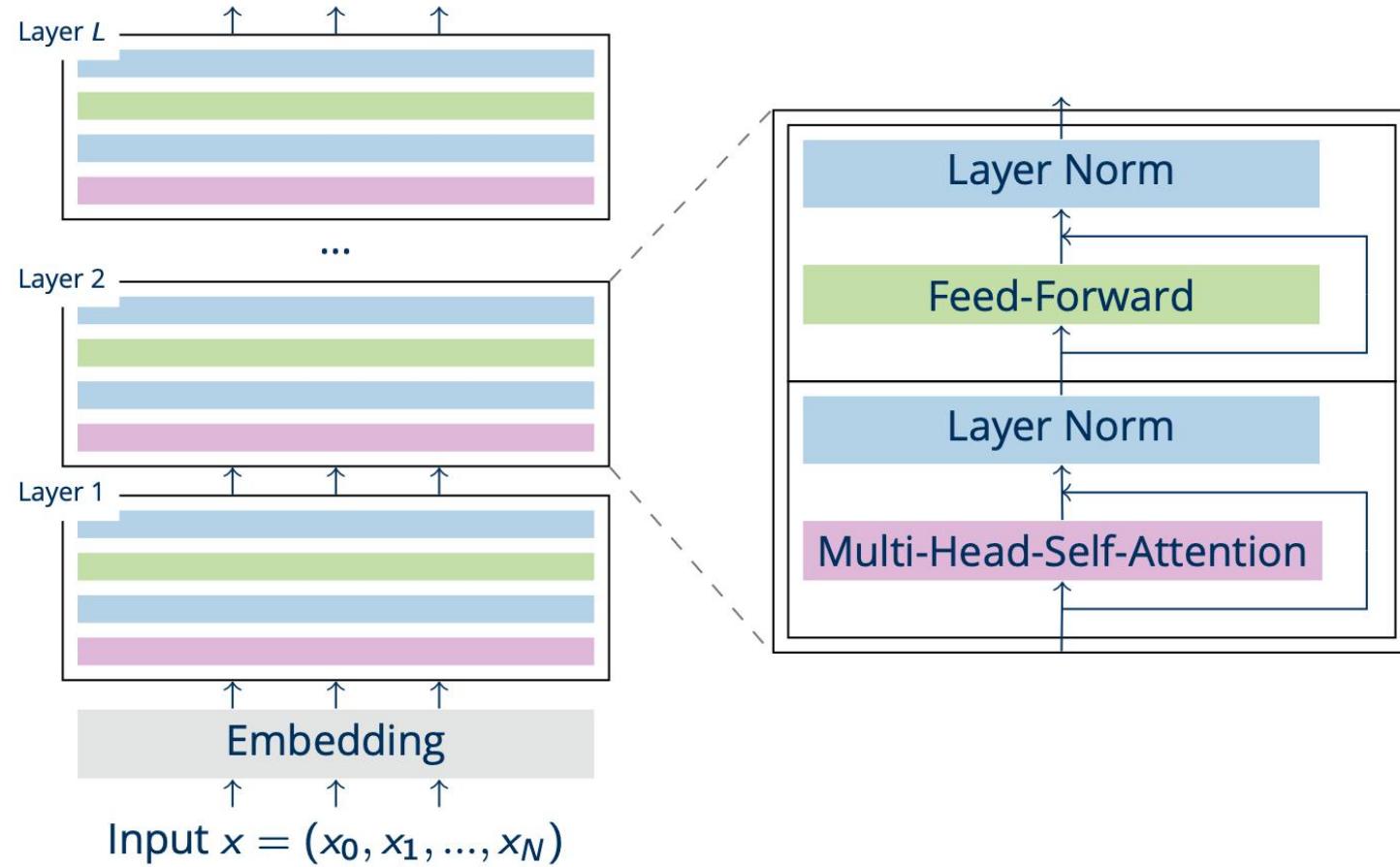
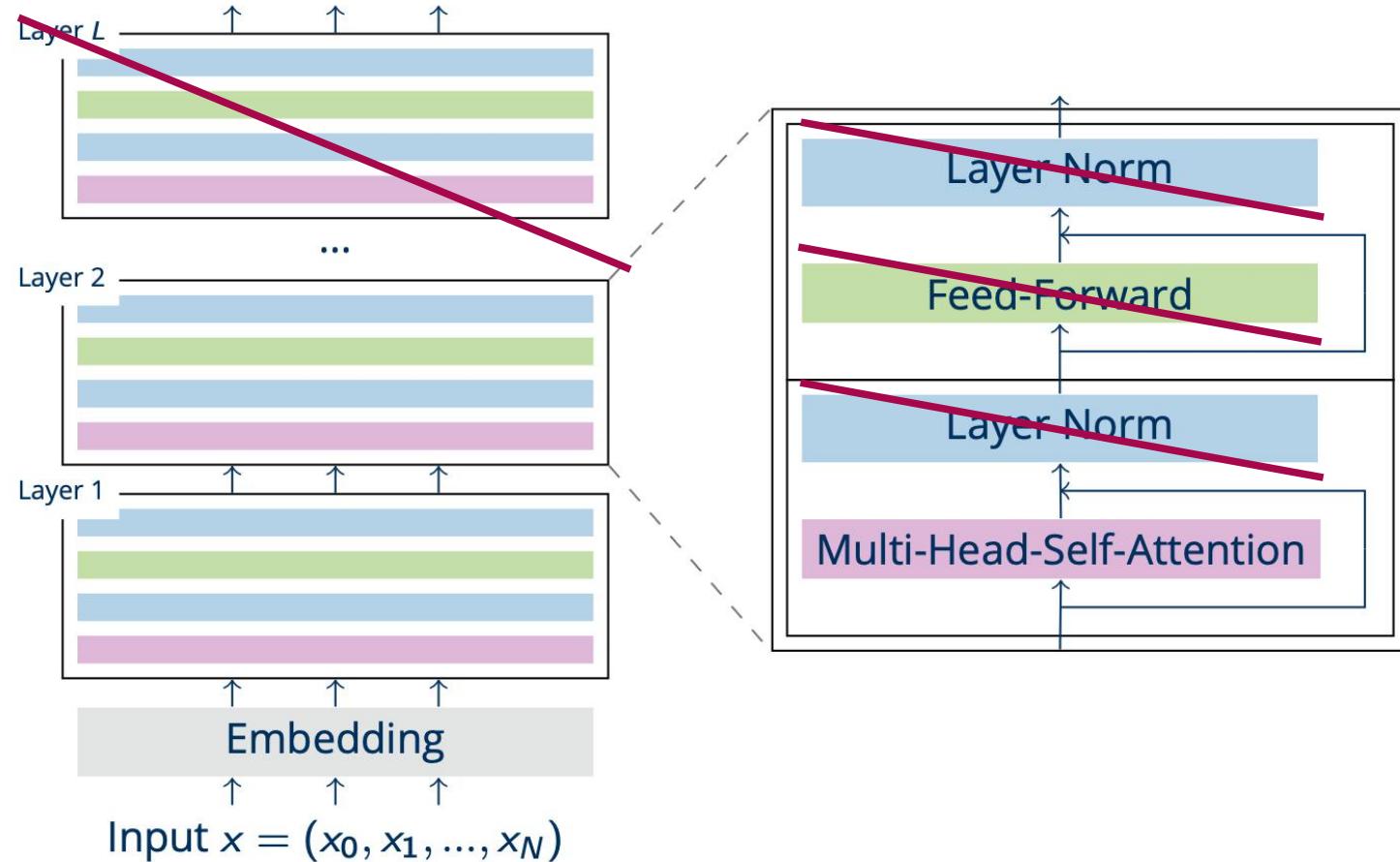


A Mathematical Framework for Transformer Circuits

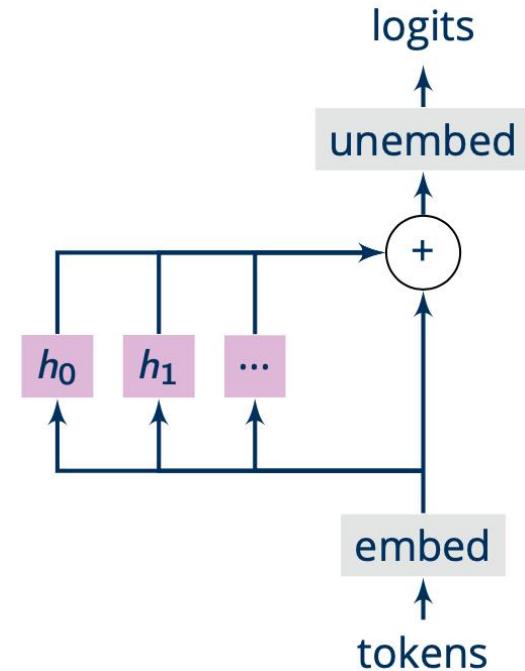
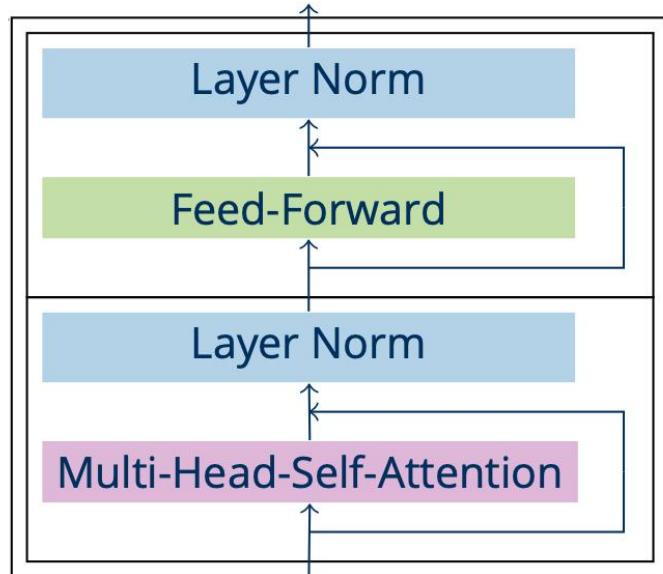
Transformer Architecture

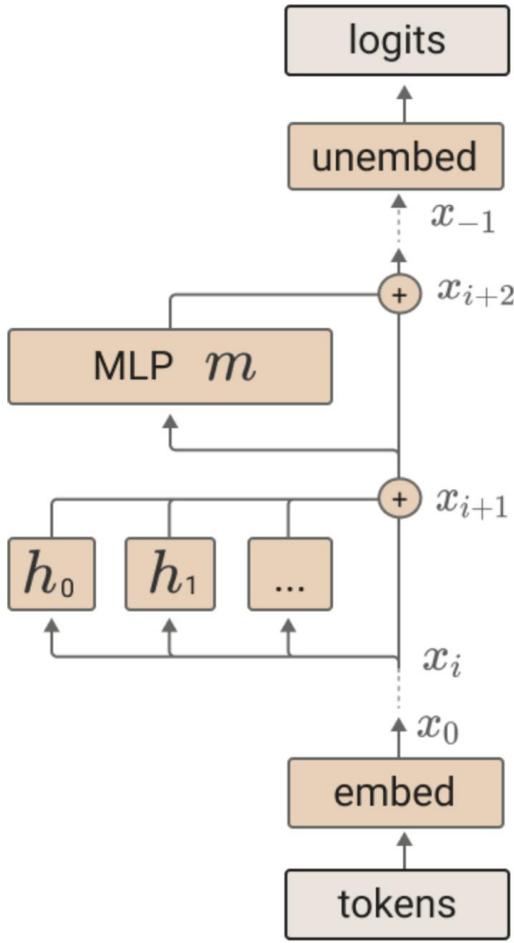


Transformer Architecture



Residual Stream





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.

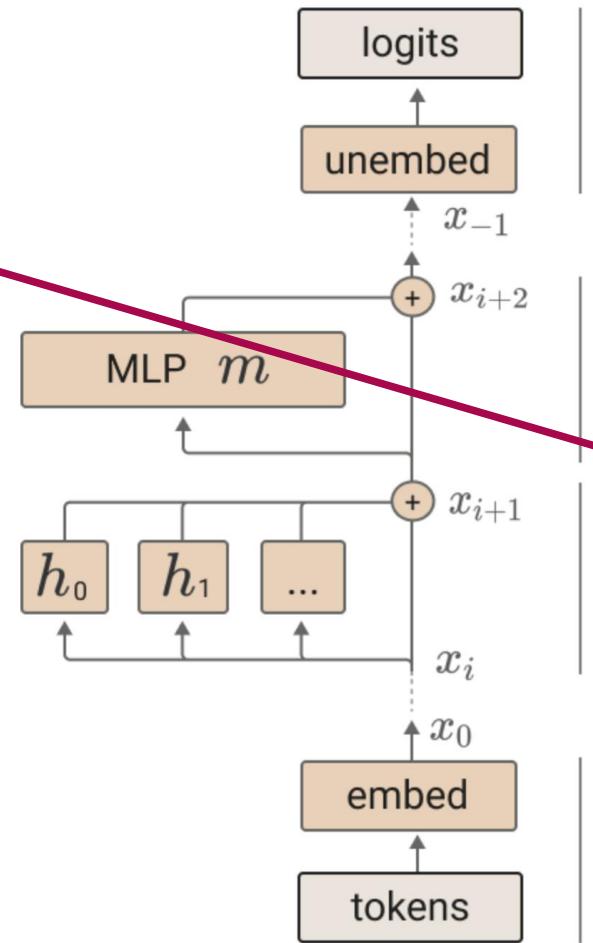
$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

Token embedding.

