

# Mechanistic Interpretability: An Introduction

12th Feb, Foundation Models Reading Group

Praveen Selvaraj

# What is Mechanistic Interpretability ?

# What is Mechanistic Interpretability ?

Reverse engineering Neural Networks. Go from trained parameters -> features, algorithms.

# What is Mechanistic Interpretability ?

Reverse engineering Neural Networks. Go from trained parameters -> features, algorithms.

“These things are totally different from us,” he says. “Sometimes I think it’s as if aliens had landed and people haven’t realized because they speak very good English.” - Hinton

# Analogy - computer programs vs neural networks



Planning documents

```
// In this supposed to happen with
// (settings.one) {
//   // remove the check
//   $unbind("appear", check);
//   if (i >= 0) $fn.appear.check;
// }
```

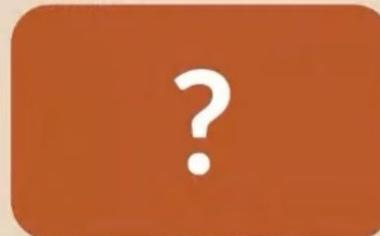
Source code

```
011101001100010101110011
1110011001100010101110011
11000100000110111010101000
110010001111110011111000
101001111010010010010010001
01101101011010000011110
011101111100101110
0001011111100101110
011111001011110
```

Compiled application



Training Objectives



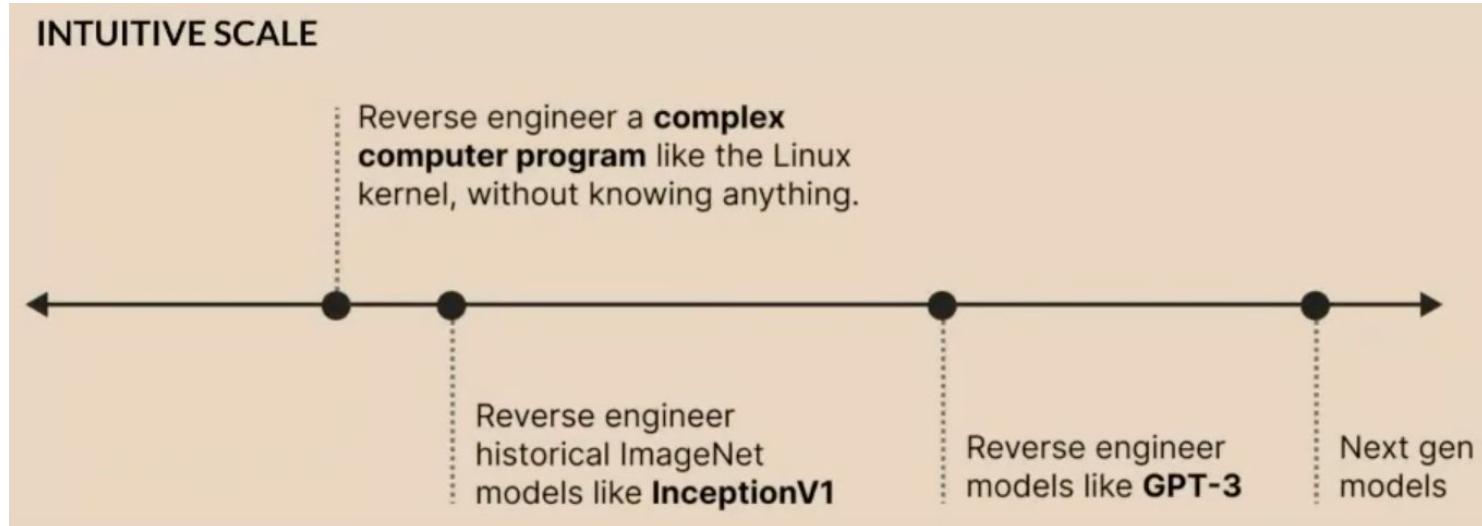
```
011101001100010101110011
1110011001100010101110011
11000100000110111010101000
110010001111110011111000
101001111010010010010010001
01101101011010000011110
011101111100101110
0001011111100101110
011111001011110
```

Trained Neural Network

# Analogy - computer programs vs neural networks

Computer programs	Neural networks
Variable	Neuron / Direction
Program state	Activations
Processor	Network architecture
Compiled binary	Network parameters
Source code	???

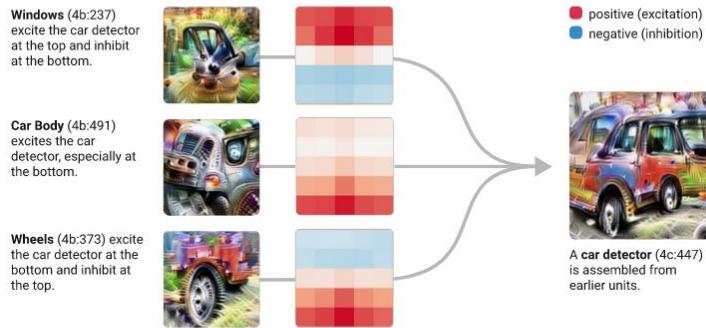
# Scale of difficulty



# Previous work

## Zoom In: An Introduction to Circuits

By studying the connections between neurons, we can find meaningful algorithms in the weights of neural networks.



AUTHORS	AFFILIATIONS	PUBLISHED	DOI
Chris Olah	OpenAI		
Nick Cammarata	OpenAI		
Ludwig Schubert	OpenAI		
Gabriel Goh	OpenAI		
Michael Petrov	OpenAI		
Shan Carter	OpenAI		

[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# The 3 Claims

- Features - fundamental unit of NNs. Correspond to ‘directions’ in activation space.

# The 3 Claims

- Features - fundamental unit of NNs. Correspond to ‘directions’ in activation space.
- Circuits - features connected by weights, form circuits.

# The 3 Claims

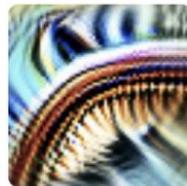
- Features - fundamental unit of NNs. Correspond to ‘directions’ in activation space.
- Circuits - features connected by weights, form circuits.
- Universality - similar features and circuits exist across models.

# The 3 Claims

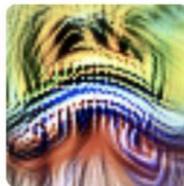
- Features - fundamental unit of NNs. Correspond to ‘directions’ in activation space.
- Circuits - features connected by weights, form circuits.
- Universality - similar features and circuits exist across models.

