

# Introduction to Agentic AI

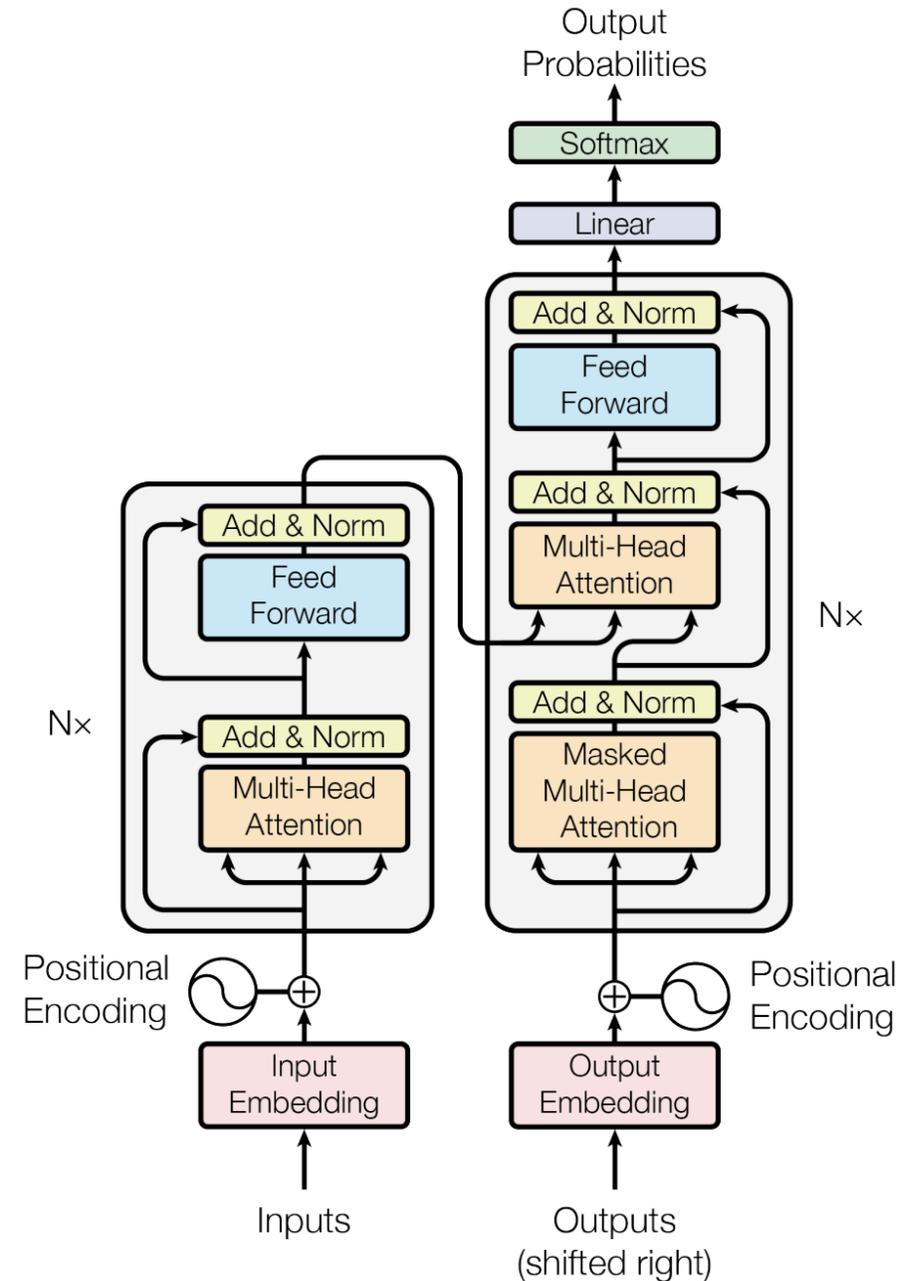
## -- Language Modeling

Instructor: Guangjing Wang

[guangjingwang@usf.edu](mailto:guangjingwang@usf.edu)

# Last Lecture

- Cross-Attention
- Self-Attention
- Transformer



# Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  [ $N \times D_{in}$ ]

Key matrix:  $\mathbf{W}_K$  [ $D_{in} \times D_{out}$ ]

Value matrix:  $\mathbf{W}_V$  [ $D_{in} \times D_{out}$ ]

Query matrix:  $\mathbf{W}_Q$  [ $D_{in} \times D_{out}$ ]

## Computation:

Queries:  $\mathbf{Q} = \mathbf{XW}_Q$  [ $N \times D_{out}$ ]

Keys:  $\mathbf{K} = \mathbf{XW}_K$  [ $N \times D_{out}$ ]

Values:  $\mathbf{V} = \mathbf{XW}_V$  [ $N \times D_{out}$ ]

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q}$  [ $N \times N$ ]

$$E_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  [ $N \times N$ ]

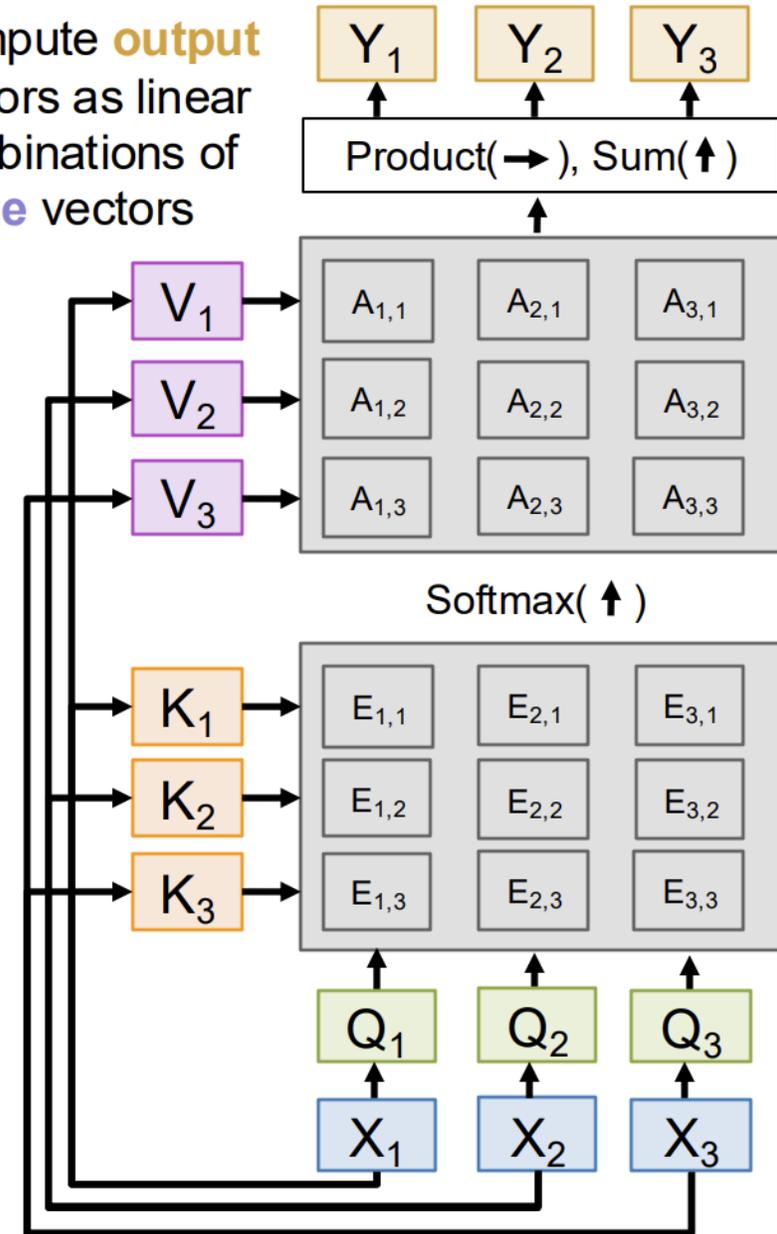
Output vector:  $\mathbf{Y} = \mathbf{AV}$  [ $N \times D_{out}$ ]

