# Introduction to Mechanistic Interpretability

TransformerLens & Induction circuits

# Transformer Architecture



**logits**

**unembed**

$x_{-1}$

The final logits are produced by applying the unembedding.

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

$x_{i+2}$

**MLP** $m$

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

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

$x_{i+1}$

$h_0$  $h_1$  ...

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

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

$x_i$

$x_0$

**embed**
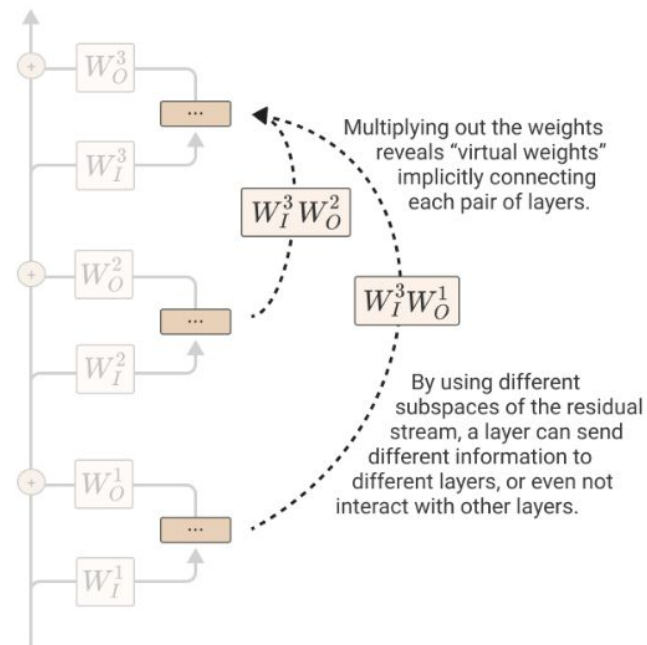
**tokens**

Token embedding.

$$x_0 = W_E t$$

One residual block

# Residual Stream
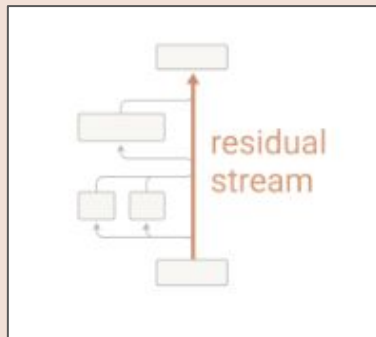
The residual stream is modified by a sequence of MLP and attention layers "reading from" and "writing to" it with linear operations.
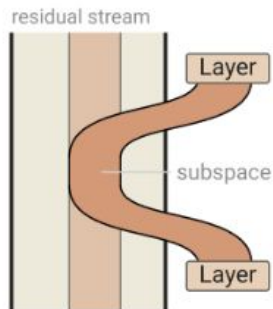
$W_O^3$

$\cdots$

$W_I^3$

Each layer **"writes"** to the residual stream by adding a linear projection of its results.

$W_O^2$

$\cdots$

$W_I^2$

Each layer **"reads"** from the residual stream with a linear projection.

$W_O^1$

$\cdots$

$W_I^1$

Because all these operations are linear, we can "multiply through" the residual stream.

$W_O^3$

$\cdots$

$W_I^3$

$W_I^3 W_O^2$

Multiplying out the weights reveals "virtual weights" implicitly connecting each pair of layers.

$W_O^2$

$\cdots$

$W_I^2$

$W_I^3 W_O^1$

$W_O^1$

$\cdots$

$W_I^1$

By using different subspaces of the residual stream, a layer can send different information to different layers, or even not interact with other layers.

# Residual Stream





The residual stream is high dimensional, and can be divided into different subspaces.

