

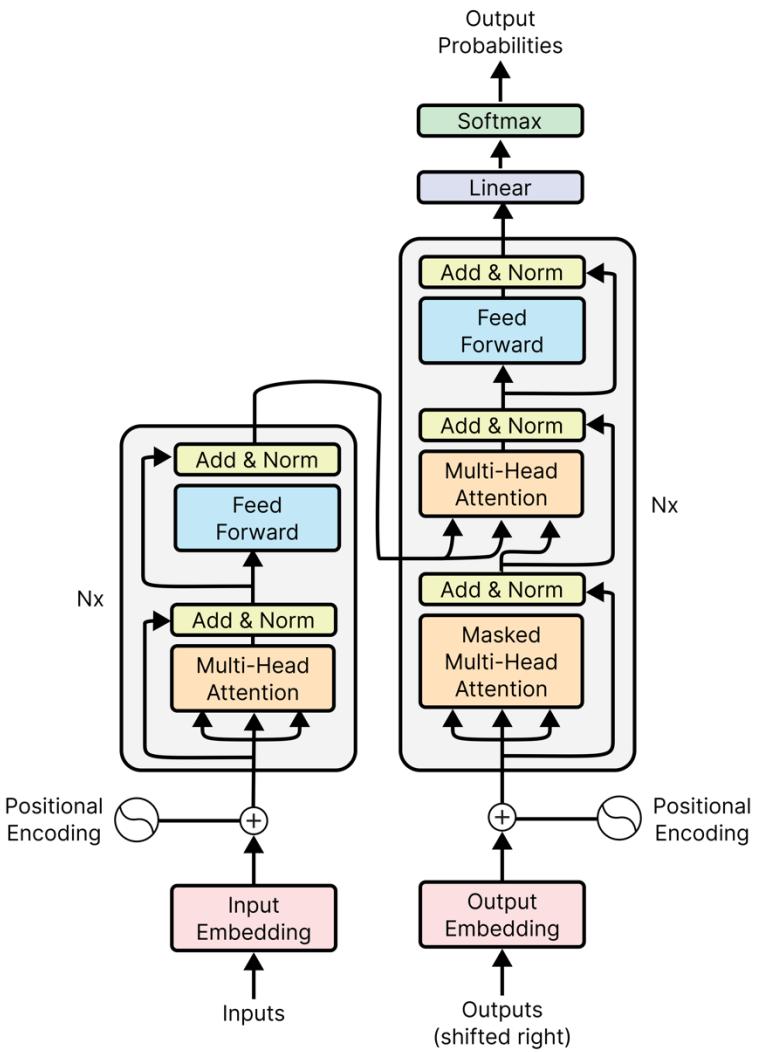
# Generative AI

2025 / 11 / 09 – Part 1

Training a Large Language Model  
*Johannes Oster*

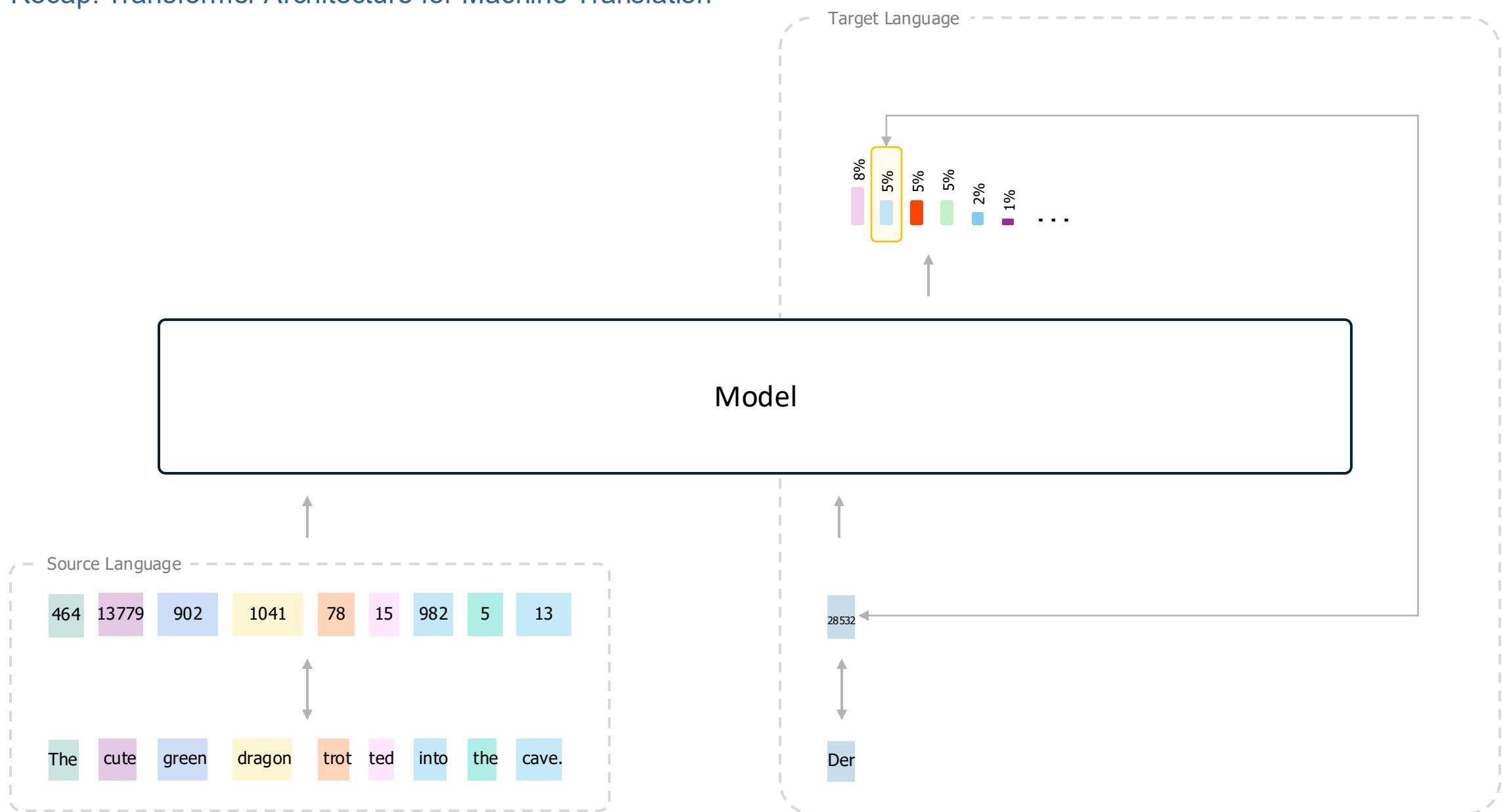
# Recap of the Transformer Architecture.

