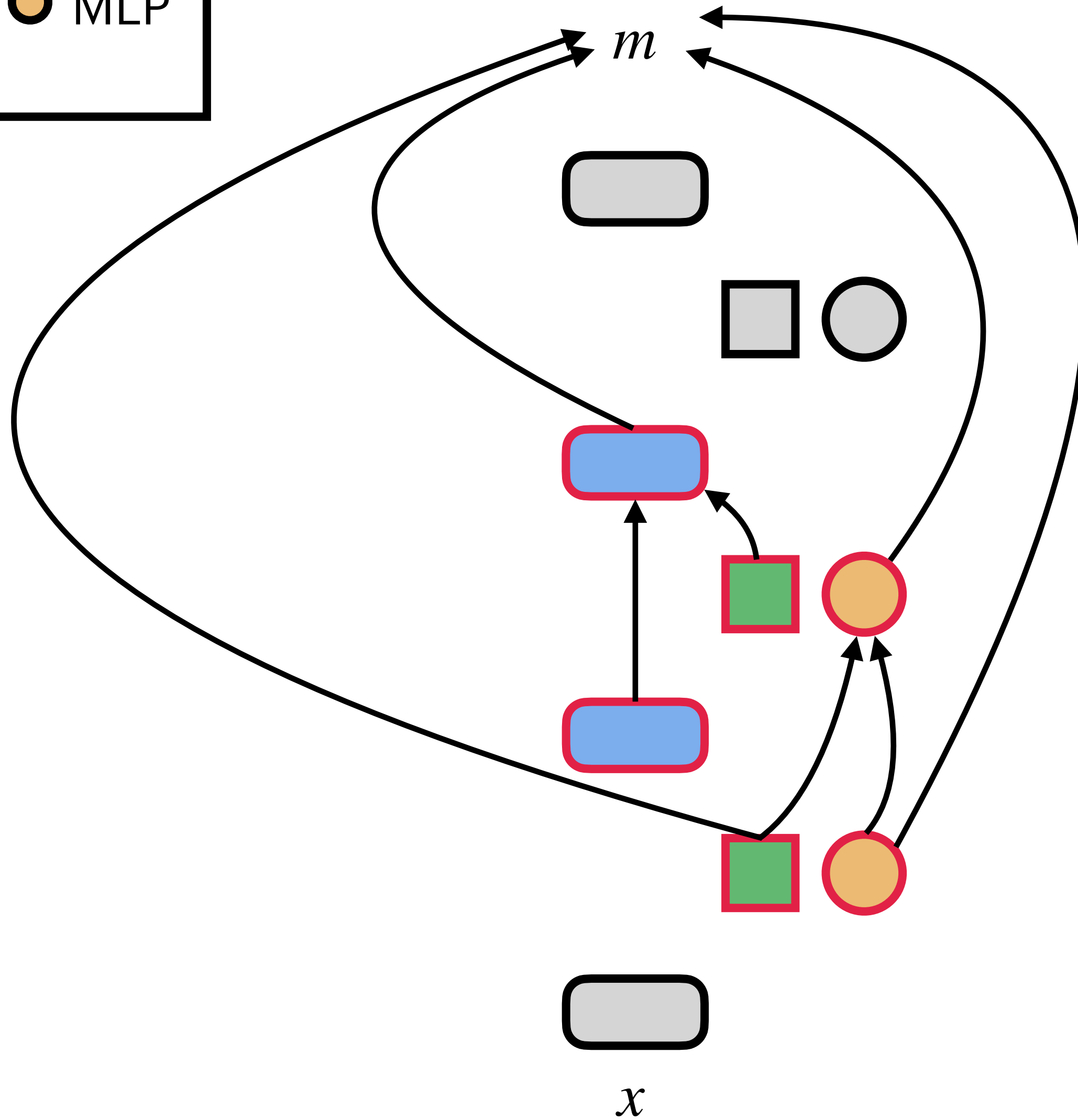
**Why** do they behave in certain ways on certain inputs?

*How can we locate and understand **unanticipated mechanisms**?*

attention · MLP · vector state

$m = \text{p}(\text{are}) - \text{p}(\text{is})$

$x = \text{The manager}$

Given a neural network, we want to know which components contribute most to the model's behavior.
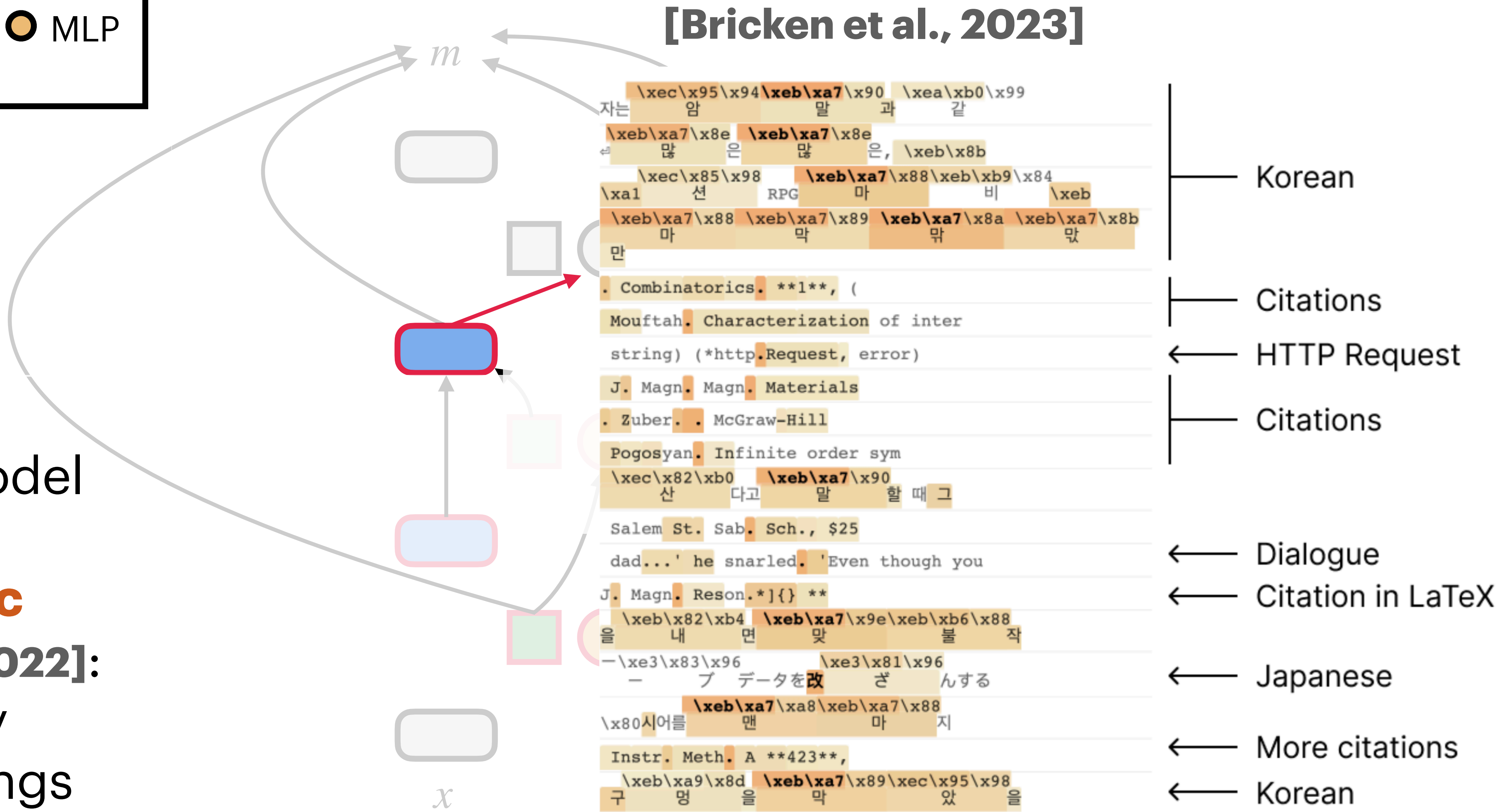
We have a **circuit**!

We have a **circuit**!

*...Now what?*

**[Bricken et al., 2023]**

Legend:
- ■ attention
- ● MLP
- ▭ vector state

Language model neurons are **polysemantic** **[Elhage et al., 2022]**: they do many unrelated things simultaneously.

