# CSCE 689-609
# Fundamentals of Large Language Models (LLMs)

Jeff Huang
jeff@cse.tamu.edu
o2lab.github.io

# The Transformer

*https://arxiv.org/pdf/1706.03762*

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

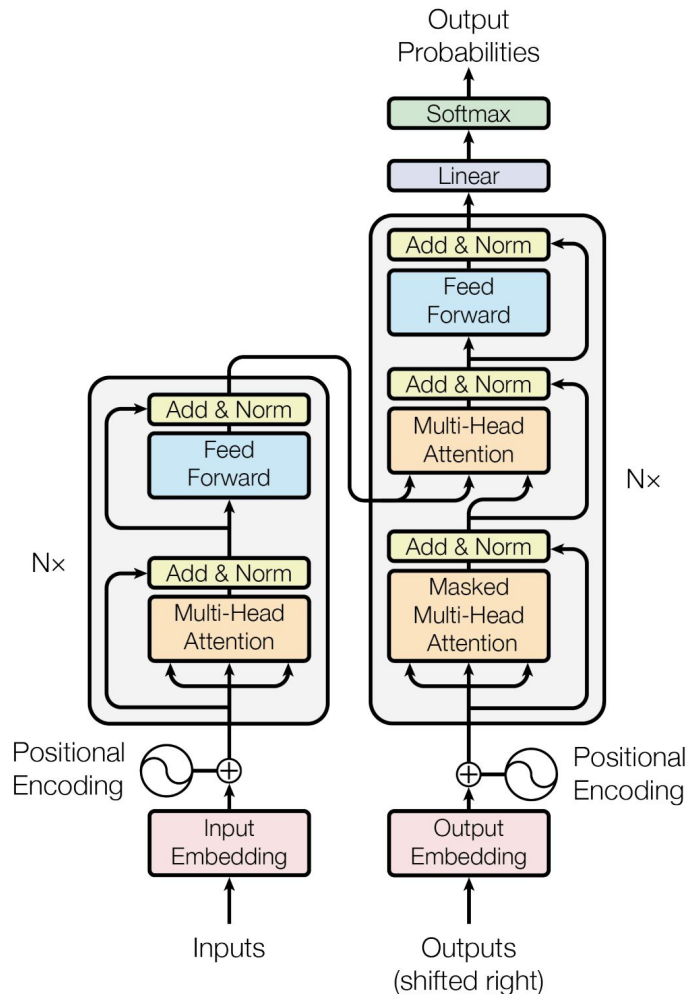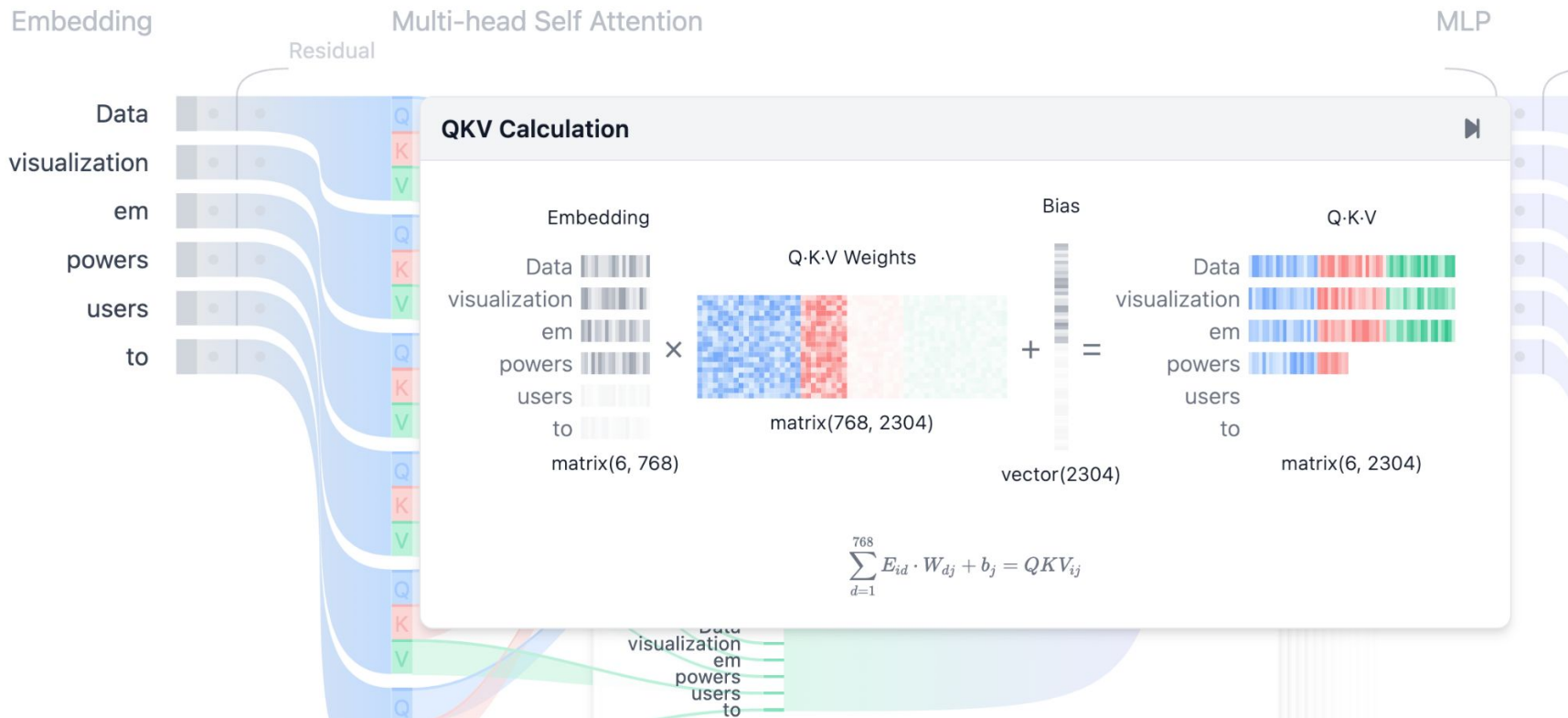Figure 1: The Transformer - model architecture.