$$x_0 = W_E t$$

One residual block

0-Layer Transformer



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

0-Layer

$$x_o = w_E \cdot t \quad T(t) = w_u x_o \Rightarrow T(t) = w_u w_E \cdot t$$

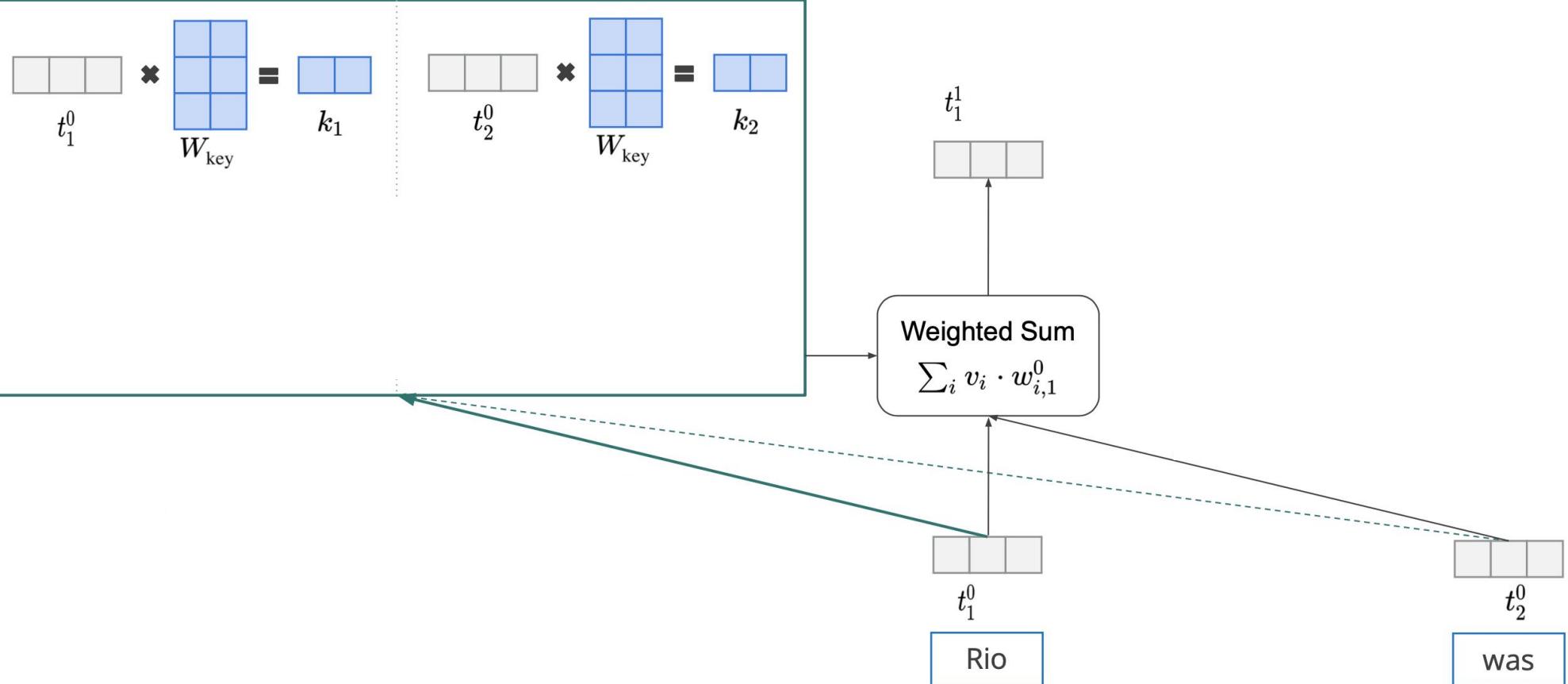
$$T("I") = \begin{pmatrix} w_u \\ \text{orange box} \\ w_u \end{pmatrix} \cdot \begin{pmatrix} \text{blue box} \end{pmatrix} = \begin{pmatrix} \text{large value for "am"} \end{pmatrix}$$

$$T(["I", "am"]) = \begin{pmatrix} w_u \\ \text{green box} \\ \text{orange box} \end{pmatrix} \cdot \begin{pmatrix} \text{blue box} & \text{blue box} \end{pmatrix} = \begin{pmatrix} \text{large value for e.g. "a"} \end{pmatrix}$$

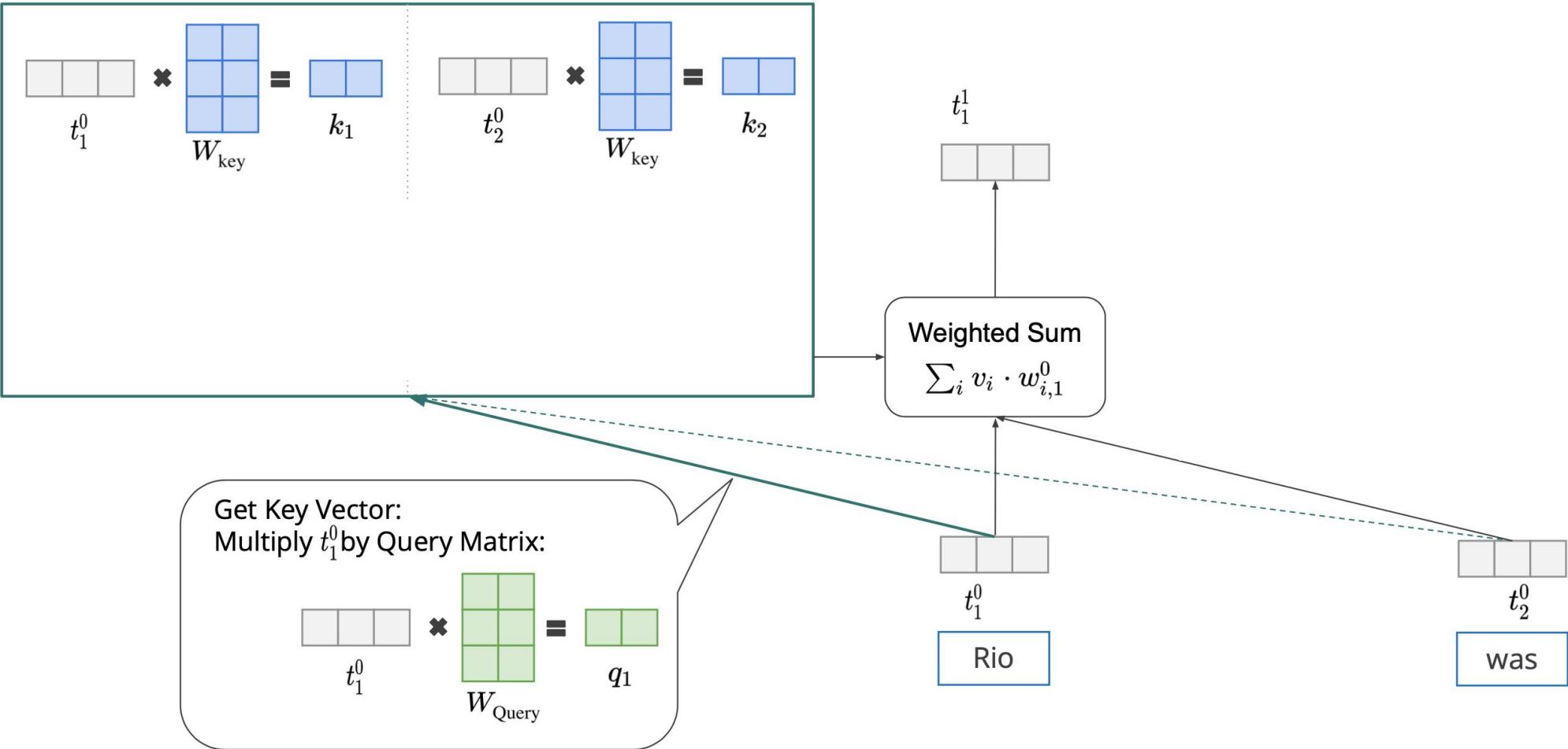
$$\begin{pmatrix} w_u \\ \text{green box} \\ \text{orange box} \end{pmatrix} \cdot \left(\begin{matrix} \leftarrow |V| \rightarrow \end{matrix} \right) = \begin{matrix} a \dots am I \dots z \\ \text{large value for "am"} \\ \sim P(a | am) \\ \text{Bigram - Probability} \end{matrix}$$

$\longleftarrow |V| \longrightarrow$

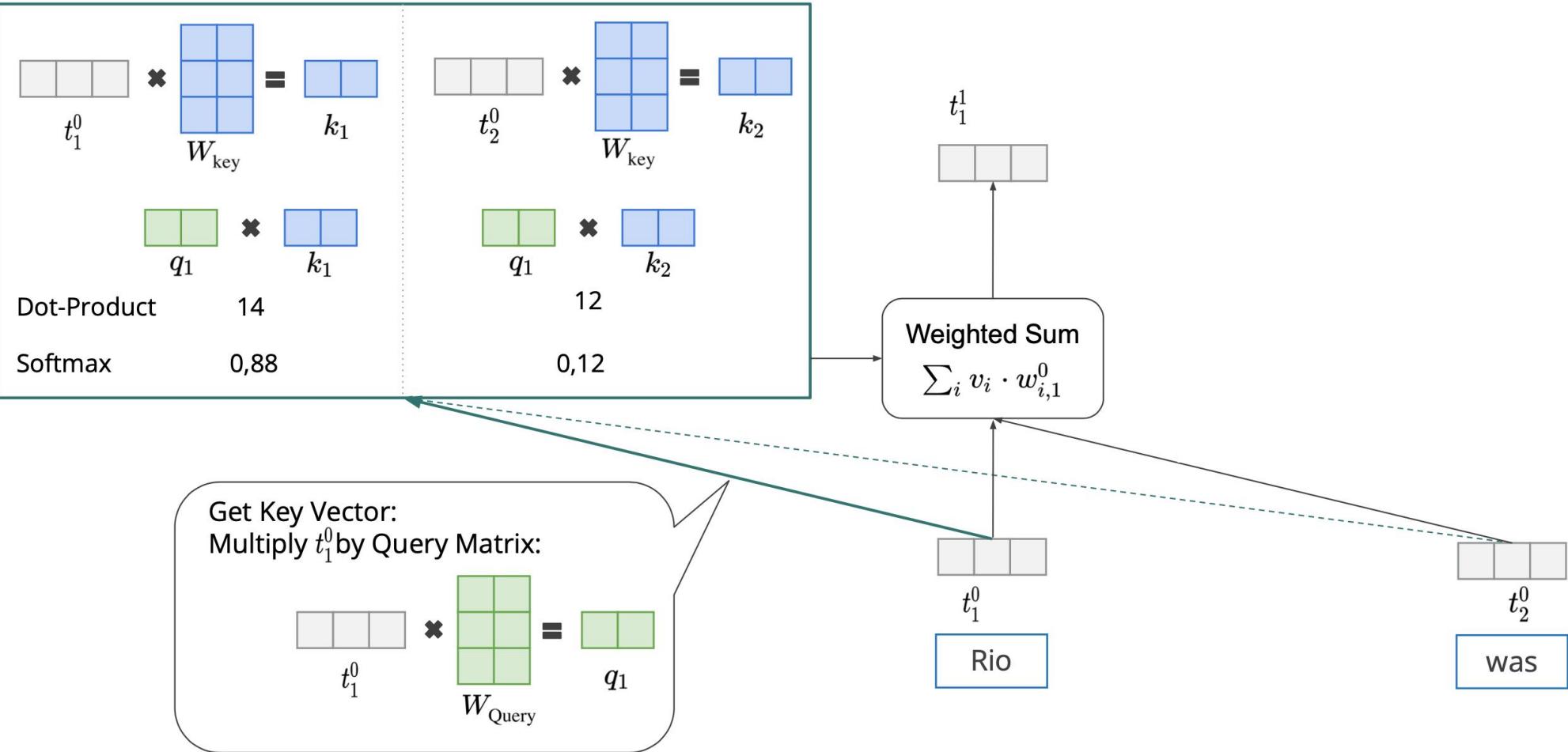
Self-Attention



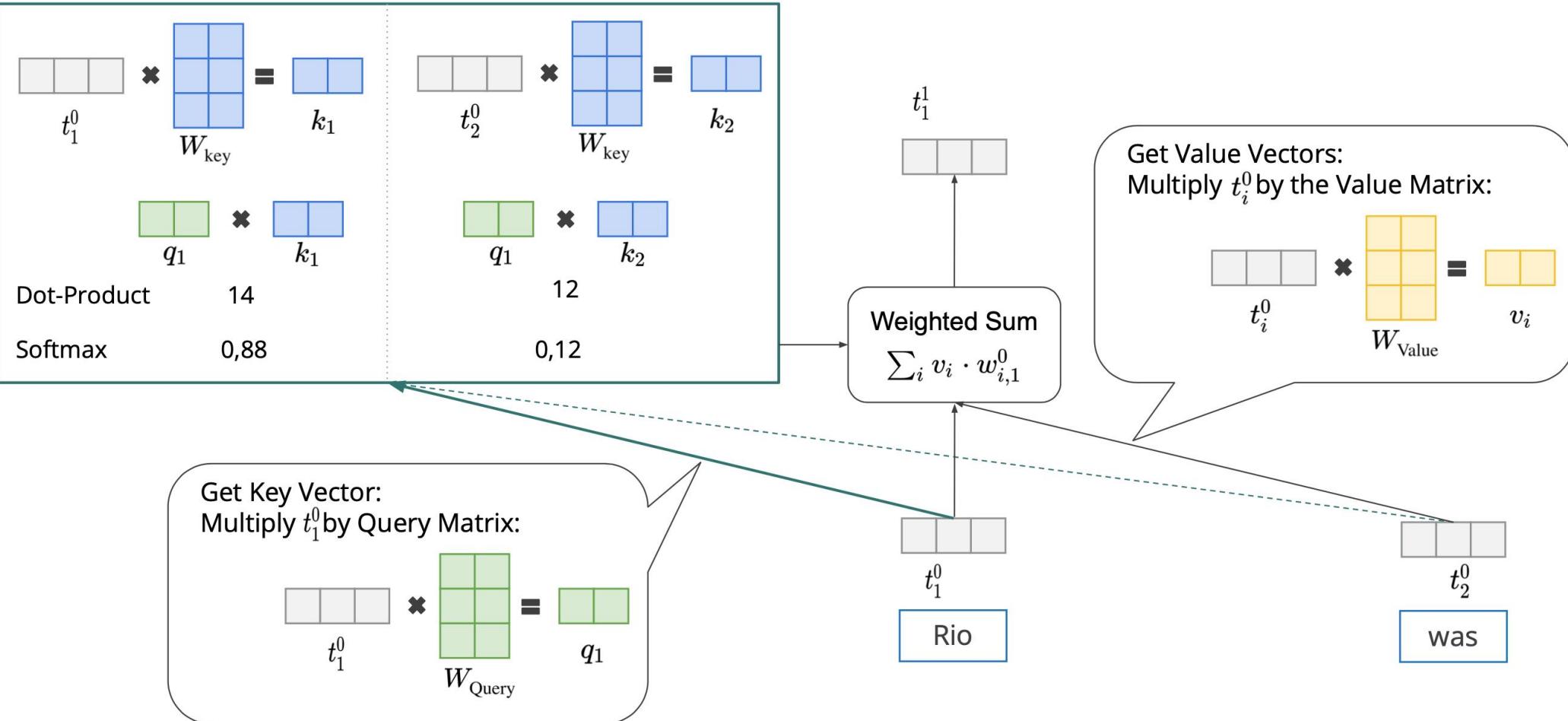
Self-Attention



Self-Attention



Self-Attention



Self-Attention

$$\begin{matrix} \square & \square & \square \\ t_1^0 & \times & \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} \\ & & W_{\text{key}} \end{matrix} = \begin{matrix} \square & \square \\ k_1 & \end{matrix}$$