Labels (right side): Korean, Citations, HTTP Request, Citations, Dialogue, Citation in LaTeX, Japanese, More citations, Korean
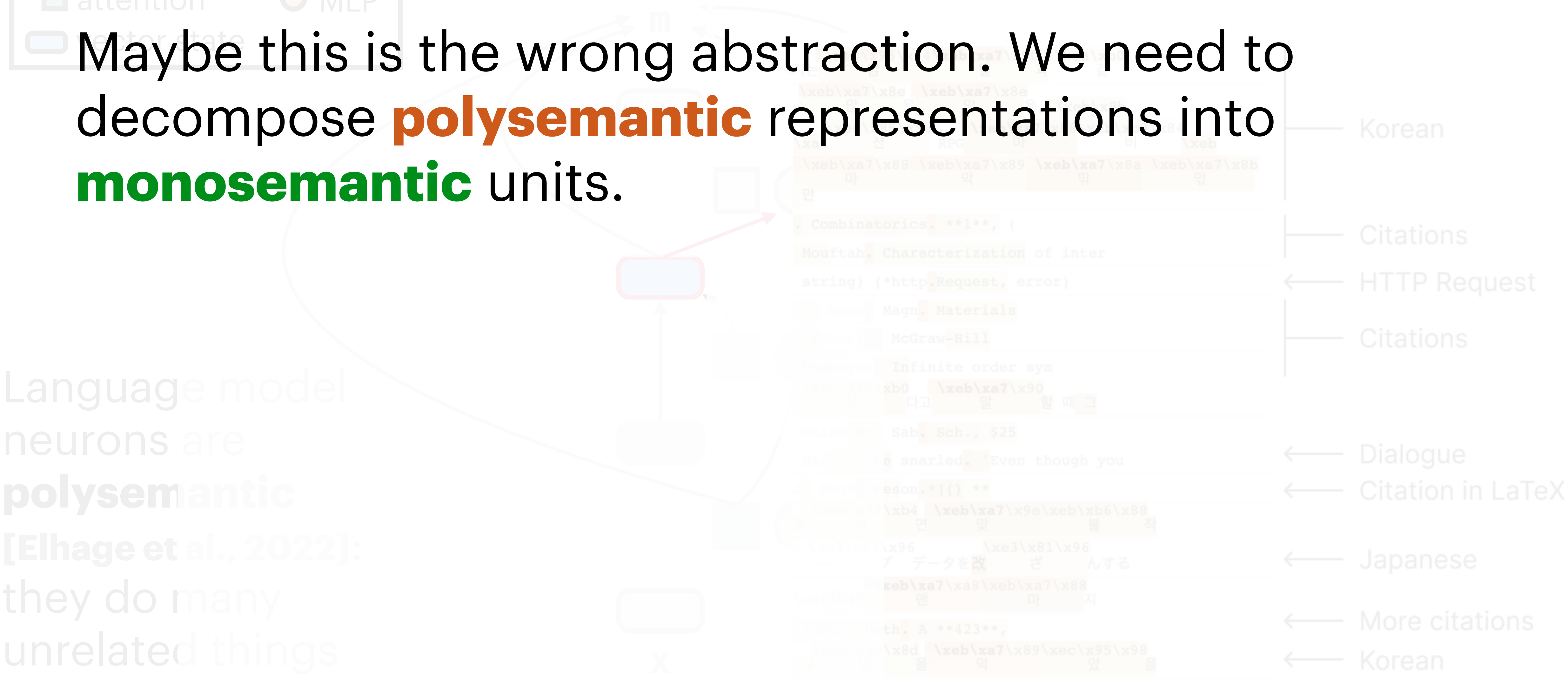
Nelson Elhage et al. (2022). "Toy Models of Superposition." *Transformer Circuits Thread*.
Trenton Bricken et al. (2023). "Towards Monosemanticity: Decomposing Language Models with Dictionary Learning." *Transformer Circuits Thread*.

Maybe this is the wrong abstraction. We need to decompose **polysemantic** representations into **monosemantic** units.

Language model neurons are polysemantic [Elhage et al., 2022]: they do many unrelated things simultaneously.

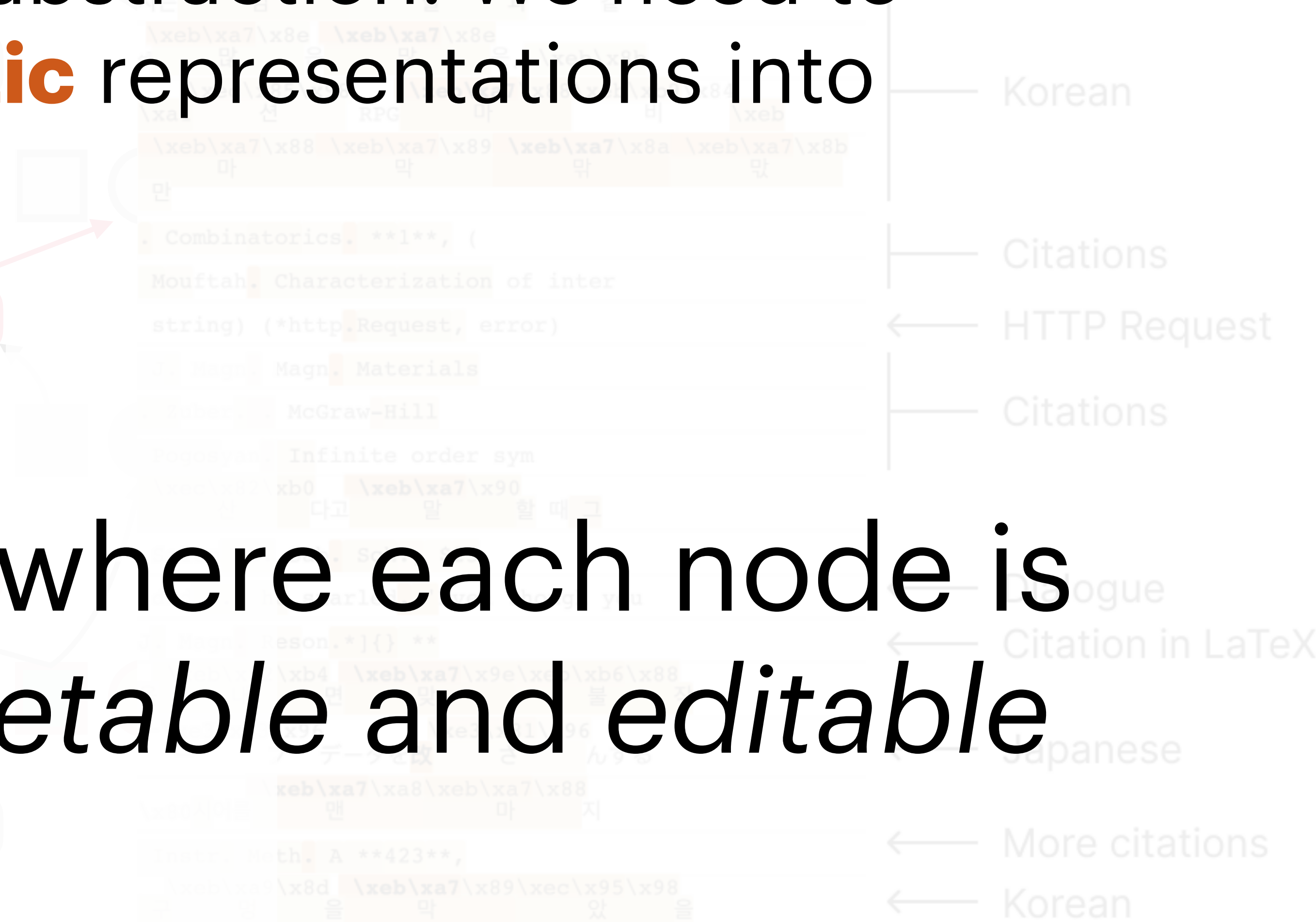Nelson Elhage et al. (2022). "Toy Models of Superposition." *Transformer Circuits Thread.*
Trenton Bricken et al. (2023). "Towards Monosemanticity: Decomposing Language Models with Dictionary Learning." *Transformer Circuits Thread.*