# Attention is all your need

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

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

# Transformer Explainer

https://poloclub.github.io/transformer-explainer/

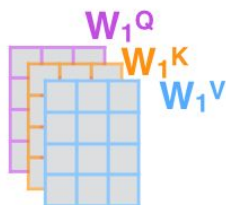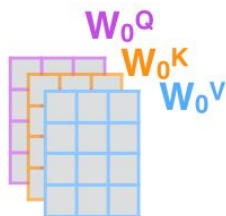# Illustrated Transformer

1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

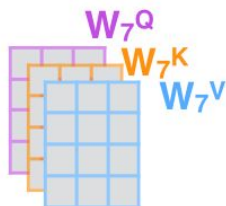5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

Thinking Machines
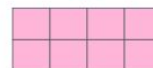
X

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

R

$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

...

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

...

$Q_7$
$K_7$
$V_7$

$Z_0$

$Z_1$

...

$Z_7$

$W^O$

Z

# A bird's view (llama2_7b)

hello

32K

**Tokenizer**

https://platform.openai.com/tokenizer

→ 15339 →

4K x 32K

**Token Embedding Table**

4K

input vector: [1.2, 0.7, 0.2, ..., 3.4, 6.8, 0.15]

world

32K

Probabilities (logits):

[0.22
0.16
0.06
0.05

aggies

...
0.01
... ]

4K x 32K

**Token output_weight Table**

**Transformer** 32 layers

block31 ← ... ← block1 ← block0

# A bird's view (llama2_7b)

# A bird's view (llama2_7b)

## Self-Attention

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

**Positional Encoding**

xi_in

4K

Q: 4K x 4K
K: 4K x 4K
V: 4K x 4K

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$W^O$: 4K x 4K
h:32

4K

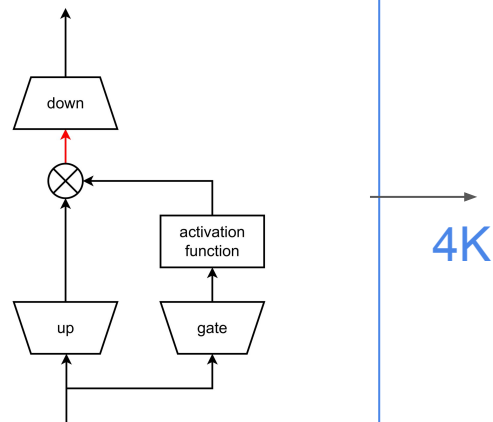# A bird's view (llama2_7b)

## FFN



$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

W2(silu(W1(x)) * W3(x))

W1(ffn_gate): 4K x 11008
W2(ffn_down): 11008 x 4K
W3(ffn_up): 4K x 11008

4K

4K

down

activation
function

up          gate

# Residual Connection

During backpropagation, gradients can become very small (vanishing gradients) or very large (exploding gradients)
- Residual connections mitigate these issues

# Layer Normalization (RMSNorm: Root Mean Square)

$$\text{RMSNorm}(x) = \frac{x}{\text{RMS}(x)} \cdot \gamma$$

where RMS is calculated as:

$$\text{RMS}(x) = \sqrt{\frac{1}{d} \sum_{i=1}^{d} x_i^2}$$

Here, $d$ is the dimensionality of the input, and $\gamma$ is a learnable scaling parameter.

