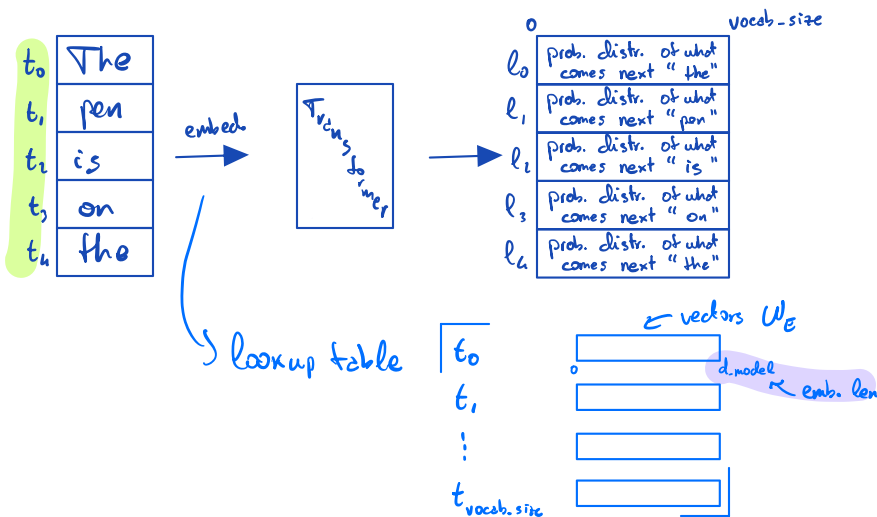


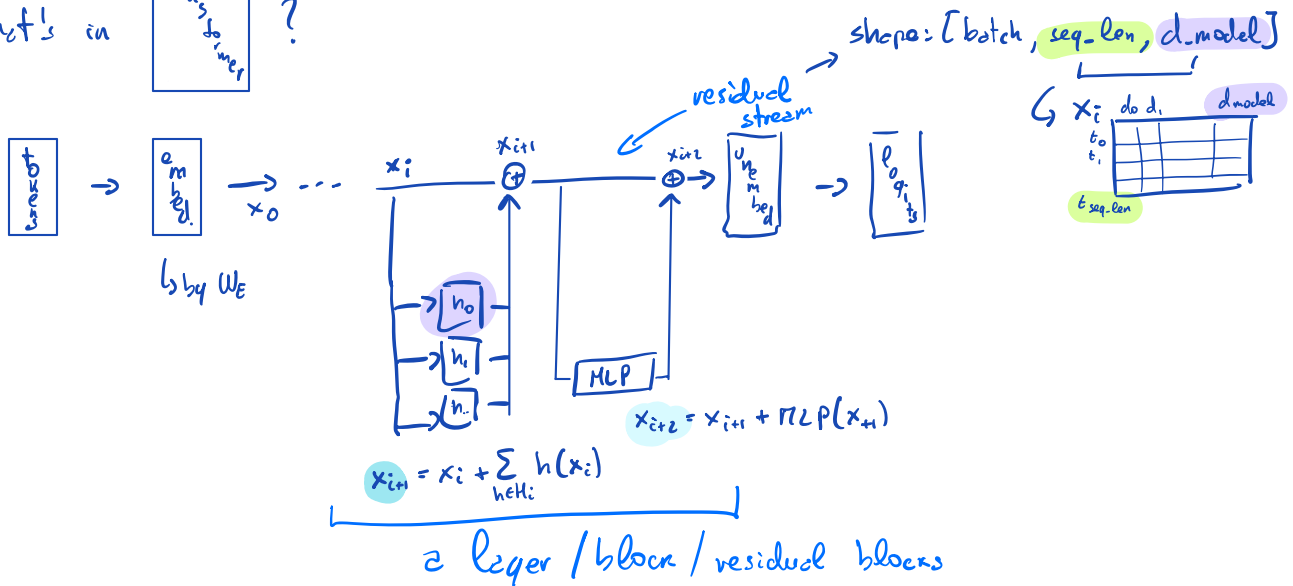
# Transformer Notes

Soft max  $x_i \rightarrow \frac{e^{x_i}}{\sum e^{x_j}}$  ← makes everything positive  
 where  $x = \begin{bmatrix} \square & \square & \square & \square \end{bmatrix}$  ← vector  
 ← everything add to 1

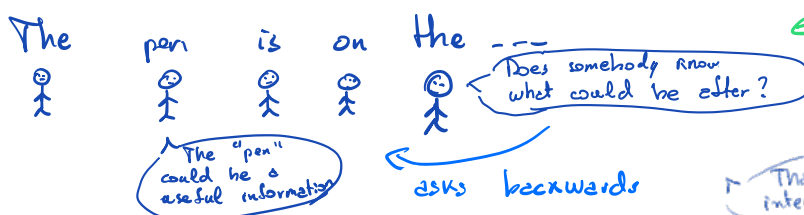
Causal attention: moves forward and predicts token  $t_n$  based only on tokens  $t_0, \dots, t_{n-1}$



What's in ?