$$\begin{matrix} \square & \square & \square \\ t_1^0 & \times & \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} \\ & & W_{\text{Query}} \end{matrix} = \begin{matrix} \square & \square \\ q_1 & \end{matrix}$$

$$\begin{matrix} \square & \square & \square \\ t_i^0 & \times & \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} \\ & & W_{\text{Value}} \end{matrix} = \begin{matrix} \square & \square \\ v_i & \end{matrix}$$

$$\alpha_{ij} = \frac{\exp e_{ij}}{\sum_{k=1}^N \exp e_{ik}}, e_{ij} = \frac{q_i \cdot k_j^T}{\sqrt{d_k}}.$$

$$c_i = \sum_{j=1}^N \alpha_{ij} \cdot v_j$$

$$q_i = x_i \cdot W_q, k_i = x_i \cdot W_k, v_i = x_i \cdot W_v$$

Multi-Head-Self-Attention

$$\text{Head}_1 = (W_q^{(1)}, W_k^{(1)}, W_v^{(1)})$$

$$\text{Head}_2 = (W_q^{(2)}, W_k^{(2)}, W_v^{(2)})$$

...

$$\text{Head}_H = (W_q^{(H)}, W_k^{(H)}, W_v^{(H)})$$

$$\text{Multi-Head-Self-Attention}(x) = [c_i^{(1)}, c_i^{(2)}, \dots, c_i^{(H)}] \cdot W_o,$$

Moving W_O into each head

$$W_O^H \begin{bmatrix} r^{h_1} \\ r^{h_2} \\ \dots \end{bmatrix} = [W_O^{h_1}, W_O^{h_2}, \dots] \cdot \begin{bmatrix} r^{h_1} \\ r^{h_2} \\ \dots \end{bmatrix} = \sum_i W_O^{h_i} r^{h_i}$$

Multi-Head-Self-Attention

$$W_O^H \begin{bmatrix} r^{h_1} \\ r^{h_2} \\ \dots \end{bmatrix} = [W_O^{h_1}, W_O^{h_2}, \dots] \cdot \begin{bmatrix} r^{h_1} \\ r^{h_2} \\ \dots \end{bmatrix} = \sum_i W_O^{h_i} r^{h_i}$$

One head:

$$\alpha_1 \cdot \boxed{\quad} + \alpha_2 \cdot \boxed{\quad} + \alpha_3 \cdot \boxed{\quad} = \boxed{\quad} = c_i$$

Results of multiple heads w_o

$$\begin{bmatrix} A & B \\ C & D \\ E & F \end{bmatrix} \cdot \begin{pmatrix} a & g \\ b & h \\ c & i \\ d & j \\ e & k \\ f & l \end{pmatrix} = \begin{bmatrix} aA + bB + cC + dD + eE + fF \\ gA + hB + iC + jD + kE + lF \end{bmatrix}$$

$$\boxed{AB} \cdot \begin{pmatrix} a & g \\ b & h \end{pmatrix} + \boxed{CD} \cdot \begin{pmatrix} c & i \\ d & j \end{pmatrix} + \boxed{EF} \cdot \begin{pmatrix} e & k \\ f & l \end{pmatrix}$$

Attention Calculation

1. Compute the value vector for each token from the residual stream ($v_i = W_V x_i$).
2. Compute the “result vector” by linearly combining value vectors according to the attention pattern ($r_i = \sum_j A_{i,j} v_j$).
3. Finally, compute the output vector of the head for each token ($h(x)_i = W_O r_i$).⁸

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

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

$$\begin{aligned}(A \otimes W_O W_V) x &= A \cdot x \cdot (W_O W_V)^T \\ &= A \cdot x \cdot W_V^T W_O^T\end{aligned}$$

Tensor Product

There are several completely equivalent ways to interpret these products. If the symbol is unfamiliar, you can pick whichever you feel most comfortable with:

- **Left-right multiplying:** Multiplying x by a tensor product $A \otimes W$ is equivalent to simultaneously left and right multiplying: $(A \otimes W)x = AxW^T$. When we add them, it is equivalent to adding the results of this multiplication: $(A_1 \otimes W_1 + A_2 \otimes W_2)x = A_1xW_1^T + A_2xW_2^T$.
- **Kronecker product:** The operations we want to perform are linear transformations on a flattened ("vectorized") version of the activation matrix x . But flattening gives us a huge vector, and we need to map our matrices to a much larger block matrix that performs operations like "multiply the elements which previously corresponded to a vector by this matrix". The correct operation to do this is the Kronecker product. So we can interpret \otimes as a Kronecker product acting on the vectorization of x , and everything works out equivalently.
- **Tensor Product:** $A \otimes W$ can be interpreted as a tensor product turning the matrices A and W into a 4D tensor. In NumPy notation, it is equivalent to `A[:, :, None, None] * W[None, None, :, :]` (although one wouldn't computationally represent them in that form). More formally, A and W are "type (1, 1)" tensors (matrices mapping vectors to vectors), and $A \otimes W$ is a "type (2, 2)" tensor (which can map matrices to matrices).

Tensor Product

$$(A \otimes w)x = Axw^T$$

in x : every row is one token vector!

$$A = \begin{pmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \end{pmatrix}$$

$$x = \begin{pmatrix} x_1 & x_4 \\ x_2 & x_5 \\ x_3 & x_6 \end{pmatrix}$$

linear comb. of
first compn.
 $A \cdot x = \begin{pmatrix} a_1 x_1 + a_2 x_2 + a_3 x_3 \\ a_4 x_1 + a_5 x_2 + a_6 x_3 \end{pmatrix}$

... of second comp.
 $\begin{pmatrix} a_1 x_4 + a_2 x_5 + a_3 x_6 \\ a_4 x_4 + a_5 x_5 + a_6 x_6 \end{pmatrix}$

$$w = \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix}$$

$$w^T = \begin{pmatrix} w_1 & w_3 \\ w_2 & w_4 \end{pmatrix}$$

Tensor Product

$$(A \otimes w)x = Axw^T$$

in x : every row is one token vector!

$$A = \begin{pmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \end{pmatrix}$$

$$x = \begin{pmatrix} x_1 & x_4 \\ x_2 & x_5 \\ x_3 & x_6 \end{pmatrix}$$

$$w = \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix}$$

linear comb. of

first compn.

$$A \cdot x = \begin{pmatrix} a_1x_1 + a_2x_2 + a_3x_3 \\ a_4x_1 + a_5x_2 + a_6x_3 \end{pmatrix}$$

... of second compn.

$$\begin{pmatrix} a_1x_4 + a_2x_5 + a_3x_6 \\ a_4x_4 + a_5x_5 + a_6x_6 \end{pmatrix}$$

$$w^T = \begin{pmatrix} w_1 & w_3 \\ w_2 & w_4 \end{pmatrix}$$

$$Ax \cdot w^T = \underbrace{(a_1x_1 + a_2x_2 + a_3x_3)w_1}_{\text{linear combination of first row vector!}} + \underbrace{(a_4x_1 + a_5x_2 + a_6x_3)w_1}_{(a_4x_1 + a_5x_2 + a_6x_3)w_1} + \underbrace{(a_1x_4 + a_2x_5 + a_3x_6)w_2}_{(a_4x_4 + a_5x_5 + a_6x_6)w_2} + \underbrace{(a_4x_4 + a_5x_5 + a_6x_6)w_2}_{(a_4x_4 + a_5x_5 + a_6x_6)w_2} + \underbrace{(a_1x_1 + a_2x_2 + a_3x_3)w_3}_{(a_4x_1 + a_5x_2 + a_6x_3)w_3} + \underbrace{(a_4x_1 + a_5x_2 + a_6x_3)w_3}_{(a_4x_1 + a_5x_2 + a_6x_3)w_3} + \underbrace{(a_1x_4 + a_2x_5 + a_3x_6)w_4}_{(a_4x_4 + a_5x_5 + a_6x_6)w_4} + \underbrace{(a_4x_4 + a_5x_5 + a_6x_6)w_4}_{(a_4x_4 + a_5x_5 + a_6x_6)w_4}$$

Attention Calculation

1. Compute the value vector for each token from the residual stream ($v_i = W_V x_i$).
2. Compute the “result vector” by linearly combining value vectors according to the attention pattern ($r_i = \sum_j A_{i,j} v_j$).
3. Finally, compute the output vector of the head for each token ($h(x)_i = W_O r_i$).⁸

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

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

$$\begin{aligned}(A \otimes W_O W_V) x &= A \cdot x \cdot (W_O W_V)^T \\ &= A \cdot x \cdot W_V^T W_O^T\end{aligned}$$

Attention Calculation

$$k_i = W_K x_i$$

$$q_i = W_Q x_i$$

$$A = \text{softmax}(q^T k)$$

$$A = \text{softmax}(x^T W_Q^T W_K x)$$

Attention Calculation

$$h(x) = (A \otimes W_O W_V) \cdot x \quad A = \text{softmax}(x^T W_Q^T W_K x)$$

Because $W_O W_V$ and $W_Q W_K$ always operate together, we like to define variables representing these combined matrices, $W_{OV} = W_O W_V$ and $W_{QK} = W_Q^T W_K$.

What happens when we compute A?

$$A = \text{Softmax} (x^T W_Q^T W_K x)$$

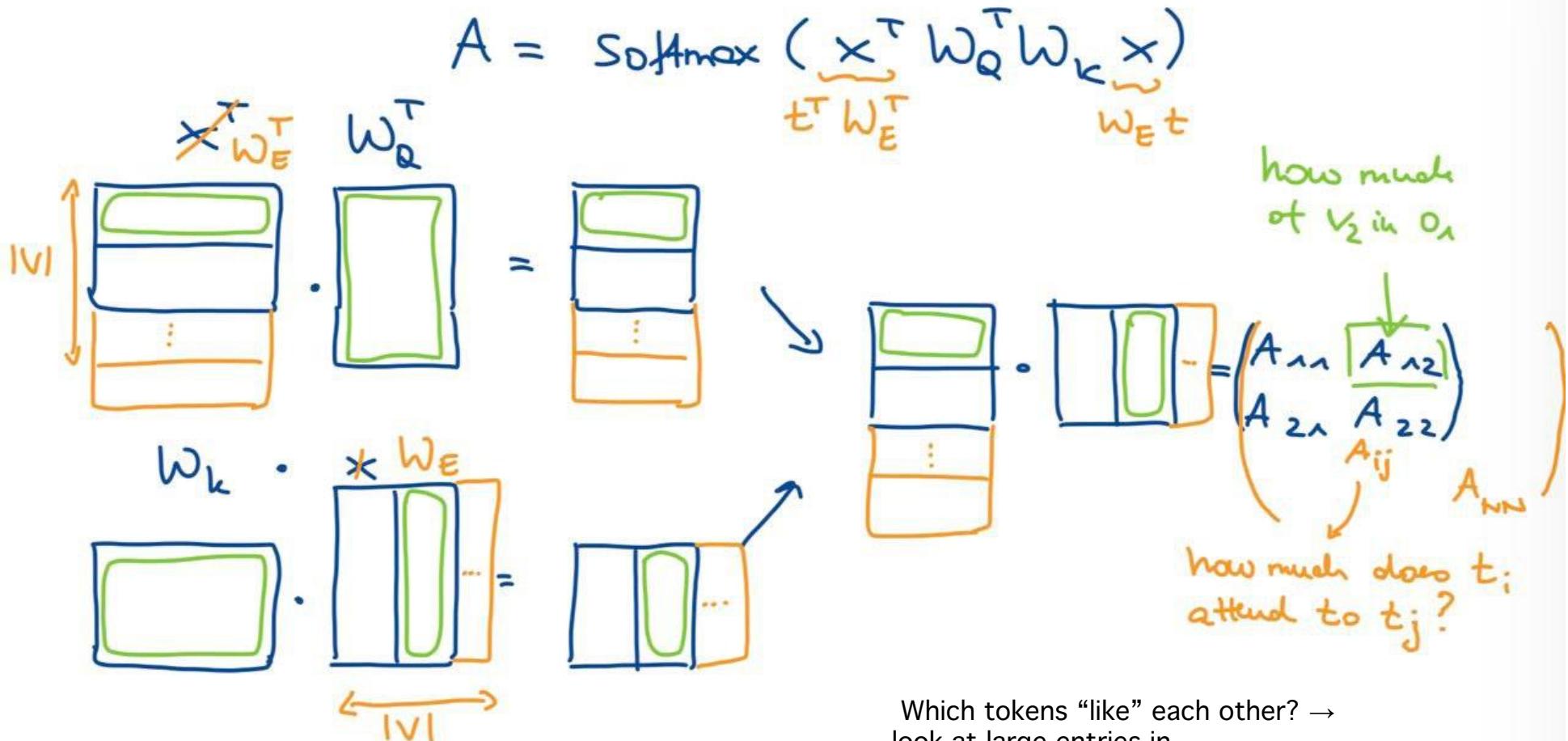
$$x^T \cdot W_Q^T = \begin{pmatrix} \text{green box} \\ \text{white box} \end{pmatrix}$$