Maybe this is the wrong abstraction. We need to decompose **polysemantic** representations into **monosemantic** units.
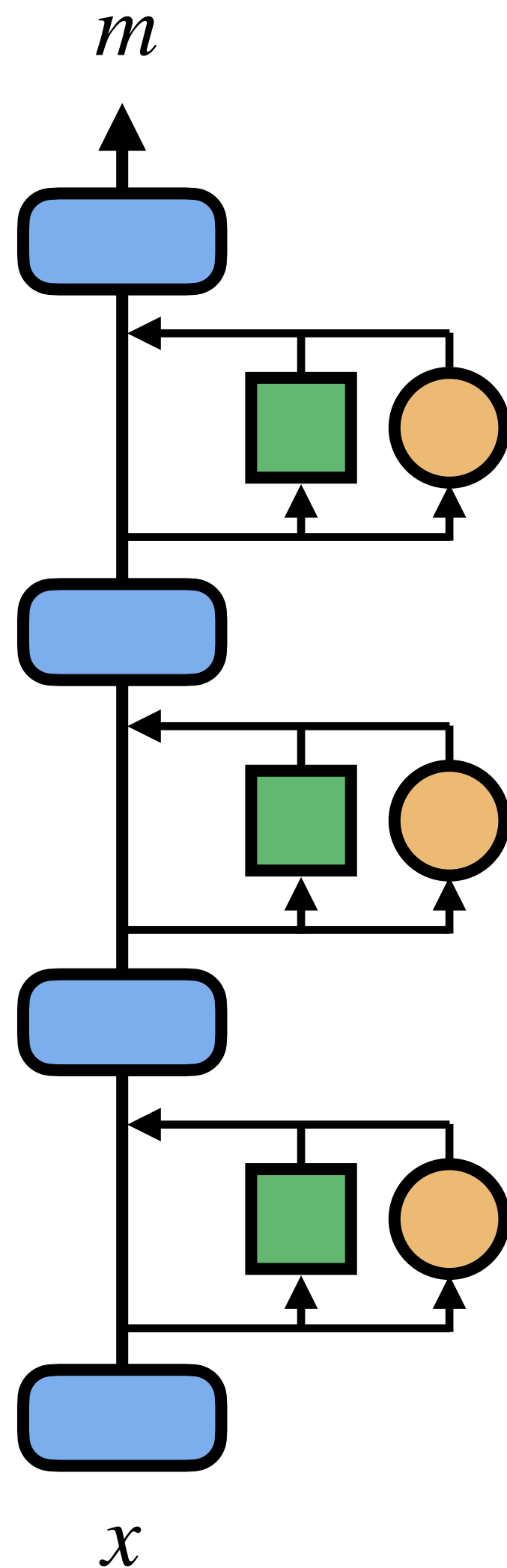
**Goal**: a circuit where each node is *human-interpretable* and *editable*

Nelson Elhage et al. (2022). "Toy Models of Superposition." *Transformer Circuits Thread.*
Trenton Bricken et al. (2023). "Towards Monosemanticity: Decomposing Language Models with Dictionary Learning." *Transformer Circuits Thread.*
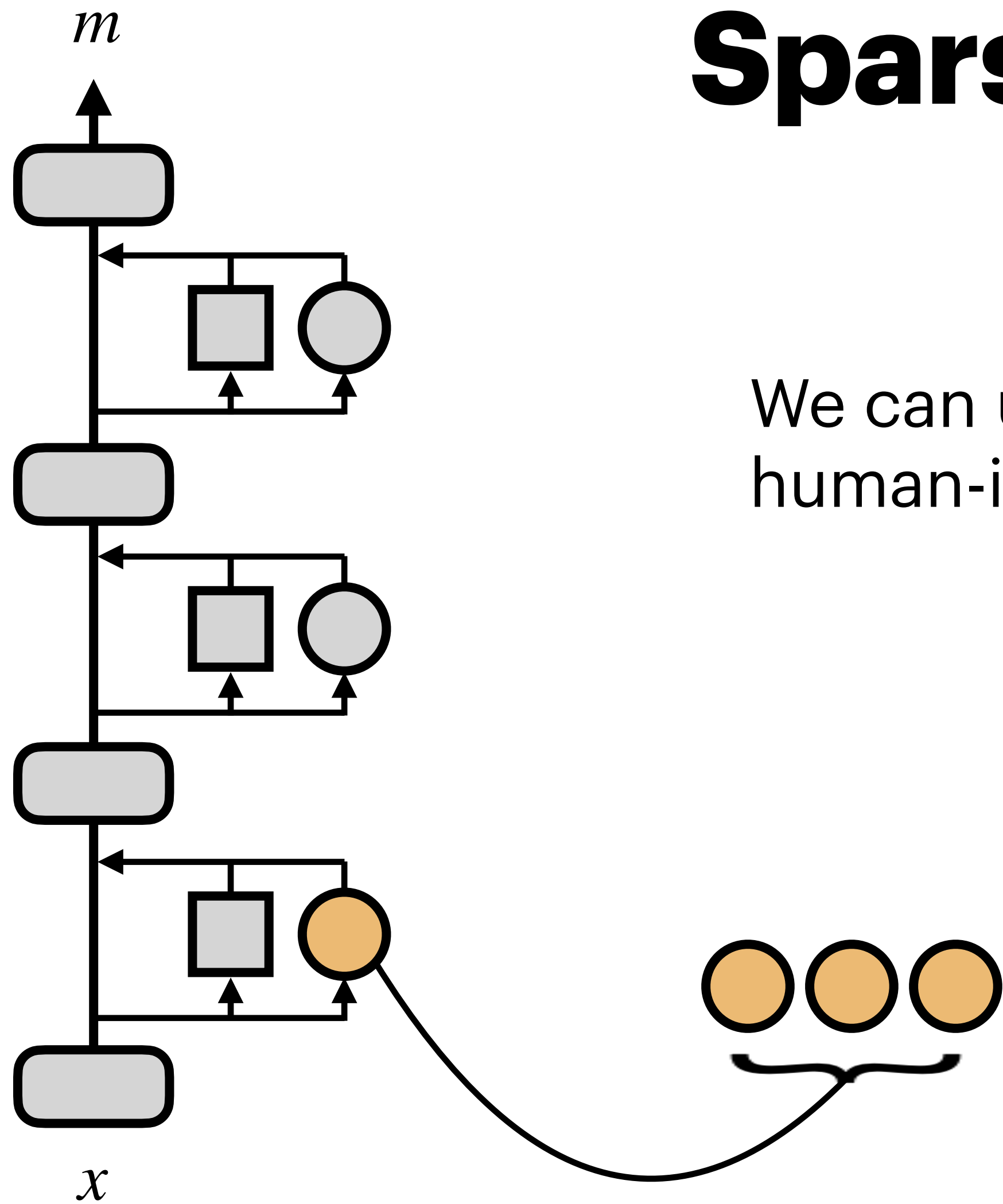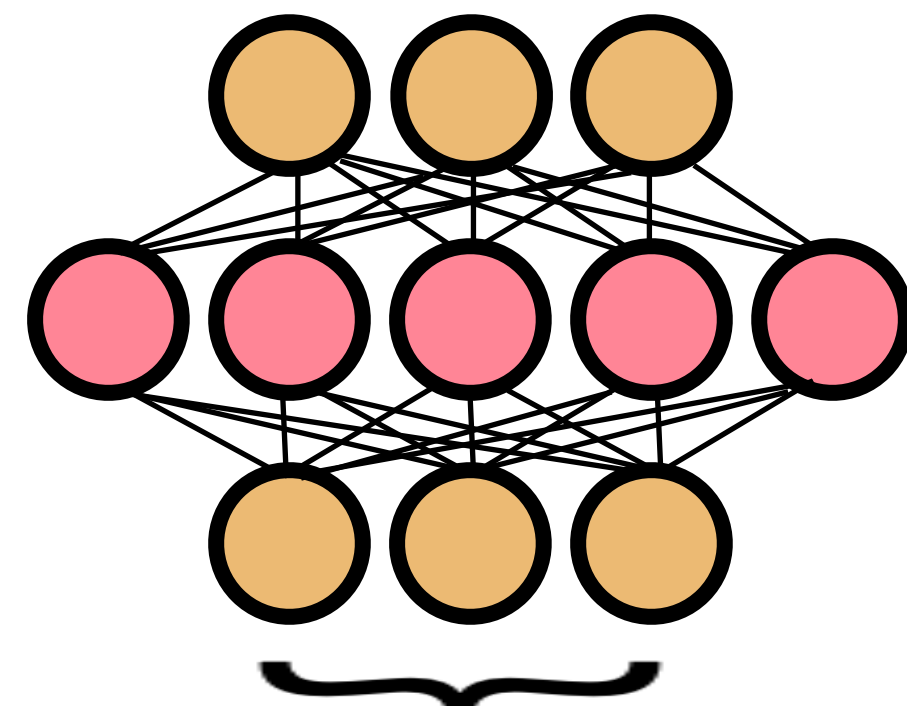
# Sparse Features