# Curve Detectors (features)

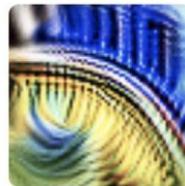
Curves



3b:379



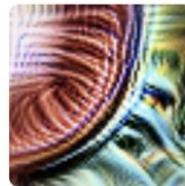
3b:406



3b:385



3b:343



3b:342



3b:388



3b:340



3b:330



3b:349



3b:324



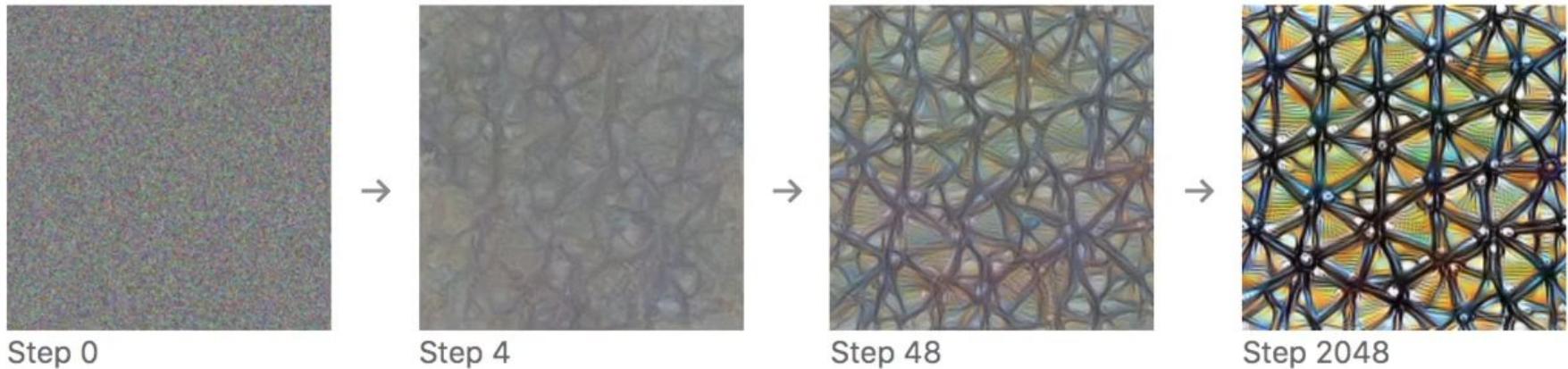
# Arguments for Curve Detectors



Arg: Feature Visualization

[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# Arguments for Curve Detectors



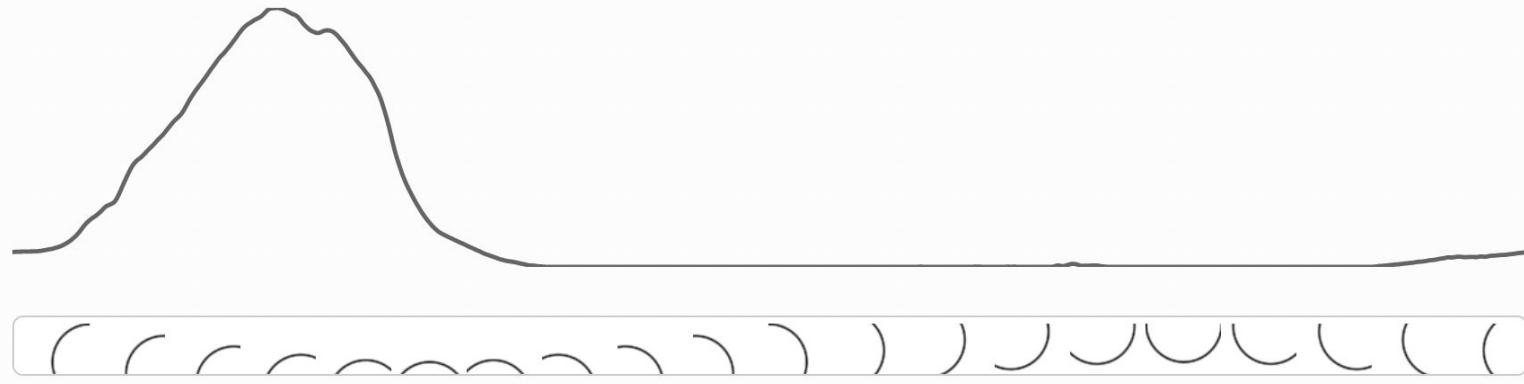
Arg: Feature Visualization

[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# Arguments for Curve Detectors



3b:379 Activations by Orientation



Arg: Feature Visualization

[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# Arguments for Curve Detectors



The images from the dataset that activate 3b:379 all contain curves that are similar to its ideal curve.

## Arg: Dataset Examples

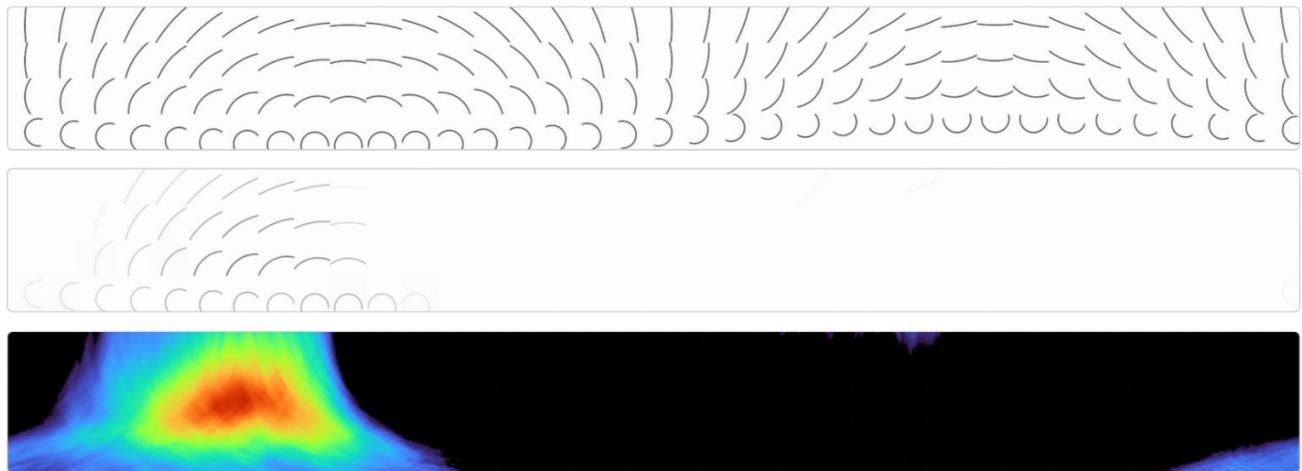
[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# Arguments for Curve Detectors

Creating synthetic stimuli, like these curves, can be helpful for probing the behavior of neurons.

We can look at which stimuli cause 3b:379 to fire.

Looking at a heatmap give us a much higher resolution view.

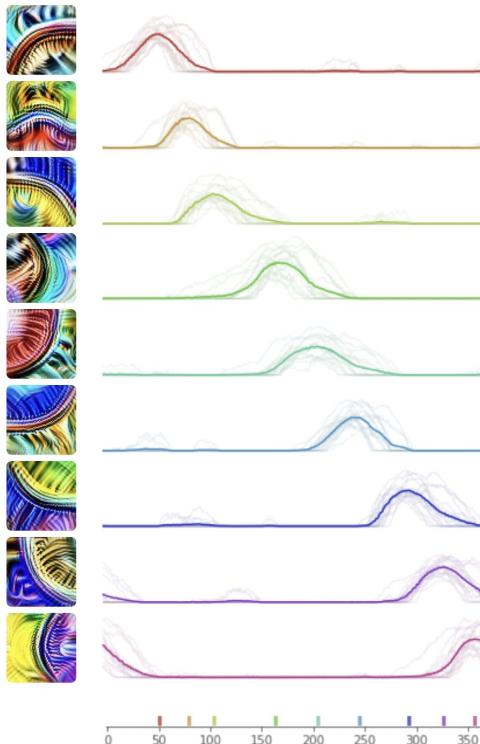


## Arg: Synthetic Examples

[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# Arguments for Curve Detectors

Arg: Joint Tuning



We collect dataset examples that maximally activate neuron. We rotate them by increments of 1 degree from 0 to 360 degrees and record activations.

The activations are shifted so that the points where each neuron responds are aligned. The curves are then averaged to create a typical response curve.



