

An Alternate Formulation of Transformers

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821



Tanmoy Chakraborty
Associate Professor, IIT Delhi
<https://tanmoychak.com/>



OpenAI introduces ChatGPT search !

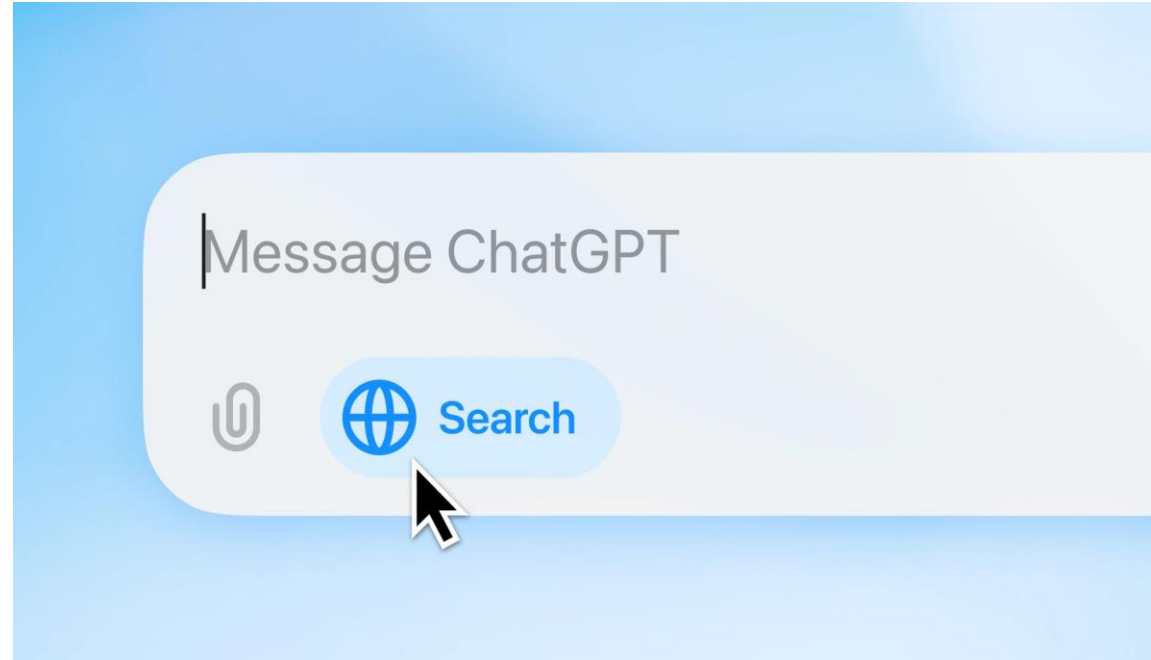
Announced on
October 31, 2024

[OpenAI Blog](#)

We can now get answers with links to relevant web sources!

The **search model** is a **fine-tuned version of GPT-4o**, **post-trained using novel synthetic data generation techniques**, including **distilling outputs from OpenAI o1-preview**.

ChatGPT search leverages third-party search providers to provide the information users are looking for.



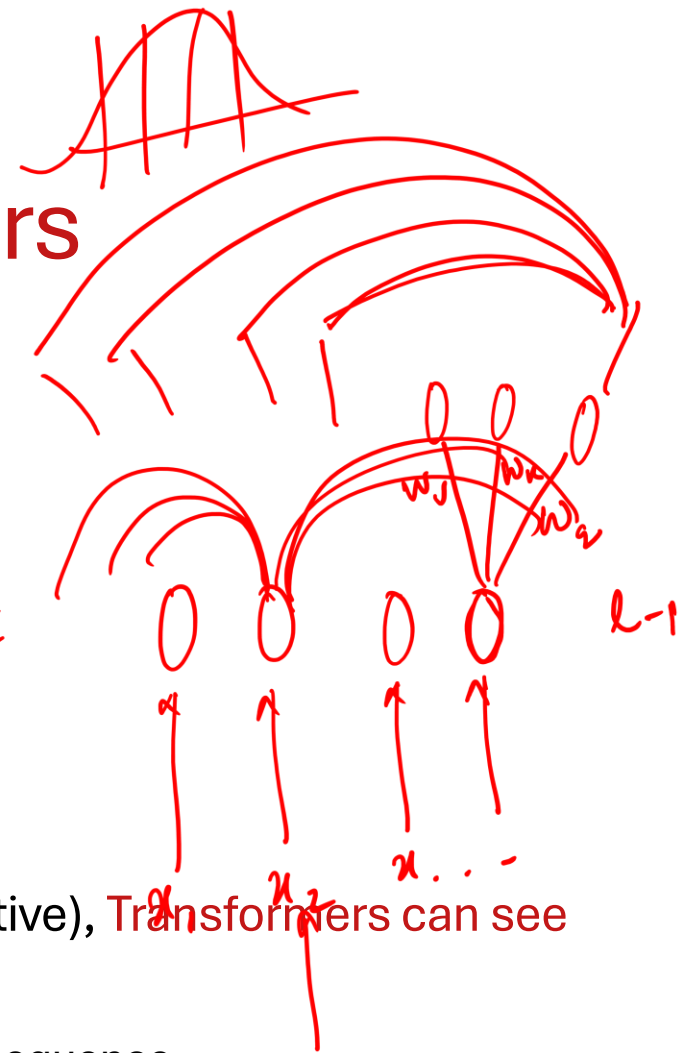
OpenAI has partnered with news and data providers to add **up-to-date information for categories like weather, stocks, sports, news, and maps**. ChatGPT will choose to **search the web based on what we ask**, or **we can manually choose to search by clicking the web search icon**.

Recall: Masked Self-Attention in Decoders

Self-Attention: Scaled dot-product attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{where, } Q = XW^Q, K = XW^K, V = XW^V$$



Problem: While training autoregressive models (with next-word-prediction objective), Transformers can see the future.

- For a current token x_i , the attention scores are computed with all tokens in the sequence including those which comes after x_i (as the whole sequence is available to us during training).