how much
of v_2 in o_1

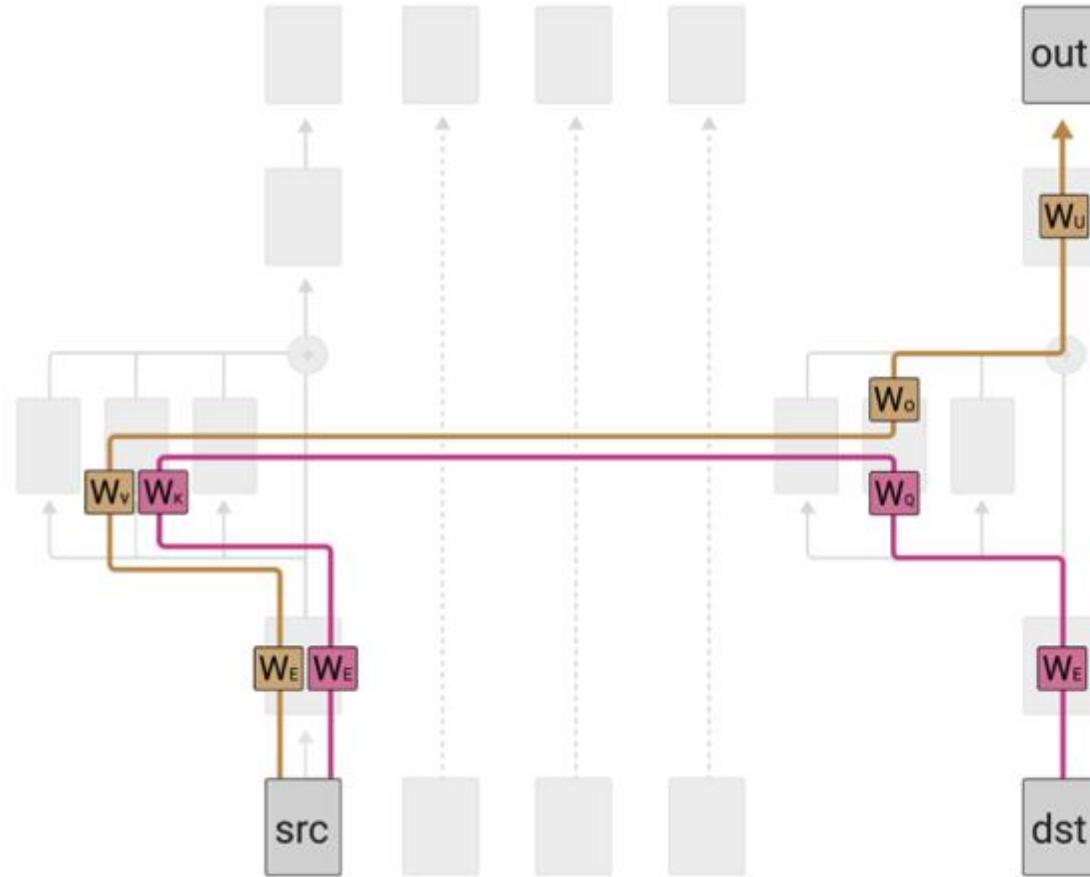
$$\begin{pmatrix} \text{green box} \\ \text{white box} \end{pmatrix} \cdot \begin{pmatrix} \text{green box} \\ \text{white box} \end{pmatrix} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

$$W_K \cdot x = \begin{pmatrix} \text{green box} \\ \text{white box} \end{pmatrix} \cdot \begin{pmatrix} \text{green box} \\ \text{white box} \end{pmatrix} = \begin{pmatrix} \text{green box} \\ \text{white box} \end{pmatrix}$$

What happens when we compute A?



1-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.

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

Interpretation?

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



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

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.

Examples

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.

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"

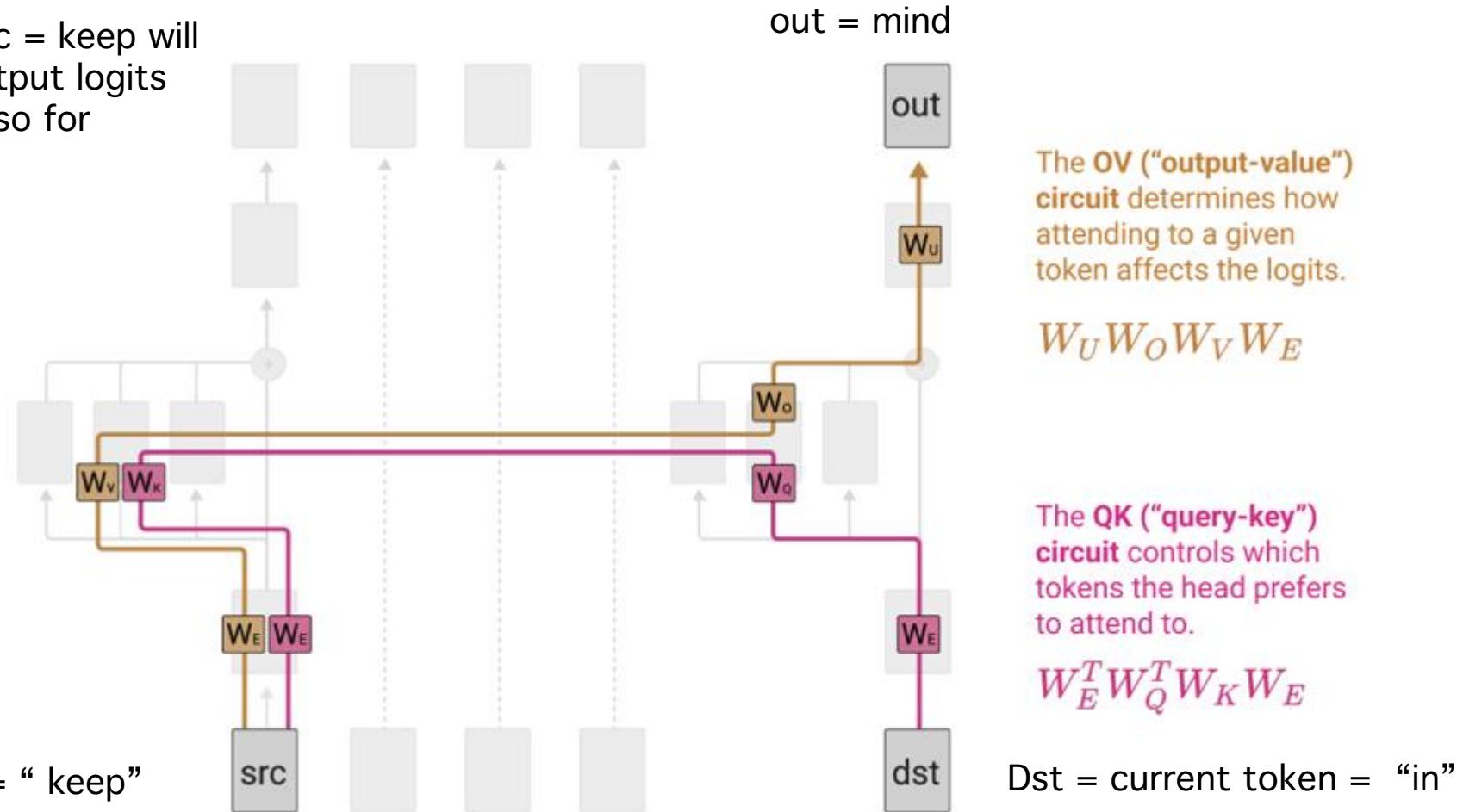
More Examples

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"

Interpretation

Attending to src = keep will increase the output logits for mind, but also for bay and wraps



“Bugs”

Limited Expressivity Can Create Bugs which Seem Strange from the Outside

Source Token	Destination Token	Out Token	“Correct” Skip Tri-grams	“Bug” Skip Tri-grams
“Pixmap”	“P”, “Q”, “P”, “p”, “U”	“ixmap”, “Canvas”, “Embed”, “grade”	“Pixmap…Pixmap”, “Pixmap…QCanvas”	“Pixmap…P Canvas ”
Source Token	Destination Token	Out Token	“Correct” Skip Tri-grams	“Bug” Skip Tri-grams
“Lloyd”	“L”, “L”, “P”, “P”, “R”, “C”	“loyd”, “alph”, “\n”, “acman”, … “atherine”	“Lloyd…Lloyd”, “Lloyd…Catherine”	“Lloyd… C loyd”, “Lloyd… L atherine”
Source Token	Destination Token	Out Token	“Correct” Skip Tri-grams	“Bug” Skip Tri-grams
“keep”	“in”, “at”, “out”, “under”, “off”	“bay”, … “mind”, … “wraps”	“keep…in mind”, “keep…at bay”, “keep…under wraps”	“keep…in bay ”, “keep…at wraps ”, “keep…under mind ”

How to find copying heads?

Look at OV-Circuit:

$$Mv_i = \lambda_i v_i \text{ with } M = W_U W_{OV} W_E$$

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

$$W_U W_O W_V W_E$$

$$m \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ \lambda_i \\ 0 \\ 0 \end{bmatrix}$$

if $\lambda_i > 1 \Rightarrow$ increase own probability

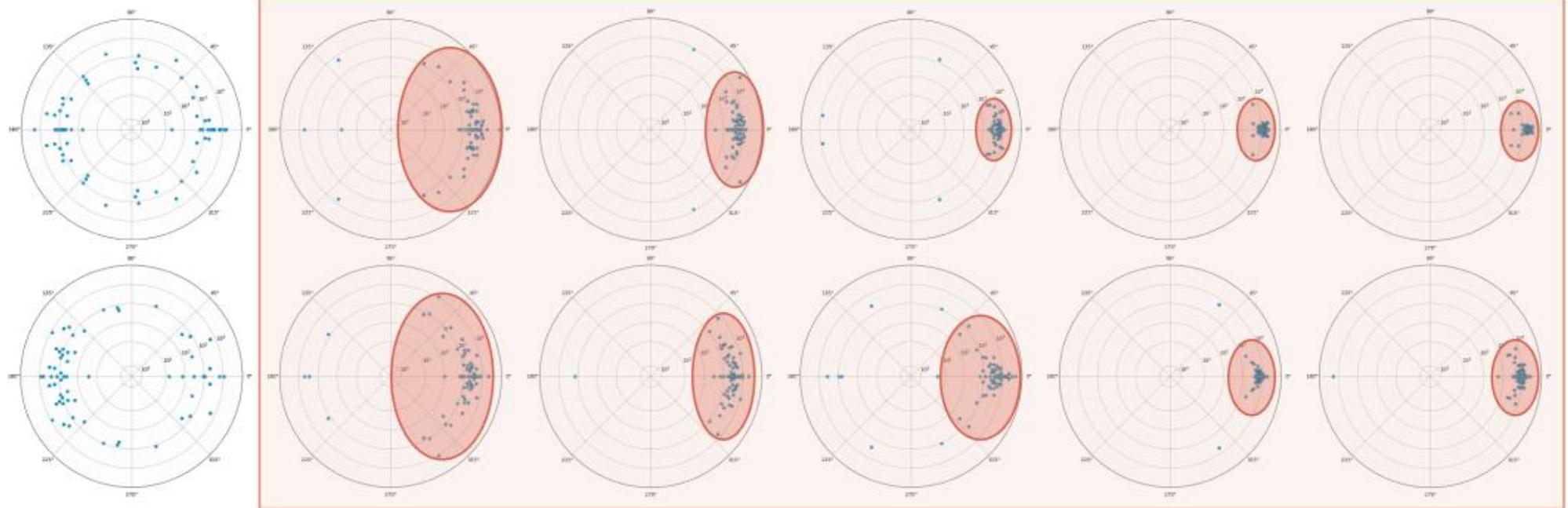
← logits = vector of values before Soft max.

one hot encoding of a token

How to find copying heads?