# Arguments for Curve Detectors



Arg: Feature Use

[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# Arguments for Curve Detectors



Arg: Handwritten Circuits

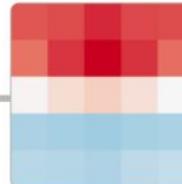
[\(2020\) Zoom In: An Intro to Circuits, Olah et. al](#)

# The 3 Claims

- Features - fundamental unit of NNs. Correspond to ‘directions’ in activation space.
- Circuits - features connected by weights, form circuits.
- Universality - similar features and circuits exist across models & tasks.

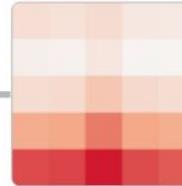
# Circuits in ConvNets

**Windows** (4b:237) excite the car detector at the top and inhibit at the bottom.

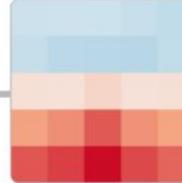


- positive (excitation)
- negative (inhibition)

**Car Body** (4b:491) excites the car detector, especially at the bottom.



**Wheels** (4b:373) excite the car detector at the bottom and inhibit at the top.



A **car detector** (4c:447) is assembled from earlier units.

# The 3 Claims

- Features - fundamental unit of NNs. Correspond to ‘directions’ in activation space.
- Circuits - features connected by weights, form circuits.
- Universality - similar features and circuits exist across models & tasks.

# Universality of Features & Circuits

Curve detectors

ALEXNET

Krizhevsky et al. [34]



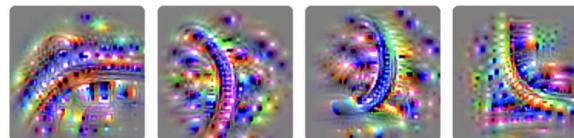
INCEPTIONV1

Szegedy et al. [26]



VGG19

Simonyan et al. [35]

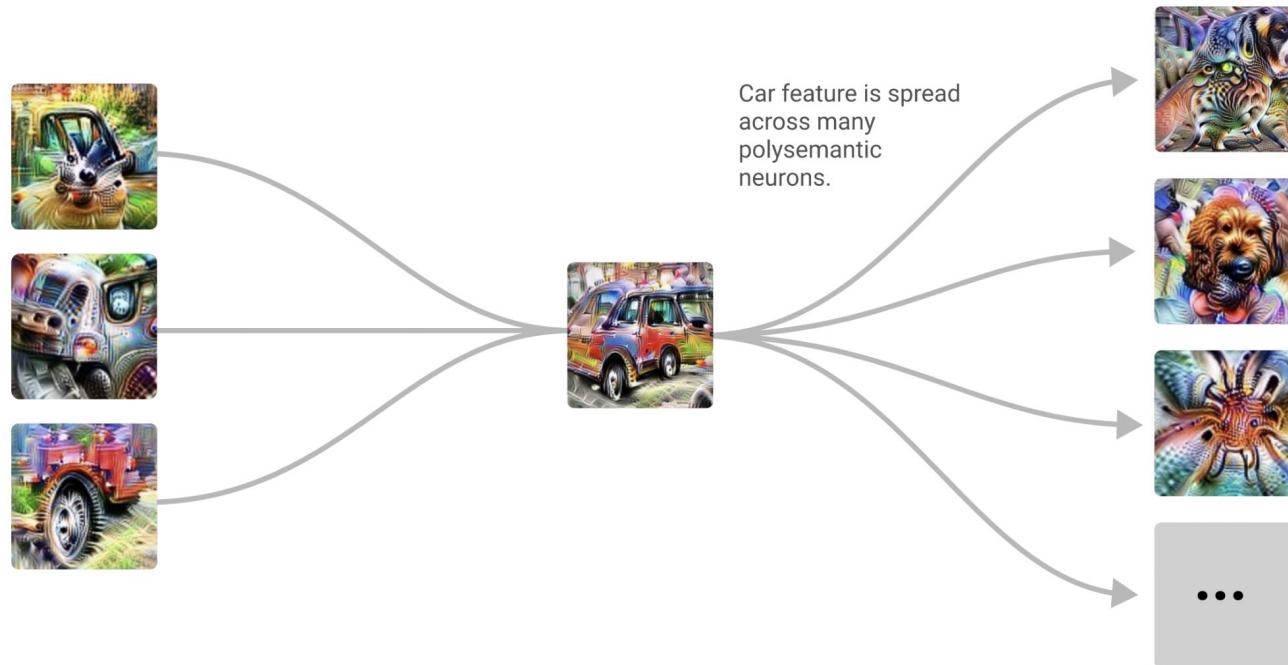


RESNETV2-50

He et al. [36]



# Major Hurdle: polysemantic neurons



# A Mathematical Framework for Transformer Circuits

---

## AUTHORS

Nelson Elhage<sup>\*†</sup>, Neel Nanda<sup>\*</sup>, Catherine Olsson<sup>\*</sup>, Tom Henighan<sup>†</sup>, Nicholas Joseph<sup>†</sup>, Ben Mann<sup>†</sup>, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, Chris Olah<sup>‡</sup>

## AFFILIATION

Anthropic

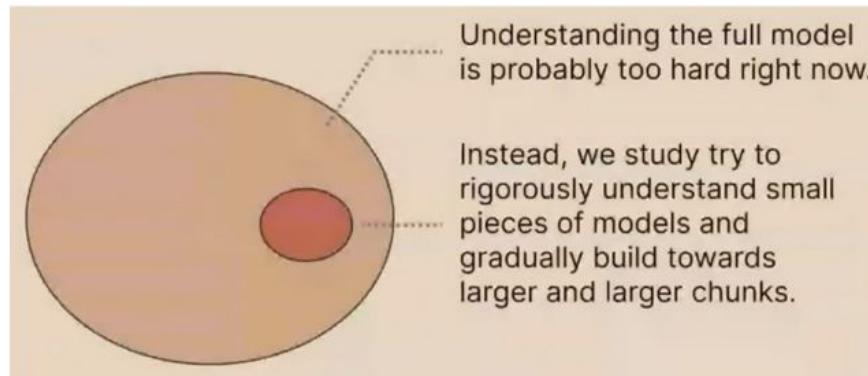
## PUBLISHED

Dec 22, 2021

\* Core Research Contributor; † Core Infrastructure Contributor; ‡ Correspondence to colah@anthropic.com;  
Author contributions statement below.