$$Y_i = \sum_j A_{ij} V_j$$

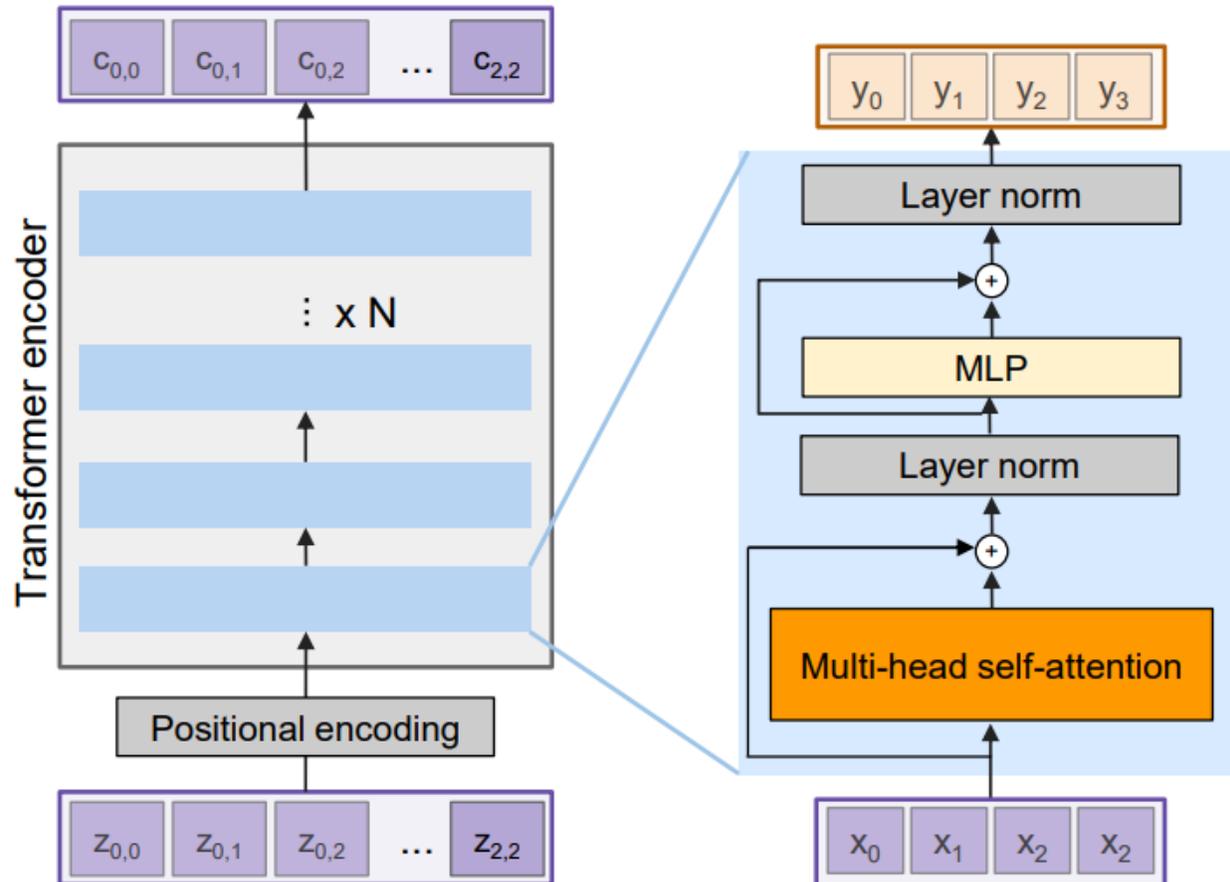
Each **input** produces one **output**, which is a mix of information from all **inputs**

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

Compute **output** vectors as linear combinations of **value** vectors



# The Transformer encoder block



## Transformer Encoder Block:

**Inputs:** Set of vectors  $x$

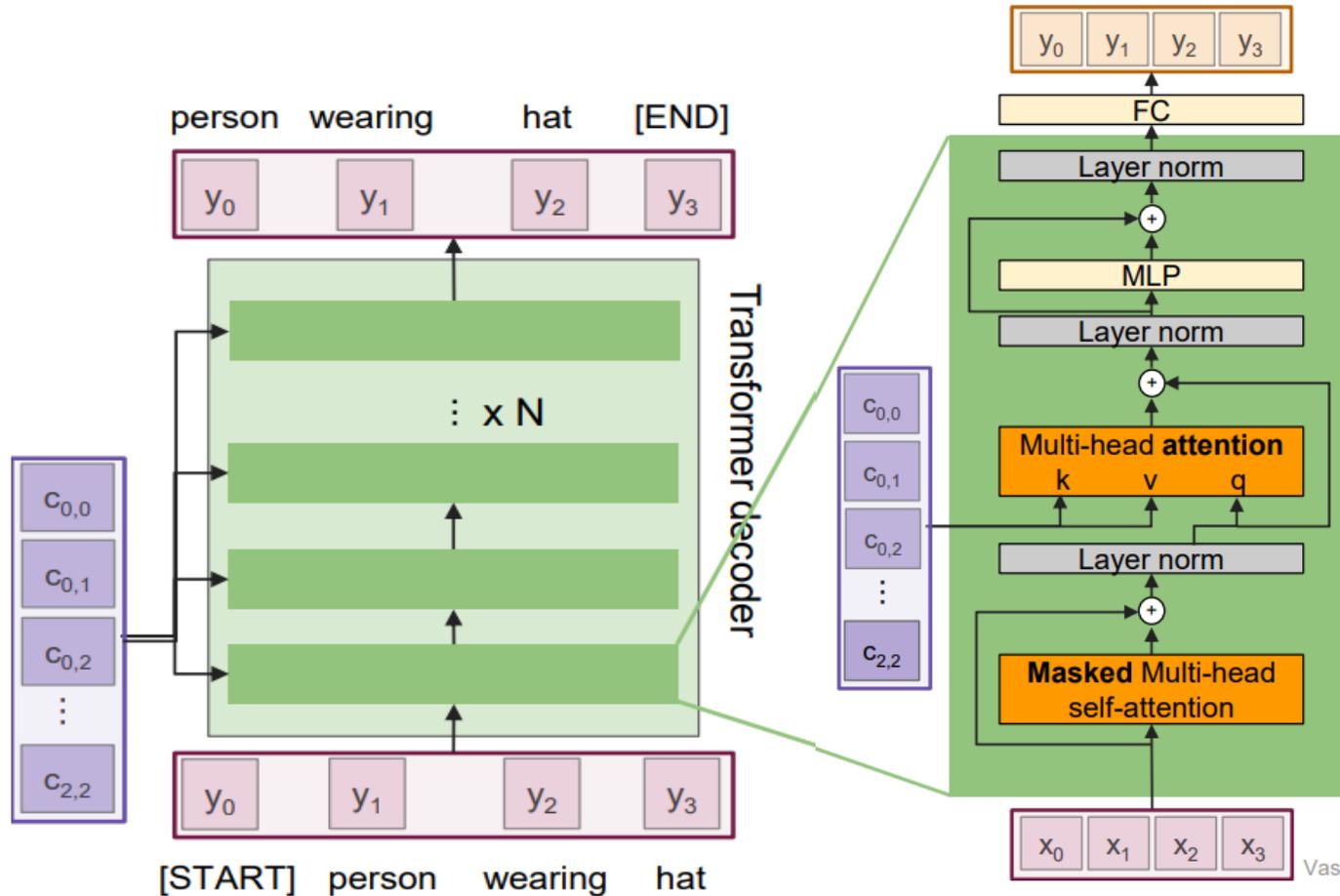
**Outputs:** Set of vectors  $y$

Self-attention is the only interaction between vectors.

Layer norm and MLP operate independently per vector.

Highly scalable, highly parallelizable, but high memory usage.

# The Transformer decoder block



## Transformer Decoder Block:

**Inputs:** Set of vectors  $\mathbf{x}$  and Set of context vectors  $\mathbf{c}$ .  
**Outputs:** Set of vectors  $\mathbf{y}$ .

Masked Self-attention only interacts with past inputs.