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

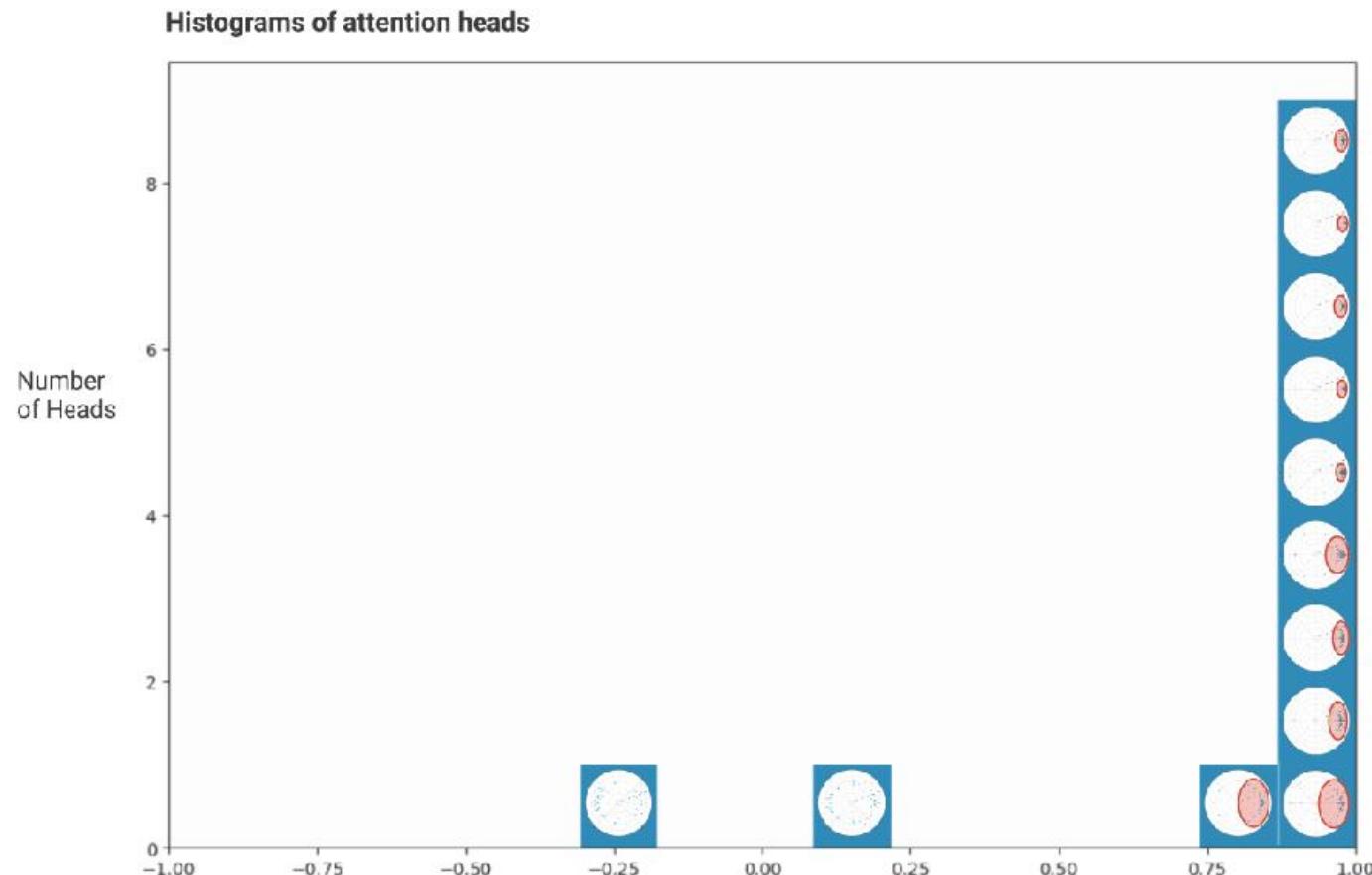
10/12 of layer 1 heads have mostly positive OV eigenvalues, and appear to significantly perform copying



← non-positive eigenvalues
not copying heads?

positive eigenvalues
copying heads? →

How to find copying heads?

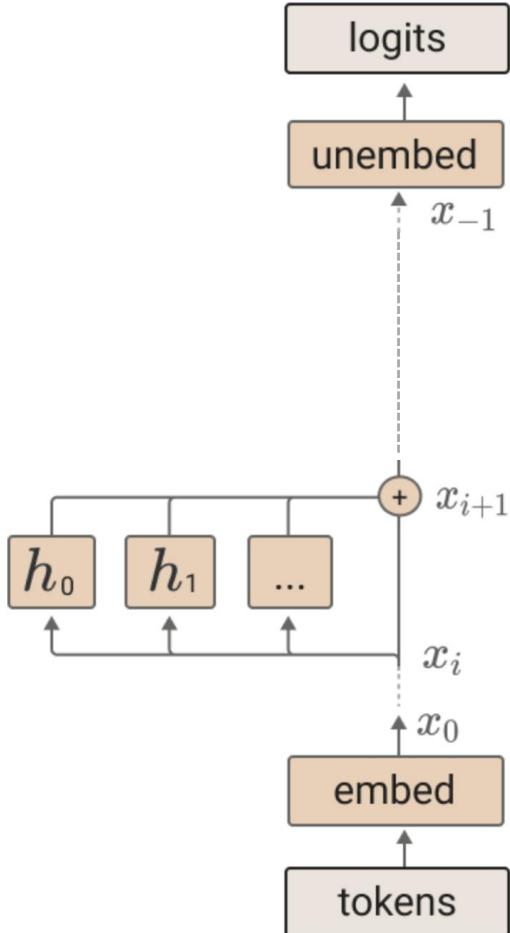


$$\frac{\sum_i \lambda_i}{\sum_i |\lambda_i|}$$

Another way to find copying heads?

- One might try to formalize "copying matrices" in other ways. One possibility is to look at the diagonal of a matrix, which describes how each token affects its own probability. As expected, entries on the diagonal are very positive-leaning.

Review Last Week



The final logits are produced by applying the unembedding.

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

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

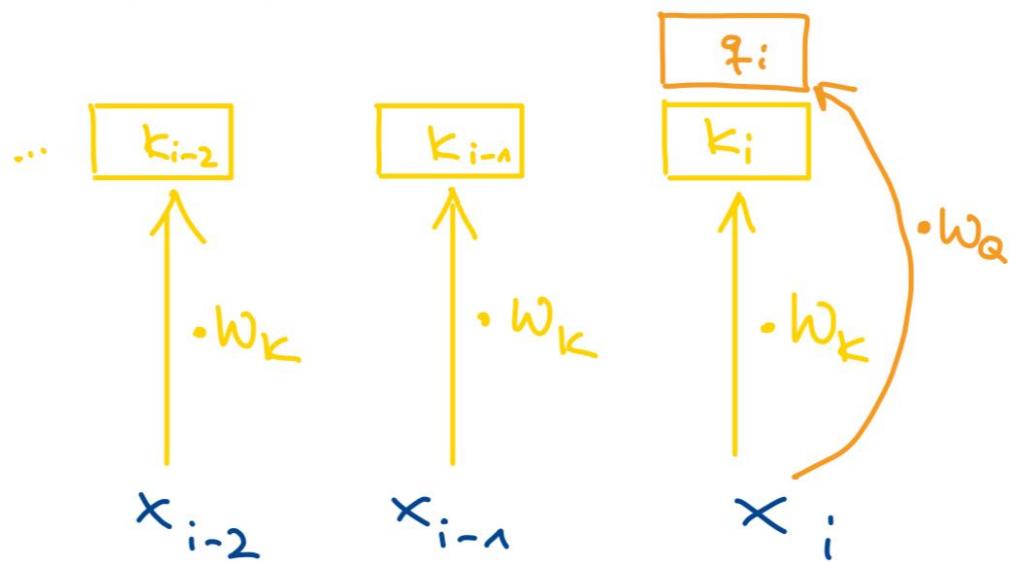
Attention Calculation

1. Compute the value vector for each token from the residual stream ($v_i = W_V x_i$).
2. Compute the “result vector” by linearly combining value vectors according to the attention pattern ($r_i = \sum_j A_{i,j} v_j$).
3. Finally, compute the output vector of the head for each token ($h(x)_i = W_O r_i$).⁸

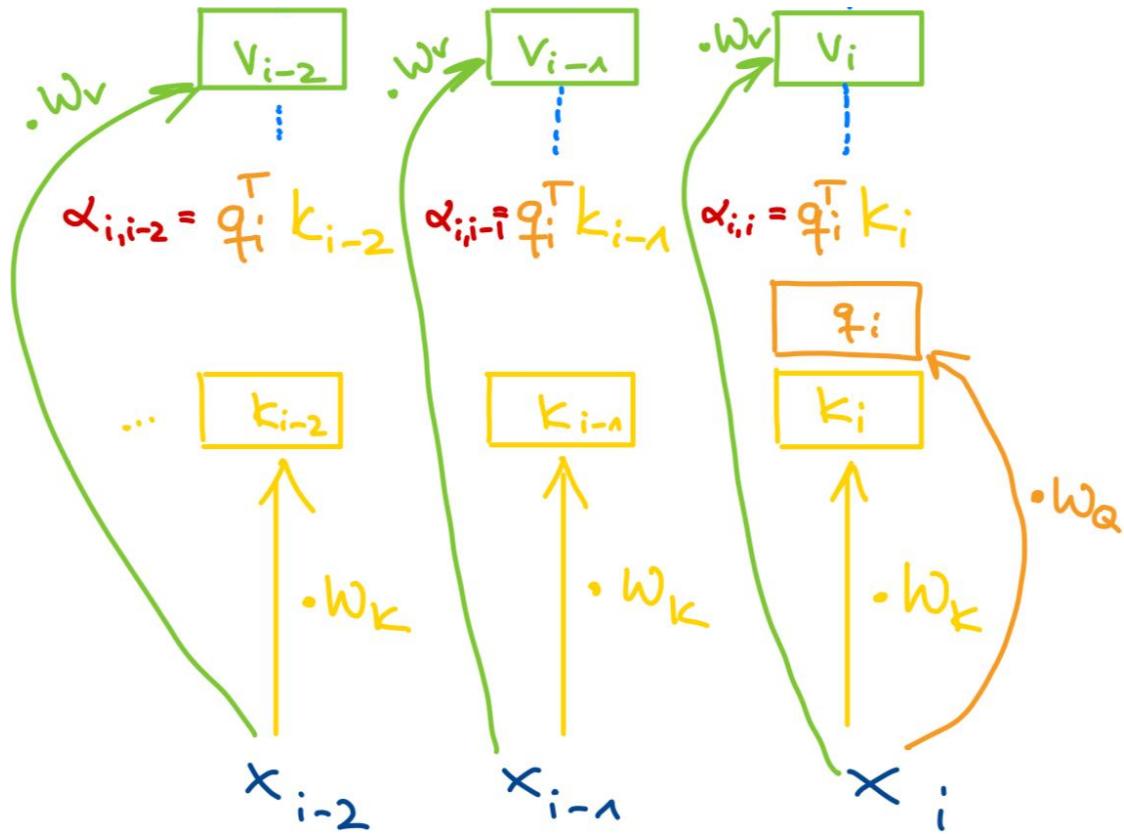
$$h(x) = (A \otimes W_O W_V) \cdot x$$

$$\left. \begin{array}{l} k_i = W_K x_i \\ q_i = W_Q x_i \\ A = \text{softmax}(q^T k) \end{array} \right\} A = \text{softmax}(x^T W_Q^T W_K x)$$

$$\alpha_{i,i-2} = q_i^T k_{i-2} \quad \alpha_{i,i-1} = q_i^T k_{i-1} \quad \alpha_{i,i} = q_i^T k_i$$