---

# (2021) A Math Framework for Transformer Circuits

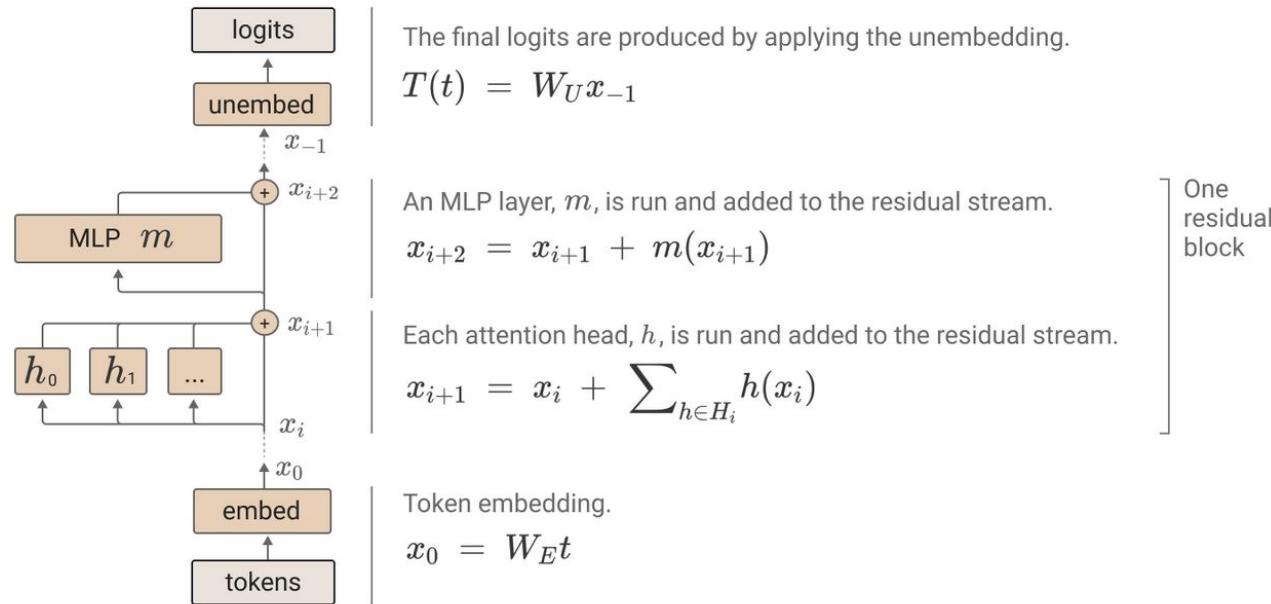


[\(2023\) Looking inside NNs, Chris Olah](#)

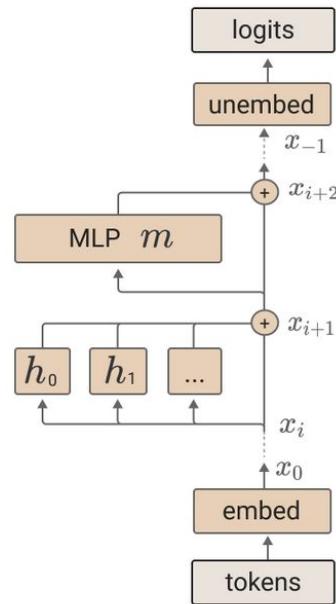
# (2021) A Math Framework for Transformer Circuits

- Zero layer transformer
- One layer attention-only transformer
- Two layer attention-only transformer

# High-level architecture



# High-level architecture



The final logits are produced by applying the unembedding.

$$T(t) = W_U x_{-1}$$

An MLP layer,  $m$ , is run and added to the residual stream.

$$x_{i+2} = x_{i+1} + m(x_{i+1})$$

Each attention head,  $h$ , is run and added to the residual stream.

$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

Token embedding.

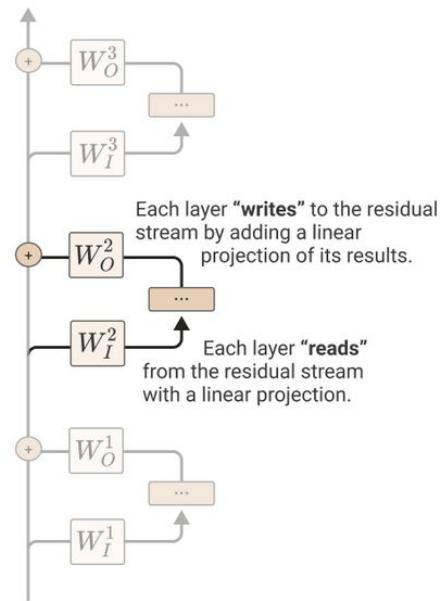
$$x_0 = W_E t$$

One residual block



# Layers reading and writing

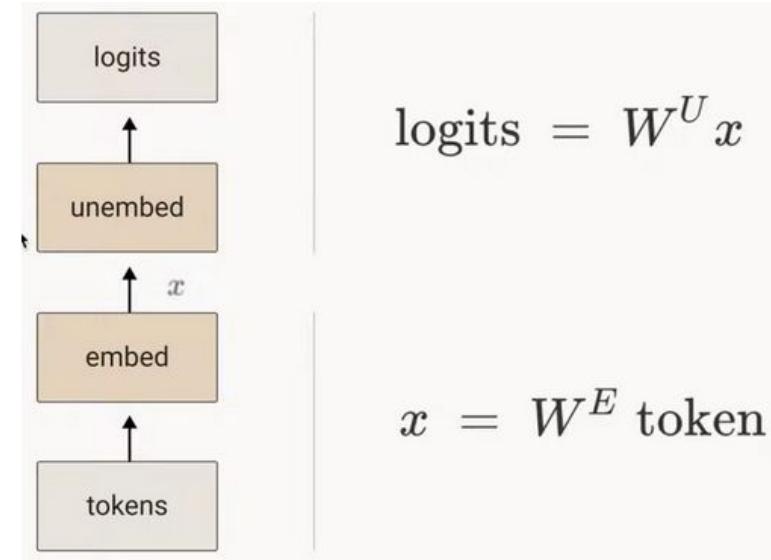
The residual stream is modified by a sequence of MLP and attention layers “reading from” and “writing to” it with linear operations.



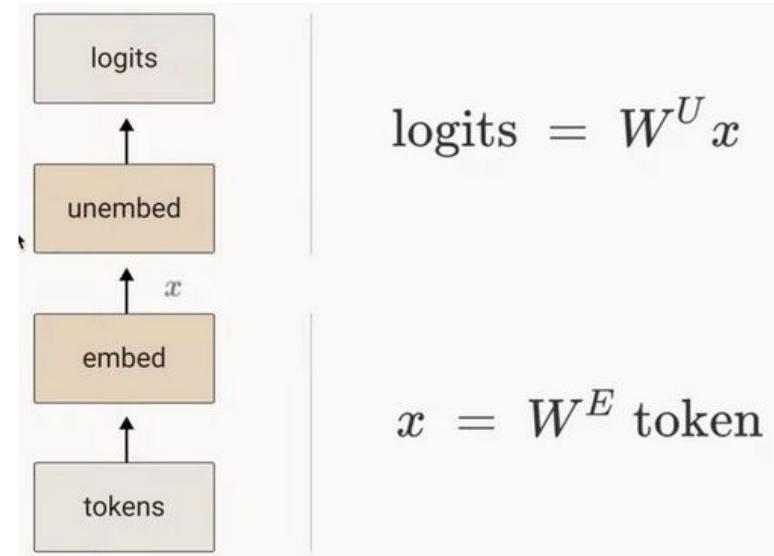
# GPT2 example

- Residual stream dimension: 768
- Attention layer dimension: 64 / 128
- MLP layer dimension: 3072
- 12 transformer blocks, 12 attention heads per block

# Zero-layer transformer



# Zero-layer transformer



Models bigram stats, for e.g. probability of seeing ‘Obama’ would be very high after seeing ‘Barrack’.

# Attention Eqn

A typical definition of an attention head might look something like:

$$h(x)_i = W_O \left( \sum_t A_{i,j} W_V x_j \right)$$

# Attention Eqn

A typical definition of an attention head might look something like:

$$h(x)_i = W_O \left( \sum_t A_{i,j} W_V x_j \right)$$

Mathematical tensors give us a cleaner way to describe this:

$$\begin{aligned} h(x) &= (\text{Id} \otimes W_O) \cdot (A \otimes \text{Id}) \cdot (\text{Id} \otimes W_V) \cdot x \\ &= (A \otimes W_O W_V) \cdot x \end{aligned}$$



Multiplication  
across positions



Multiplication of vector  
at each position

# Tensor product rewrite

$$h(x) = (A \otimes W_O W_V) \cdot x$$

---

$A$  mixes across tokens while  
 $W_O W_V$  acts on each vector  
independently.