Multi-head attention block is NOT self-attention. It attends over encoder outputs.

Highly scalable, highly parallelizable, but high memory usage.

Vaswani et al, "Attention is all you need", NeurIPS 2017

# Transformer Variant: DeepSeek V3

**RMSNorm: root mean square layer normalization**

$$y_i = \frac{x_i}{\text{RMS}(x)} * \gamma_i, \quad \text{where} \quad \text{RMS}(x) = \sqrt{\epsilon + \frac{1}{n} \sum_{i=1}^n x_i^2}$$

**Query Vector**

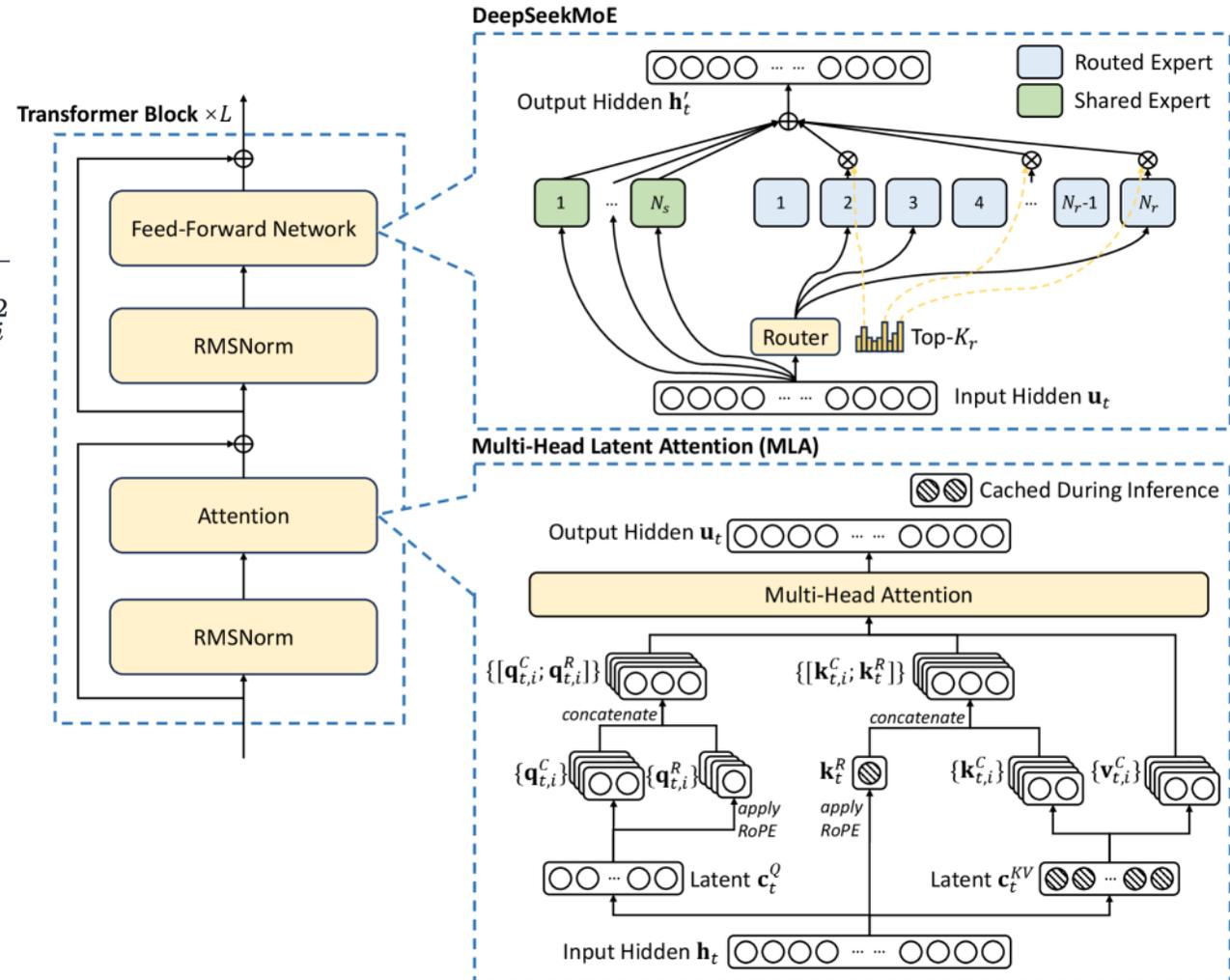
$$\mathbf{c}_t^Q = W^{DQ} \mathbf{h}_t,$$

$$[\mathbf{q}_{t,1}^C; \mathbf{q}_{t,2}^C; \dots; \mathbf{q}_{t,n_h}^C] = \mathbf{q}_t^C = W^{UQ} \mathbf{c}_t^Q,$$

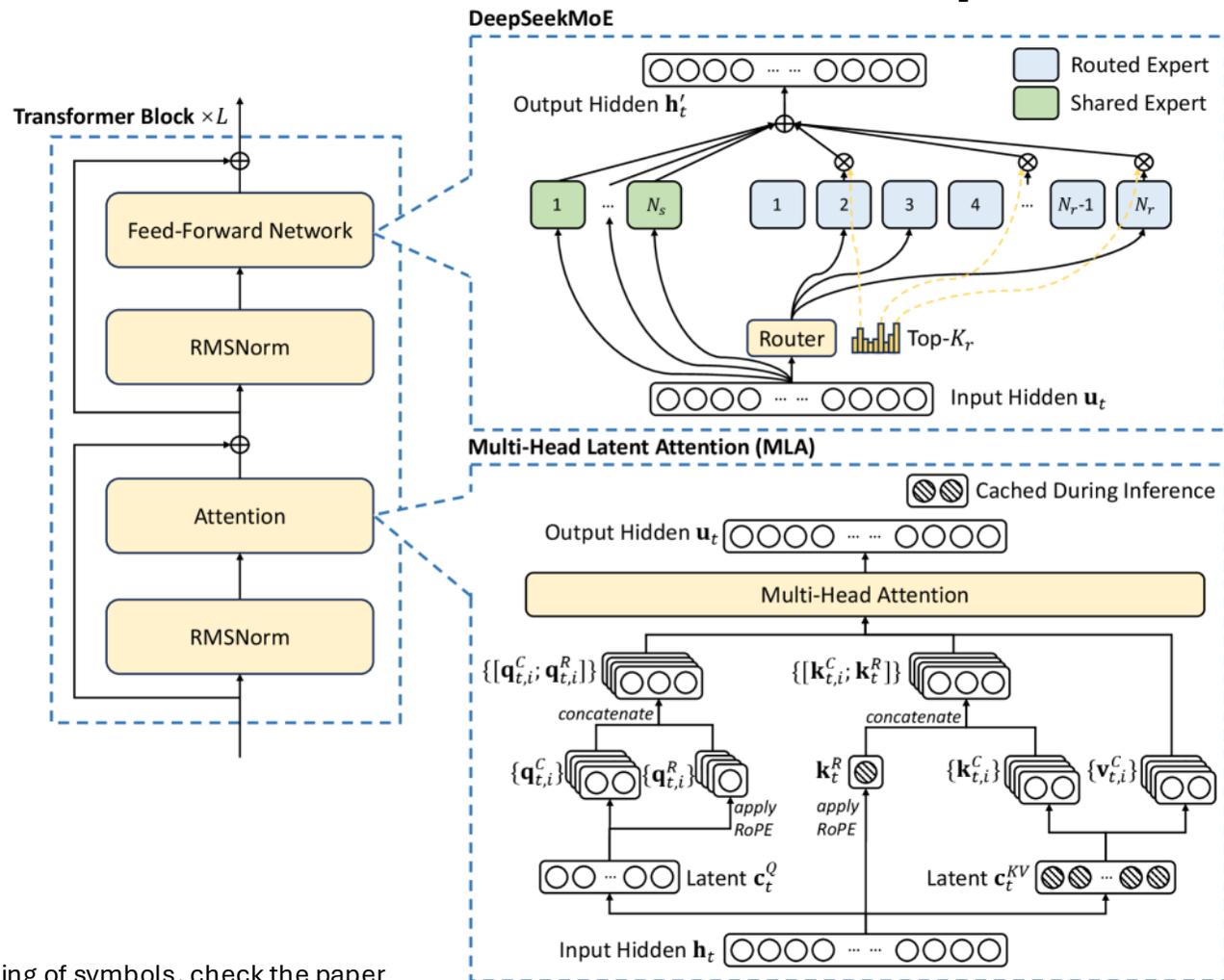
$$[\mathbf{q}_{t,1}^R; \mathbf{q}_{t,2}^R; \dots; \mathbf{q}_{t,n_h}^R] = \mathbf{q}_t^R = \text{RoPE}(W^{QR} \mathbf{c}_t^Q),$$

$$\mathbf{q}_{t,i} = [\mathbf{q}_{t,i}^C; \mathbf{q}_{t,i}^R],$$

For the meaning of symbols, please check the paper  
<https://arxiv.org/html/2412.19437v1>