Query-Calculation
Key-Calculation

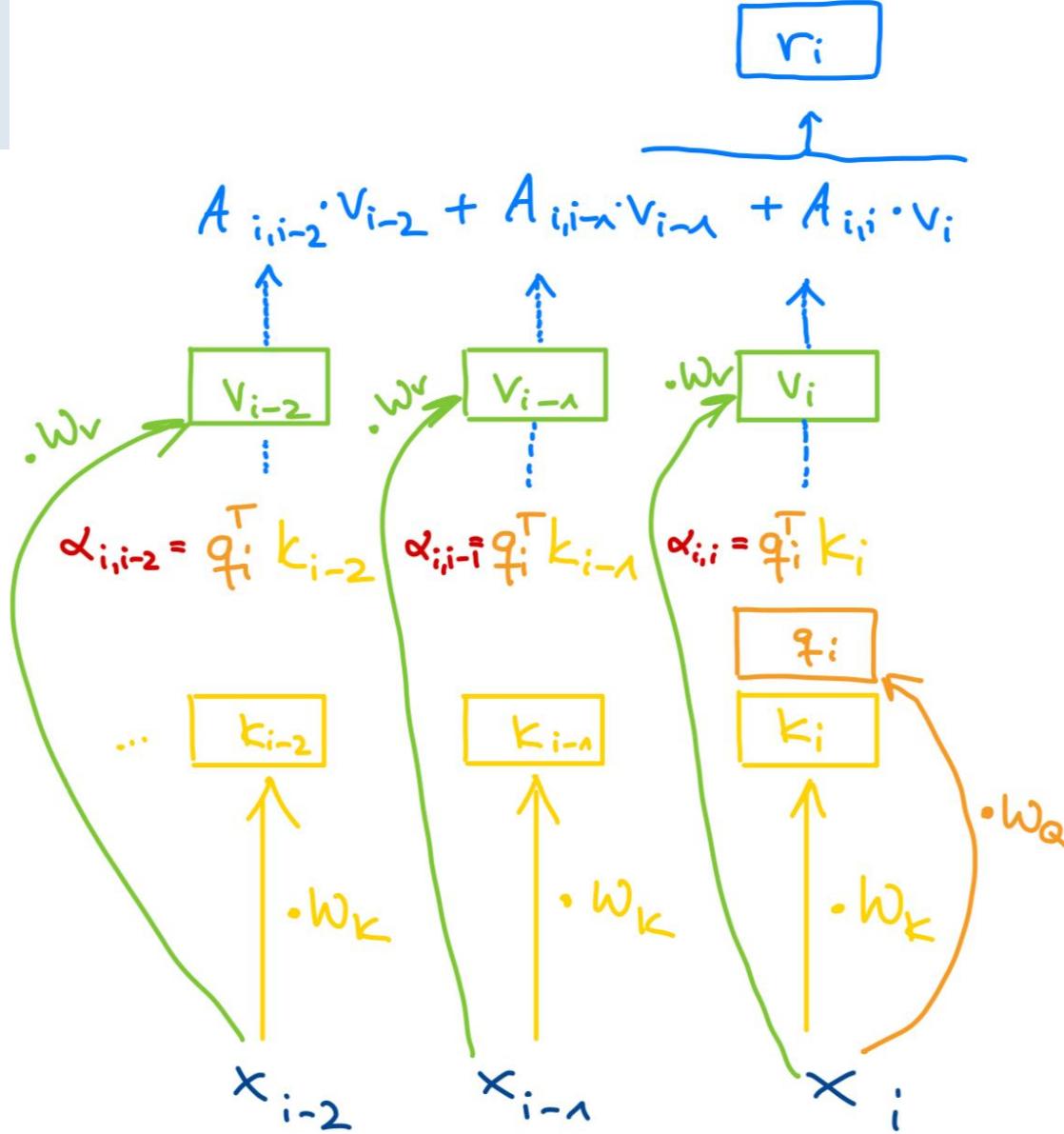


Value - Calculation

Query - Calculation

Key - Calculation

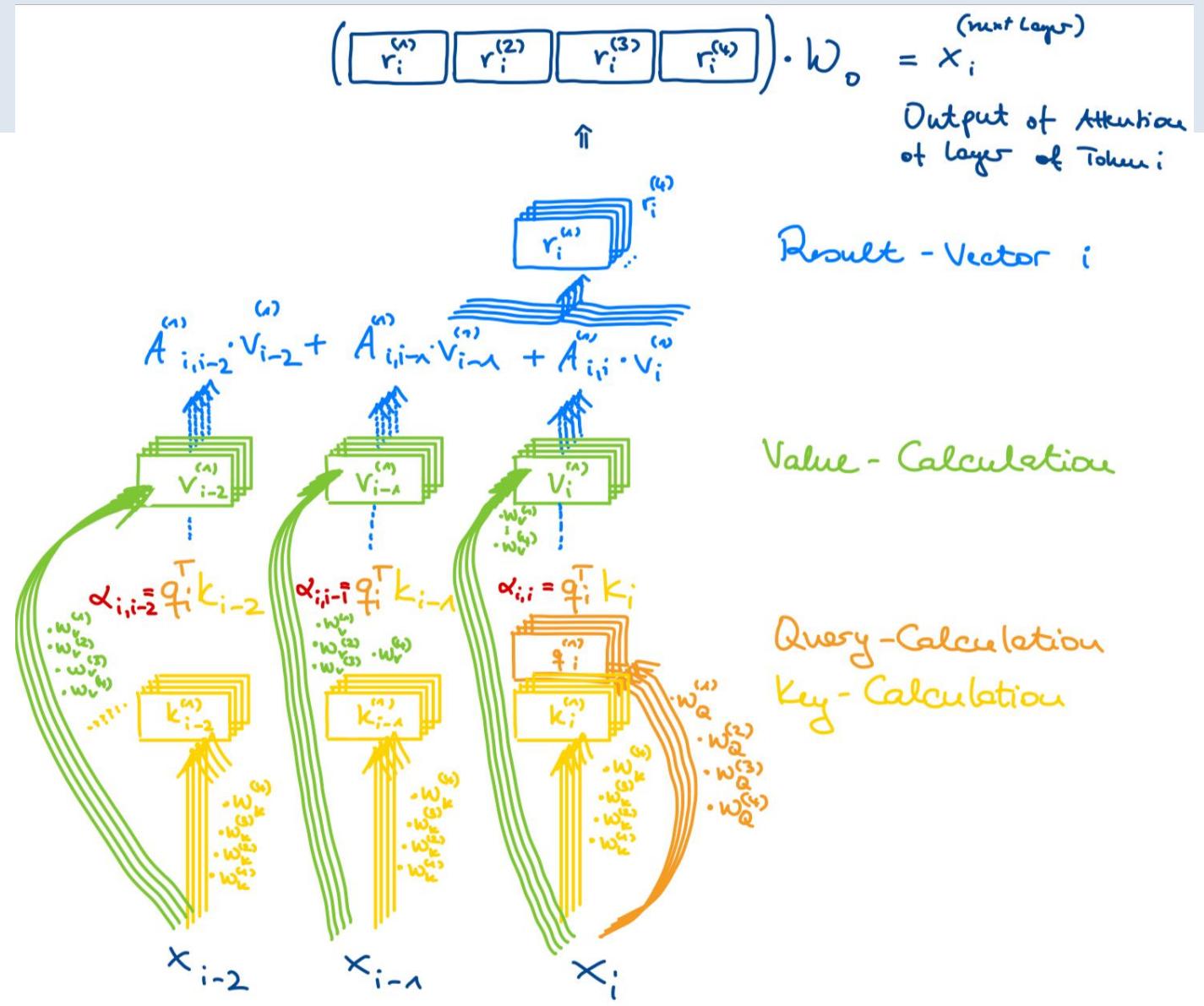
Result - Vector i



Value - Calculation

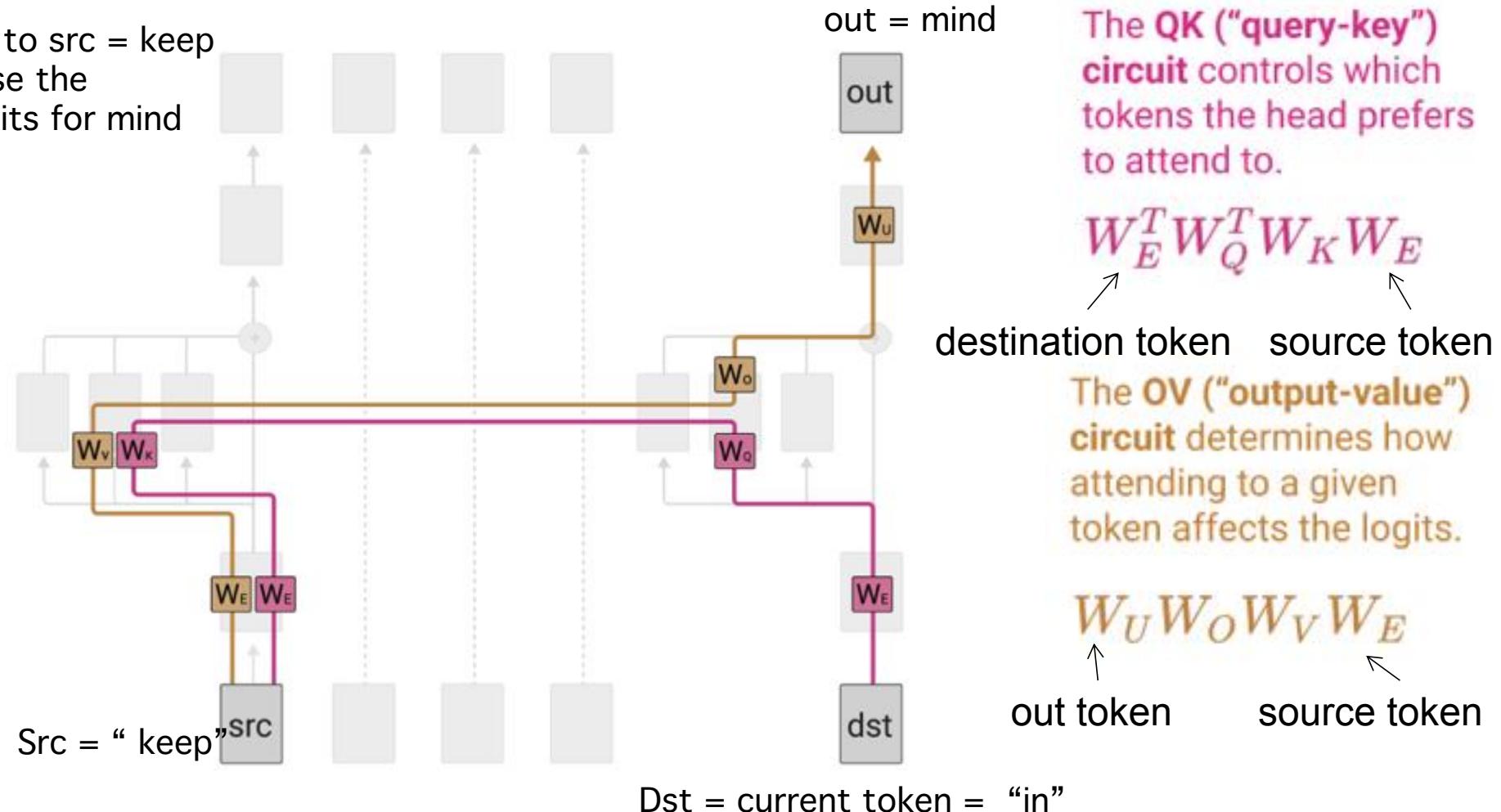
Query - Calculation
Key - Calculation

Multiple Heads parallel



Review Layer 1

Attending to src = keep
will increase the
output logits for mind



2-Layers Transformer

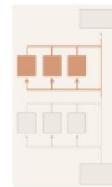
$$T = \underbrace{\text{Id} \otimes W_U}_{\text{---}} + \left(\text{Id} + \sum_{h \in H_2} A^h \otimes W_{OV}^h \right) + \left(\text{Id} + \sum_{h \in H_1} A^h \otimes W_{OV}^h \right) + \underbrace{\text{Id} \otimes W_E}_{\text{---}}$$

The second **attention layer** has multiple attention heads which add into the residual stream

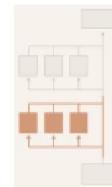
The first **attention layer** has multiple attention heads which add into the residual stream

2-Layers Transformer

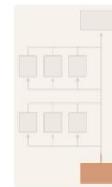
$$T = \underbrace{\text{Id} \otimes W_U}_{\text{---}} \cdot \left(\text{Id} + \sum_{h \in H_2} A^h \otimes W_{OV}^h \right) \cdot \left(\text{Id} + \sum_{h \in H_1} A^h \otimes W_{OV}^h \right) \cdot \underbrace{\text{Id} \otimes W_E}_{\text{---}}$$



The second **attention layer** has multiple attention heads which add into the residual stream



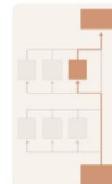
The first **attention layer** has multiple attention heads which add into the residual stream



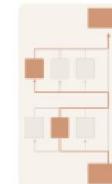
$$= \underbrace{\text{Id} \otimes W_U W_E}_{\text{---}} + \sum_{h \in H_1 \cup H_2} A^h \otimes (W_U W_{OV}^h W_E) + \sum_{h_2 \in H_2} \sum_{h_1 \in H_1} (A^{h_2} A^{h_1}) \otimes (W_U W_{OV}^{h_2} W_{OV}^{h_1} W_E)$$



“Direct path” term contributes to 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.



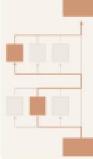
The **virtual attention head** terms correspond to V-composition of attention heads. They function a lot like individual attention heads, with their own attention patterns (the composition of the heads patterns) and own OV matrix.

2-Layers Transformer

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


"Direct path" term contributes to 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 V-composition of attention heads. They function a lot like individual attention heads, with their own attention patterns (the composition of the heads patterns) and own OV matrix.

Same as for 0 layer Same as for 1 layer New!

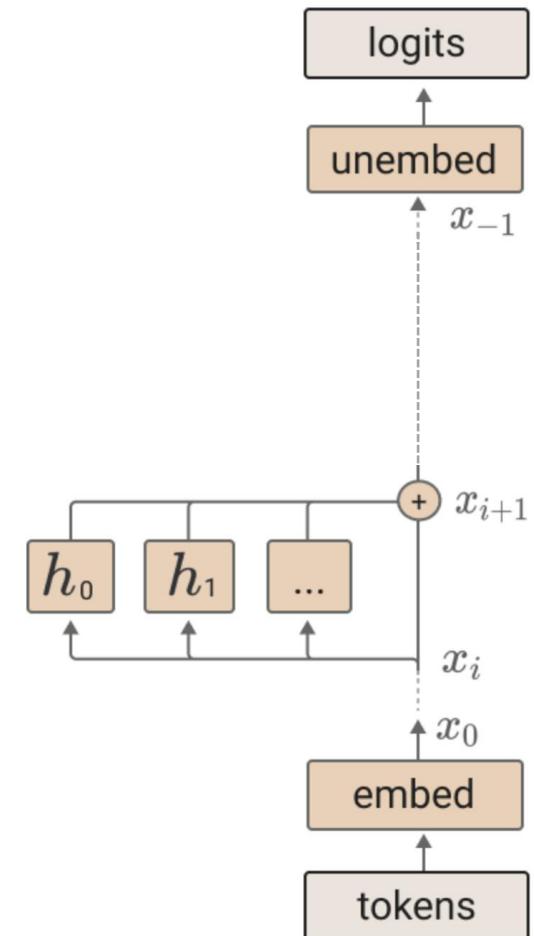
Calculating Attention

$$A^h = \text{softmax}^*(t^T \cdot C_{QK}^h t)$$

$$C_{QK}^{h \in H_1} = x_0^T W_{QK}^h x_0 = W_E^T W_{QK}^h W_E$$

$$C_{QK}^{h \in H_2} = x_1^T W_{QK}^h x_1$$

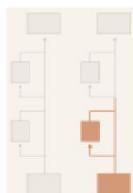
x_1 is residual stream after layer 1 attention



Calculating Attention Pattern in 2nd Layer

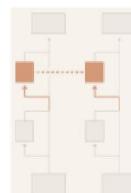
$$C_{QK}^{h \in H_2} = x_1^T W_{QK}^h x_1 \quad h(x) = (A \otimes W_O W_V) \cdot x \quad A = \text{softmax}(x^T W_Q^T W_K x)$$

$$= \left(\text{Id} \otimes \text{Id} \otimes W_E^T + \sum_{h_q \in H_1} A^{h_q} \otimes \text{Id} \otimes (W_{OV}^{h_q} W_E)^T \right)$$



The “**query side**” **residual stream** at the start of the second layer contains both the layer 1 direct path and layer 1 attention heads. All terms are of the form $\dots \otimes \text{Id} \otimes \dots$ because they don’t move key information.

$$\cdot \text{Id} \otimes \text{Id} \otimes W_{QK}^h \cdot$$



W_{QK} of the second layer head combines both sides into attention scores.

$$\left(\text{Id} \otimes \text{Id} \otimes W_E + \sum_{h_k \in H_1} \text{Id} \otimes A^{h_k} \otimes W_{OV}^{h_k} W_E \right)$$

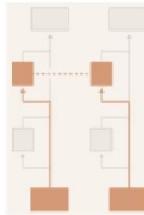


The “**key side**” **residual stream** at the start of the second layer contains both the layer 1 direct path and attention heads. All terms are of the form $\text{Id} \otimes \dots$ because they don’t move query information.

$$C_{QK}^{h \in H_2} = x_1^T W_{QK}^h x_1$$

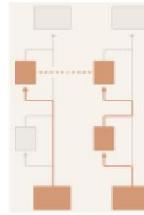
$$x_1 = x_0 + \sum_{h \in H_1} (A^h \otimes W_O W_V) x_0 = W_E \cdot t + \sum_{h \in H_1} (A^h \otimes W_O W_V) W_E \cdot t$$

$$= \frac{\text{Id} \otimes \text{Id} \otimes (W_E^T W_{QK}^h W_E)}{+}$$



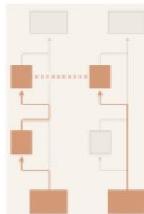
The no composition term. Both first layer follows the direct path on both the key and query side.

$$+ \sum_{h_q \in H_1} A^{h_q} \otimes \text{Id} \otimes \left(W_E^T W_{OV}^{h_q T} W_{QK}^h W_E \right)$$

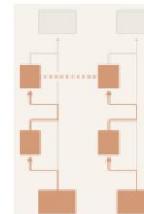


These terms correspond to pure **Q-composition**. A previous attention head is used to generate the query side, but the key side is the first layer direct path.

$$+ \sum_{h_k \in H_1} \text{Id} \otimes A^{h_k} \otimes \left(W_E^T W_{QK}^h W_{OV}^{h_k} W_E \right) + \sum_{h_q \in H_1} \sum_{h_k \in H_1} A^{h_q} \otimes A^{h_k} \otimes \left(W_E^T W_{OV}^{h_q T} W_{QK}^h W_{OV}^{h_k} W_E \right)$$



These terms correspond to pure **K-composition**. A previous attention head is used to generate part of the key, but the query side is the first layer direct path.