- $A$  controls which token's information is moved from and to

# Tensor product rewrite

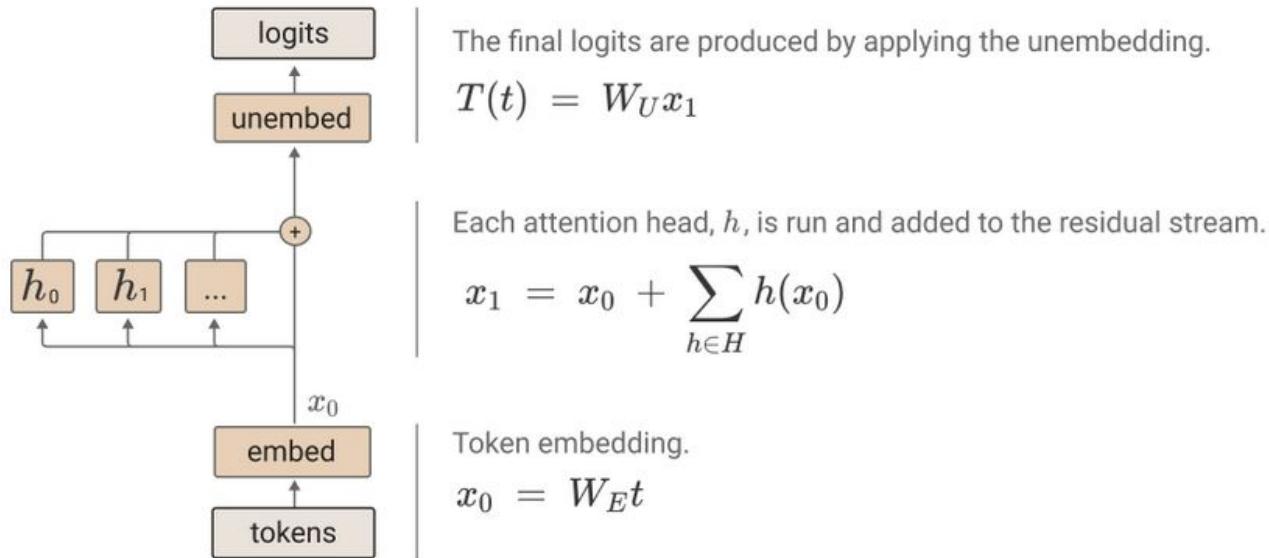
$$h(x) = (A \otimes W_O W_V) \cdot x$$

---

$A$  mixes across tokens while  
 $W_O W_V$  acts on each vector  
independently.

- $A$  controls which token's information is moved from and to
- $W_O W_V$  controls what information is read from a source token and written to the destination token

# One-layer transformer



# One-layer transformer

$$T = \underbrace{\text{Id} \otimes W_U W_E}_{\text{---}} + \sum_{h \in H} A^h \otimes (W_U W_{OV}^h W_E)$$

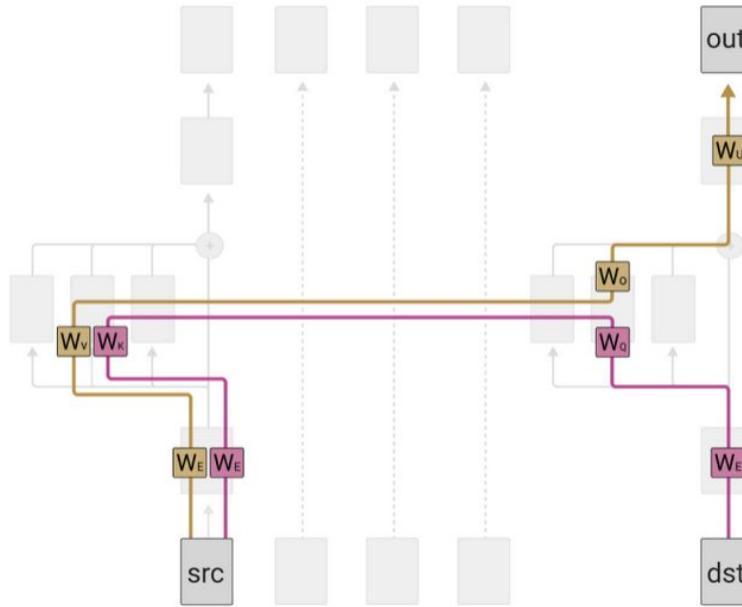


"Direct path" term contributes to bigram statistics.



The **attention head** terms describe the effects of attention heads in linking input tokens to logits.  $A^h$  describes which tokens are attended to while  $W_U W_{OV}^h W_E$  describes how each token changes the logits if attended to.

# One-layer transformer



The **OV** (“output-value”) circuit determines how attending to a given token affects the logits.

$$W_U W_O W_V W_E$$

The **QK** (“query-key”) circuit controls which tokens the head prefers to attend to.

$$W_E^T W_Q^T W_K W_E$$

# Skip Tri-grams

Some examples of large entries QK/OV circuit

Source Token	Destination Token	Out Token	Example Skip Tri-grams
"perfect"	"are", "looks", "is", "provides"	"perfect", "super", "absolute", "pure"	"perfect... are perfect", "perfect... looks super"
"large"	"contains", "using", "specify", "contain"	"large", "small", "very", "huge"	"large... using large", "large... contains small"
"two"	"One", "\n", "has", \r\n, "One"	"two", "three", "four", "five", "one"	"two... One two", "two... has three"
"lambda"	"\$\\"", "}{\"", "+\"", "({\"", "\${\""	"lambda", "sorted", "lambda", "operator"	"lambda... \$\lambda", "lambda... +\lambda"
"nbsp"	"&", "\&", "&", >&, "="	"nbsp", "01", "gt", "00012", "nbs", "quot"	"nbsp...  ", "nbsp... > "
"Great"	"The", "The", "the", "contains", "/"	"Great", "great", "poor", "Every"	"Great... The Great", "Great... the great"

In the above example, we fix a given source token and look at the largest corresponding QK entries (the destination token) and largest corresponding OV entries (the out token). The source token is selected to show interesting behavior, but the destination and out token are the top entries unless entries are explicitly skipped with an ellipsis; they are colored by the intensity of their value in the matrix.