# Transformer Variant: DeepSeek V3



## Value Vector

$$[\mathbf{v}_{t,1}^C; \mathbf{v}_{t,2}^C; \dots; \mathbf{v}_{t,n_h}^C] = \mathbf{v}_t^C = W^{UV} \mathbf{c}_t^{KV},$$

## Key Vector

$$\mathbf{c}_t^{KV} = W^{DKV} \mathbf{h}_t,$$

$$[\mathbf{k}_{t,1}^C; \mathbf{k}_{t,2}^C; \dots; \mathbf{k}_{t,n_h}^C] = \mathbf{k}_t^C = W^{UK} \mathbf{c}_t^{KV},$$

$$\mathbf{k}_t^R = \text{RoPE}(W^{KR} \mathbf{h}_t),$$

$$\mathbf{k}_{t,i} = [\mathbf{k}_{t,i}^C; \mathbf{k}_t^R],$$

For the meaning of symbols, check the paper  
<https://arxiv.org/html/2412.19437v1>

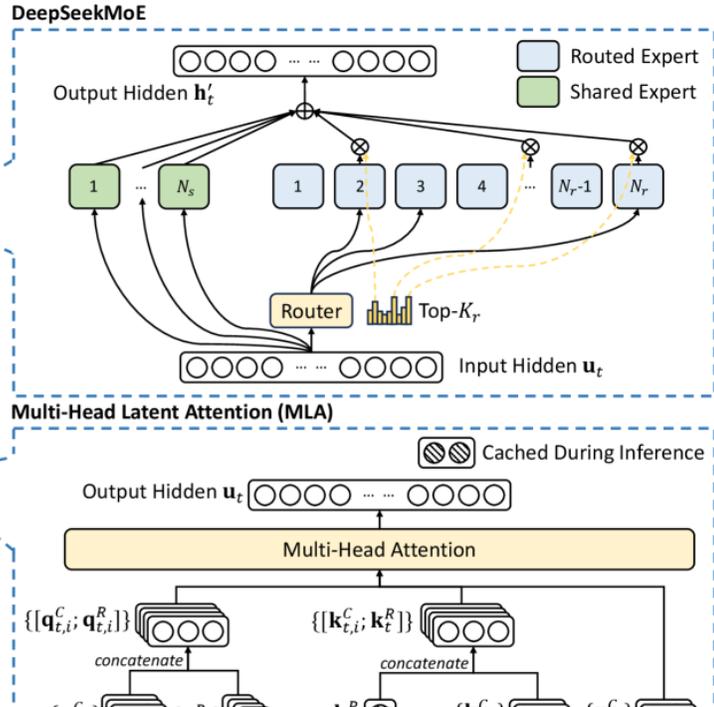
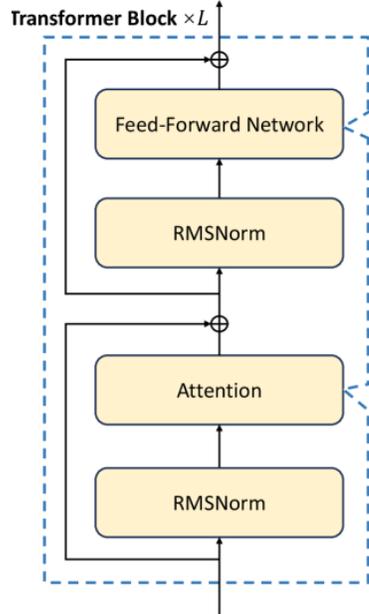
# DeepSeekMoE: Mixture of Expert

$$\mathbf{h}'_t = \mathbf{u}_t + \sum_{i=1}^{N_s} \text{FFN}_i^{(s)}(\mathbf{u}_t) + \sum_{i=1}^{N_r} g_{i,t} \text{FFN}_i^{(r)}(\mathbf{u}_t),$$

$$g_{i,t} = \frac{g'_{i,t}}{\sum_{j=1}^{N_r} g'_{j,t}},$$

$$g'_{i,t} = \begin{cases} s_{i,t}, & s_{i,t} \in \text{Topk}(\{s_{j,t} | 1 \leq j \leq N_r\}, K_r), \\ 0, & \text{otherwise,} \end{cases}$$

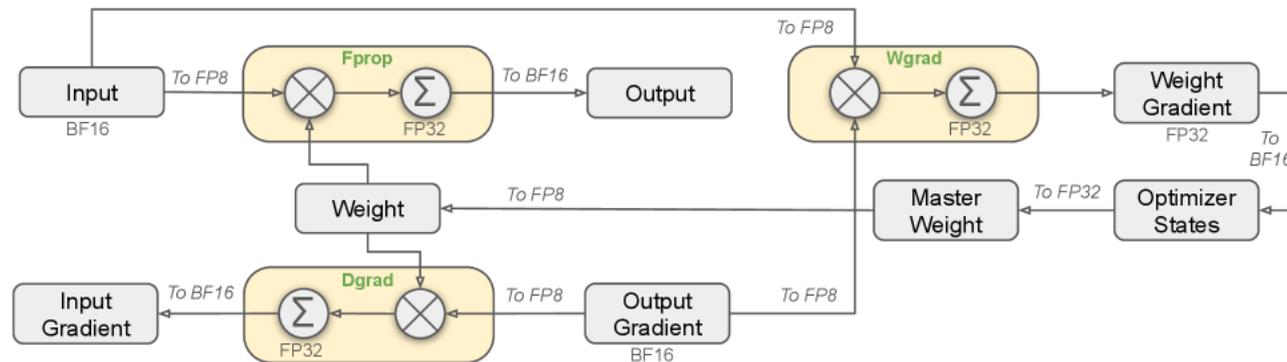
$$s_{i,t} = \text{Sigmoid}(\mathbf{u}_t^T \mathbf{e}_i),$$



where  $N_s$  and  $N_r$  denote the numbers of shared experts and routed experts, respectively;  $\text{FFN}_i^{(s)}(\cdot)$  and  $\text{FFN}_i^{(r)}(\cdot)$  denote the  $i$ -th shared expert and the  $i$ -th routed expert, respectively;  $K_r$  denotes the number of activated routed experts;  $g_{i,t}$  is the gating value for the  $i$ -th expert;  $s_{i,t}$  is the token-to-expert affinity;  $\mathbf{e}_i$  is the centroid vector of the  $i$ -th routed expert; and  $\text{Topk}(\cdot, K)$  denotes the set comprising  $K$  highest scores among the affinity scores calculated for the  $t$ -th token and all routed experts. Slightly different from DeepSeek-V2, DeepSeek-V3 uses the sigmoid function to compute the affinity

# Pre-training of DeepSeek-V3 (671B)

- 2048 NVIDIA H800 GPUs for 3.7 days (180K H800 GPU hours) **per trillion token**
  - GPUs connected by NVLink and NVSwitch within nodes
  - InfiniBand (IB) interconnects are utilized to facilitate communications across nodes
- Low-precision training: FP8 (8-bit floating point).



- Pre-train DeepSeek-V3 on 14.8 trillion diverse and high-quality tokens.
- Total training cost \$ 5.576 million, if H800 GPU is \$2 per GPU hour.