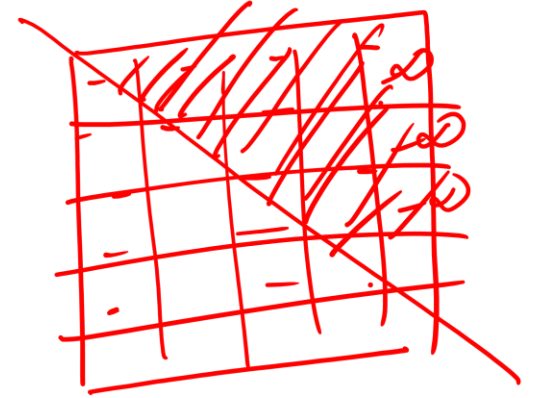
Solution: Masking



Recall: Masked Self-Attention in Decoders

Masking: 'Masked' scaled dot-product attention

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} + M \right) V$$



where, masking matrix M is defined as:

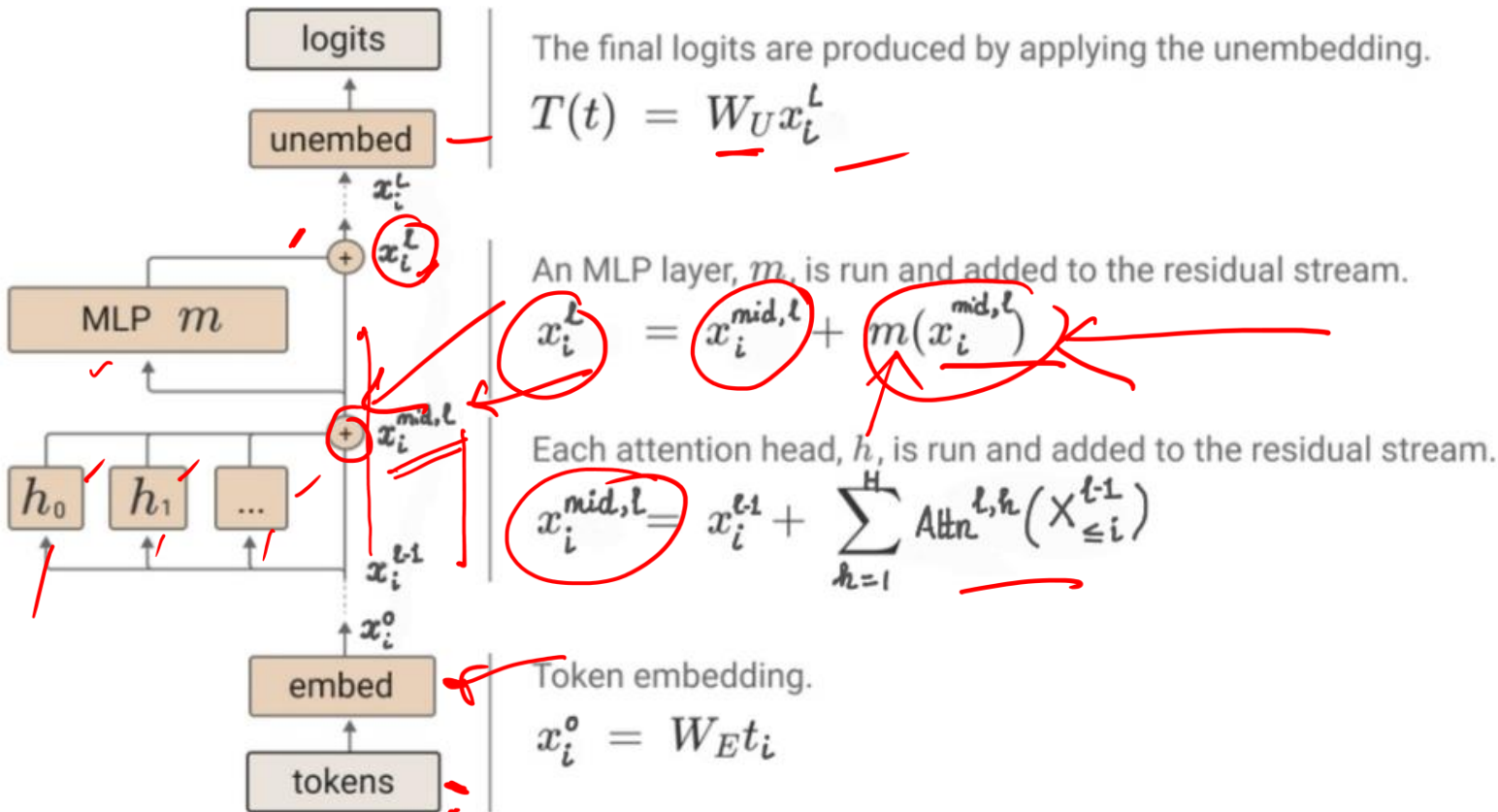
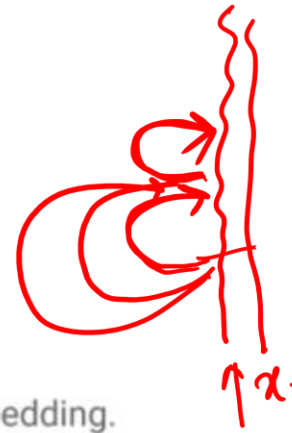
$$M_{ij} = \begin{cases} 0 & \text{if } j \leq i \\ -\infty & \text{if } j > i \end{cases}$$

For future tokens, the attention scores become zero after applying softmax [$\text{softmax}(-\infty) = 0$].

- Effectively, **after masking**, the **query is the current token x_i** , and the **keys and values come from the tokens before it, including itself (i.e., $x_j, j \leq i$)**.



Residual Stream Perspective



The final logits are produced by applying the unembedding.

$$T(t) = W_U x_i^l$$

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

$$x_i^l = x_i^{mid,l} + m(x_i^{mid,l})$$

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

$$x_i^{mid,l} = x_i^{l-1} + \sum_{h=1}^H \text{Attn}^{l,h}(X_{\leq i}^{l-1})$$

Token embedding.