**Question: Size of each component in llama2_7b?**

# Question: Size of each component in llama2_7b?

Total: 25.1G (type f32: 32-bit floating point)
Token embedding: 4K x 32K x 4 bytes = 0.5 GB
Each block:
    K/Q/V: 4K x 4K x 4 bytes x 3 = 192 MB
    attn_out: 4K x 4K x 4 bytes = 64 MB
    attn_norm: 4K x 4 bytes = 16 KB
    ffn_gate: 4K x 11008 x 4 bytes = 172 MB
    ffn_down: 11008 x 4K x 4 bytes = 172 MB
    ffn_up: 4K x 11008 x 4 bytes = 172 MB
    ffn_norm: 4K x 4 bytes = 16KB
Block x 32 = 772 MB x 32 = 24.1 GB
Output_weight: 4K x 32K x 4 bytes = 0.5 GB

# Question: Why QKV?

# Question: Why QKV?

**Retrieval**

**Database**

| Key 1 | Value 1 |
|-------|---------|
| Key 2 | Value 2 |
| Key 3 | Value 3 |
| Key 4 | Value 4 |
| Key 5 | Value 5 |

**Query**

Value 5

**Output**

$$attention(q, k, v) = \sum_i similarity(q, k_i) * v_i$$

1. It measures the similarity between the query and each key
2. This similarity returns a weight for each key
3. Output is the weighted combination of all values

# QKV

**Scaled dot-product attention**



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

## Step 1

| Q | | K$^T$ | | QK$^T$ | | V | | Attention |
|---|---|---|---|---|---|---|---|---|
| Query Token 1 | x | Key Token 1 | = | Q$_1$.K$_1$ | | Value Token 1 | x | Token 1 |
| (1, emb_size) | | (emb_size, 1) | | (1, 1) | | (1, emb_size) | | (1, emb_size) |

☐ Values that will be masked

*Zoom-in! (simplified without Scale and Softmax)*

# QKV

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



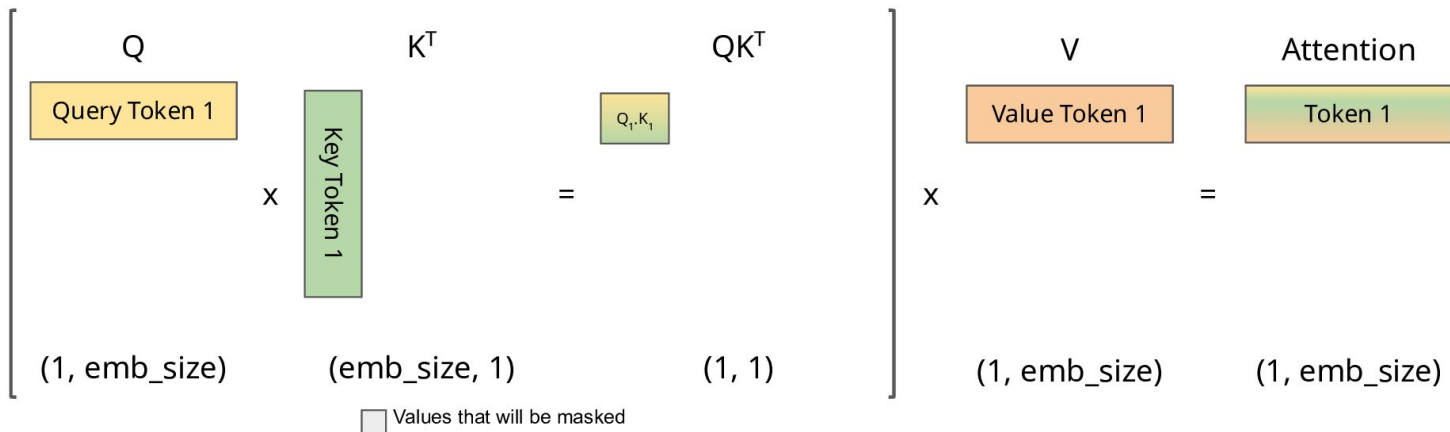Zoom-in! (simplified without Scale and Softmax)