# Skip Tri-grams

More examples of large entries QK/OV circuit

Source Token	Destination Token	Out Token	Example Skip Tri-grams
"indy"	"C", "C", "V", "V", "R", "c"	"indy", "obby", "INDY", "loyd"	"indy... Cindy", "indy... CINDY"
" Pike"	"P", "P", "V", "Sp", "V", "R"	"ike", "ikes", "ishing", "owler"	" Pike... Pike", " Pike... Spikes"
" Ralph"	"R", "R", "P", "P", "V", "r"	"alph", "ALPH", "obby", "erald"	" Ralph... Ralph", " Ralph... RALPH"
" Lloyd"	"L", "L", "P", "P", "R", "C"	"loyd", "alph", "\n ", "acman", ... "atherine"	" Lloyd... Lloyd", " Lloyd... Catherine"
" Pixmap"	"P", "Q", "P", "p", "U"	"ixmap", "Canvas", "Embed", "grade"	" Pixmap... Pixmap", " Pixmap... QCanvas"

# General patterns

## Primitive In-Context Learning Patterns

---

**[b]...[a] → [b]**

[ two]...[ One] → [ two]

[ perfect]...[ are] → [ perfect]

[nbsp]...[ &] → [nbsp]

[lambda]...[ \$\\] → [lambda]

**[b]...[a] → [b']**

[ two]...[ has] → [ three]

[ perfect]...[ looks] → [ super]

[nbsp]...[ &] → [gt]

[lambda]...[ \$\\] → [operator]

**[ab]...[a] → [b]**

[Ralph]...[ R] → [alph]

[Pike]...[ P] → [ike]

[Pixmap]...[ P] → [ixmap]

[ Lloyd]...[ L] → [loyd]

**[ab]...[a] → [b']**

[Ralph]...[ R] → [ALPH]

[Pike]...[ P] → [ikes]

# Detecting copying heads

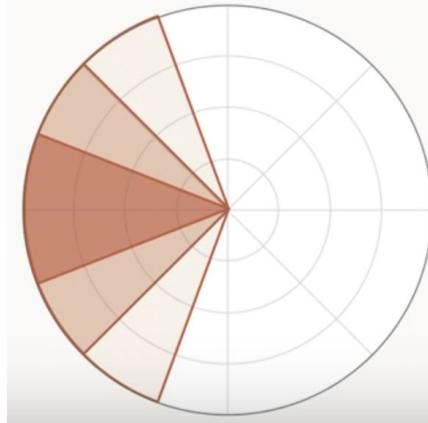
**Eigenvalues** and **eigenvectors** are a useful tool when one has maps from the same vector space onto itself

$$Wv = \lambda v$$

# Detecting copying heads

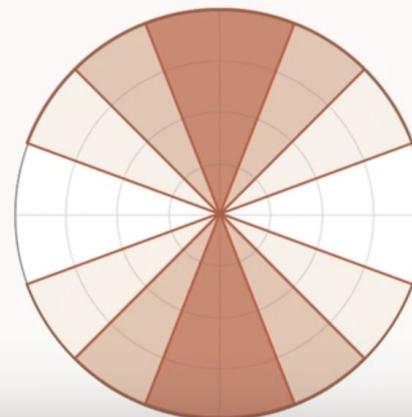
## Negative: Anti-Copying

Negative OV circuit eigenvalues mean that some tokens decrease the probability of the same token being the output.



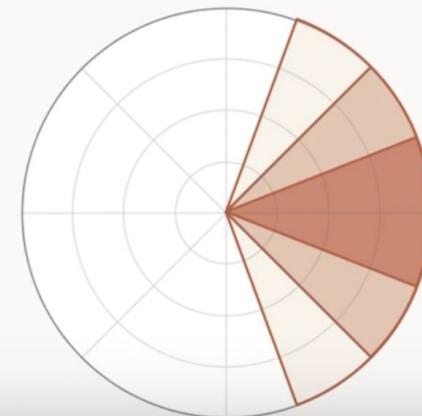
## Imaginary: Different Tokens

Imaginary OV circuit eigenvalues mean that some tokens affect the probability of different tokens in the output.



## Positive: Copying

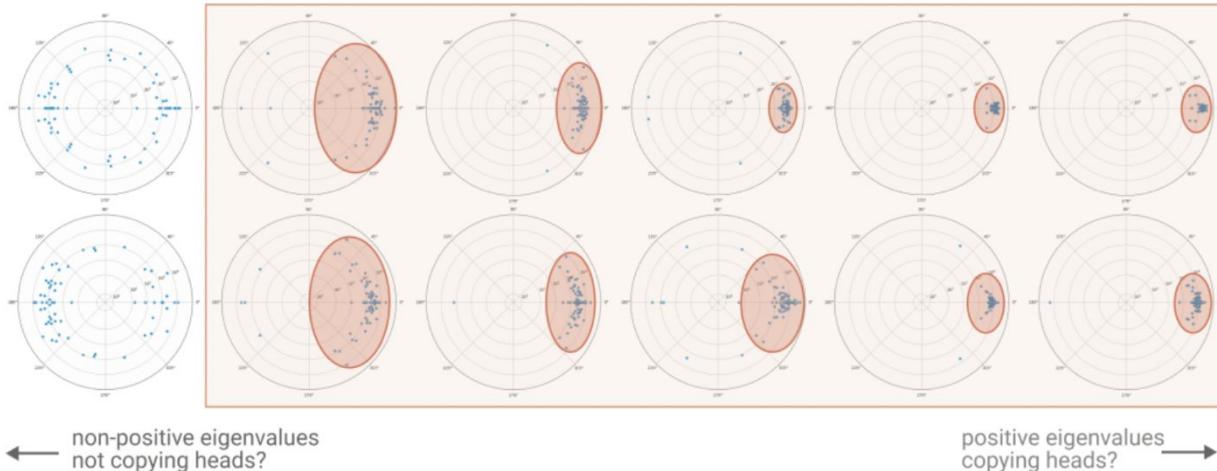
Positive OV circuit eigenvalues mean that some tokens increase the probability of the same token being the output.



# Detecting copying heads

Eigenvalue analysis of **first layer** attention head OV circuits

10/12 of layer 1 heads have mostly positive OV eigenvalues, and appear to significantly perform copying

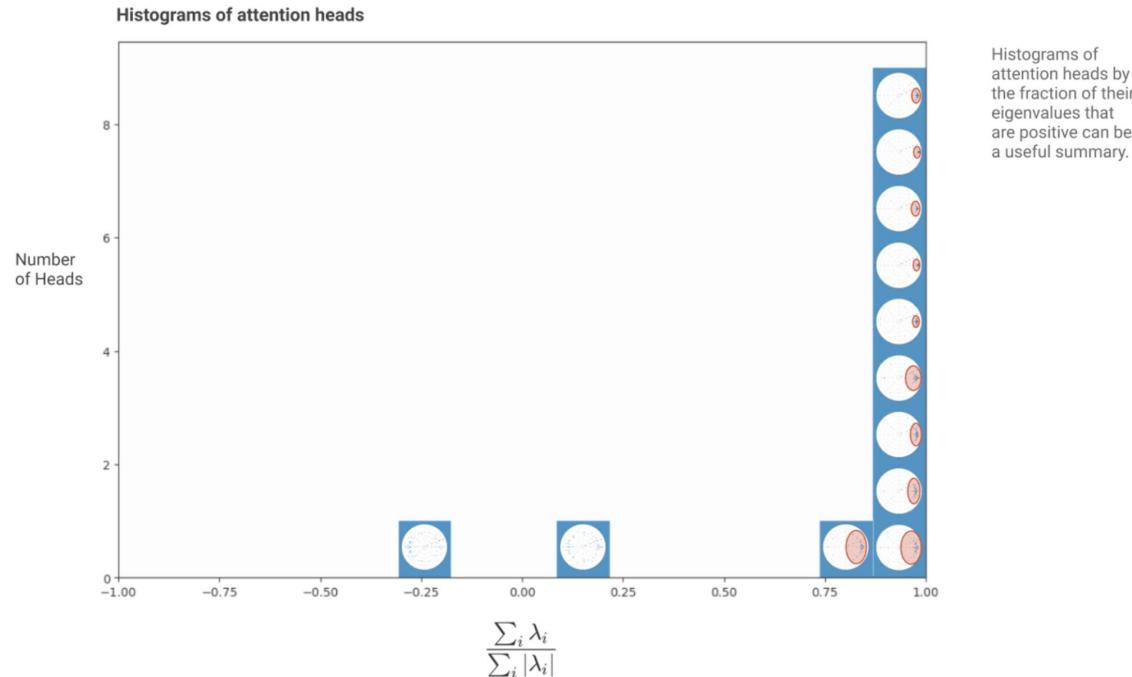


We use a **log scale** to represent magnitude, since it varies by many orders of magnitude.