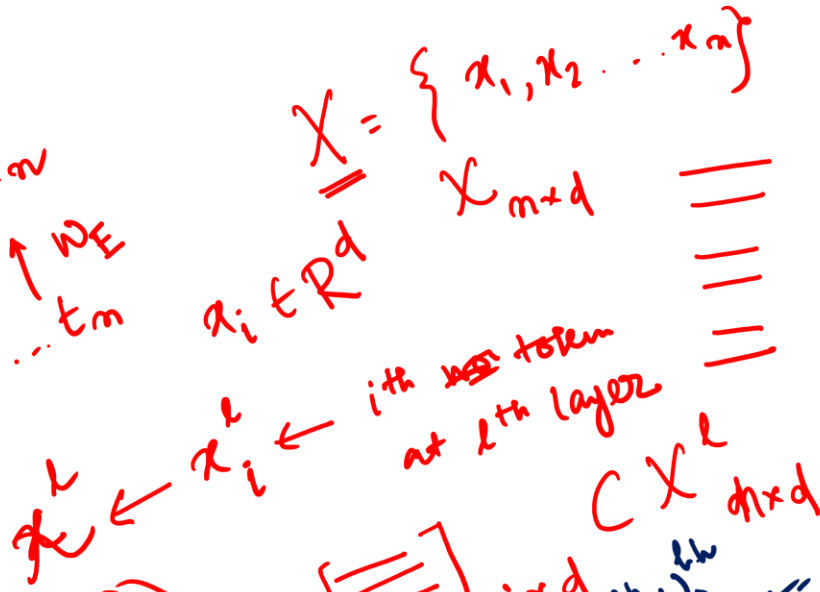
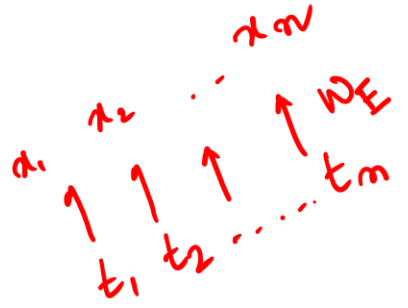
$$x_i^0 = W_{E} t_i$$

One residual block

- Each input embedding gets updated via vector additions from the attention and feed-forward blocks producing **residual stream states** (or intermediate representations).
- The final layer residual stream state is then projected into the vocabulary space via the unembedding matrix $W_U \in R^{d \times |V|}$ and normalized via the *softmax*.

Elhage, et al., A Mathematical Framework for Transformer Circuits





$W_{ov}^{l,h} = W_v^{l,h} W_o^{l,h}$
Output Value Circuit

$$Attn^{l,h}(x_{\leq i}^{l-1}) = \sum_{j \leq i} a_{i,j}^{l,h} (x_j^{l-1} W_v^{l,h}) W_o^{l,h} a_{i,j}^{l,h}$$

Diagram illustrating the attention mechanism. The input $x_{\leq i}^{l-1}$ (dimension $1 \times d$) is multiplied by $W_v^{l,h}$ (dimension $d \times d_h$) to produce $(x_j^{l-1} W_v^{l,h})$ (dimension $1 \times d_h$). This is then multiplied by $W_o^{l,h}$ (dimension $d_h \times d$) to produce the output $a_{i,j}^{l,h}$ (dimension $1 \times d$).

Attention Block

$$a_i^{l,h} = \text{Softmax} \left(x_i^{l-1} W_q^{l,h} (x_{\leq i}^{l-1} W_k^{l,h})^T \right)$$

Diagram illustrating the attention mechanism. The input x_i^{l-1} (dimension $1 \times d$) is multiplied by $W_q^{l,h}$ (dimension $d \times d_h$) to produce $x_i^{l-1} W_q^{l,h}$ (dimension $1 \times d_h$). The input $x_{\leq i}^{l-1}$ (dimension $i \times d$) is multiplied by $W_k^{l,h}$ (dimension $d \times d_h$) to produce $(x_{\leq i}^{l-1} W_k^{l,h})^T$ (dimension $d_h \times i$). The dot product of these two vectors is $(1 \times d_h) \cdot (d_h \times i) = 1 \times i$. The result is passed through a Softmax function to produce the attention score $a_i^{l,h}$ (dimension 1×1).

$$= \left(x_i^{l-1} W_q^{l,h} W_k^{l,h} x_{\leq i}^{l-1} \right) \left(W_q^{l,h} W_k^{l,h} \right)^T$$

Diagram illustrating the attention mechanism. The input x_i^{l-1} (dimension $1 \times d$) is multiplied by $W_q^{l,h}$ (dimension $d \times d_h$) to produce $x_i^{l-1} W_q^{l,h}$ (dimension $1 \times d_h$). The input $x_{\leq i}^{l-1}$ (dimension $i \times d$) is multiplied by $W_k^{l,h}$ (dimension $d \times d_h$) to produce $(x_{\leq i}^{l-1} W_k^{l,h})^T$ (dimension $d_h \times i$). The dot product of these two vectors is $(1 \times d_h) \cdot (d_h \times i) = 1 \times i$. The result is passed through a Softmax function to produce the attention score $a_i^{l,h}$ (dimension 1×1).

$$Attn^l(x_{\leq i}^{l-1}) = \sum_{k=1}^{\#} Attn^{l,h}(x_{\leq i}^{l-1})$$

Diagram illustrating the attention mechanism. The input $x_{\leq i}^{l-1}$ (dimension $i \times d$) is processed by the attention block to produce the output x_i^l (dimension $1 \times d$).



Re-writing the Masked Self-Attention Equation

Now let's re-write the masked attention equation for a current token x_i .

- Assume that we are considering the **attention head h** of **layer l** .
- Let's denote the matrix with the **output hidden representation from layer k of previous tokens $x_j, j \leq i$** as $X_{\leq i}^k$.

Thus, for calculating attention scores for **attention head h** of **layer l** , input to the attention sub-layer is the **output representation from the previous layer $l-1$** .

• **Query:** $x_i^{l-1} W_Q^{l,h}$

• **Keys:** $X_{\leq i}^{l-1} W_K^{l,h}$

$$a_i^{l,h} = \text{softmax} \left(\frac{\begin{array}{|c|c|} \hline \text{Query vector} & \\ \hline x_i^{l-1} W_Q^{l,h} & (X_{\leq i}^{l-1} W_K^{l,h})^\top \\ \hline \end{array}}{\sqrt{d_k}} \right)$$

Key vector

Ferrando et al., A Primer on the Inner Workings of Transformer-based Language Models





QK Circuit

$$\begin{aligned} \mathbf{a}_i^{l,h} &= \text{softmax} \left(\frac{\begin{array}{|c|c|} \hline \mathbf{x}_i^{l-1} \mathbf{W}_Q^{l,h} & (\mathbf{X}_{\leq i}^{l-1} \mathbf{W}_K^{l,h})^\top \\ \hline \end{array}}{\sqrt{d_k}} \right) \\ &= \text{softmax} \left(\frac{\mathbf{x}_i^{l-1} \mathbf{W}_{QK}^h \mathbf{X}_{\leq i}^{l-1 \top}}{\sqrt{d_k}} \right), \end{aligned}$$

Diagram annotations: A red arrow labeled "Query vector" points to the $\mathbf{x}_i^{l-1} \mathbf{W}_Q^{l,h}$ term. A green arrow labeled "Key vector" points to the $(\mathbf{X}_{\leq i}^{l-1} \mathbf{W}_K^{l,h})^\top$ term. A red circle highlights the \mathbf{W}_{QK}^h term in the simplified equation, with a red arrow pointing to the text below.