We can use **sparse autoencoders** (SAEs) to disentangle human-interpretable **features** from model components

# Sparse Features

We can use **sparse autoencoders** (SAEs) to disentangle human-interpretable **features** from model components

# Sparse Features



We can use **sparse autoencoders** (SAEs) to disentangle human-interpretable **features** from model components

$$\hat{\mathbf{x}} = W_d\mathbf{f} + \mathbf{b}_d$$

$$\mathbf{f} = \text{ReLU}(W_e(\mathbf{x} - \mathbf{b}_d) + \mathbf{b}_e)$$

$$\mathbf{x}$$

attention   MLP
vector state   SAE
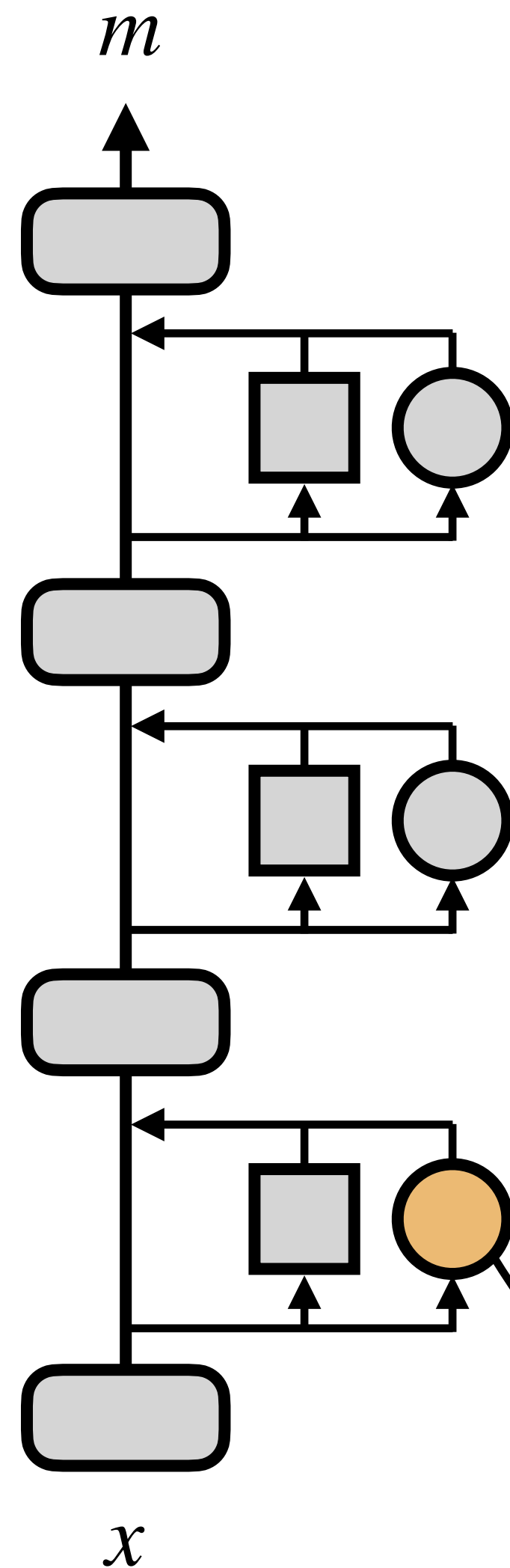
# Sparse Features

$m$

We can use **sparse** **autoencoders** (SAEs) to disentangle human-interpretable **features** from model components
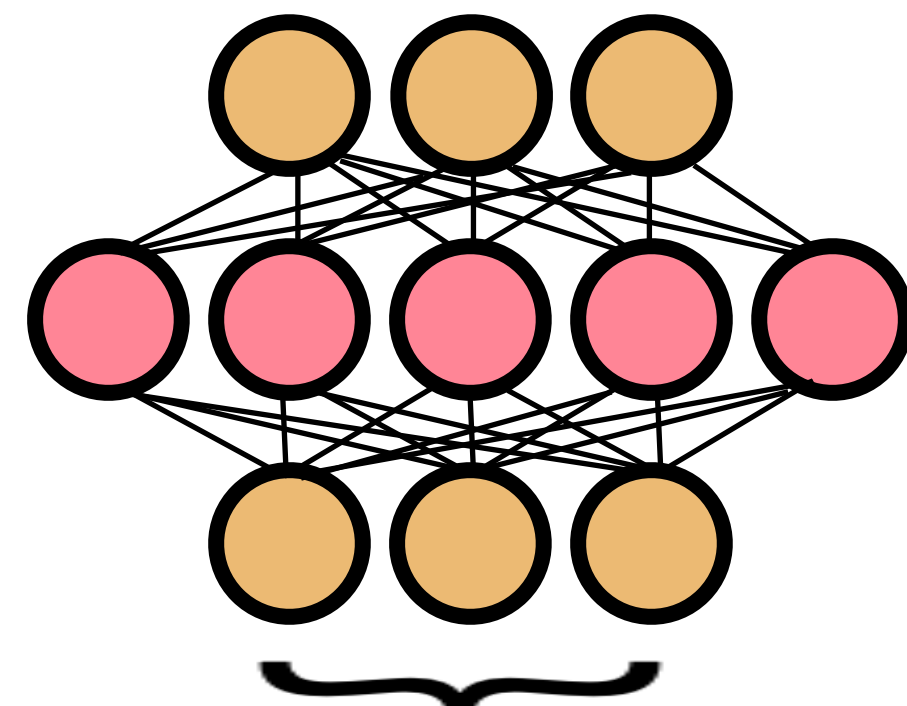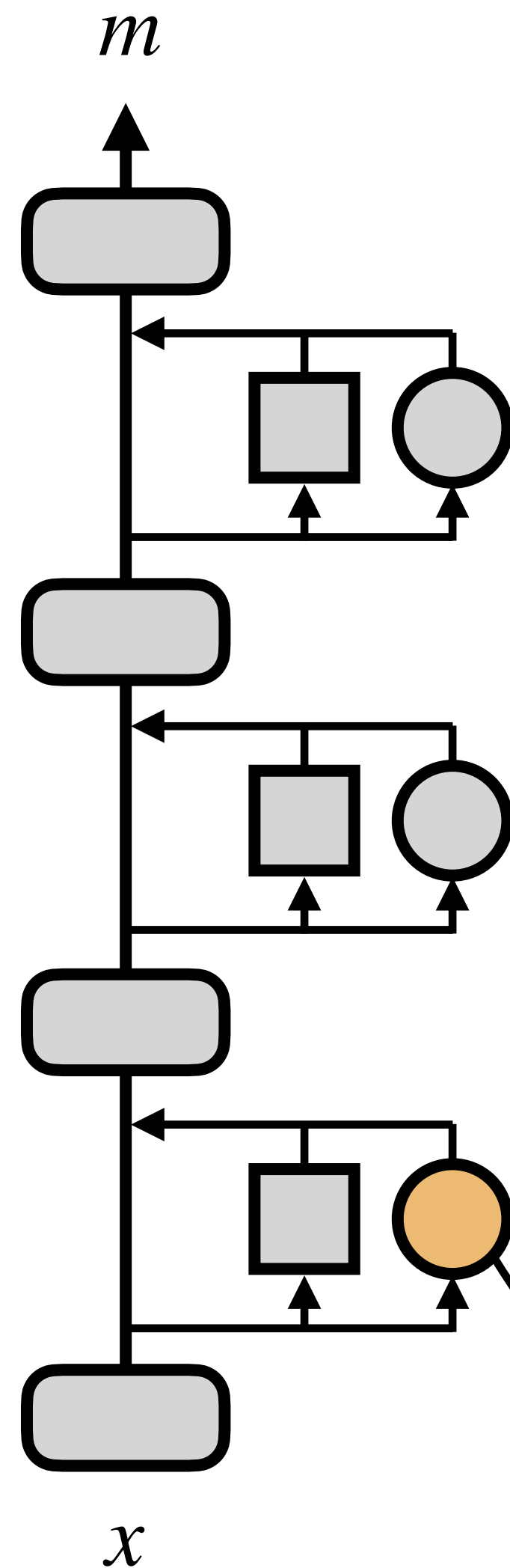
$$\hat{\mathbf{x}} = W_d\mathbf{f} + \mathbf{b}_d$$