# QKV

Scaled dot-product attention

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

## Step 3



| Q | | $K^T$ | | $QK^T$ | | V | | Attention |
|---|---|---|---|---|---|---|---|---|
| (3, emb_size) | | (emb_size, 3) | | (3, 3) | | (3, emb_size) | | (3, emb_size) |

☐ Values that will be masked

Zoom-in! (simplified without Scale and Softmax)

# KV Cache

https://medium.com/@joaolages/kv-caching-explained-276520203249

**Step 3**

**Without cache**

| Q | | $K^T$ | | $QK^T$ | | V | | Attention |
|---|---|---|---|---|---|---|---|---|

Q:
- Query Token 1
- Query Token 2
- Query Token 3

$K^T$:
- Key Token 1
- Key Token 2
- Key Token 3

$QK^T$:

| $Q_1.K_1$ | $Q_1.K_2$ | $Q_1.K_3$ |
| $Q_2.K_1$ | $Q_2.K_2$ | $Q_2.K_3$ |
| $Q_3.K_1$ | $Q_3.K_2$ | $Q_3.K_3$ |

V:
- Value Token 1
- Value Token 2
- Value Token 3

Attention:
- Token 1
- Token 2
- Token 3

(3, emb_size)   (emb_size, 3)   (3, 3)   (3, emb_size)   (3, emb_size)

**With cache**

Q:
- Query Token 3

$K^T$:
- Key Token 1
- Key Token 2
- Key Token 3

$QK^T$:

| $Q_3.K_1$ | $Q_3.K_2$ | $Q_3.K_3$ |

V:
- Value Token 1
- Value Token 2
- Value Token 3

Attention:
- Token 3

(1, emb_size)   (emb_size, 3)   (1, 3)   (3, emb_size)   (1, emb_size)

☐ Values that will be masked   ☐ Values that will be taken from cache

# Multihead Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

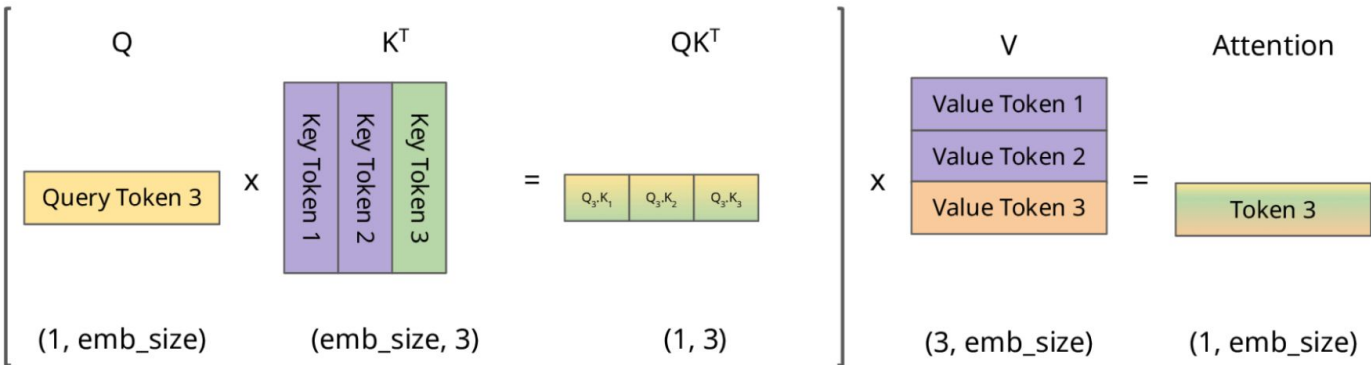Where the projections are parameter matrices $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ and $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$.

In this work we employ $h = 8$ parallel attention layers, or heads. For each of these we use $d_k = d_v = d_{\text{model}}/h = 64$. Due to the reduced dimension of each head, the total computational cost is similar to that of single-head attention with full dimensionality.