QK (query-key) circuit: $\mathbf{W}_{QK}^h = \mathbf{W}_Q^h \mathbf{W}_K^{h \top}$

- QK circuits are responsible for reading from the **residual stream**.

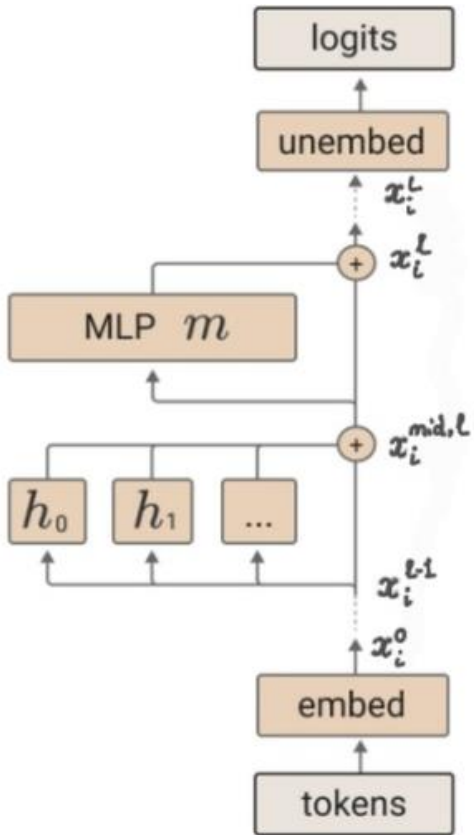
Let's now look at the residual stream

Ferrando et al., A Primer on the Inner Workings of Transformer-based Language Models





Residual Stream Perspective



The final logits are produced by applying the unembedding.

$$T(t) = W_U x_i^l$$

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

$$x_i^l = x_i^{mid,l} + m(x_i^{mid,l})$$

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

$$x_i^{mid,l} = x_i^{l-1} + \sum_{h=1}^H \text{Attn}^{l,h}(X_{\leq i}^{l-1})$$

Token embedding.

$$x_i^0 = W_{E} t_i$$

One residual block

- Each input embedding gets updated via vector additions from the attention and feed-forward blocks producing **residual stream states** (or intermediate representations).
- The final layer residual stream state is then projected into the vocabulary space via the unembedding matrix $W_U \in R^{d \times |V|}$ and normalized via the *softmax*.

Elhage, et al., A Mathematical Framework for Transformer Circuits



Combining the Output of Multiple Attention Heads

$$\begin{aligned}
 \mathbf{a}_i^{l,h} &= \text{softmax} \left(\frac{\overbrace{\mathbf{x}_i^{l-1} \mathbf{W}_Q^{l,h}}^{\text{Query vector}} \underbrace{(\mathbf{X}_{\leq i}^{l-1} \mathbf{W}_K^{l,h})^\top}_{\text{Key vector}}}{\sqrt{d_k}} \right) \\
 &= \text{softmax} \left(\frac{\mathbf{x}_i^{l-1} \mathbf{W}_{QK}^{l,h} \mathbf{X}_{\leq i}^{l-1 \top}}{\sqrt{d_k}} \right),
 \end{aligned}$$

$$\begin{aligned}
 \text{Attn}^{l,h}(\mathbf{X}_{\leq i}^{l-1}) &= \sum_{j \leq i} a_{i,j}^{l,h} \underbrace{\mathbf{x}_j^{l-1} \mathbf{W}_V^{l,h}}_{\text{Value vector}} \mathbf{W}_O^{l,h} \\
 &= \sum_{j \leq i} a_{i,j}^{l,h} \mathbf{x}_j^{l-1} \mathbf{W}_{OV}^{l,h},
 \end{aligned}$$

Ferrando et al., A Primer on the Inner Workings of Transformer-based Language Models



OV Circuit

$$\begin{aligned}\text{Attn}^{l,h}(\mathbf{X}_{\leq i}^{l-1}) &= \sum_{j \leq i} a_{i,j}^{l,h} \mathbf{x}_j^{l-1} \mathbf{W}_V^{l,h} \mathbf{W}_O^{l,h} \\ &= \sum_{j \leq i} a_{i,j}^{l,h} \mathbf{x}_j^{l-1} \mathbf{W}_{OV}^{l,h}\end{aligned}$$

Value vector

OV (output-value) circuit: $W_{OV}^{l,h} = W_V^{l,h} W_O^{l,h}$

- OV circuits are responsible for writing to the **residual stream**.

Ferrando et al., A Primer on the Inner Workings of Transformer-based Language Models



Attention Block Output

The attention block output is the **sum of individual attention heads**, which is subsequently added back into the residual stream.

$$\begin{aligned}\text{Attn}^{l,h}(\mathbf{X}_{\leq i}^{l-1}) &= \sum_{j \leq i} a_{i,j}^{l,h} \mathbf{x}_j^{l-1} \mathbf{W}_V^{l,h} \mathbf{W}_O^{l,h} \\ &= \sum_{j \leq i} a_{i,j}^{l,h} \mathbf{x}_j^{l-1} \mathbf{W}_{OV}^{l,h},\end{aligned}$$

Value vector

$$\begin{aligned}a_i^{l,h} &= \text{softmax} \left(\frac{\mathbf{x}_i^{l-1} \mathbf{W}_Q^{l,h} (\mathbf{X}_{\leq i}^{l-1} \mathbf{W}_K^{l,h})^\top}{\sqrt{d_k}} \right) \\ &= \text{softmax} \left(\frac{\mathbf{x}_i^{l-1} \mathbf{W}_{QK}^{l,h} \mathbf{X}_{\leq i}^{l-1 \top}}{\sqrt{d_k}} \right),\end{aligned}$$