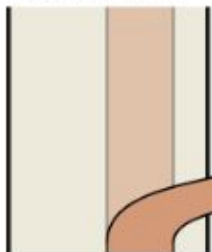residual stream

subspace

Layer

Layer

Layers can interact by writing to and reading from the same or overlapping subspaces. If they write to and read from disjoint subspaces, they won't interact. Typically the spaces only partially overlap.
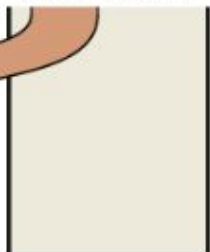
residual stream

Layer

Layers can delete information from the residual stream by reading in a subspace and then writing the negative verison.

# Attention Heads as Information Movement
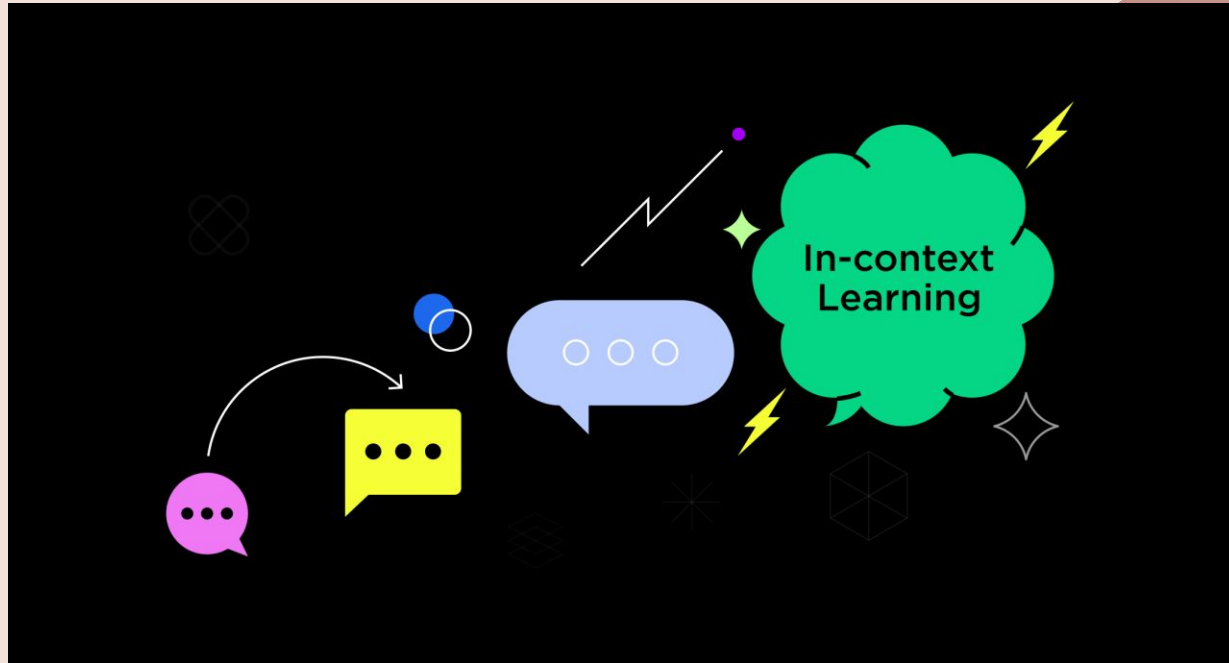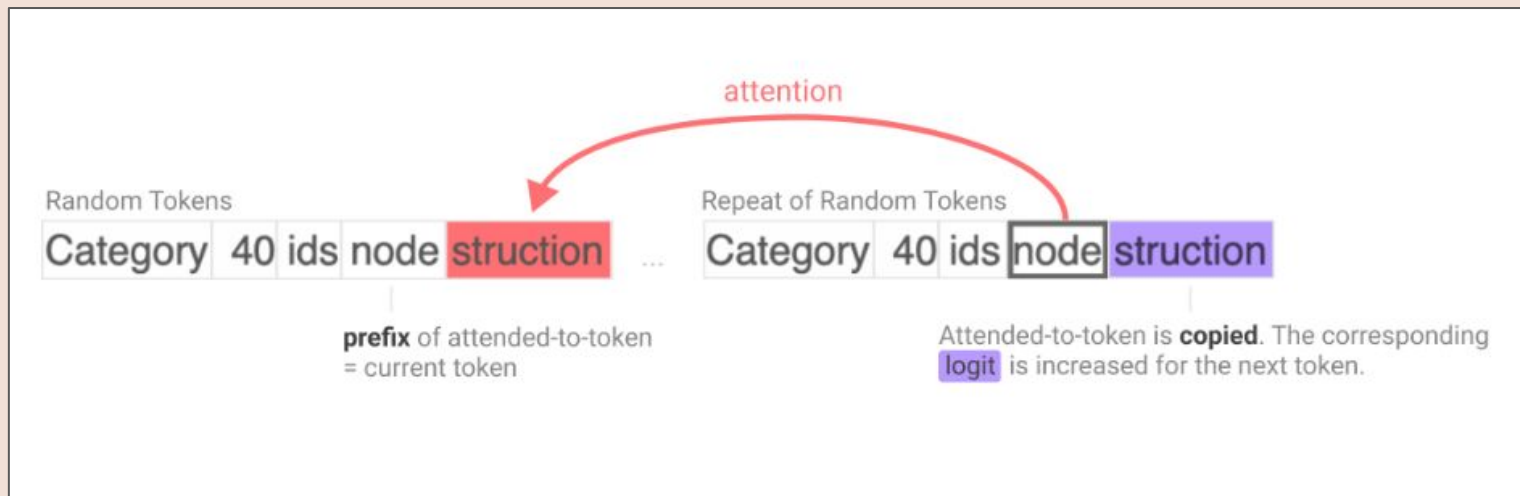