With reference to 'Attention Is All You Need' by A. Vaswani et al.<sup>1</sup>

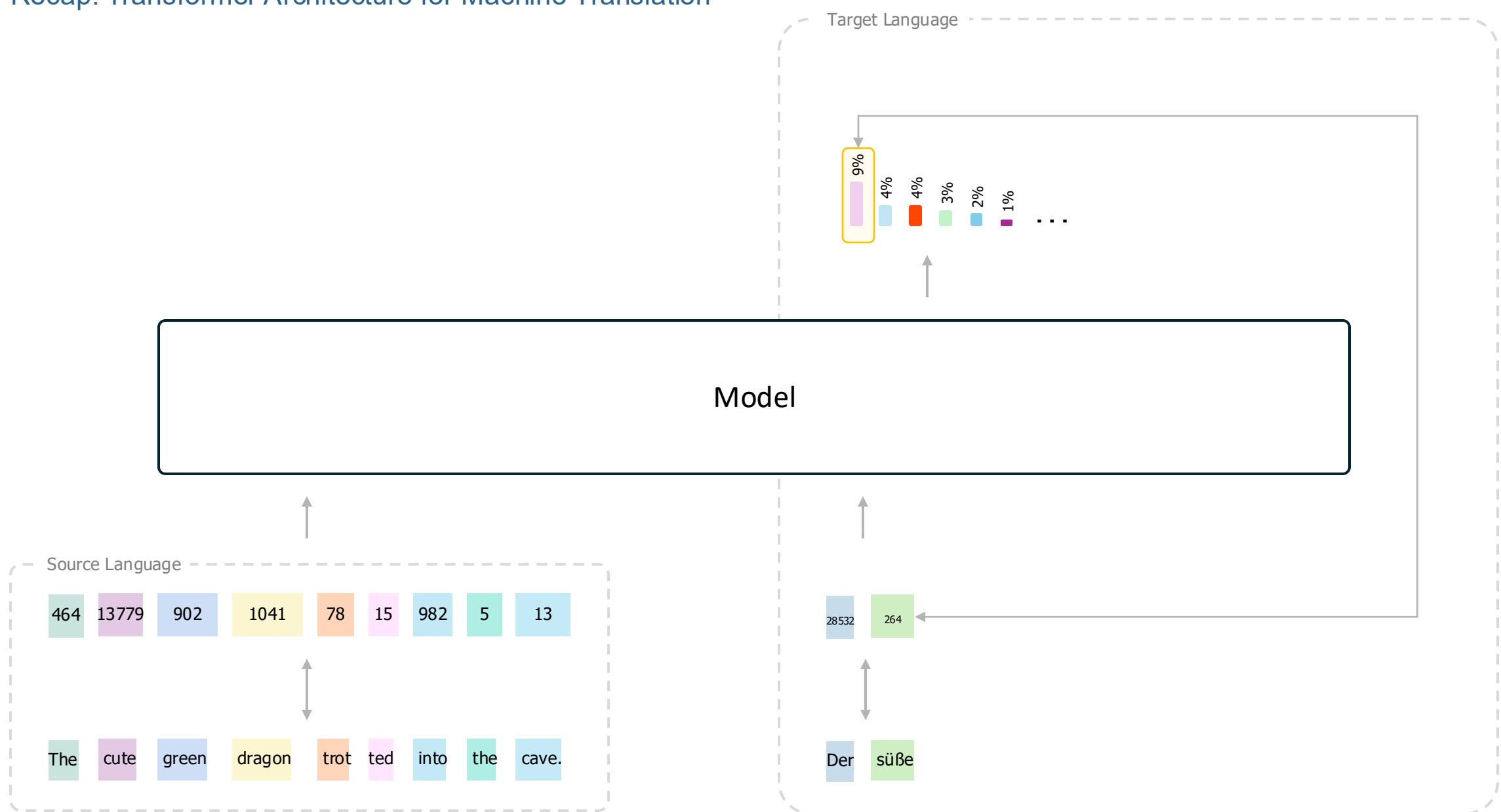


1. [https://papers.nips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fb053c1c4a845