Query vector

Key vector

$$\begin{aligned}\text{Attn}^l(\mathbf{X}_{\leq i}^{l-1}) &= \sum_{h=1}^H \text{Attn}^{l,h}(\mathbf{X}_{\leq i}^{l-1}) \\ \mathbf{x}_i^{\text{mid},l} &= \mathbf{x}_i^{l-1} + \text{Attn}^l(\mathbf{X}_{\leq i}^{l-1}).\end{aligned}$$

Ferrando et al., A Primer on the Inner Workings of Transformer-based Language Models

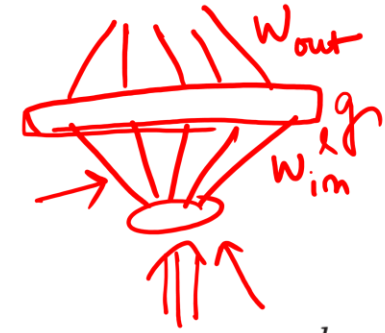




Feed-Forward Network (FFN)

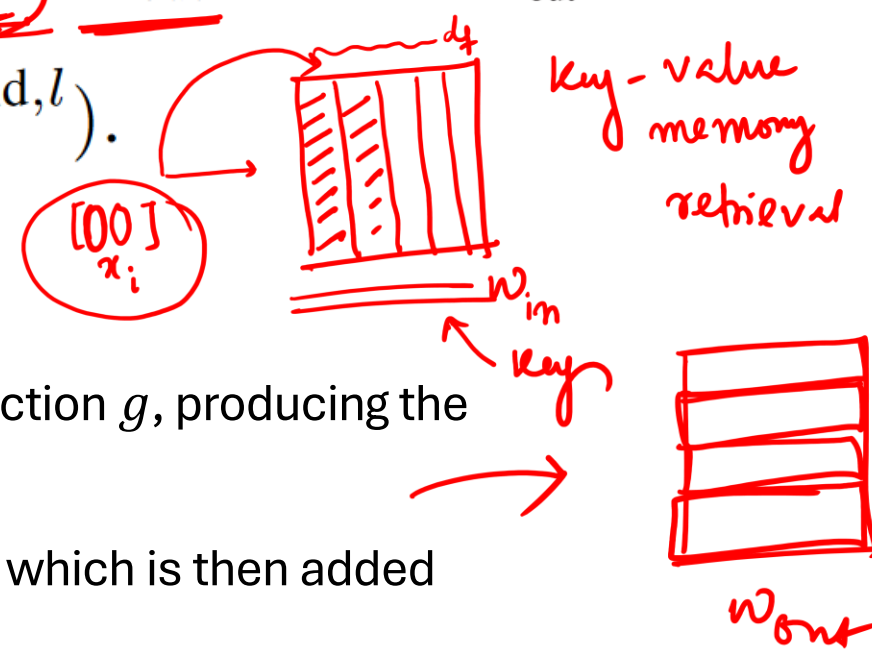
$$\text{FFN}^l(\mathbf{x}_i^{\text{mid},l}) = g(\mathbf{x}_i^{\text{mid},l} \mathbf{W}_{\text{in}}^l) \mathbf{W}_{\text{out}}^l$$

$$\mathbf{x}_i^l = \mathbf{x}_i^{\text{mid},l} + \text{FFN}^l(\mathbf{x}_i^{\text{mid},l})$$