$$\mathbf{f} = \text{ReLU}(W_e(\mathbf{x} - \mathbf{b}_d) + \mathbf{b}_e)$$

$$\mathbf{x}$$

$$L = \sqrt{\text{MSE}(\mathbf{x}, \hat{\mathbf{x}})} + \lambda\|\mathbf{f}\|_1$$

$$\mathbf{x} = \hat{\mathbf{x}} + \epsilon$$

$x$

- ■ attention
- ▭ vector state
- ● MLP
- ● SAE

# Sparse Features

obau, the daughter of Ratu Sir George

office by a homeless woman named Lois Lang.

Benedict debate. But she has some thoughts on

of these creative women, the reader gets

"Ma'am?" "You

the physician who examined her body was unable to

you hear her towards the end what

Norma and Sherryl suggest that there was

*Words related to women*

goal of our research program on innate immune sensors

4. His research interests include bioinformatics

.K.'s group are funded by the

Dave Lovinger?s Laboratory, investigates the

a Hungarian mathematician who works as a professor at

in the Kalluri laboratory, where both tumor

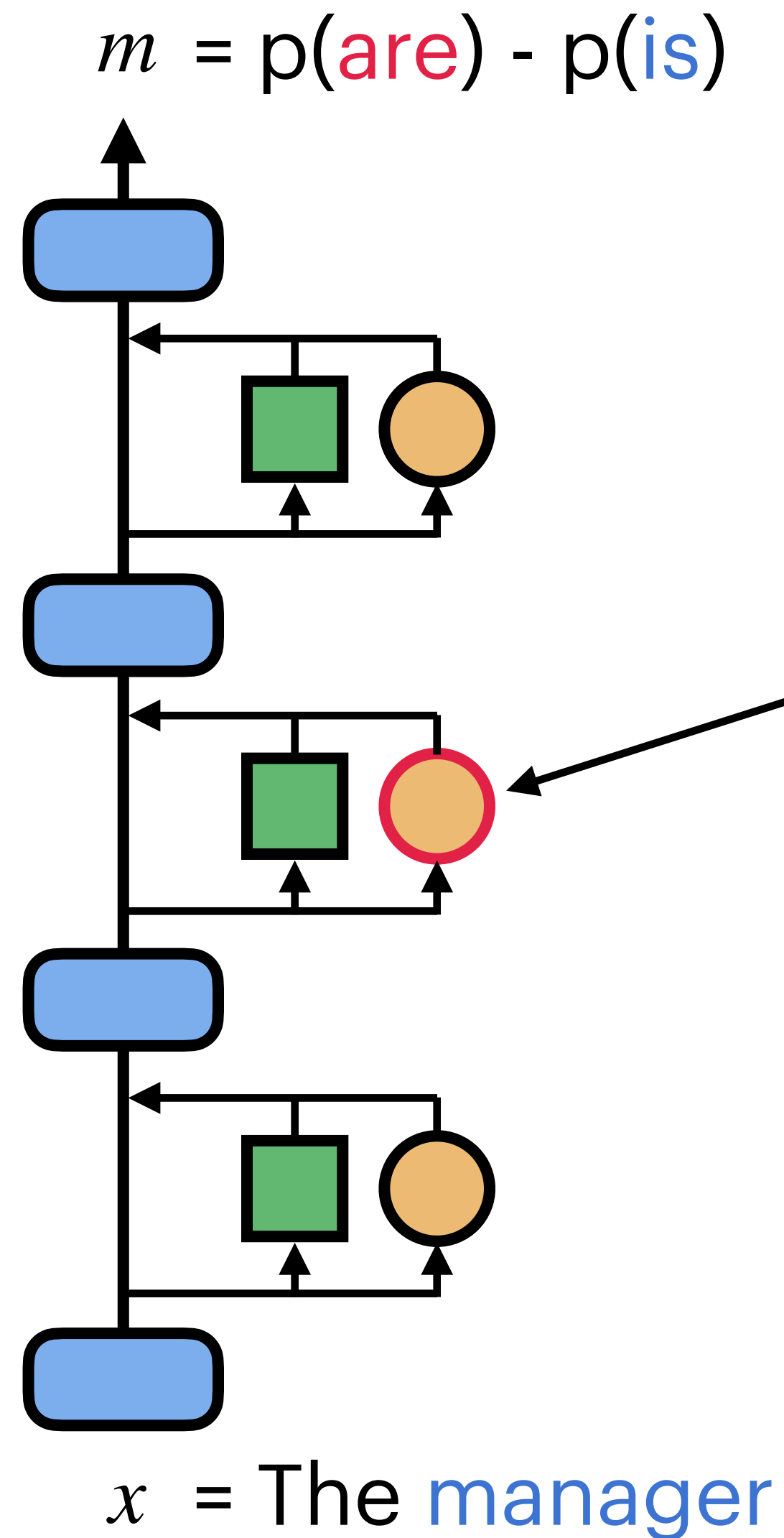the Human Cognitive Neuroscience Unit at Northumbria University

Murakami's research team, which received a

*Passages related to academia, research*

# Language Models and SAEs

- **Models:** Pythia (70M), Gemma 2 (2B)

- **SAEs:** GemmaScope **[Lieberum et al, 2024]**, or trained by us on model activations given documents from The Pile

- SAE features are interpreted using activations and logits from The Pile

Tom Lieberum et al. (2024). "Gemma Scope: Open Sparse Autoencoders Everywhere All at Once On Gemma 2." *BlackboxNLP*.

**Activation Patching**

$m$ = p(are) - p(is)

$x'$ = The ~~manager~~ managers

do(swap-number): Set $a$ to what it would have been if the subject in $x$ were the opposite number