aa-Paper.pdf](https://papers.nips.cc/paper_files/paper/2017/file/3f5ee243547dee91fb053c1c4a845aa-Paper.pdf)

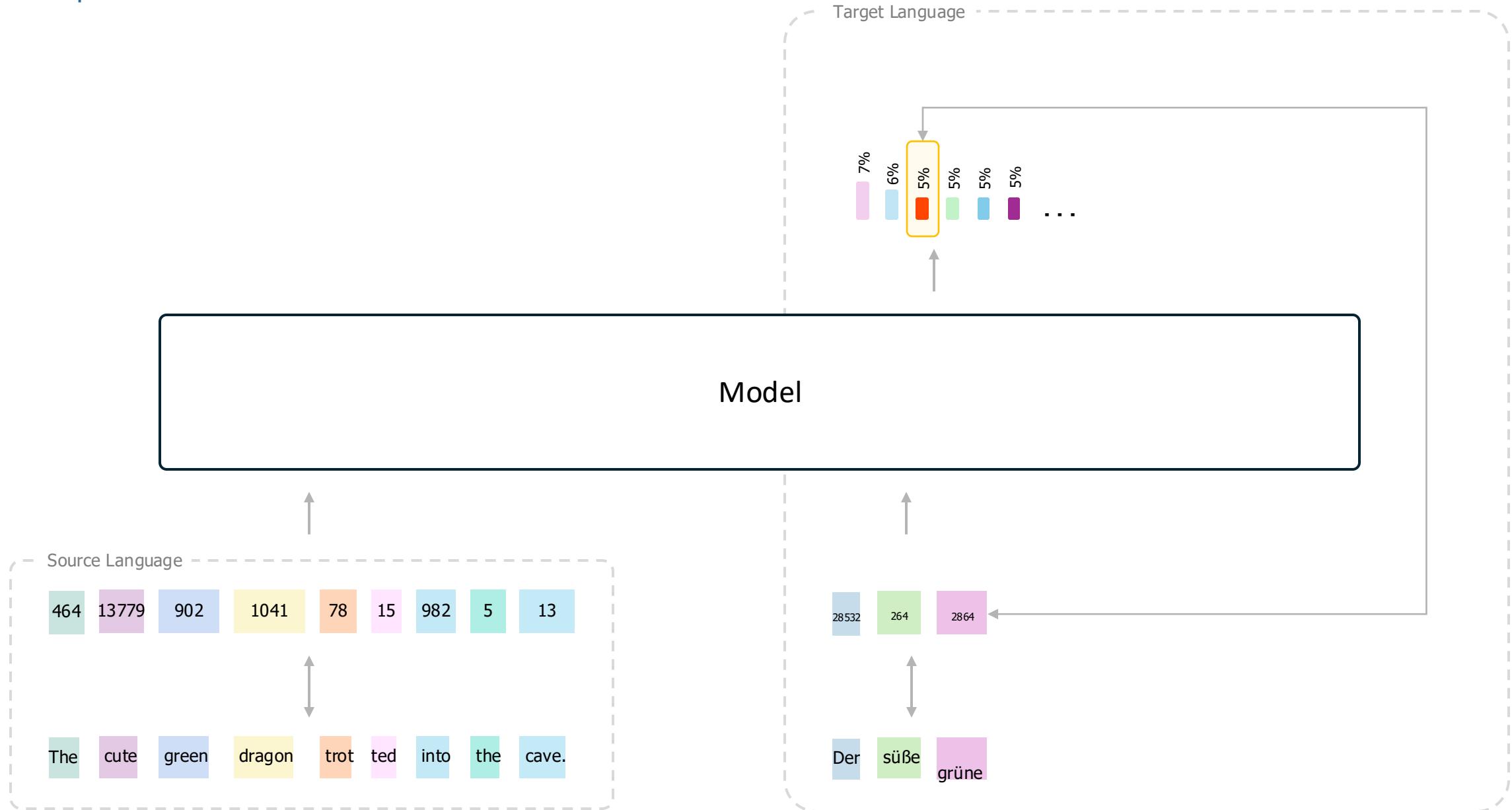
## Recap: Transformer Architecture for Machine Translation