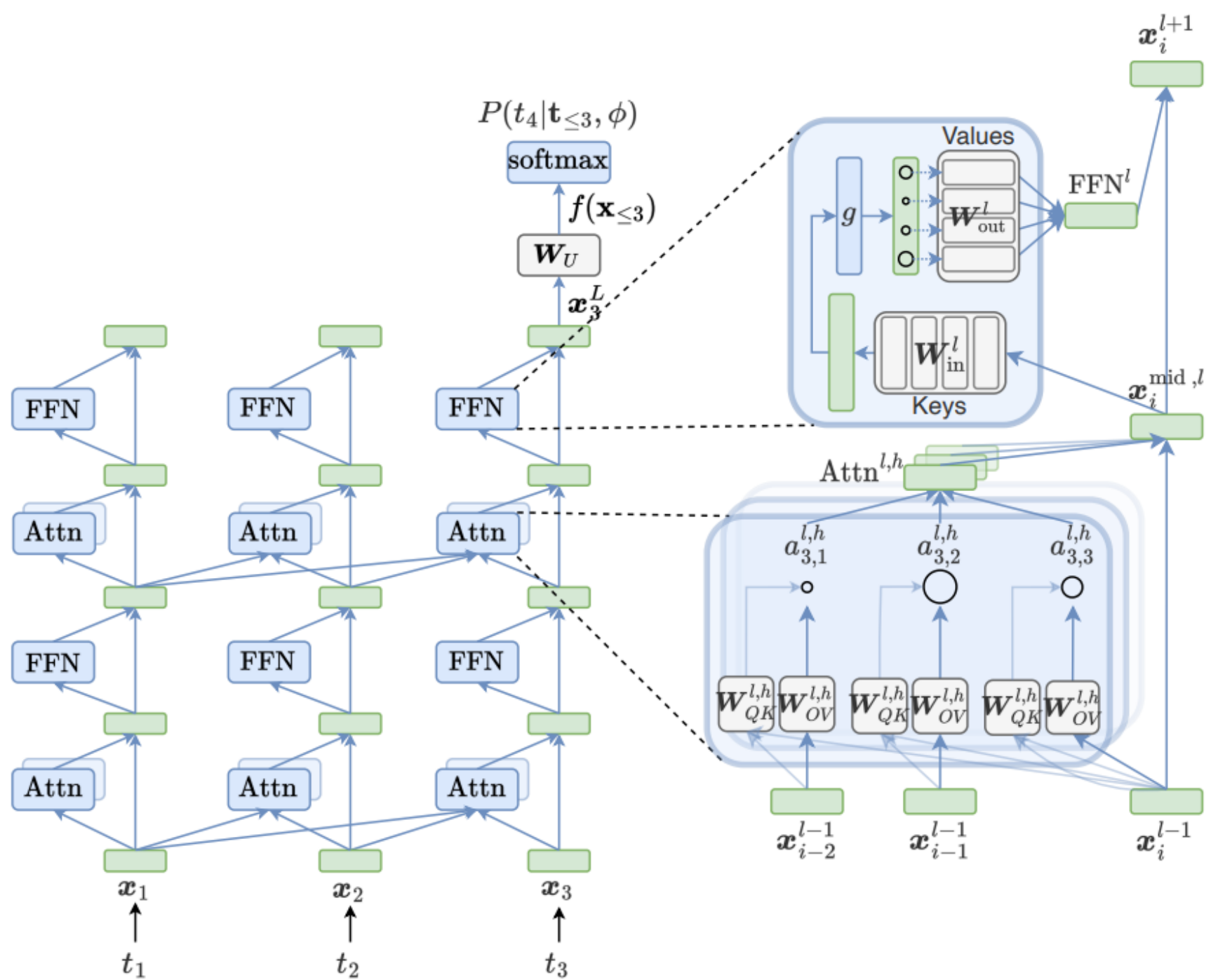


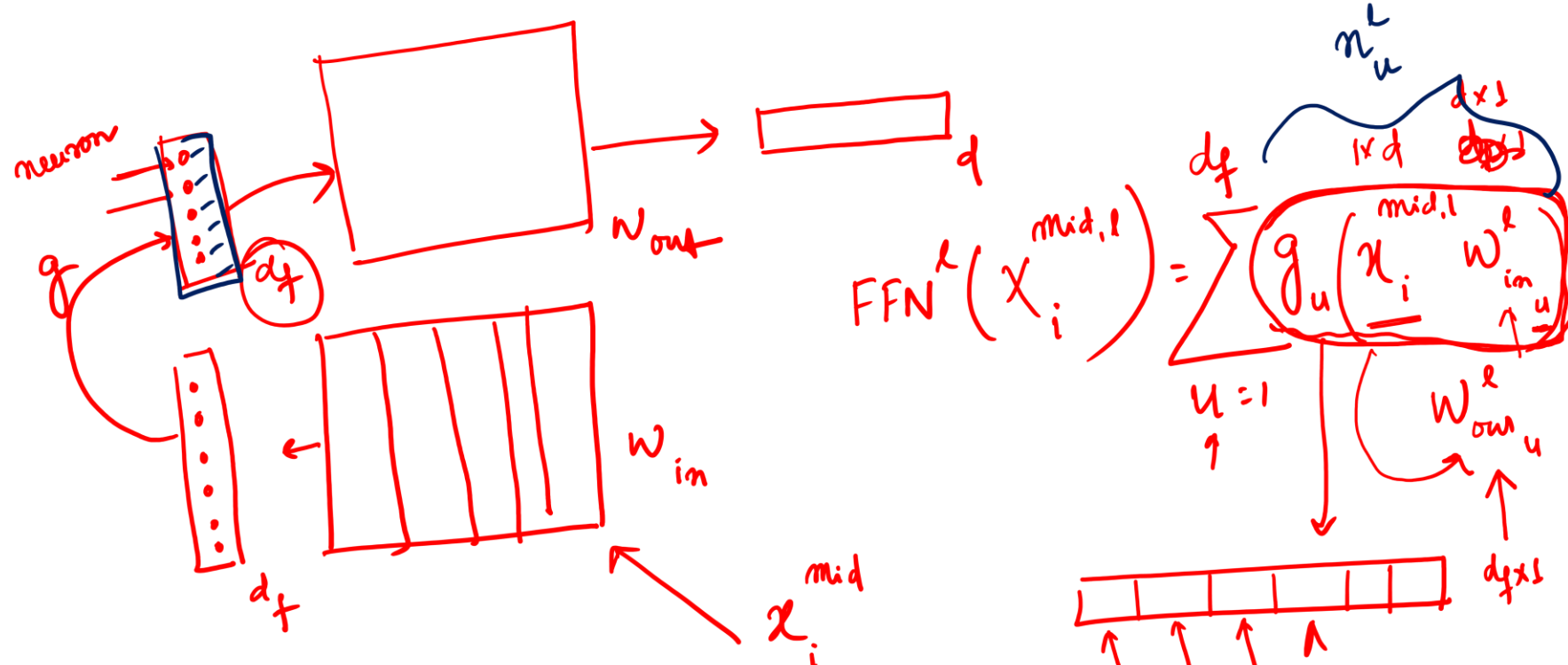
$$\mathbf{W}_{\text{in}}^l \in \mathbb{R}^{d \times d_{\text{ffn}}}$$

$$\mathbf{W}_{\text{out}}^l \in \mathbb{R}^{d_{\text{ffn}} \times d}$$



- \mathbf{W}_{in}^l reads from the residual stream state $\mathbf{x}_i^{\text{mid},l}$.
- Its result is passed through an element-wise non-linear activation function g , producing the neuron activations.
- These get transformed by $\mathbf{W}_{\text{out}}^l$ to produce the output $\text{FFN}^l(\mathbf{x}_i^{\text{mid},l})$, which is then added back to the residual stream





$$FFN^l(x_i^{mid,l}) = g_u(x_i^{mid,l}, W_{in}^l, u)$$

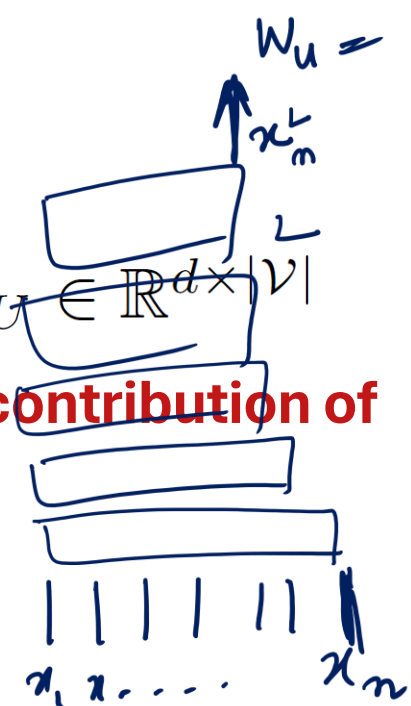
n_u^l :
 $n_u^l \in \mathbb{R}^{d_f}$
 vector of neuron activation



Prediction as a Sum of Component Outputs

- Prediction head of a Transformer consists of an unembedding matrix: $\mathbf{W}_U \in \mathbb{R}^{d \times |\mathcal{V}|}$