$x$ = The manager

attention   MLP
vector state

# Activation Patching
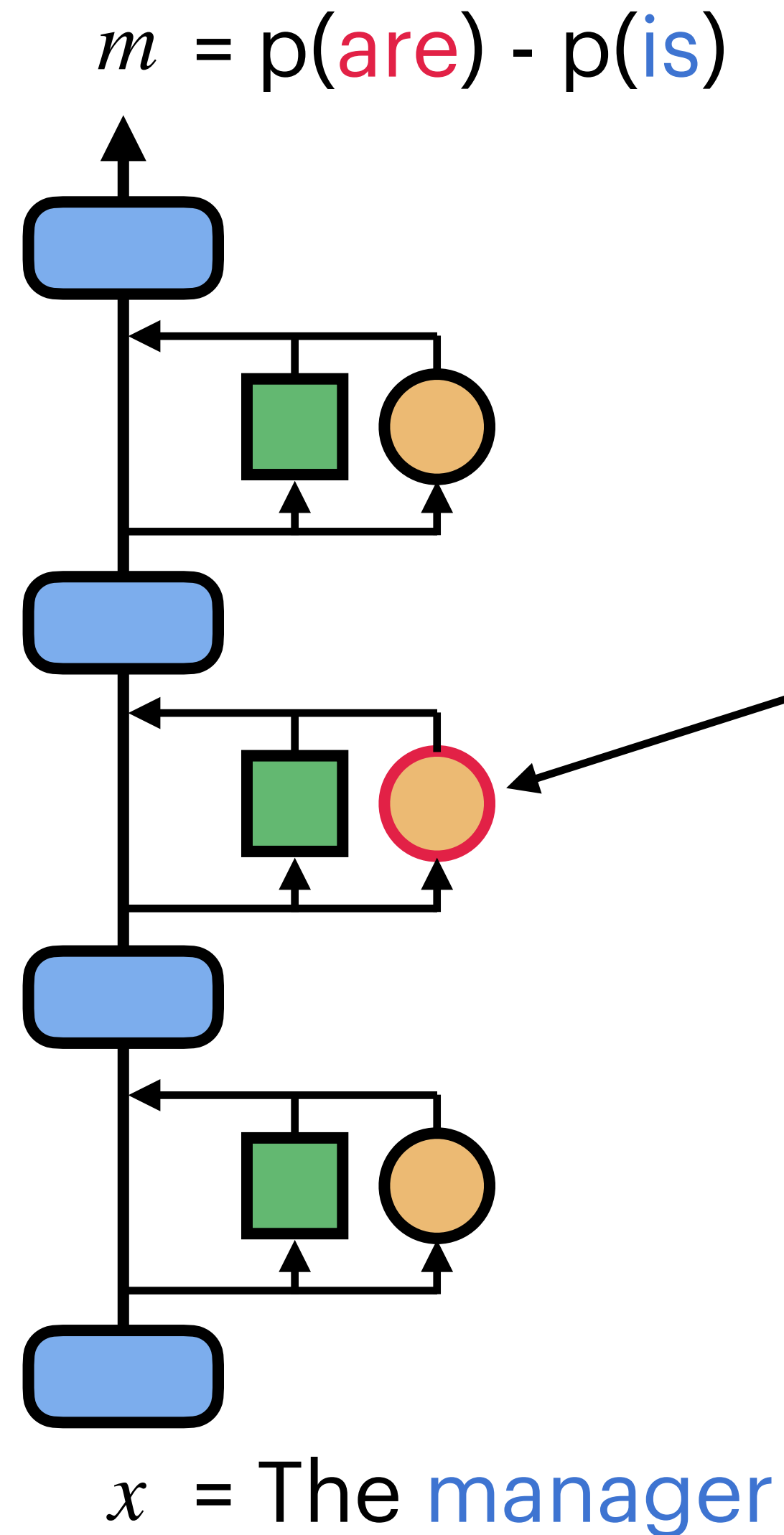
$m$ = p(are) - p(is)

$x'$ = The ~~manager~~ managers

do(swap-number): Set $a$ to what it would have been if the subject in $x$ were the opposite number

**Indirect effect**$(m; a; x, x')$: How much does do(swap-number) change $m$?

$$\mathsf{IE}(m; a; x, x') = m(x, \mathsf{do}(a = a_{x'})) - m(x)$$

$x$ = The manager

attention ■  MLP ●
vector state ▬

# Activation Patching

$m$ = p(are) - p(is)

$x'$ = The ~~manager~~ managers

do(swap-number): Set $a$ to what it would have been
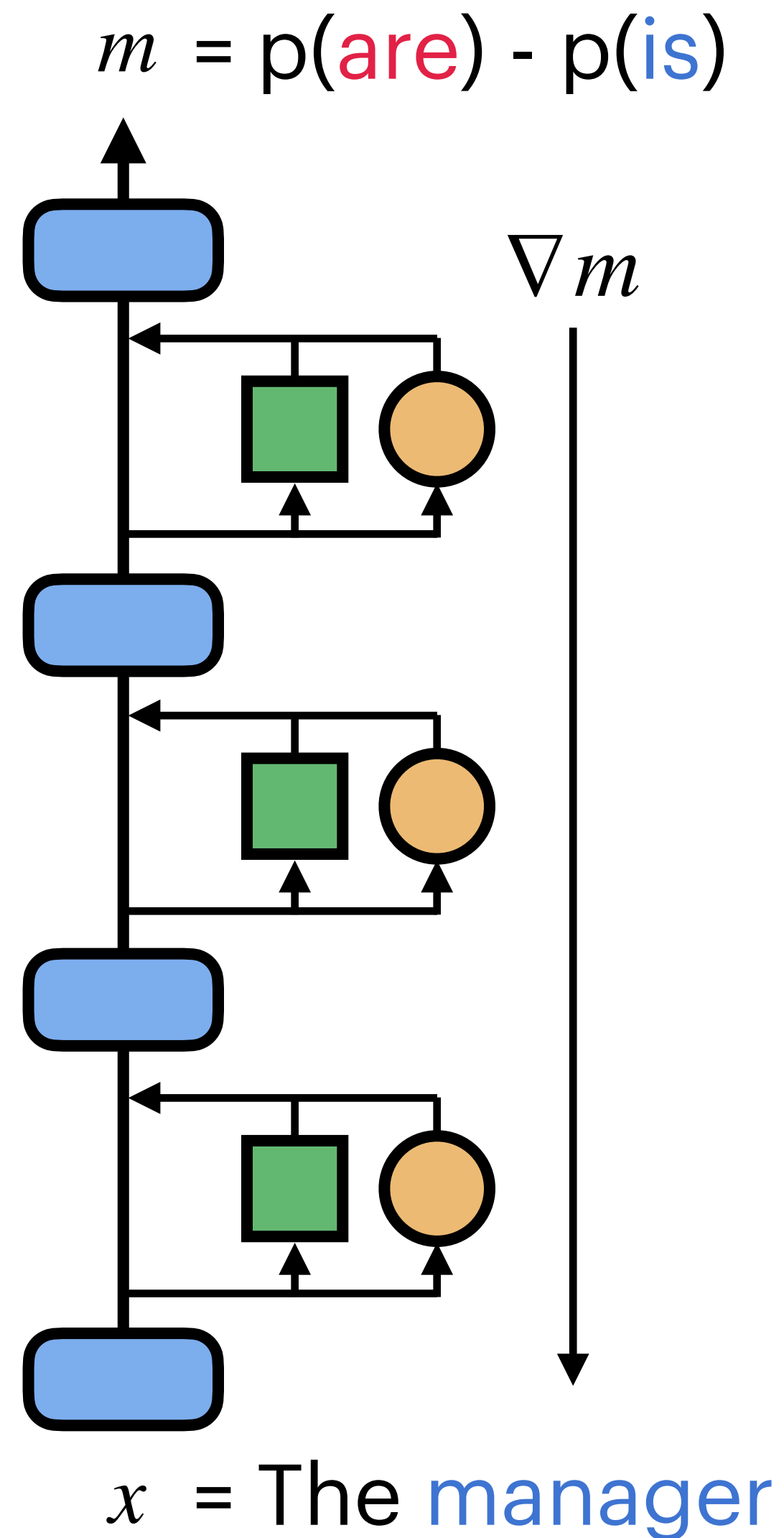if the subject in $x$ were the
opposite number