These terms are interactions between both **Q-composition** and **K-composition**. A previous attention head is used to generate the query and key sides.

Which L2 head uses information of which previous L1 Head?

The above diagram shows Q-, K-, and V-Composition between attention heads in the first and second layer. That is, how much does the query, key or value vector of a second layer head read in information from a given first layer head? This is measured by looking at the Frobenius norm of the product of the relevant matrices, divided by the norms of the individual matrices. For Q-Composition, $\|W_{QK}^{h_2T} W_{OV}^{h_1}\|_F / (\|W_{QK}^{h_2T}\|_F \|W_{OV}^{h_1}\|_F)$, for K-Composition $\|W_{QK}^{h_2} W_{OV}^{h_1}\|_F / (\|W_{QK}^{h_2}\|_F \|W_{OV}^{h_1}\|_F)$, for V-Composition $\|W_{OV}^{h_2} W_{OV}^{h_1}\|_F / (\|W_{OV}^{h_2}\|_F \|W_{OV}^{h_1}\|_F)$. By default, we subtract off the empirical expected amount for random matrices of the same shapes (most attention heads have a much smaller composition than random matrices). In practice, for this model, there is only significant K-composition, and only with one layer 0 head.

Q-Composition

$$\frac{\|W_{QK}^{h_2T} W_{OV}^{h_1}\|_F}{\|W_{QK}^{h_2T}\|_F \|W_{OV}^{h_1}\|_F}$$

K-Composition

$$\frac{\|W_{QK}^{h_2} W_{OV}^{h_1}\|_F}{\|W_{QK}^{h_2}\|_F \|W_{OV}^{h_1}\|_F}$$

V-Composition

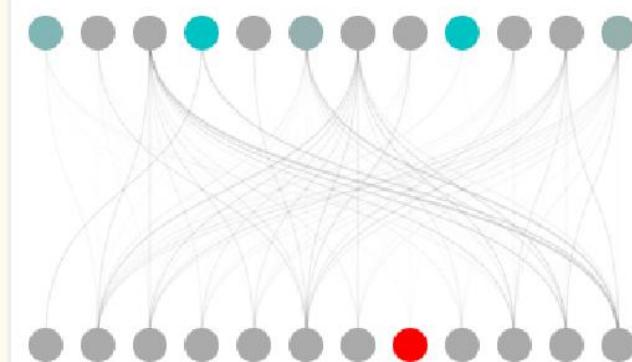
$$\frac{\|W_{OV}^{h_2} W_{OV}^{h_1}\|_F}{\|W_{OV}^{h_2}\|_F \|W_{OV}^{h_1}\|_F}$$

$$\|A\|_F = \sqrt{\sum_i^m \sum_j^n |a_{ij}|^2}$$

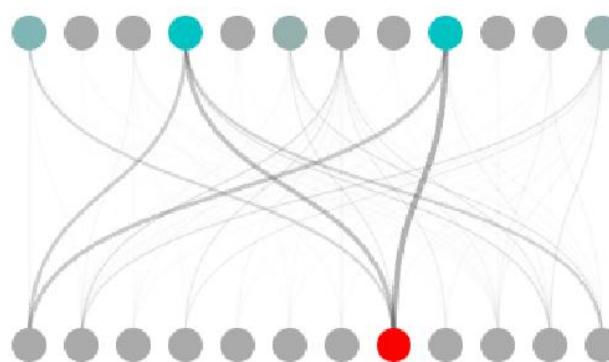
Which L2 head uses information of which previous L1 Head?

CORRECTED PLOT (BASELINE SUBTRACTED)

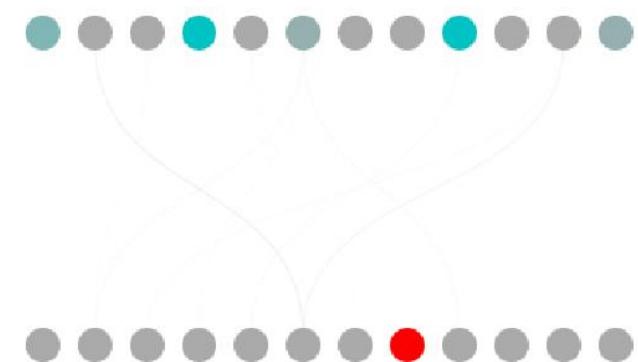
Q-Composition



K-Composition



V-Composition



- Induction Heads
- Previous Token Heads

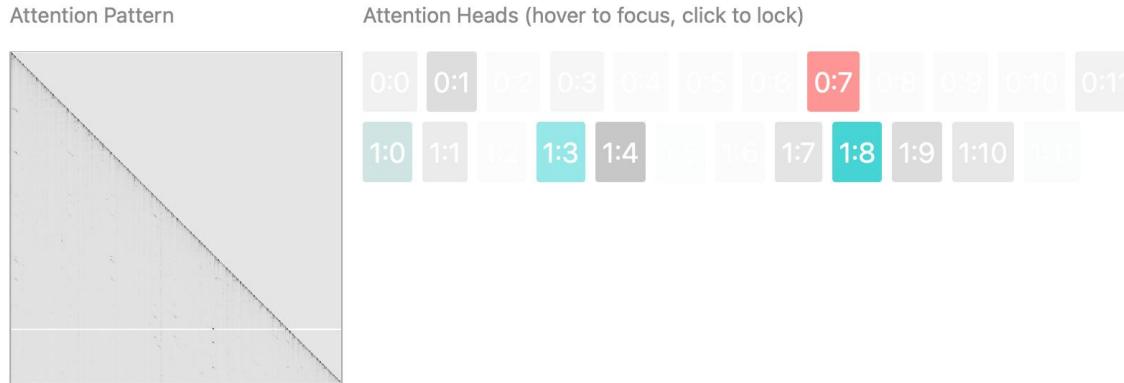
Note: Line widths between original and corrected plots are not directly comparable.

Which L2 head uses information of which previous L1 Head?

The above diagram shows the *value-weighted attention pattern* for various attention heads; that is, the attention patterns with attention weights scaled by the norm of the value vector at the source position $\|v_{src}^h\|$. You can think of the value-weighted attention pattern as showing "how big a vector is moved from each position." (This approach was also recently introduced by Kobayashi *et al.* [16].) This is especially useful because attention heads will sometimes use certain tokens as a kind of default or resting position when there isn't a token that matches what they're looking for; the value vector at these default positions will be small, and so the value weighted pattern is more informative.

The interface allows one to isolate attention heads, shows the overall attention pattern, and allows one to explore the attention for individual tokens. Attention heads involved in K-composition are colored using the same scheme as above. We suggest trying to isolate these heads.

When destination = “Pot”, which tokens are the source?

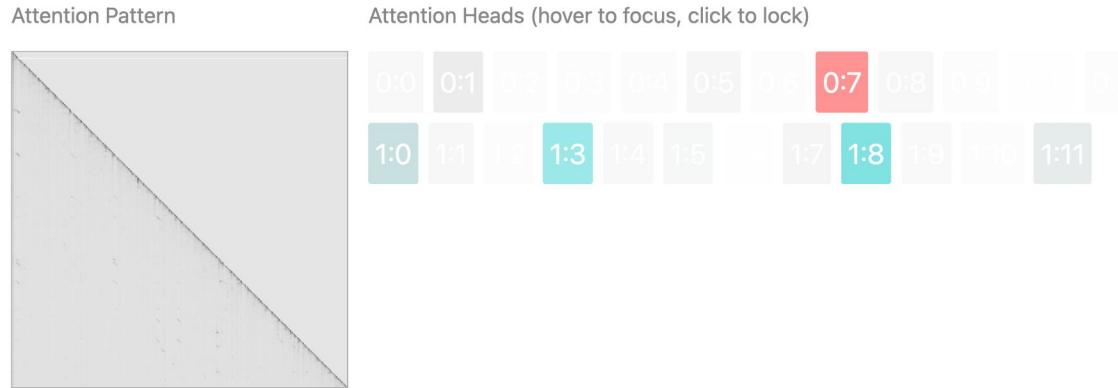


Tokens (hover to focus, click to lock)

source ← destination (destination attention when none selected) ↴

<START>Mr and Mrs Dursley, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs Dursley was thin and blonde and had nearly twice the usual amount of neck, which came in very useful as she spent so much of her time craning over garden fences, spying on the neighbours. The Dursleys had a small son called Dudley and in their opinion there was no finer boy anywhere. The Dursleys had everything they wanted, but they also had a secret, and their greatest fear was that somebody would discover it. They didn't think they could bear it if anyone found out about the Potters. Mrs Potter was Mrs Dursley's sister, but they hadn't met for several years; in fact, Mrs Dursley pretended she didn't have a sister, because her sister and her good-for-nothing husband were as undursleyish as it was possible to be. The Dursleys shuddered to think what the neighbours would say if the Potters arrived in the street. The Dursleys knew that the Potters had a small son, too, but they had never even seen him. This boy was another good reason for keeping the Potters away; they didn't want Dudley mixing with a child like that.

When reading which token does the model attend to ‘urs’ -> everytime ‘D’ is the current tokens, but also when predicting the following token of ‘urs’ (but here we use more Head 0:7)



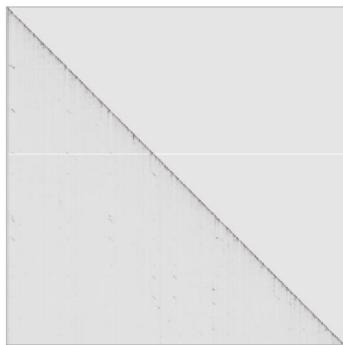
Tokens (hover to focus, click to lock)

source → destination (source attention when none selected)

<START> Mr and Mrs Dursley, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs Dursley was thin and blonde and had nearly twice the usual amount of neck, which came in very useful as she spent so much of her time craning over garden fences, spying on the neighbours. The Dursleys had a small son called Dudley and in their opinion there was no finer boy anywhere. The Dursleys had everything they wanted, but they also had a secret, and their greatest fear was that somebody would discover it. They didn't think they could bear it if anyone found out about the Potters. Mrs Potter was Mrs Dursley's sister, but they hadn't met for several years; in fact, Mrs Dursley pretended she didn't have a sister, because her sister and her good-for-nothing husband were as unDursleyish as it was possible to be. The Dursleys shuddered to think what the neighbours would say if the Potters arrived in the street. The Dursleys knew that the Potters had a small son, too, but they had never even seen him. This boy was another good reason for keeping the Potters away; they didn't want Dudley mixing with a child like that.

“When is “urs”
the source
token?

Attention Pattern



Attention Heads (hover to focus, click to lock)

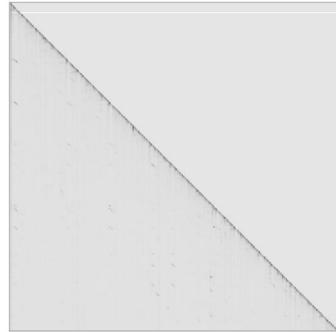


Tokens (hover to focus, click to lock)

source → destination (source attention when none selected) ▼

<START>Mr and Mrs Dursley, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs Dursley was thin and blonde and had nearly twice the usual amount of neck, which came in very useful as she spent so much of her time craning over garden fences, spying on the neighbours. The Dursleys had a small son called Dudley and in their opinion there was no finer boy anywhere. The Dursleys had everything they wanted, but they also had a secret, and their greatest fear was that somebody would discover it. They didn't think they could bear it if anyone found out about the Potters. Mrs Potter was Mrs Dursley's sister, but they hadn't met for several years; in fact, Mrs Dursley pretended she didn't have a sister, because her sister and her good-for-nothing husband were as unDursleyish as it was possible to be. The Dursleys shuddered to think what the neighbours would say if the Potters arrived in the street. The Dursleys knew that the Potters had a small son, too, but they had never even seen him. This boy was another good reason for keeping the Potters away; they didn't want Dudley mixing with a child like that.