token A residual stream

token B residual stream

Attention heads copy information from the residual stream of one token to the residual stream of another. They typically write to a different subspace than they read from.

# In-Context Learning

# Induction Heads

Induction heads are implemented by a circuit consisting of
a pair of attention heads in different layers that work together
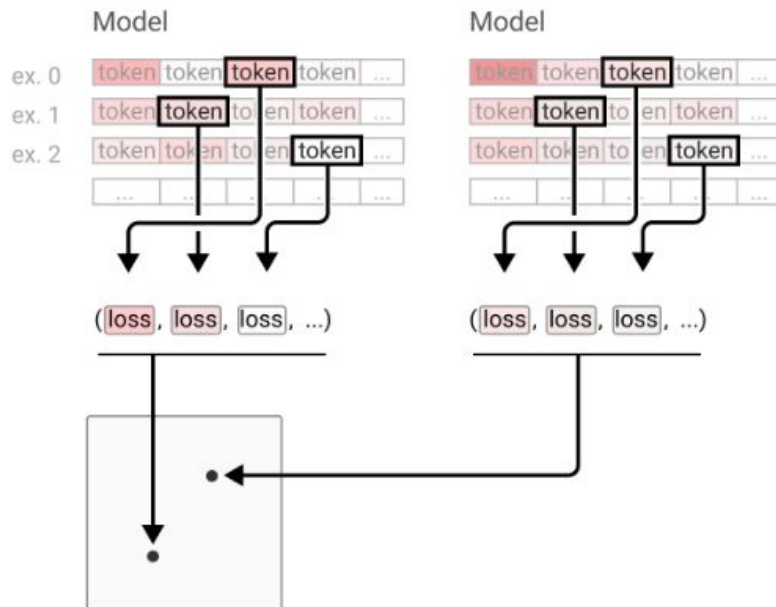to copy or complete patterns

# ICL Score

**In-context learning score:** the loss of the 500th token in the context minus the average loss of the 50th token in the context, averaged over dataset examples



**Step 1:** Run each model / snapshot over the same set of multiple dataset examples, collecting one token's loss per example.

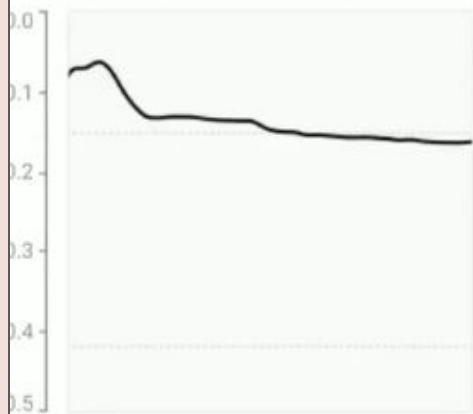**Step 2:** For each sample, extract the loss of a consistent token. Combine these to make a vector of losses per model / snapshot.