We can rearrange the traditional forward pass formulation to separate the **contribution of each model component to the output logits**:



$$\begin{aligned}
 \underline{f(\mathbf{x})} &= \mathbf{x}_n^L \mathbf{W}_U \\
 &= \left(\sum_{l=1}^L \sum_{h=1}^H \text{Attn}^{l,h}(\mathbf{X}_{\leq n}^{l-1}) + \sum_{l=1}^L \text{FFN}^l(\mathbf{x}_n^{\text{mid},l}) + \mathbf{x}_n \right) \mathbf{W}_U \\
 &= \sum_{l=1}^L \sum_{h=1}^H \text{Attn}^{l,h}(\mathbf{X}_{\leq n}^{l-1}) \mathbf{W}_U + \sum_{l=1}^L \text{FFN}^l(\mathbf{x}_n^{\text{mid},l}) \mathbf{W}_U + \mathbf{x}_n \mathbf{W}_U.
 \end{aligned}$$

Attention head logits update FFN logits update

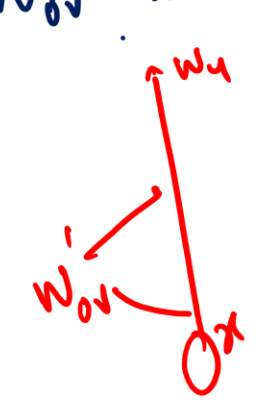
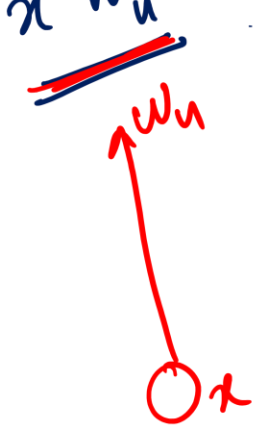
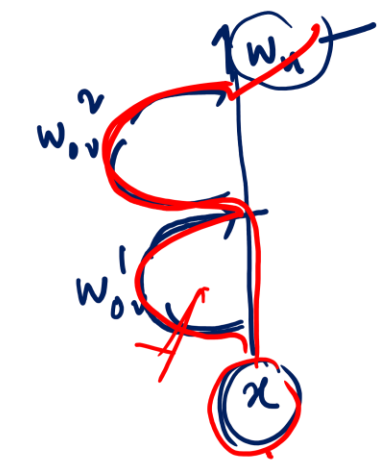
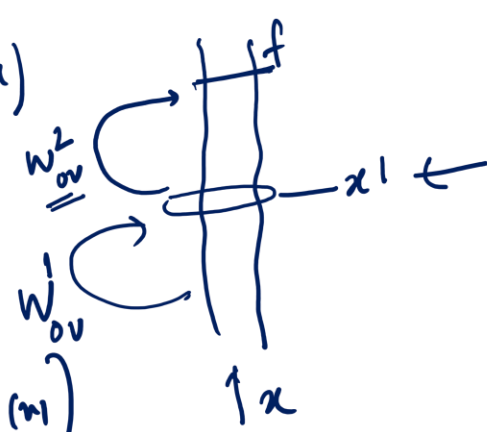
$$\underline{f(x)} = \underline{x' + W_{ov}^2(x')}$$

$$= x + W_{ov}^1(x) + W_{ov}^2(x + W_{ov}^1(x))$$

$$= (x + W_{ov}^1(x) + W_{ov}^2(x) + W_{ov}^1 W_{ov}^2(x)) W_u$$

$$= \underline{x W_u} + x W_{ov}^1 W_u + x W_{ov}^2 W_u + x W_{ov}^1 W_{ov}^2 W_u$$

$$x' = x + W_{ov}^1(x)$$



Prediction as an Ensemble of Shallow Networks

- Residual networks work as ensembles of shallow networks, where **each subnetwork defines a path in the computational graph**.

Consider a two-layer attention-only Transformer, where each attention head is composed just by an OV matrix:

$$f(x) = x^1 + \mathbf{W}_{OV}^2(x^1), \text{ with } x^1 = x + \mathbf{W}_{OV}^1(x)$$

We can decompose the forward pass as:

$$f(x) = \overbrace{x\mathbf{W}_U}^{\text{Direct path}} + \underbrace{x\mathbf{W}_{OV}^1\mathbf{W}_U}_{\text{Full OV circuits}} + \underbrace{x\mathbf{W}_{OV}^1\mathbf{W}_{OV}^2\mathbf{W}_U}_{\text{Virtual attention heads (V-composition)}} + \underbrace{x\mathbf{W}_{OV}^2\mathbf{W}_U}_{\text{Full OV circuits}}.$$