Attention Pattern



Attention Heads (hover to focus, click to lock)



Tokens (hover to focus, click to lock)

source → destination (source attention when none selected) ▾

<START> Mr and Mrs Dursley, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs Dursley was thin and blonde and had nearly twice the usual amount of neck, which came in very useful as she spent so much of her time craning over garden fences, spying on the neighbours. The Dursleys had a small son called Dudley and in their opinion there was no finer boy anywhere. The Dursleys had everything they wanted, but they also had a secret, and their greatest fear was that somebody would discover it. They didn't think they could bear it if anyone found out about the Potters. Mrs Potter was Mrs Dursley's sister, but they hadn't met for several years; in fact, Mrs Dursley pretended she didn't have a sister, because her sister and her good-for-nothing husband were as unDursleyish as it was possible to be. The Dursleys shuddered to think what the neighbours would say if the Potters arrived in the street. The Dursleys knew that the Potters had a small son, too, but they had never even seen him. This boy was another good reason for keeping the Potters away; they didn't want Dudley mixing with a child like that.

If you look carefully, you'll notice that the aqua colored "induction heads" often attend back to previous instances of the token which will come next. We'll investigate this more in the next section.



<START>Mr and Mrs Dursley, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs Dursley was thin and blonde and had nearly twice the usual amount of neck, which came in very useful as she spent so much of her time craning over garden fences, spying on the neighbours. The Dursleys had a small son called Dudley and in their opinion there was no finer boy anywhere. The Dursleys had everything they wanted, but they also had a secret, and their greatest fear was that somebody would discover it. They didn't think they could bear it if anyone found out about the Potters. Mrs Potter was Mrs Dursley's sister, but they hadn't met for several years; in fact, Mrs Dursley pretended she didn't have a sister, because her sister and her good-for-nothing husband were as unDursleyish as it was possible to be. The Dursleys shuddered to think what the neighbours would say if the Potters arrived in the street. The Dursleys knew that the Potters had a small son, too, but they had never seen him. This was another thing they liked about living in the Dursley street.

Induction Heads

HOW INDUCTION HEADS WORK

The central trick to induction heads is that the key is computed from tokens shifted one token back.²⁰ The query searches for "similar" key vectors, but because keys are shifted, finds the next token.

The diagram shows three rows of text, each containing the words "out about the Potters. Mrs Potter was ... neighbours would say if the".

- Top row:** The word "Potters" is highlighted in blue. A red bracket labeled "attention pattern moves information" spans from the first "Potters" to the second "Potters". A purple arrow labeled "logit effect" points from the second "Potters" to the third "Potters".
- Middle row:** The word "Potters" is highlighted in red. A green bracket labeled "key" spans from the first "Potters" to the second "Potters". A green arrow labeled "query" points from the second "Potters" to the third "Potters".
- Bottom row:** The word "Potters" is highlighted in green. A blue bracket labeled "query" spans from the first "Potters" to the second "Potters". A blue arrow labeled "logit effect" points from the second "Potters" to the third "Potters".

Induction Heads

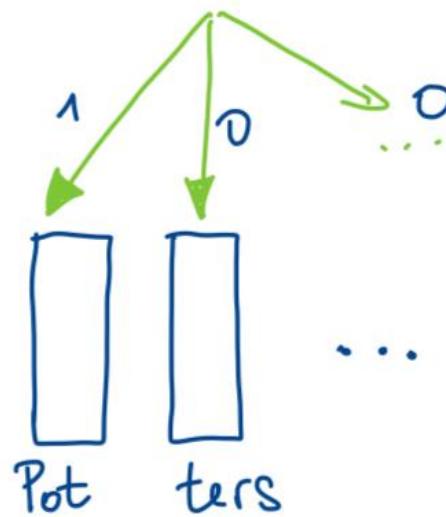
The minimal way to create an induction head is to use K-composition with a previous token head to shift the key vector forward one token. This creates a term of the form $\text{Id} \otimes A^{h-1} \otimes W$ in the QK-circuit (where A^{h-1} denotes an attention pattern attending to the previous token). If W matches cases where the tokens are the same — the QK version of a "copying matrix" — then this term will increase attention scores when the previous token before the source position is the same as the destination token.

Induction Heads

Previous
Token
Head

$$x_{1,t_\text{ns}} = x_{0,\text{Pot}}$$

$$x_{1,t_\text{ns}} = 1 \cdot x_{0,\text{Pot}} + 0 \cdot x_{0,\text{ters}} + \dots$$



Layer 1

Induction Heads

Induction Head

K-COMP. =

only W_K
reads information
from h_1 only,
(W_V, W_Q read information
from W_E)

Previous Token Head

$$x_{2,Pot} = \sum \alpha_i \cdot v_i = 0 \cdot \dots + 1 \cdot v_{ters} + 0 \cdot \dots$$

$$k = x_{1,ters} = x_{0,Pot}$$

$$q^T k = 1 !$$

$$q = x_{0,Pot}$$

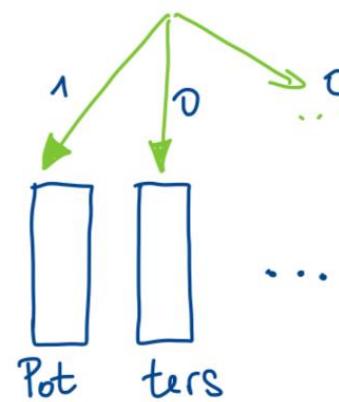
$$x_{0,ters} \cdot w_v$$

Layer 2

$$\text{softmax}^* W_E^T W_Q^T W_K W_E \approx 1$$

$$x_{1,ters} = x_{0,Pot}$$

$$x_{1,ters} = 1 \cdot x_{0,Pot} + 0 \cdot x_{0,ters} + \dots$$



Layer 1

Induction Heads



$$W_U \cdot W_{OV}^{h_2} \cdot x_{0, \text{ters}} \Rightarrow \text{"ters"}$$

Output

Copy OV-Circuit:

$$W_U \cdot W_{OV}^{h_2} \cdot W_E \approx 1$$

--> simply enhances the logit of the token that "goes in"



K-Comp. =
only W_K
reads information
from h_1 only,
(W_V, W_Q read information
from W_E)

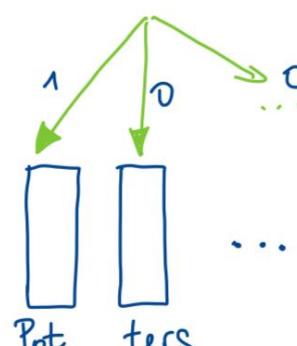


$$x_{2, \text{pot}} = \sum \alpha_i \cdot v_i = 0 \dots + 1 \cdot v_{\text{ters}} + 0 \dots$$

$$k = x_{1, \text{ters}} = x_{0, \text{pot}}$$

$q^T k = 1 !$

$$x_{1, \text{ters}} = 1 \cdot x_{0, \text{pot}} + 0 \cdot x_{0, \text{ters}} + \dots$$



$$\text{softmax}^* W_E^T W_Q^T W_K W_E \approx 1$$

Layer 2

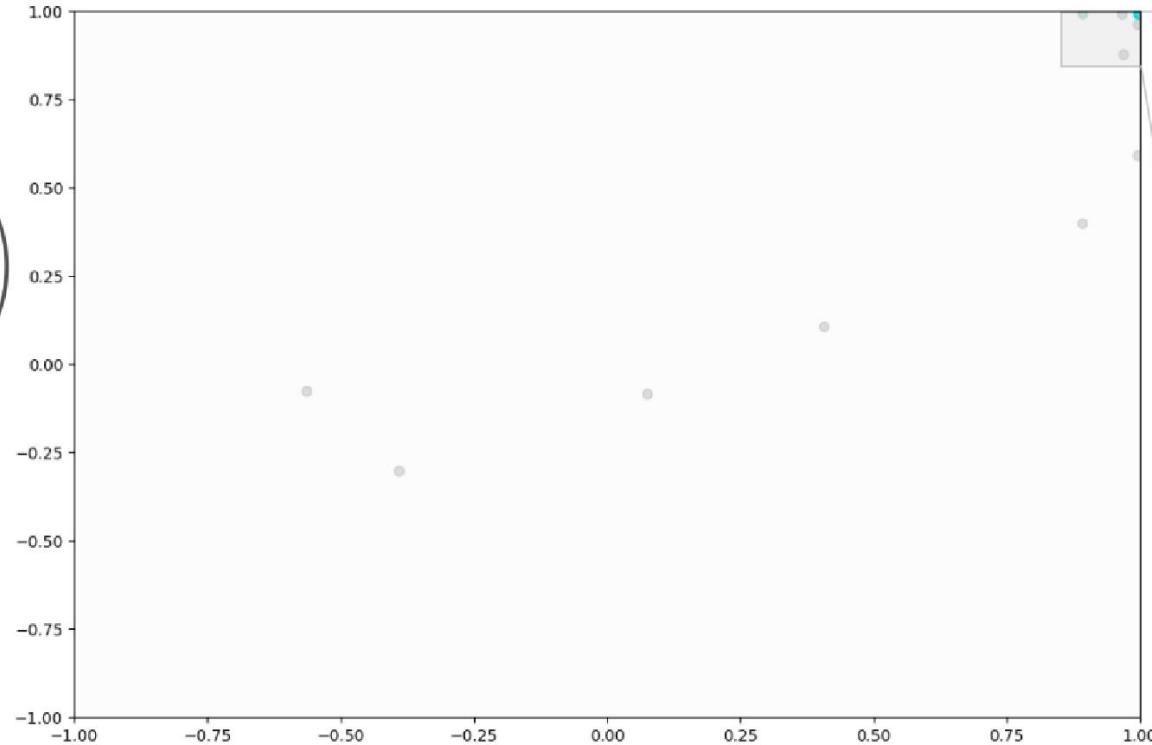
Layer 1

Find these Heads

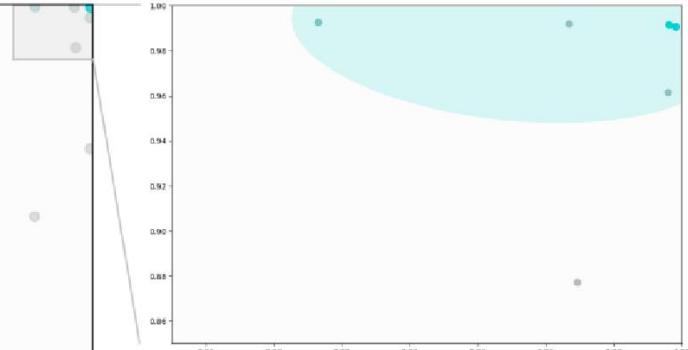
Eigenvalue positivity of W in the QK term
 $\text{Id} \otimes A^{h-1} \otimes W$

$$\left(\frac{\sum \lambda_i}{\sum |\lambda_i|} \text{ of } W_E^T W_{QK}^h W_{OV}^{h-1} W_E \right)$$

Second Layer Attention Heads by OV and QK Eigenvalue Positivity



All the induction heads are at the extreme of OV and QK positivity

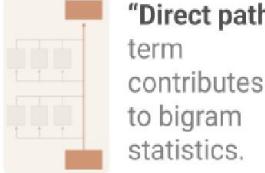


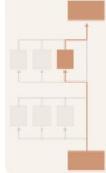
Eigenvalue positivity of the OV circuit

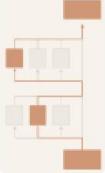
$$\left(\frac{\sum \lambda_i}{\sum |\lambda_i|} \text{ of } W_U W_{OV}^h W_E \right)$$

Which terms are actually important?

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

 "Direct path" term contributes to 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 V-composition of attention heads. They function a lot like individual attention heads, with their own attention patterns (the composition of the heads patterns) and own OV matrix.

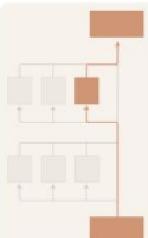
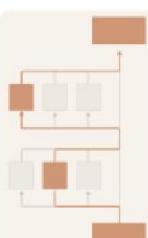
ALGORITHM FOR MEASURING MARGINAL LOSS REDUCTION OF N TH ORDER TERMS

Step 1: Run model, save all attention patterns.

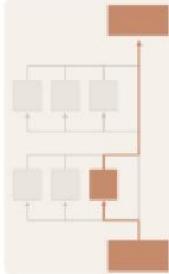
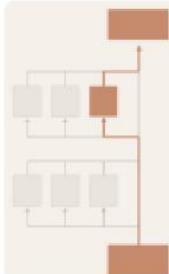
Step 2: Run model, forcing all attention patterns to be the version you recorded, and instead of adding attention head outputs to residual stream, save the output, and then replace it with a zero tensor of the same shape. Record resulting loss.

Step n : Run model, forcing all attention patterns to be the version you recorded, and instead of adding attention head outputs to residual stream, save the output, and replace it with the value you saved for this head last time. Record resulting loss.

Which terms are actually important?

Type	Example	Equation	Marginal Loss Reduction
direct path order 0		$W_U W_E$	- 1.8 nats relative to uniform predictions -1.8 nats/term (- 1.8 nats / 1 term)
individual attention head order 1		$A^h \otimes (W_U W_{OV}^h W_E)$	- 5.2 nats relative to only using direct path -0.2 nats/term (5.2 nats / 24 terms)
virtual attention head order 2		$(A^{h_2} A^{h_1}) \otimes (W_U W_{OV}^{h_2} W_{OV}^{h_1} W_E)$	- 0.3 nats relative to only using above -0.002 nats/term (0.3 nats / 144 terms)

Individual Attention Heads

Type	Example	Marginal Loss Reduction	Notes
Layer 1 Attention Heads		- 0.05 nats -0.004 nats/head relative to direct path + layer 2 - 1.3 nats -0.1 nats/head relative to direct path only	Relatively small effect, but keep in mind these heads also contribute to layer 2 QK circuits.
Layer 2 Attention Heads		- 4.0 nats -0.3 nats/head relative to direct path + layer 1 - 5.2 nats -0.4 nats/head relative to direct path only	We'll focus on these. Much larger effect. These heads are a lot more sophisticated than the layer 1 heads, since they can use layer 1 heads in their QK circuits.