- Residual stream: where model "stores" information (and remembers)
- Attention?
  - ↳ Moves info from prior positions to current tokens



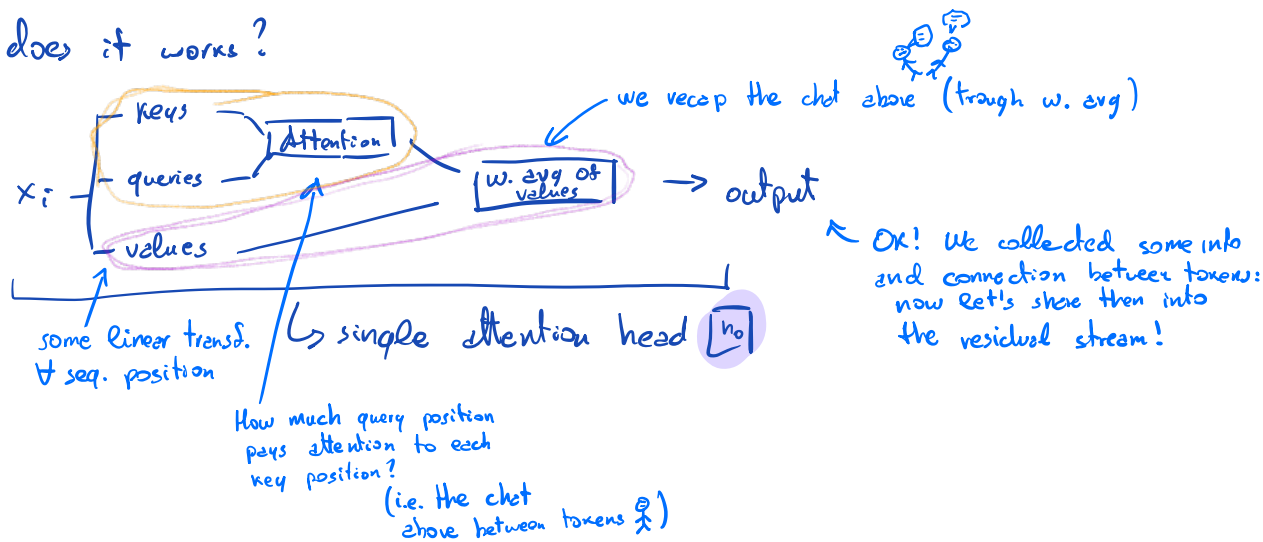
← This chat compose the residual stream. It's where everybody communicate!

↳ QK circuit

↳ That's interesting

→ OV circuit

How does it work?



## MLP?

↳ simple NN with (usually) a hidden layer  $h \times$  bigger than  $d_{\text{model}}$

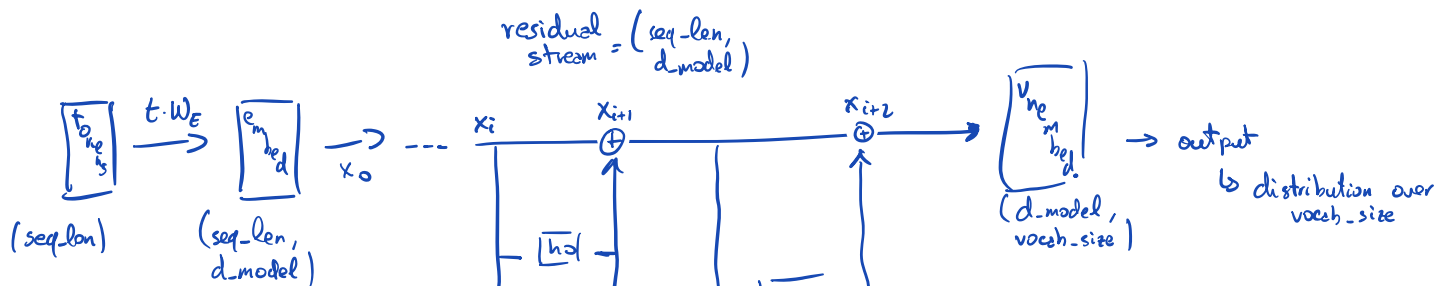
In math?  $\text{MLP} := f(x^T W^{\text{in}}) W^{\text{out}} = \sum_{i=1}^{h \times d_{\text{model}}} f(x^T W_{:,i}^{\text{in}}) W_{i,:}^{\text{out}}$