# Multihead Attention (Reduced KV Heads)

Llama 2:
llama.attention.head_count         = 32
llama.attention.head_count_kv      = 32

Llama 3.1:
llama.attention.head_count         = 32
llama.attention.head_count_kv      **= 8**
**grouped query attention (GQA)**

# Grouped Query Attention (GQA)

Paper: https://arxiv.org/pdf/2305.13245



Difference between MHA, GQA, and MQA (Source — https://arxiv.org/pdf/2305.13245.pdf)

Blog: https://towardsdatascience.com/demystifying-gqa-grouped-query-attention-3fb97b678e4a

# softmax

$$\text{softmax}(x_i) = \frac{\exp(x_i - \max(x))}{\sum_{j=0}^{\text{size}-1} \exp(x_j - \max(x))}$$

# Positional Encoding (RoPE: rotary positional encoding)

$$\text{RoPE}(x_{2i}, x_{2i+1}, p) = \begin{pmatrix} x_{2i} \cdot \cos(p \cdot \omega_i) - x_{2i+1} \cdot \sin(p \cdot \omega_i) \\ x_{2i} \cdot \sin(p \cdot \omega_i) + x_{2i+1} \cdot \cos(p \cdot \omega_i) \end{pmatrix}$$

The frequency for the $i$-th dimension is calculated as:

$$\omega_i = \frac{1}{10000^{\frac{2i}{d}}}$$

where $d$ is the dimension of the embedding vector.

Paper: Extending Context Window of Large Language Models via Positional Interpolation

# Let's write one!

Meta's Llama 2 7B model as an example

https://github.com/karpathy/llama2.c

https://github.com/karpathy/llama2.c/blob/master/run.c#L249-L354

```
248        // forward all the layers
249        for(unsigned long long l = 0; l < p->n_layers; l++) {
250
251            // attention rmsnorm
252            rmsnorm(s->xb, x, w->rms_att_weight + l*dim, dim);
253
254            // key and value point to the kv cache
255            int loff = l * p->seq_len * kv_dim; // kv cache layer offset for convenience
256            s->k = s->key_cache + loff + pos * kv_dim;
257            s->v = s->value_cache + loff + pos * kv_dim;
258
259            // qkv matmuls for this position
260            matmul(s->q, s->xb, w->wq + l*dim*dim, dim, dim);
261            matmul(s->k, s->xb, w->wk + l*dim*kv_dim, dim, kv_dim);
262            matmul(s->v, s->xb, w->wv + l*dim*kv_dim, dim, kv_dim);
263
264            // RoPE relative positional encoding: complex-valued rotate q and k in each head
265            for (int i = 0; i < dim; i+=2) {
```

```
tensor   0 :        token_embd.weight q4_K    [ 4096, 32000,   1,   1]
tensor   1 :      blk.0.attn_norm.weight f32   [ 4096,   1,   1,   1]
tensor   2 :      blk.0.ffn_down.weight q6_K   [ 11008, 4096,   1,   1]
tensor   3 :      blk.0.ffn_gate.weight q4_K   [ 4096, 11008,   1,   1]
tensor   4 :       blk.0.ffn_up.weight q4_K   [ 4096, 11008,   1,   1]
tensor   5 :      blk.0.ffn_norm.weight f32   [ 4096,   1,   1,   1]
tensor   6 :       blk.0.attn_k.weight q4_K   [ 4096, 4096,   1,   1]
tensor   7 :   blk.0.attn_output.weight q4_K   [ 4096, 4096,   1,   1]
tensor   8 :       blk.0.attn_q.weight q4_K   [ 4096, 4096,   1,   1]
tensor   9 :       blk.0.attn_v.weight q6_K   [ 4096, 4096,   1,   1]
…
tensor 281 :      blk.31.attn_norm.weight f32   [ 4096,   1,   1,   1]
tensor 282 :      blk.31.ffn_down.weight q6_K   [ 11008, 4096,   1,   1]
tensor 283 :      blk.31.ffn_gate.weight q4_K   [ 4096, 11008,   1,   1]
tensor 284 :       blk.31.ffn_up.weight q4_K   [ 4096, 11008,   1,   1]
tensor 285 :      blk.31.ffn_norm.weight f32   [ 4096,   1,   1,   1]
tensor 286 :       blk.31.attn_k.weight q4_K   [ 4096, 4096,   1,   1]
tensor 287 :   blk.31.attn_output.weight q4_K   [ 4096, 4096,   1,   1]
tensor 288 :       blk.31.attn_q.weight q4_K   [ 4096, 4096,   1,   1]
tensor 289 :       blk.31.attn_v.weight q6_K   [ 4096, 4096,   1,   1]
tensor 290 :       output_norm.weight f32   [ 4096,   1,   1,   1]
```

**Question: How many matmuls to generate a token?**

# Important Notes

- Read:
    - llama2.c: https://github.com/karpathy/llama2.c
    - SGLang https://github.com/sgl-project/sglang
- Due
    - HW0 (Saturday)
        - Don Knuth: https://cs.stanford.edu/~knuth/chatGPT20.txt