# This Lecture

- Masked language modeling
- Causal language modeling
- Sequence to sequence modeling
- LLM inference process
- Evaluation and Optimization Inference

# Language Modeling

- LM: probability distribution over sequence of tokens

| Step | Token ( $w_i$ ) | Context (History) | Conditional Probability              | Estimated Value |
|------|-----------------|-------------------|--------------------------------------|-----------------|
| 1    | The             | [Start]           | $P(\text{"The"})$                    | 0.1             |
| 2    | cat             | The               | $P(\text{"cat"} \text{"The"})$       | 0.05            |
| 3    | sat             | The cat           | $P(\text{"sat"} \text{"The cat"})$   | 0.2             |
| 4    | .               | The cat sat       | $P(\text{"."} \text{"The cat sat"})$ | 0.8             |

$$P(\text{"The cat sat."}) = 0.10 \times 0.05 \times 0.20 \times 0.80$$

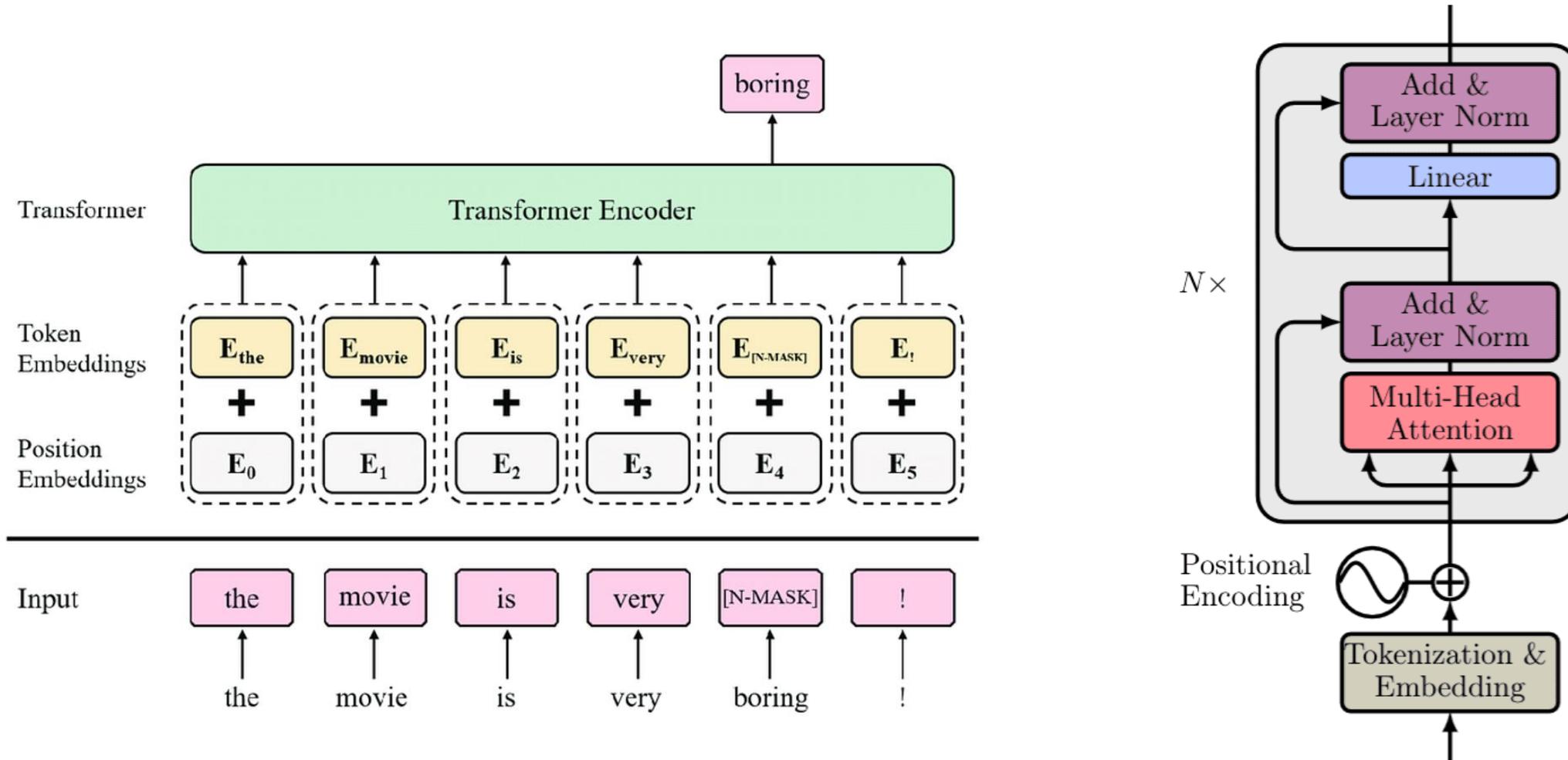
- **Autoregressive** language models:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$$

# Masked Language Modeling

- Language models work by being trained to **predict the probability of a word** given **the context of surrounding words**.
- **Masked language modeling (MLM)**: randomly masks some tokens in the input and trains the model to predict the original tokens based on the surrounding context.
  - Bidirectional context (looking at words both before and after the masked word)
  - Model example: BERT
  - Tasks such as classification, named entity recognition, and question answering