Annotations:   
 -  $f$ : linear transformations (the neurons)   
 -  $W^{\text{in}}$ :  $d_{\text{model}} \times h$    
 -  $W^{\text{out}}$ :  $h \times d_{\text{model}}$    
 - residual stream

Is this detected feature relevant? (e.g., the subject is a cat, seems important? Probably yes, let me write it to the residual stream)

By training a transformer we train layers to understand connection and what's important in a sentence!

I really like to look at matrix dimensions (shape) so:



For GPT-2:

$d_{\text{model}} = 768$

$\text{vocab\_size} = 50257$

$n_{\text{heads}} = 12$

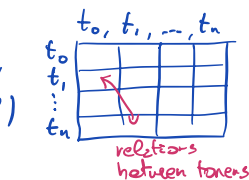
$n_{\text{layers}} = 12$

$d_{\text{heads}} = d_{\text{model}} / n_{\text{heads}} = 64$

$d_{\text{mlp}} = 4 \cdot d_{\text{model}} = 3072$

$Q, K, V = (\text{seq\_len}, n_{\text{heads}}, d_{\text{model}}/n_{\text{heads}})$

$\text{Attentions} = (n_{\text{heads}}, \text{seq\_len}, \text{seq\_len})$



# Implementation:

Steps:  $[batch, seq-len, d-model]$   $[batch, seq-len, n-heads, d-head]$

1. Linear map residual stream to queries, keys and values

Math:  $[n-heads, d-model, d-head]$

$$queries: x^T W_Q + bias_Q$$

$$keys: x W_K^T + bias_K$$

$$values: x^T W_V + bias_V$$

Pseudo-code:

```
queries = einsum (
    x, W_Q,
    "b s d_model, n-heads d_model d-head -> b s n-heads d-heads"
) + bias_Q
```

keys: same as above but with  $W_K$  and  $bias_K$

values: same as above but with  $W_V$  and  $bias_V$

$[batch, n-heads, s_Q, s_K]$   $[batch, s_{(Q,K)}, n-heads, d-heads]$

2. Get attention scores using queries and keys

Math:

$$scores = x^T W_Q W_K^T x$$

Pseudo-code:

```
scores = einsum (
    queries, keys,
    "b s_Q n-heads d-heads, b s_K n-heads d-heads -> b n-heads s_Q s_K"
)
```

$[s_Q, s_K]$  where Upper Triangular filled w/ 0

3. Scale it (by  $\sqrt{d-head}$ ), mask and apply softmax (along keys dimension)

Math:

$$A = \text{softmax} \left( \frac{\text{mask}(scores)}{\sqrt{d-head}} \right)$$

Pseudo-code:

$$A = \text{softmax}(\text{mask}(scores) / \text{sqrt}(d-head), \text{dim} = -1)$$

$[batch, n-heads, s_Q, s_K]$

$[batch, s_V = s_K, n-head, d-head]$

get a prob. distr. for each query

$[batch, n-heads, s_Q, s_K]$

4. Apply linear map from source A to destination tokens  $x W_V$  getting the w. avg of values

Math:

$$z = A \cdot x W_V$$

Pseudo-code:

```
z = einsum (
    A, values,
    "b n-heads s_Q s_K, b s_K n-heads d-head -> b s_V n-head d-head"
)
```

$[batch, s, d-model]$

5. Get the output OV summing over heads (or concatenation of heads) + bias

Math:

$[n-heads, d-heads, d-model]$   $[b, s, n-heads, d-heads]$

$$OV = z \cdot W_O + bias_O$$

Pseudo-code

$OV = \text{einsum} ($

$z, W_0$

$\xrightarrow{\text{implicit connect or position sum}}$

$"b_{s_v} \text{ n-heads d-heads, n-heads d-heads d-model} \rightarrow b_{s_v} \text{ d-model}"$   
 $) + \text{bias}_v$

Note that there are two circuits:  $QK$  and  $OV$  circuits