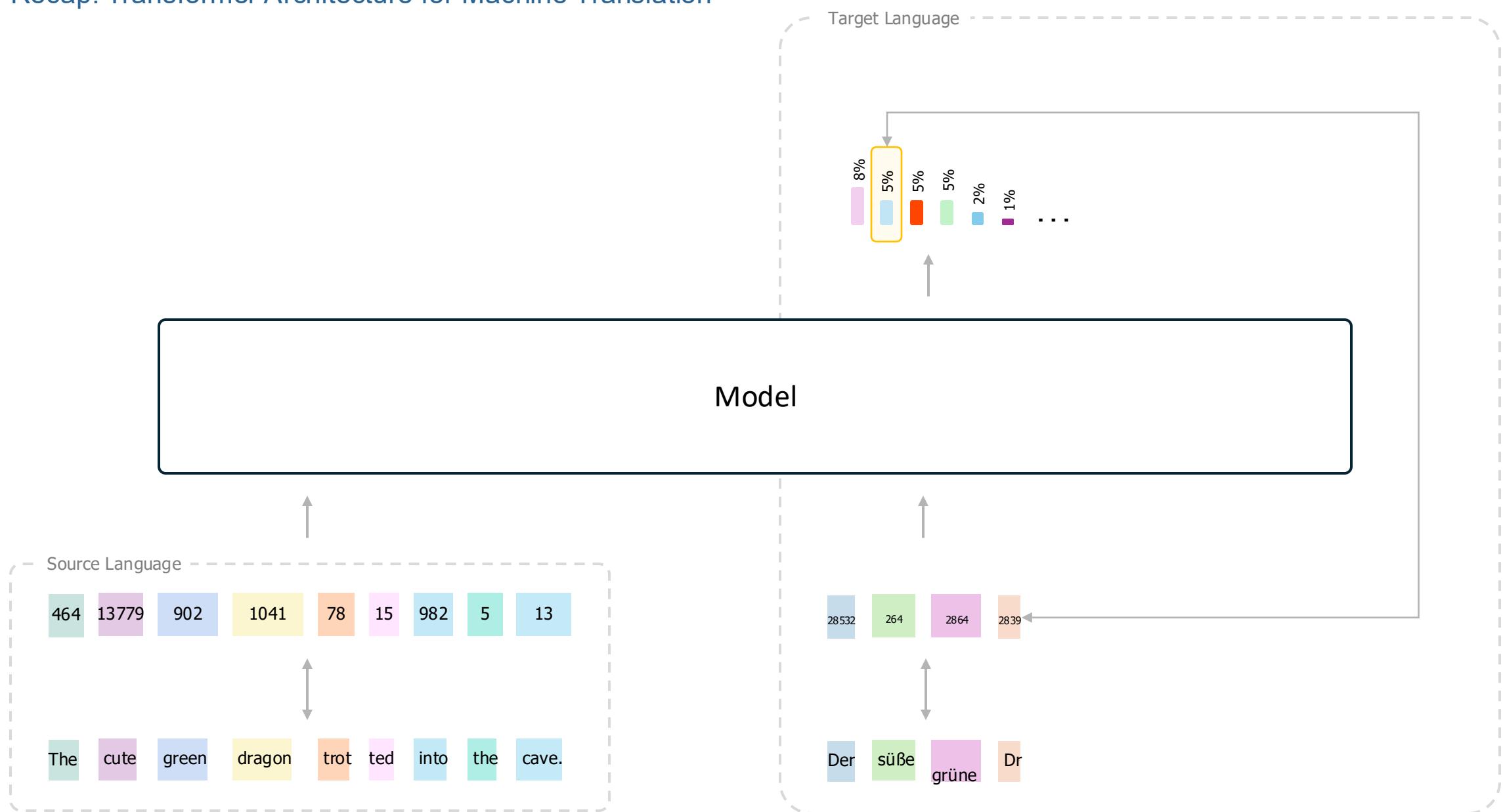
## Recap: Transformer Architecture for Machine Translation