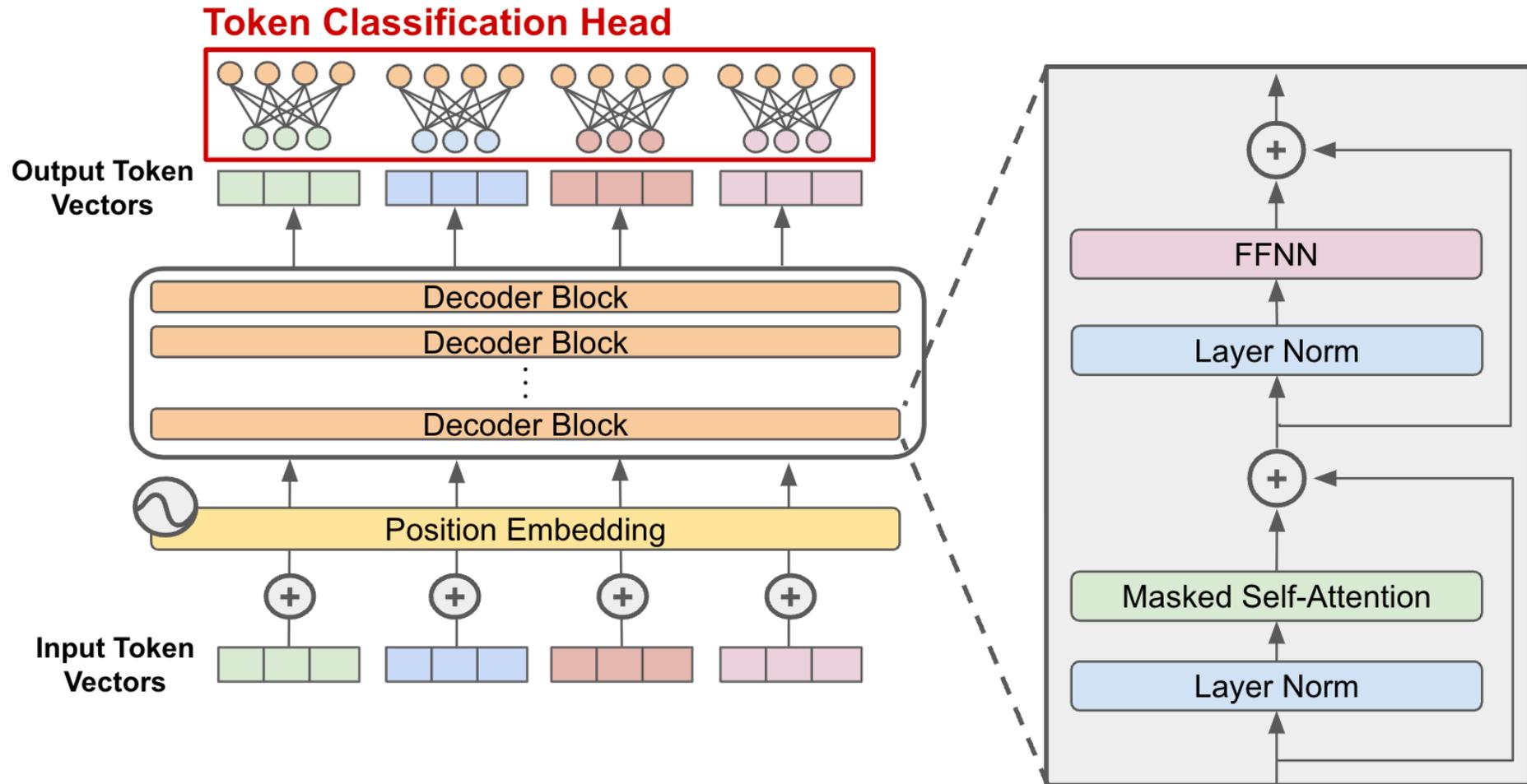
# Masked Language Modeling



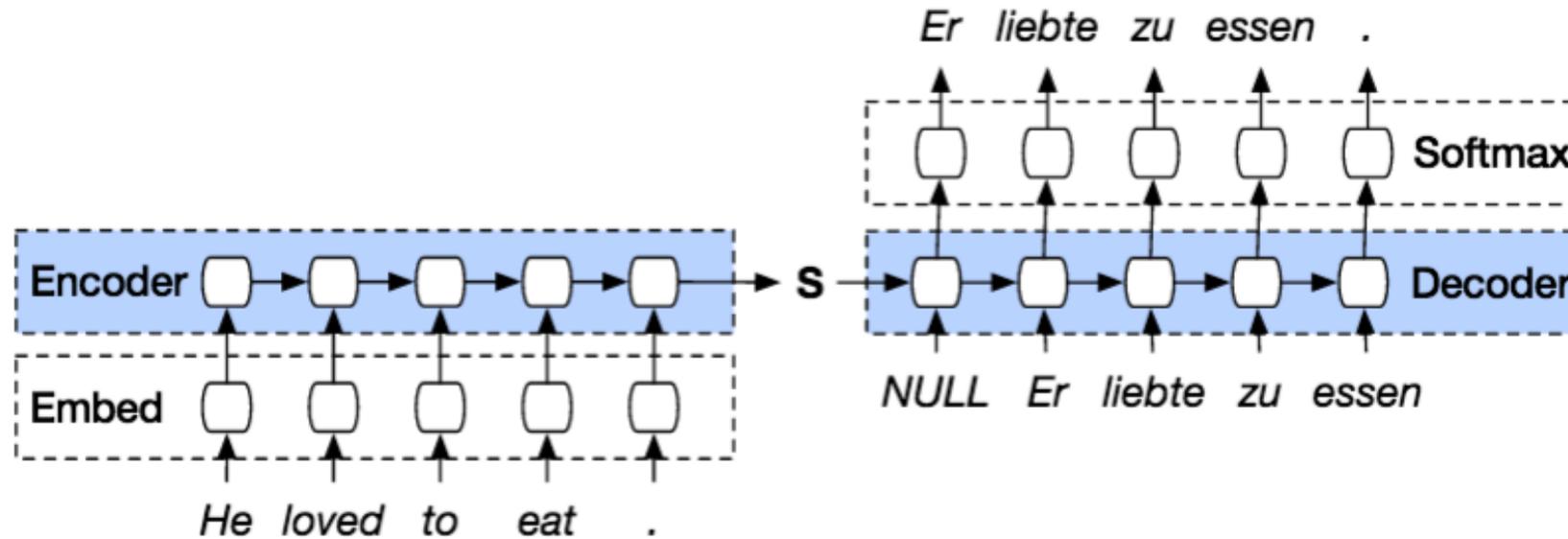
# Causal Language Modeling

- **Causal language modeling (CLM):** predicts the next token based on all previous tokens in the sequence.
- The model can only use context from the left (previous tokens) to predict the next token.
  - Model example: GPT, Llama
  - They can complete sentences, write essays, or generate code based on a prompt.

# Causal Language Modeling



# Sequence to Sequence Modeling



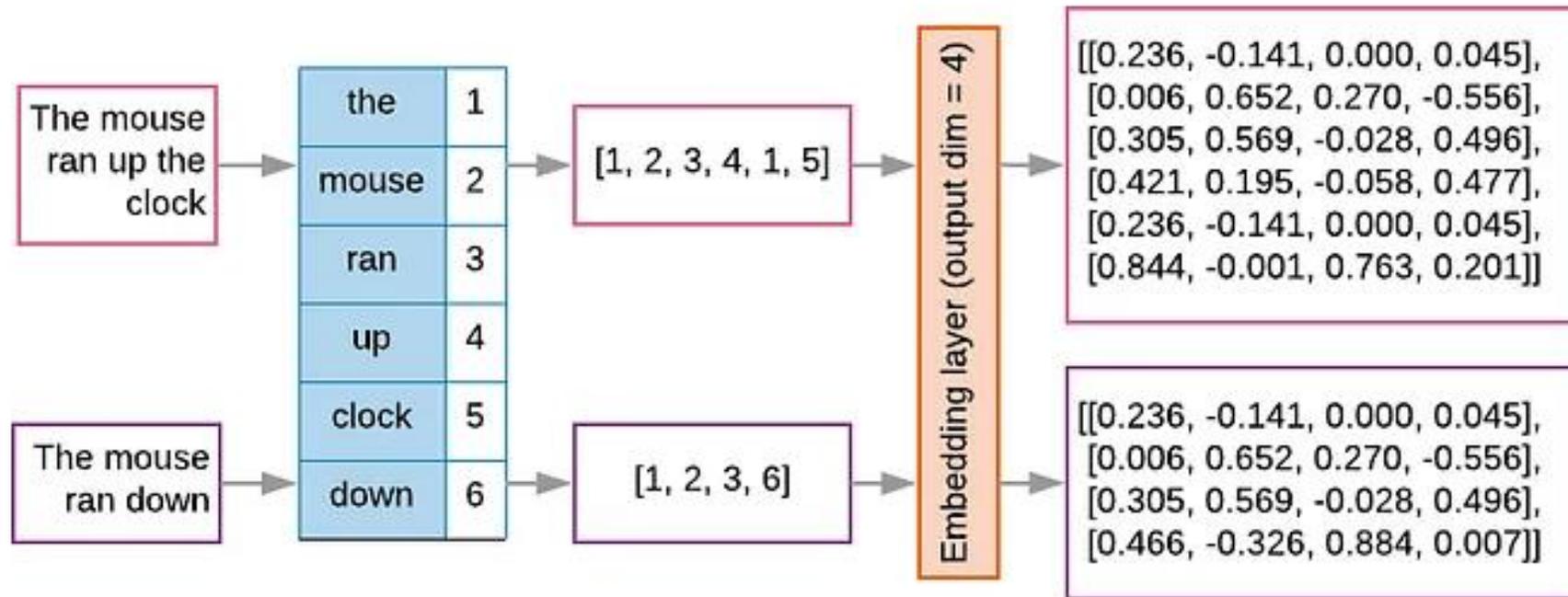
Models such as T5, BART use an encoder to understand the input and a decoder to generate output. Tasks such as translation, summarization, and question answering