**Step 3:** The vectors are jointly reduced with principal component analysis to project them into a shared 2D space.

Change architecture to promote induction heads => phase change happens earlier

Elapsed Training Tokens

0    2.5e9    5.0e9    7.5e9    1e10

phase change

phase change occurs earlier than in baseline

**ABLATION TO "IN-CONTEXT LEARNING" SCORE**

contributes to in-context learning

0.15
0.10
0.05
0.00
-0.05
-0.10
-0.15

one-layer model
no change

models with more than one layer
have a phase change

The attention heads which increase in-context learning are almost entirely induction heads They form as in-context learning increases and drive its increase.

Heads are colored by type:
previous token;
induction; other.

# Subspaces in the Residual Stream & Embedding

## Layer 0 attention head

Summary: each token looks one position backwards, and gets the information about which token preceded it.

I am "D" is converted to I follow "D" and moved to the destination token, ready to be read by second layer.

I am "urs"
at posn=1
I follow "D"

at posn=1 looks one token back, i.e. searches for at posn=0

(this requires a rotation)

value → output

key    query

I am "D"
at posn=0

I am "urs"
at posn=1

I am "D"
at posn=n

I am "urs"
at posn=n+1

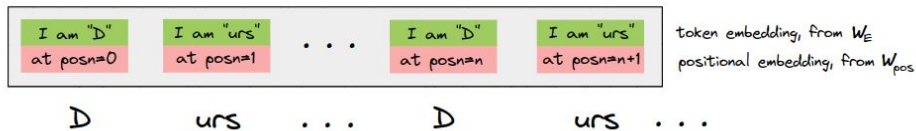D          urs          . . .          D          urs          . . .

### QK circuit

Query vector is $\begin{bmatrix} \text{I am "urs"} \\ \text{at posn=1} \end{bmatrix}^T W_Q =$ "I'm looking for posn 1 - 1 = 0"

Key vector is $\begin{bmatrix} \text{I am "D"} \\ \text{at posn=0} \end{bmatrix}^T W_K =$ "I'm at posn 0"

i.e. we subtract one from the current position (this requires a rotation[0])

The attention score (the dot product of the key and query vectors) is large, because the key is a good match for the query.

### OV circuit

The OV circuit reads token information, and moves it to a different subspace (so it doesn't erase the token embedding information which is already stored at the destination token).

We have:

$\begin{bmatrix} \text{I am "D"} \\ \text{at posn=0} \end{bmatrix}^T W_V W_O =$ I follow "D"

which gets added to the residual stream for the first "urs" token.

# Layer 0 Attention Head Mathematically

Mathematically, the attention scores are:

$$\begin{bmatrix} \text{I am "X"} \\ \text{at posn=i} \end{bmatrix}^{T} W_Q W_K^{T} \begin{bmatrix} \text{I am "Y"} \\ \text{at posn=j} \end{bmatrix} \approx \begin{bmatrix} \text{at posn=i} \end{bmatrix}^{T} W_Q W_K^{T} \begin{bmatrix} \text{at posn=j} \end{bmatrix} = \begin{cases} \text{large} & \text{if} \quad j = i - 1 \\ \text{small} & \text{if else} \end{cases}$$

How can we prove this?

Define the "full QK circuit" $W_{pos} W_{QK} W_{pos}^{T}$

The (i, j) th element of this matrix is $\begin{bmatrix} \text{at posn=i} \end{bmatrix}^{T} W_Q W_K^{T} \begin{bmatrix} \text{at posn=j} \end{bmatrix}$, so we can verify this.

# Layer 1 Attention Head

## Layer 1 attention head

Summary: the 2nd "D" token looks back for the token following the 1st "D" (which is "urs"), and uses that as its prediction.