**Eigenvalue distribution for randomly initialized weights.** Note that the mostly – and in some cases, entirely– positive eigenvalues we observe are very different from what we randomly expect.

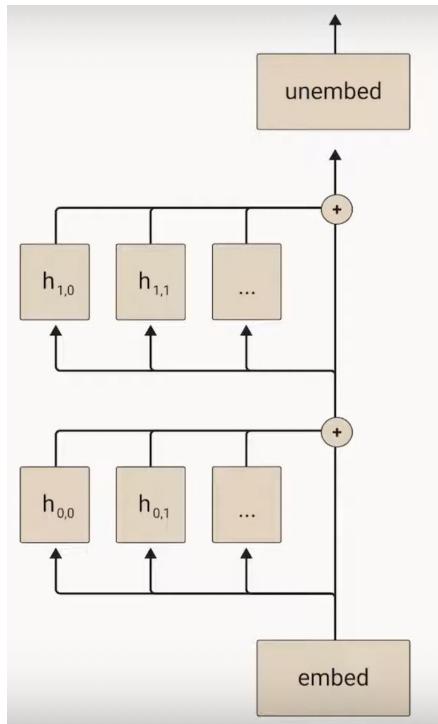


# Detecting copying heads



[\(2021\) A Mathematical Framework for Transformer Circuits](#)

# Two-layer transformer



$$\begin{aligned} T = & W_U \cdot \left( \text{Id} + \sum_{h \in H_2} A^h \otimes W_{OV}^h \right) \\ & \cdot \left( \text{Id} + \sum_{h \in H_1} A^h \otimes W_{OV}^h \right) \cdot W_E \end{aligned}$$

## Two-layer transformer

Recall that tensor products obey the “mixed product” identity:

$$(A_2 \otimes W_2) \cdot (A_1 \otimes W_1) = (A_2 \cdot A_1) \otimes (W_2 \cdot W_1)$$

# Two-layer transformer

$$T = W_U W_E + \sum_{h \in H_1 \cup H_2} A^h \otimes (W_U W_{OV}^h W_E) + \sum_{h_2 \in H_2} \sum_{h_1 \in H_1} (A^{h_2} A^{h_1}) \otimes (W_U W_{OV}^{h_2} W_{OV}^{h_1} W_E)$$

**“Direct path”**  
term tends  
to represent  
bigram  
statistics.

The **individual attention head** terms  
describe the effects of individual  
attention heads in linking input  
tokens to logits, similar to those we  
saw in the one layer model.

The **virtual attention head** terms correspond to  
compositions of attention heads, but function a  
lot like normal attention heads. They have their  
own attention patterns (the composition of the  
heads patterns) and own OV circuits.

# Composition of heads

Example 1

Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the

 Present Token

 Attention

Example 2

the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they

# Composition of heads - induction

Induction Head - Example 1

Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the  
Mr and Mrs Dursley , of ... such nonsense. Mr Dursley was the

- Present Token
- Attention
- Logit Effect

Induction Head - Example 2

the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they

# Composition of heads - induction

out about the Potters. Mrs Potter was ... neighbours would say if the **Potters** arrived in

**attention pattern moves information**

out about the Potters. Mrs Potter was ... neighbours would say if the **Potters** arrived in

key

out about the **Potters**. Mrs Potter was ... neighbours would say if the **Potters** arrived in

query

logit effect

Mr and Mrs Dursley, of number ... with such nonsense. Mr **Dursley** was the

**attention pattern moves information**

Mr and Mrs **Dursley**, of number ... with such nonsense. Mr **Dursley** was the

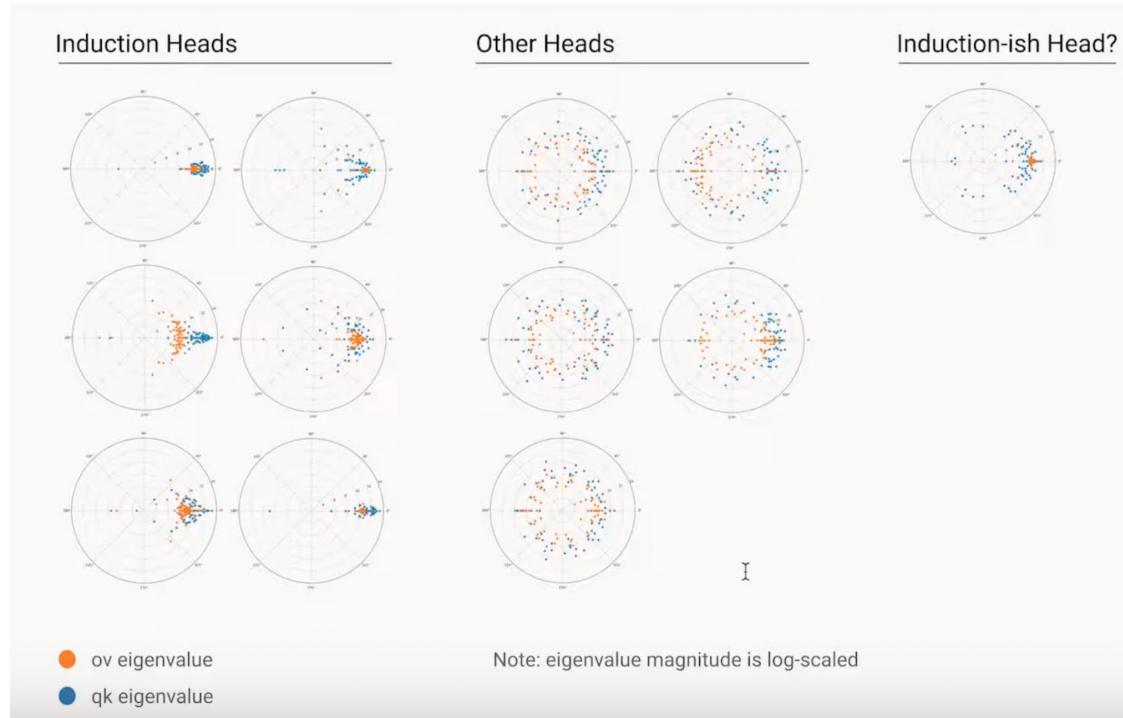
key

Mr and **Mrs Dursley**, of number ... with such nonsense. **Mr Dursley** was the

query

logit effect

# Detecting induction heads



[\(2021\) A Mathematical Framework for Transformer Circuits](#)