# Prefill Phase in Inference Process:

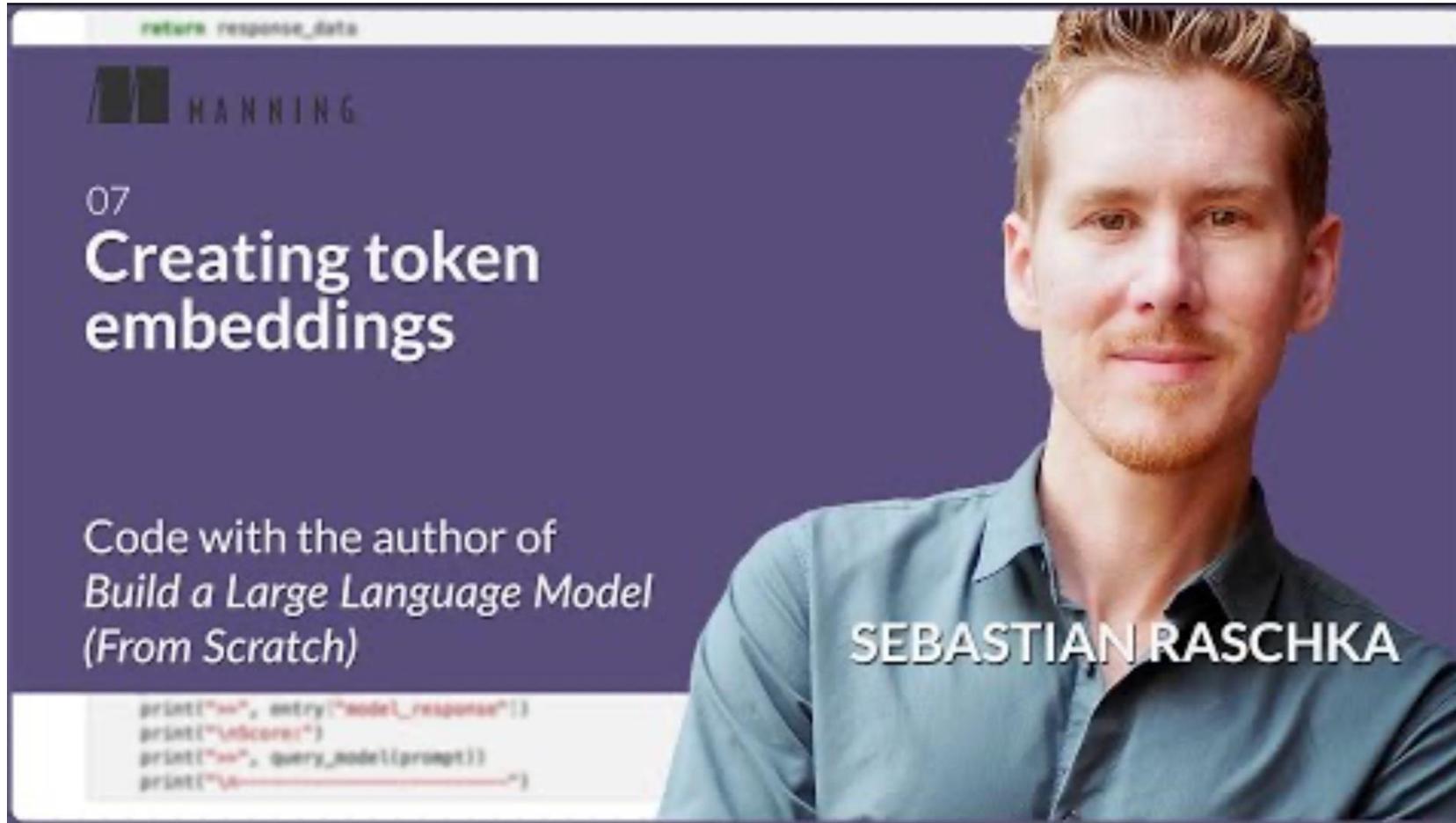
- 1. Tokenization:** Converting the input text into tokens (think of these as the basic building blocks the model understands)
- 2. Embedding Conversion:** Transforming these tokens into numerical representations that capture their meaning
- 3. Initial Processing:** Running these embeddings through the model's neural networks to create a rich understanding of the context

# Tokenization

- Splitting the input into words, subwords, or symbols (like punctuation) that are called **tokens**



# Creating Token Embedding During Model Training



[https://www.youtube.com/watch?v=x2tYF\\_DHpgc](https://www.youtube.com/watch?v=x2tYF_DHpgc)

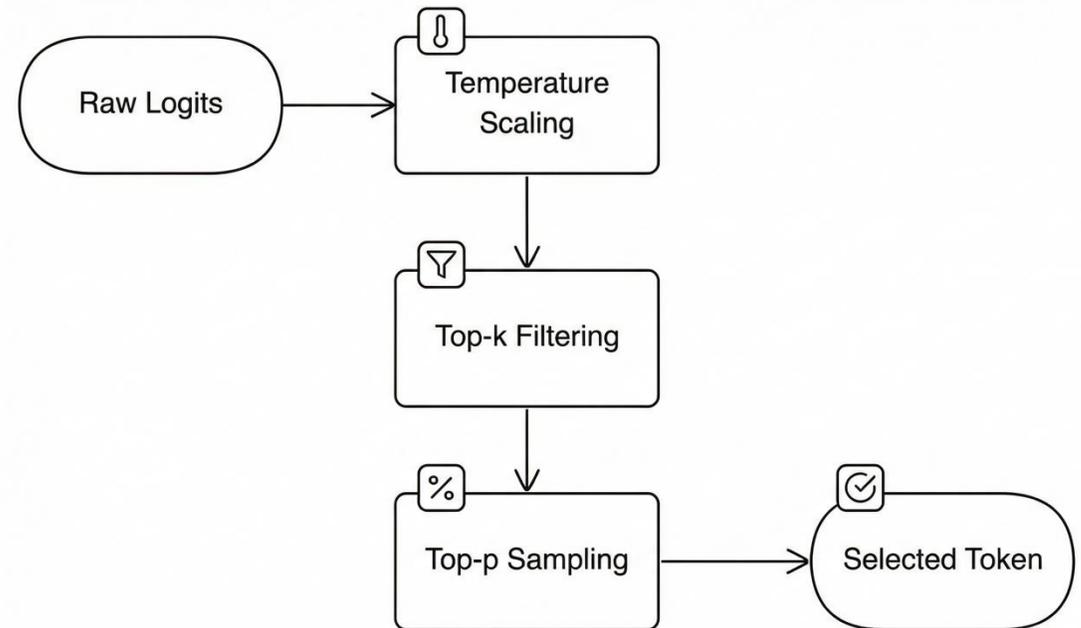
<https://github.com/rasbt/LLMs-from-scratch>

# Decoding Phase in Inference Process:

- 1. Attention Computation:** Looking back at all previous tokens to understand context
- 2. Probability Calculation:** Determining the likelihood of each possible next token
- 3. Token Selection:** Choosing the next token based on these probabilities
- 4. Continuation Check:** Deciding whether to continue or stop generation

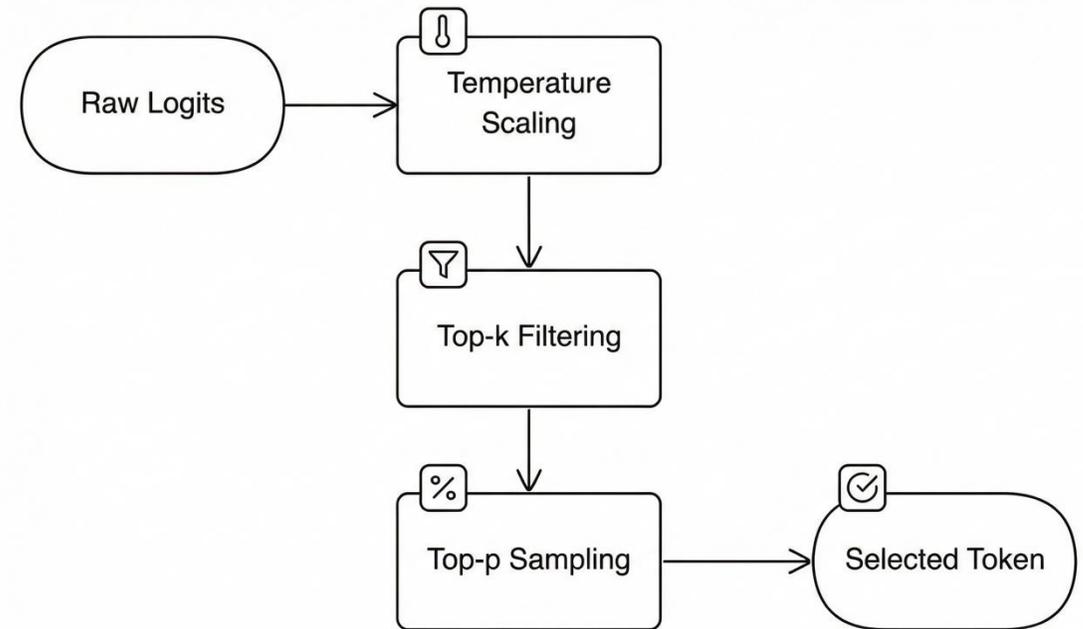
# Token Selection: Probabilities to Token Choices

- **Raw Logits:** Think of these as the model's initial gut feelings about each possible next word.
- **Temperature Control:** Like a creativity dial - higher settings ( $>1.0$ ) make choices more random and creative, lower settings ( $<1.0$ ) make them more focused and deterministic.