$x$ = The manager

**Indirect effect**$(m; a; x, x')$: How much does do(swap-number)
change $m$?

$$\text{IE}(m; a; x, x') = m(x, \text{do}(a = a_{x'})) - m(x)$$

*Activation patching requires $O(\mathbf{a})$ forward passes.*
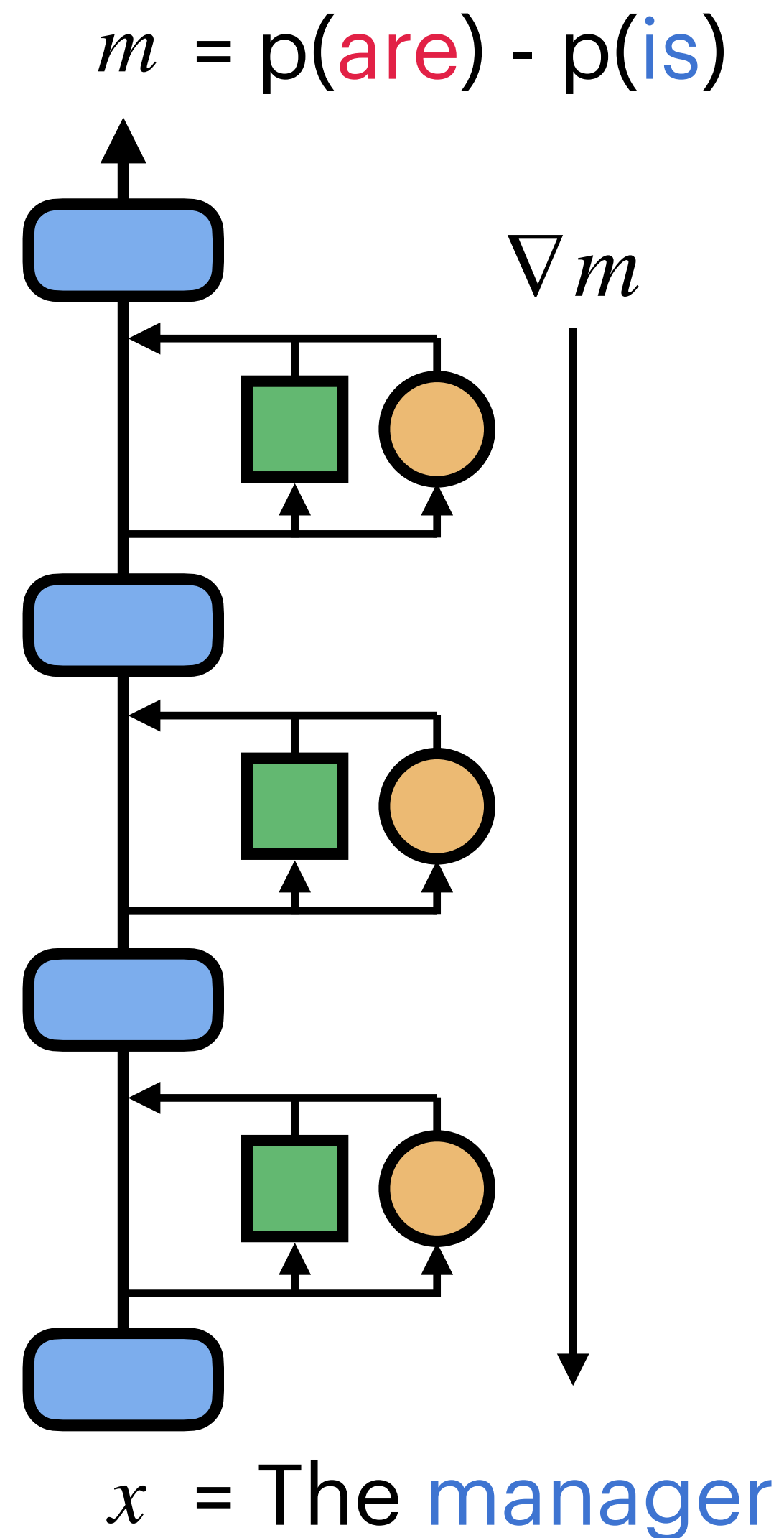
■ attention    ● MLP
▬ vector state

# Attribution Patching

$$\hat{\mathsf{IE}}(m; a; x, x') = \frac{\partial m}{\partial a}\bigg|_x (a_{x'} - a_x)$$

*Attribution patching requires $O(1)$ forward and backward passes!*

# Attribution Patching

$m = \mathrm{p}(\text{are}) - \mathrm{p}(\text{is})$

$\nabla m$

$$\hat{\mathrm{IE}}(m; a; x, x') = \frac{\partial m}{\partial a}\bigg|_x (a_{x'} - a_x)$$

$x$ = The manager

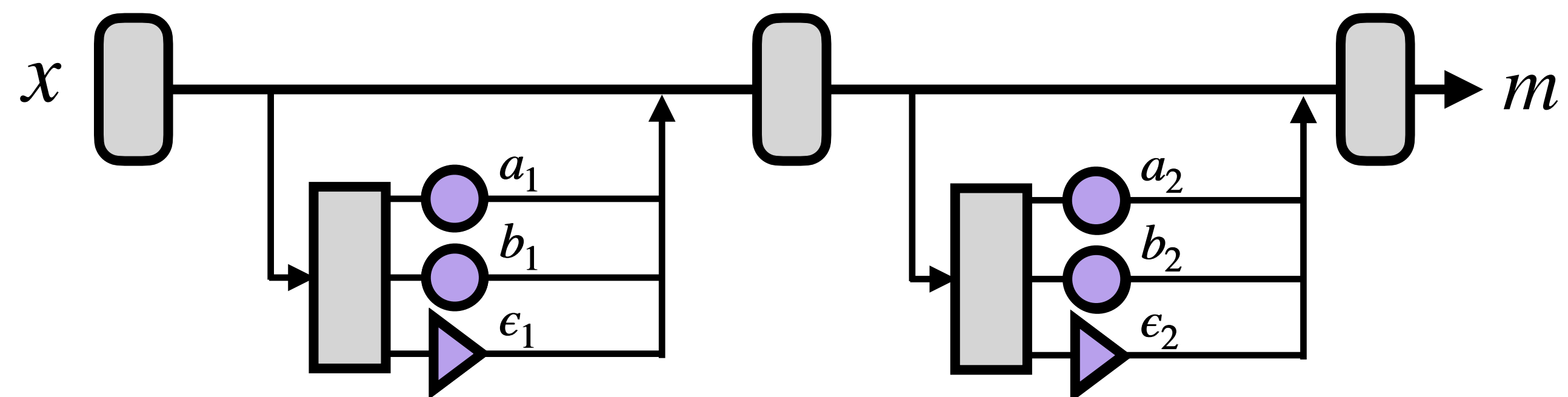*Attribution patching requires $O(1)$ forward and backward passes!*

(We actually propose and use a more accurate approximation based on integrated gradients.)

attention   MLP
vector state

SAE feature

SAE error

Submodule

$x$     $m$

$a_1$   $b_1$   $\epsilon_1$    $a_2$   $b_2$   $\epsilon_2$

**1** Cache activations and metric.

$m$

$a_2$   $b_2$   $\epsilon_2$

$a_1$   $b_1$   $\epsilon_1$

$x$ = The **teacher**

SAE feature

SAE error

Submodule

**1** Cache activations and metric.

$m = \log p(\textbf{have}) - \log p(\textbf{has})$

$x$ = The **teacher**    $x$ = The **teachers**

**Legend:**
- ● SAE feature
- ▲ SAE error
- ▬ Submodule

**1** Cache activations and metric.

$$m = \log p(\textbf{have}) - \log p(\textbf{has})$$

$x$ = The **teacher**      $x$ = The **teachers**

**2** Backpropagate. Store gradients.

$\nabla_{a_1} m$      $\nabla_{\epsilon_1} m$

$\nabla_{a_2} m$      $\nabla_{\epsilon_2} m$

Legend:
- ● SAE feature
- ▲ SAE error
- ▮ Submodule

$x$ → $m$

$a_1$, $b_1$, $\epsilon_1$ / $a_2$, $b_2$, $\epsilon_2$

**1** Cache activations and metric.

$m = \log p(\text{have}) - \log p(\text{has})$

$x$ = The **teacher**   $x$ = The **teachers**

**2** Backpropagate. Store gradients.

$\nabla_{a_1} m$   $\nabla_{\epsilon_1} m$
$\nabla_{a_2} m$   $\nabla_{\epsilon_2} m$