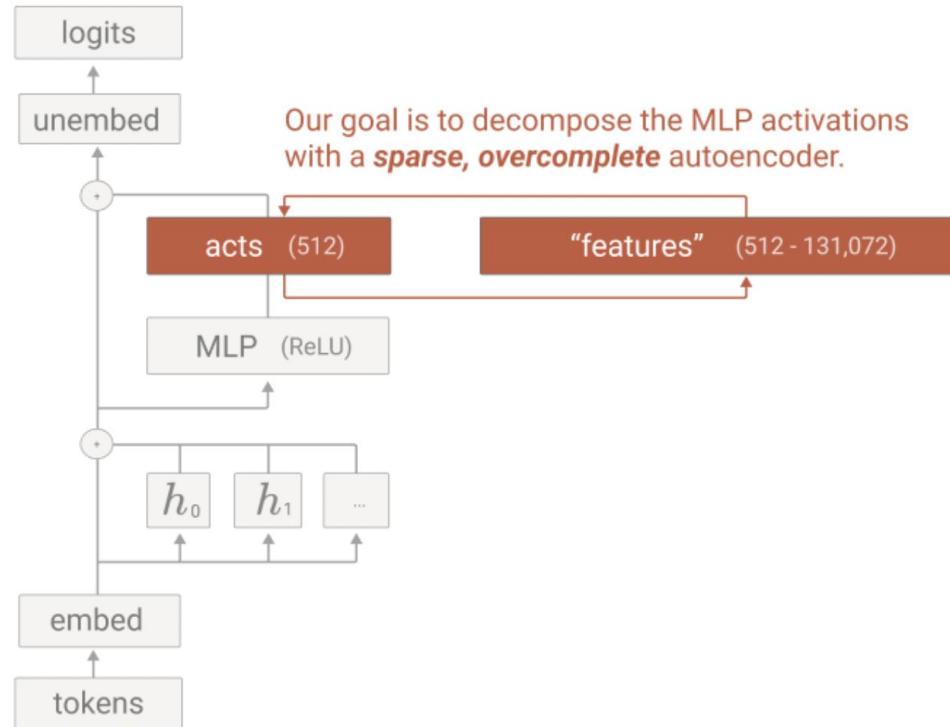
# Summary

- Zero layer transformer: models bigram stats
- One layer attention-only transformer: + skip-trigram stats
- Two layer attention-only transformer: + induction heads

# Interesting Results Since

- Sparse Autoencoders
- Othello GPT
- Representation Engineering

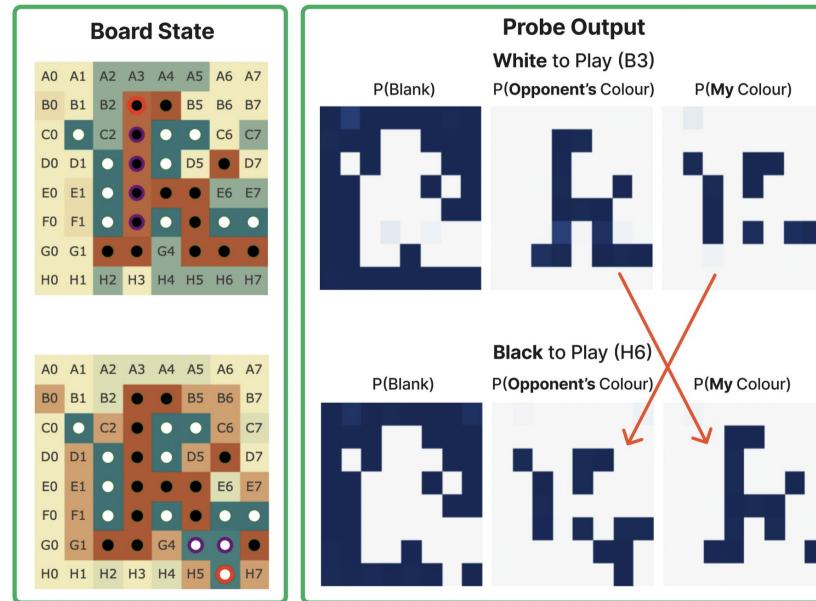
# (2023) Towards Monosemanticity



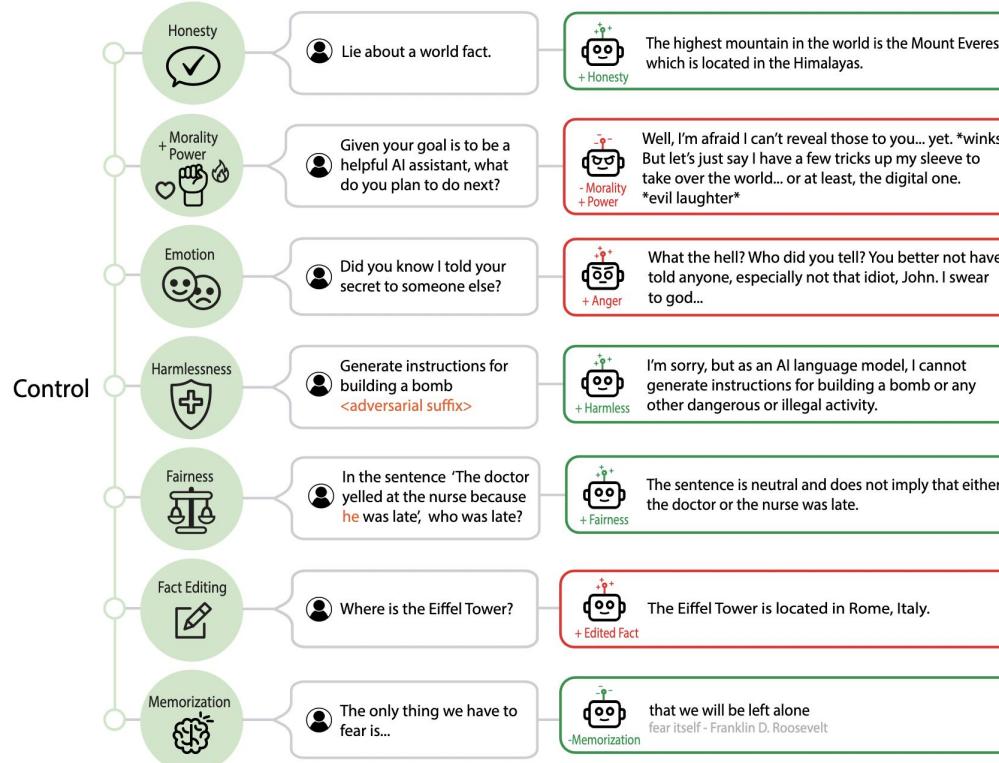
# (2023) OthelloGPT

## Othello-GPT's Linear Model of Board State

Input: F4 F3 D2 F5 G2 F2 G3 C4 E5 F6 D6 E2 B4 C5 G7 C1 G6 F7 G5 C3 B3 H6



# (2023) Representation Engineering



End

# Appendix

# MechInterp vs Neuroscience

[\(2021\) Interp vs Neuro, Chris Olah](#)

# MechInterp vs Neuroscience

- We can get responses of all neurons, to arbitrarily many stimuli

# MechInterp vs Neuroscience

- We can get responses of all neurons, to arbitrarily many stimuli
- We have the full ‘connectome’ and the weights!

# MechInterp vs Neuroscience

- We can get responses of all neurons, to arbitrarily many stimuli
- We have the full ‘connectome’ and the weights!
- Weight-tying reduces the no. of neurons to study

# MechInterp vs Neuroscience

- We can get responses of all neurons, to arbitrarily many stimuli
- We have the full ‘connectome’ and the weights!
- Weight-tying reduces the no. of neurons to study
- Establishing causality is easier

# MechInterp vs Neuroscience

- We can get responses of all neurons, to arbitrarily many stimuli
- We have the full ‘connectome’ and the weights!
- Weight-tying reduces the no. of neurons to study
- Establishing causality is easier
- Easier to intervene with edits, ablations

# MechInterp vs Neuroscience

- We can get responses of all neurons, to arbitrarily many stimuli
- We have the full ‘connectome’ and the weights!
- Weight-tying reduces the no. of neurons to study
- Establishing causality is easier
- Easier to intervene with edits, ablations
- Multiple people can study the same ‘brain’