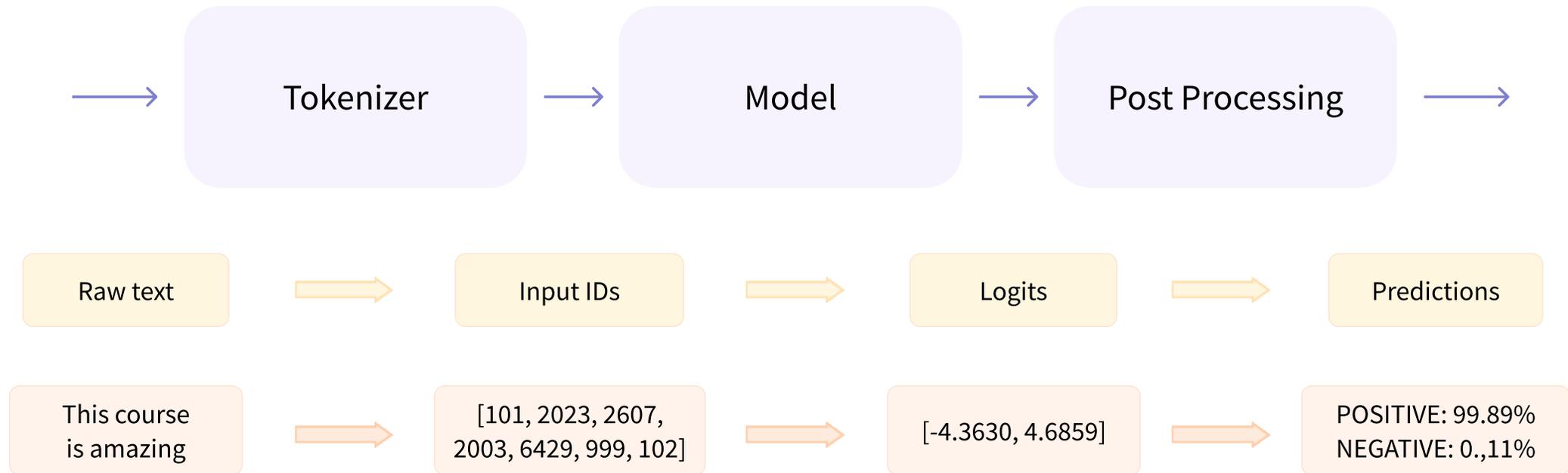


# Token Selection: Probabilities to Token Choices

- **Top-p (Nucleus) Sampling:** Instead of considering all possible words, we only look at the most likely ones that add up to our chosen probability threshold (e.g., top 90%).
- **Top-k Filtering:** An alternative approach where we only consider the k most likely next words.



# Full LLM Pipeline in Inference



<https://huggingface.co/learn/llm-course/chapter2/2>

# Key Metrics

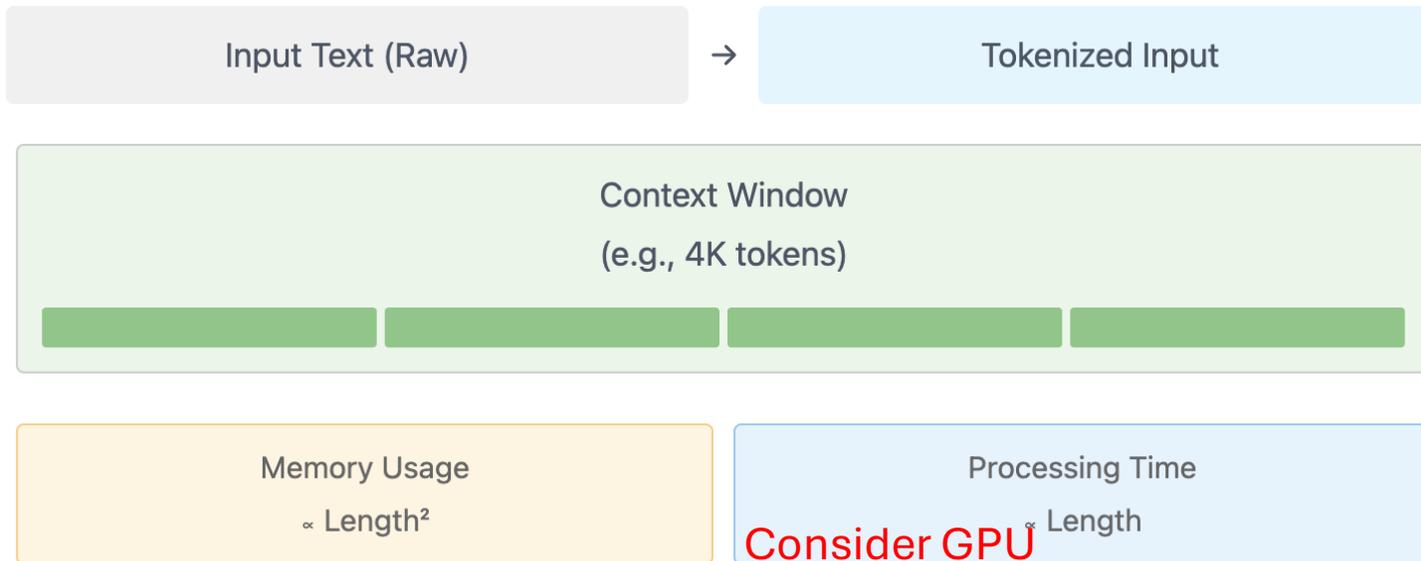
- **Time to First Token (TTFT):** How quickly can you get the first response? This is crucial for user experience and is primarily affected by the prefill phase.
- **Time Per Output Token (TPOT):** How fast can you generate subsequent tokens? This determines the overall generation speed.
- **Throughput:** How many requests can you handle simultaneously? This affects scaling and cost efficiency.
- **VRAM Usage:** How much GPU memory do you need? This often becomes the primary constraint in real-world applications.

# Attention and Context Length

- The attention mechanism is what gives LLMs their ability to understand context and focus on relevant information.
  - **Remember the similarity score and attention weights in last lecture?**
  - When predicting the next word, not every word in a sentence carries equal weight.
- The context length refers to the maximum number of tokens (words or parts of words) that the LLM can process at once.
  - <https://arxiv.org/html/2412.19437v1#S4> (Long Context Extension)
- The context is the LLM working memory, limited by **model architecture and size, computational resources, and the complexity of input and the desired output.**

# Context Length Management in Inference

- **Memory Usage:** Grows quadratically with context length
- **Processing Speed:** Decreases linearly with longer contexts
- **Resource Allocation:** Requires careful balancing of VRAM usage



Consider GPU  
limited bandwidth

**One solution: KV Cache;  
e.g., PagedAttention  
<https://arxiv.org/abs/2309.06180>**

# Building Large Language Models



<https://www.youtube.com/watch?v=7xTGNNLPyMI>

# Another Suggested Video

---



---

<https://www.youtube.com/watch?v=9vM4p9NN0Ts>

# References

- <https://huggingface.co/learn/llm-course/chapter2/1>
- <https://arxiv.org/html/2412.19437v1>
- [https://huggingface.co/spaces/hesamation/primer-llm-embedding?section=what\\_are\\_embeddings?](https://huggingface.co/spaces/hesamation/primer-llm-embedding?section=what_are_embeddings?)