I am "urs" is converted into "urs" is next and moved to the destination token, ready to be decoded.

I am "D"
at posn=n
"urs" is next

value → output

key   query

I am "D" searches for tokens containing I follow "D"

I am "urs"
at posn=1
I follow "D"

K          Q

I am "D"        I am "urs"                    I am "D"        I am "urs"
at posn=0       at posn=1                     at posn=n       at posn=n+1

D              urs          . . .             D              urs          . . .

## QK circuit

The query is  I am "D" at posn=n $W_Q$ = "I'm looking for a token following "D" "

The key is  I am "urs" at posn=1 I follow "D" $W_K$ = "I follow "D""

The attention score (the dot product of the key and query vectors) is large, because the key is a good match for the query.

## OV circuit

The OV circuit reads token information, and coverts it into something, which will result in a prediction of "urs".

We have:

I am "urs" at posn=1 I follow "D" $W_V W_O$ = "urs" is next

which gets added to the residual stream of the second "D" token.

# Layer 1 Attention Head Mathematically
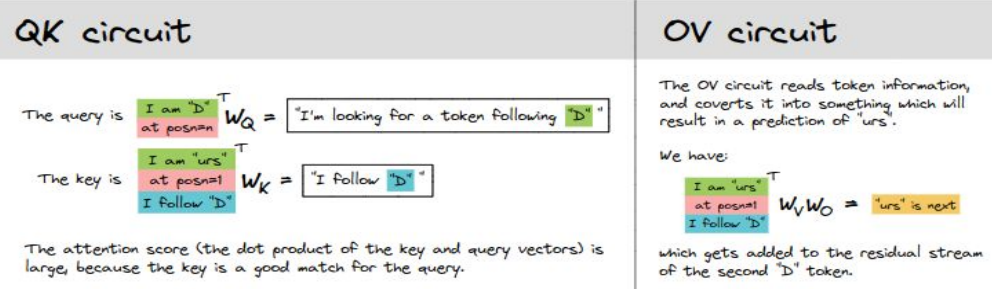
Mathematically, the attention scores are:

$$\begin{bmatrix} \text{I am "x"} \\ \text{at pos}n_{;}i \\ \text{I follow "p"} \end{bmatrix}^{\top} W_Q W_K^{\top} \begin{bmatrix} \text{I am "Y"} \\ \text{at pos}n_{;}j \\ \text{I follow "Q"} \end{bmatrix} \approx \boxed{\text{I am "x"}}^{\top} W_Q W_K^{\top} \boxed{\text{I follow "Q"}}$$

$$= \left( \boxed{\text{I am "x"}}^{\top} W_Q \right) \cdot \left( \boxed{\text{I follow "Q"}}^{\top} W_K \right) = \begin{cases} \text{large} & \text{if } \quad "x" = "Q" \\ \text{small} & \text{if else} \end{cases}$$

"I am looking for something which follows "x""          "I follow "Q""

How can we prove this?

Define the "full K-composition circuit" $W_E W_{QK}^{1} (W_{OV}^{0})^{\top} W_E^{\top}$

If x and Q are tokens, then the (x, Q) th element of this matrix is:

$$"x"^{\top} W_E W_{QK}^{1} (W_{OV}^{0})^{\top} W_E^{\top} "Q" = \boxed{\text{I am "x"}}^{\top} W_{QK}^{1} (W_{OV}^{0})^{\top} \boxed{\text{I am "Q"}}$$

$$= \boxed{\text{I am "x"}}^{\top} W_{QK}^{1} \boxed{\text{I follow "Q"}}$$

Superscript 0 means the head in layer 0, and 1 to means the head in layer 1.

So we just need to verify that this is $\begin{cases} \text{large} & \text{if } \quad "x" = "Q" \\ \text{small} & \text{if else} \end{cases}$

# quAIdditch

In the context of transformer models, which specific circuit is responsible for writing information into the residual stream that helps subsequent layers make predictions about token sequences?

OV (Output-Value) circuit