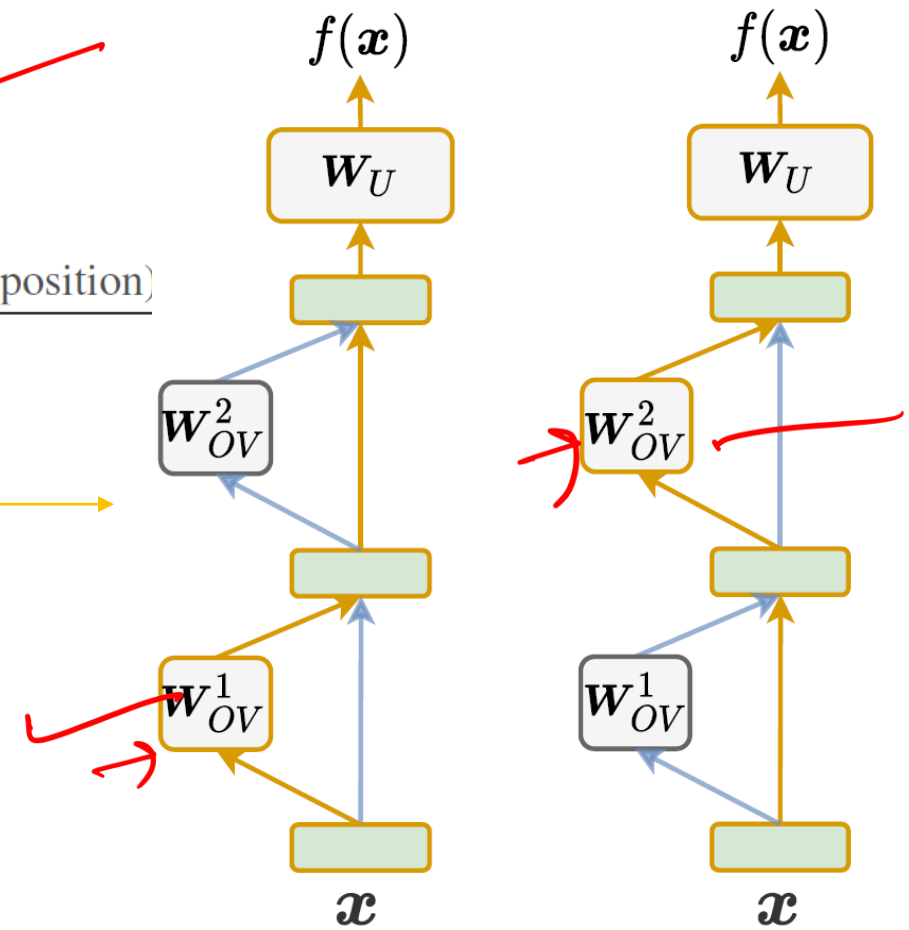


Prediction as an Ensemble of Shallow Networks

$$f(\mathbf{x}) = \overbrace{\mathbf{x}W_U}^{\text{Direct path}} + \underbrace{\mathbf{x}W_{OV}^1W_U + \mathbf{x}W_{OV}^1W_{OV}^2W_U + \mathbf{x}W_{OV}^2W_U}_{\text{Full OV circuits}}.$$

Virtual attention heads (V-composition)

- These terms depict paths traversing a single OV circuit, and are named **full OV circuits**
- The contribution of each OV circuit towards the output logit of the next token to be predicted.

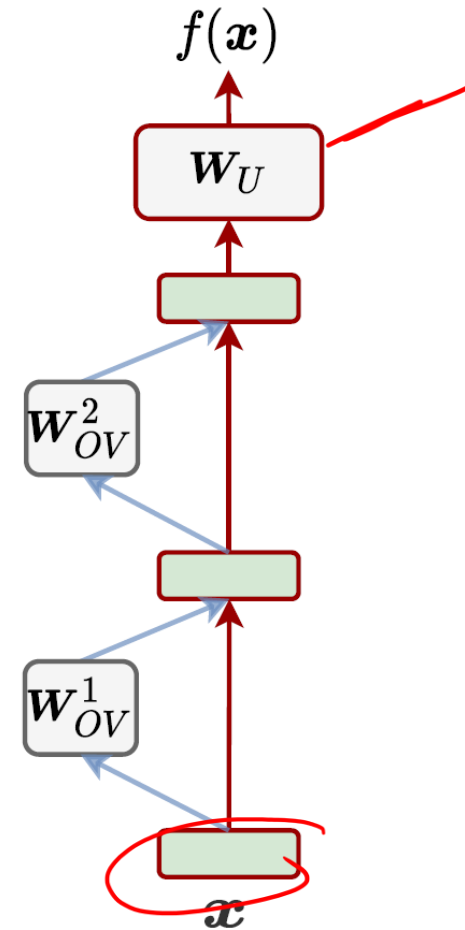


Prediction as an Ensemble of Shallow Networks

$$f(x) = \overbrace{xW_U}^{\text{Direct path}} + \overbrace{xW_{OV}^1 W_U}^{\text{Full OV circuits}} + \overbrace{xW_{OV}^1 W_{OV}^2 W_U}^{\text{Full OV circuits}} + \overbrace{xW_{OV}^2 W_U}^{\text{Full OV circuits}}.$$

Virtual attention heads (V-compositor)

- This term links the input embedding to the unembedding matrix and is referred to as the **direct path**
- It shows the contribution of the input embedding towards the output logit of the next token to be predicted.



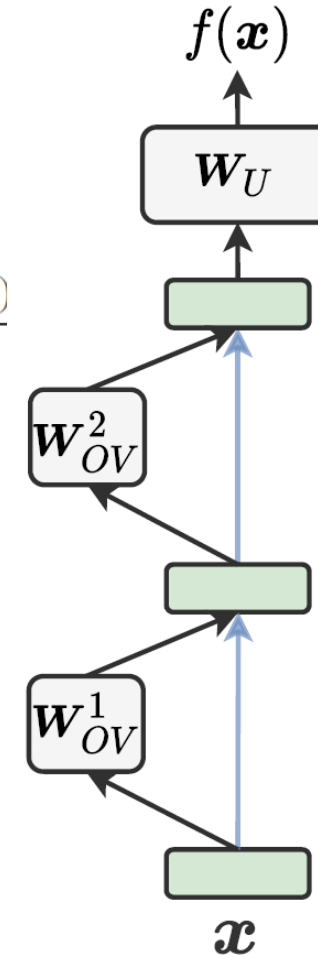
Prediction as an Ensemble of Shallow Networks

$$f(\mathbf{x}) = \overbrace{\mathbf{x}W_U}^{\text{Direct path}} + \underbrace{\mathbf{x}W_{OV}^1W_U + \mathbf{x}W_{OV}^1W_{OV}^2W_U + \mathbf{x}W_{OV}^2W_U}_{\text{Full OV circuits}}$$

Virtual attention heads (V-composition)

- This term depicts the path involving both attention heads, and is referred to as **virtual attention heads doing V-composition**
- This is called ‘composition’ since the sequential writing and reading of the two heads is seen as OV matrices composing together.
 - The amount of composition can be measured as:

$$\frac{\|W_{OV}^1W_{OV}^2\|_F}{\|W_{OV}^1\|_F\|W_{OV}^2\|_F}$$



Prediction as an Ensemble of Shallow Networks

- In full Transformer models, **Q-composition** and **K-composition**, i.e., compositions of W_Q and W_K with the W_{OV} output of previous layers, can also be found.
- Such decomposition enables us to **localize the inputs or model components** responsible for a particular prediction.



Why Do We Need Such a Formulation?

- Better understand the information flow within Transformer-based LLMs.
- Reveals how each layer incrementally transforms token representations.
 - Shows how attention heads and FFNs contribute to language modeling.
- Breaking down the contributions of individual circuits allows us to interpret which aspects of the model influence specific predictions.

Thus, through this formulation, the behavior of attention heads, the interaction between tokens, and the role of the residual stream can be explored more clearly.