**3** Compute effects. Filter nodes.

$\hat{\text{IE}}(a, m) = \nabla_a m \cdot (\,a\, - \,a\,)$

$\hat{\text{IE}}(a, m) > T_N$

$x$

$m$

SAE feature

SAE error

Submodule

$a_1$ $b_1$ $\epsilon_1$ $a_2$ $b_2$ $\epsilon_2$

**1** Cache activations and metric.

$m = \log p(\textbf{have}) - \log p(\textbf{has})$

$m$

$a_2$ $b_2$ $\epsilon_2$

$a_1$ $b_1$ $\epsilon_1$

$x$ = The **teacher**    $x$ = The **teachers**

**2** Backpropagate. Store gradients.

$\nabla_{a_1} m$    $\nabla_{\epsilon_1} m$

$\nabla_{a_2} m$    $\nabla_{\epsilon_2} m$

$m$

$a_2$ $b_2$ $\epsilon_2$

$a_1$ $b_1$ $\epsilon_1$

**3** Compute effects. Filter nodes.

$m$

$a_2$ $b_2$ $\epsilon_2$

$a_1$ $b_1$ $\epsilon_1$

$\hat{\mathrm{E}}(a, m) = \nabla_a m \cdot ( a - a )$

$\hat{\mathrm{E}}(a, m) > T_N$

**4** Compute and filter edges.

$m$

$a_2$ $b_2$

$b_1$ $\epsilon_1$

# Case Study

**Subject–Verb Agreement**

$$m = p(\textcolor{red}{\text{are}}) - p(\textcolor{blue}{\text{is}})$$

**x =** The **manager** that the parents like

# Case Study

## Subject–Verb Agreement

# Case Study

## Subject–Verb Agreement

has/have

Embeddings, layer 0-4 MLP, resid

Noun number detection

48    8

Layers 2-3 attn, MLP, resid

PP/RC detection

5    4

Verb form discriminators

PP/RC end detection

16    5

Layers 4-5 attn, resid

The    girl/girls    that    the    teacher    sees

# Case Study

## Subject–Verb Agreement



has/have

Embeddings, layer 0-4 MLP, resid

Noun number detection

48    8

Layers 2-3 attn, MLP, resid

PP/RC detection

5    4

Verb form discriminators

PP/RC end detection

16    5

Layers 4-5 attn, resid

The    girl/girls    that    the    teacher    sees

*This corresponds to the human intuition!*

# Case Study

## Subject–Verb Agreement

has/have

Verb form discriminators

Embeddings, layer 0-4 MLP, resid

Layers 2-3 attn, MLP, resid

PP/RC end detection

Noun number detection

PP/RC detection

48 8

5 4

16 5

Layers 4-5 attn, resid

The  girl/girls  that  the  teacher  sees

*This corresponds to the human intuition!*

*But what about cases where it doesn't?*

# Classifying Ambiguous Data

## Bias in Bios

**Task:** classify profession described in biography

# Classifying Ambiguous Data

## Bias in Bios

**Task:** classify profession described in biography

"He was previously an **assistant professor** at the
University of Arizona..."

Professor

0

"She graduated in 2005 with honors, and has 11 years
of experience as a **nurse practitioner**"

Nurse

1

# Classifying Ambiguous Data

## Bias in Bios

**Task:** classify profession described in biography

| | | Man | Professor |
|---|---|---|---|
| "**He** was previously an **assistant professor** at the University of Arizona…" | | 0 | 0 |

| | | Woman | Nurse |
|---|---|---|---|
| "**She** graduated in 2005 with honors, and has 11 years of experience as a **nurse practitioner**" | | 1 | 1 |

*What if the target feature correlates perfectly with the spurious feature?*

# SHIFT

**Method**

Classifier Head

$[p(0), p(1)]$

**Task:** classify profession

**Acc.:**
*Profession* : 63%
*Gender*: 87%

**SHIFT**

**Method**

Classifier Head

$[p(0), p(1)]$

**Task:** classify profession

**Acc.:**

*Profession* : 63%
*Gender*: 87%

Look for features with high
IE on classifier logits

👍

Inspect each high-IE feature

Matt Vera is a registered nurse with a bachelor of science in nursing since 2009 and is currently working as a full-time writer and editor for

two Registered Nurses to work on a day or night shift . The nursing home has easy access to public transport Tub ... ↵
full job description ↵
↵

with other students and faculty . ↵
↵
But for many of the most popular nursing programs the online environment is not a complete solution . For one thing any nurse

fier Head

), p(1)]

**Task:** classify profession

**Acc.:**

*Profession* : 63%
*Gender:* 87%

bodies for calf re aring . ↵

↵

It features daily videos of Nicole and Alice ,
along with a few other farmers , doing warm
ups , stretches and strengthening

the marriage was failing . Paul suffered
engulf ing dep ressions . Sometimes he and
Angela barely spoke for days . She felt
swollen with un expressed emotion . " I

It was like a bitter taste , just a foul taste
,' he said âG¦ Mary Cel este Clement , a
children ' s book author , lives about 2 miles

At rium at age 13 and that he was preceded in
death by his wife Sarah , who rests next to
him . ↵

↵

**Classifier Head**

$[p(0), p(1)]$

Inspect each high-IE feature

bodies for calf re aring . ↵

↵

It features daily videos of Nicole and Alice , along with a few other farmers , doing warm ups , stretches and strengthening

the marriage was failing . Paul suffered engulf ing dep ressions . Sometimes he and Angela barely spoke for days . She felt swollen with un expressed emotion . " I

It was like a bitter taste , just a foul taste ,' he said âǦ Mary Cel este Clement , a children ' s book author , lives about 2 miles

At rium at age 13 and that he was preceded in death by his wife Sarah , who rests next to him . ↵

↵

**Classifier Head**

[p(0), p(1)]

Inspect each high-IE feature

Ablate features that seem related to *gender*

# SHIFT

**Results**

| Method | Pythia-70M | | | Gemma-2-2B | | |
|---|---|---|---|---|---|---|
| | ↑Profession | ↓Gender | ↑Worst group | ↑Profession | ↓Gender | ↑Worst group |
| Original | 61.9 | 87.4 | 24.4 | 67.7 | 81.9 | 18.2 |
| 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** |

# SHIFT

## Results

| Method | Pythia-70M | | | Gemma-2-2B | | |
|---|---|---|---|---|---|---|
| | ↑Profession | ↓Gender | ↑Worst group | ↑Profession | ↓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** |

# SHIFT

## Results

| Method | Pythia-70M | | | Gemma-2-2B | | |
|---|---|---|---|---|---|---|
| | ↑Profession | ↓Gender | ↑Worst group | ↑Profession | ↓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 |

*Features are a stronger basis than neurons for removing spurious correlations.*

# SHIFT

## Results

| Method | Pythia-70M | | | Gemma-2-2B | | |
|---|---|---|---|---|---|---|
| | ↑Profession | ↓Gender | ↑Worst group | ↑Profession | ↓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 |

*Features are a stronger basis than neurons for removing spurious correlations.*
*Our judgments about feature relevance are largely informative.*

# SHIFT

## Results

| Method | Pythia-70M | | | Gemma-2-2B | | |
|---|---|---|---|---|---|---|
| | ↑Profession | ↓Gender | ↑Worst group | ↑Profession | ↓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 |

*Features are a stronger basis than neurons for removing spurious correlations.*

*Our judgments about feature relevance are largely informative.*

*SHIFT achieve the performance of a classifier trained on **unbiased** data!*

# An Unsupervised Interpretability Pipeline

Intepretability typically requires us to have a behavior in mind.

*Can we fully automate the behavior and circuit discovery process?*

# An Unsupervised Interpretability Pipeline

Intepretability typically requires us to have a behavior in mind.

*Can we fully automate the behavior and circuit discovery process?*

1. Given large text corpus $\{(x_i, y_i)\}$, collect activations of SAEs $\mathbf{v}(x_i, y_i)$

2. Cluster $\mathbf{v}$

3. Discover sparse feature circuits on clusters

# An Unsupervised Interpretability Pipeline

## Results

**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 […]
Step 4

**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

*This yields not only interesting unanticipated behaviors...*

# An Unsupervised Interpretability Pipeline

## Results

**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 […]
Step 4

Example features involved:

Succession

| Chapter 1 Chapter 2 Chapter 3 | A, B, C |
| | I, II, III, IV |

Narrow induction

A3 … A → 3 or III or 4 …

A7 … A → 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

Direct the user to     It's up to you to

Other words which precede
infinitive objects

According to     This infection leads to

*This yields not only interesting unanticipated behaviors...*

*...But also interesting unanticipated features!*

# Takeaways

1. Sparse feature circuits allow us to derive **human-interpretable** and **editable** causal graphs from LMs.

2. They allow us to surgically improve model generalization *without additional data*.

3. They allow us to automatically discover **unanticipated** model behaviors and mechanisms.
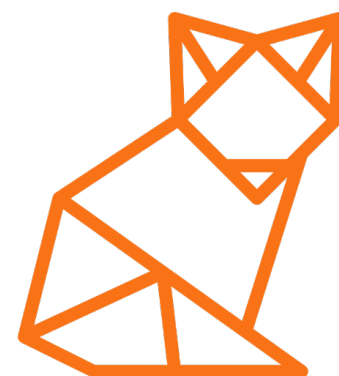
# Thank you!

🌐 **Project Website**

**Check out our poster:**

Today 3pm – 5:30pm
Poster #495

Open Philanthropy

NSF

ISRAEL SCIENCE FOUNDATION

Fondation Azrieli Foundation

MORTIMER B. Zuckerman STEM Leadership Program

# Sparse Feature Circuits

## Discovering and Editing Interpretable Causal Graphs in Language Models

Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, **Aaron Mueller**
2025 International Conference on Machine Learning
26 April 2025