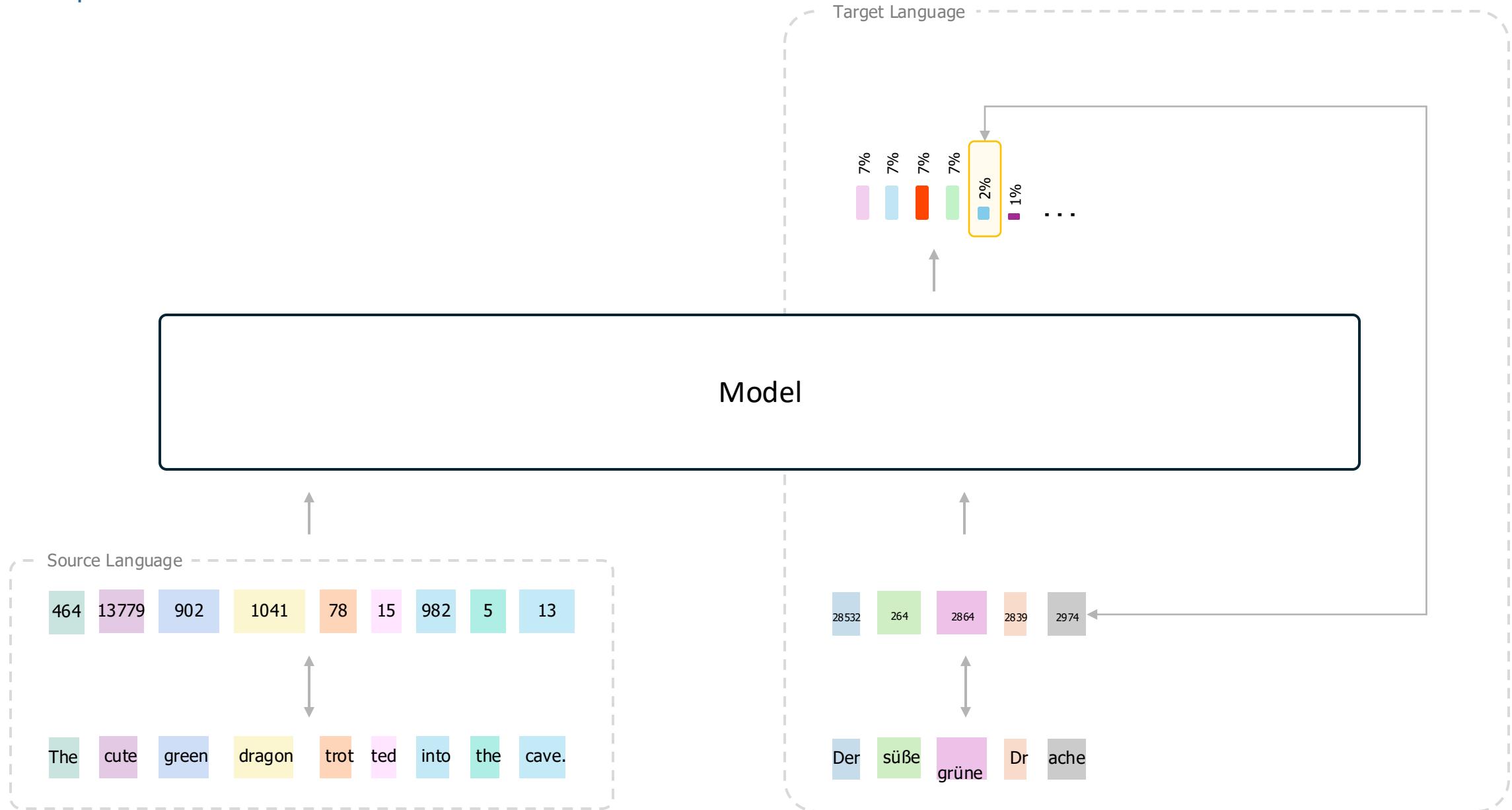
## Recap: Transformer Architecture for Machine Translation



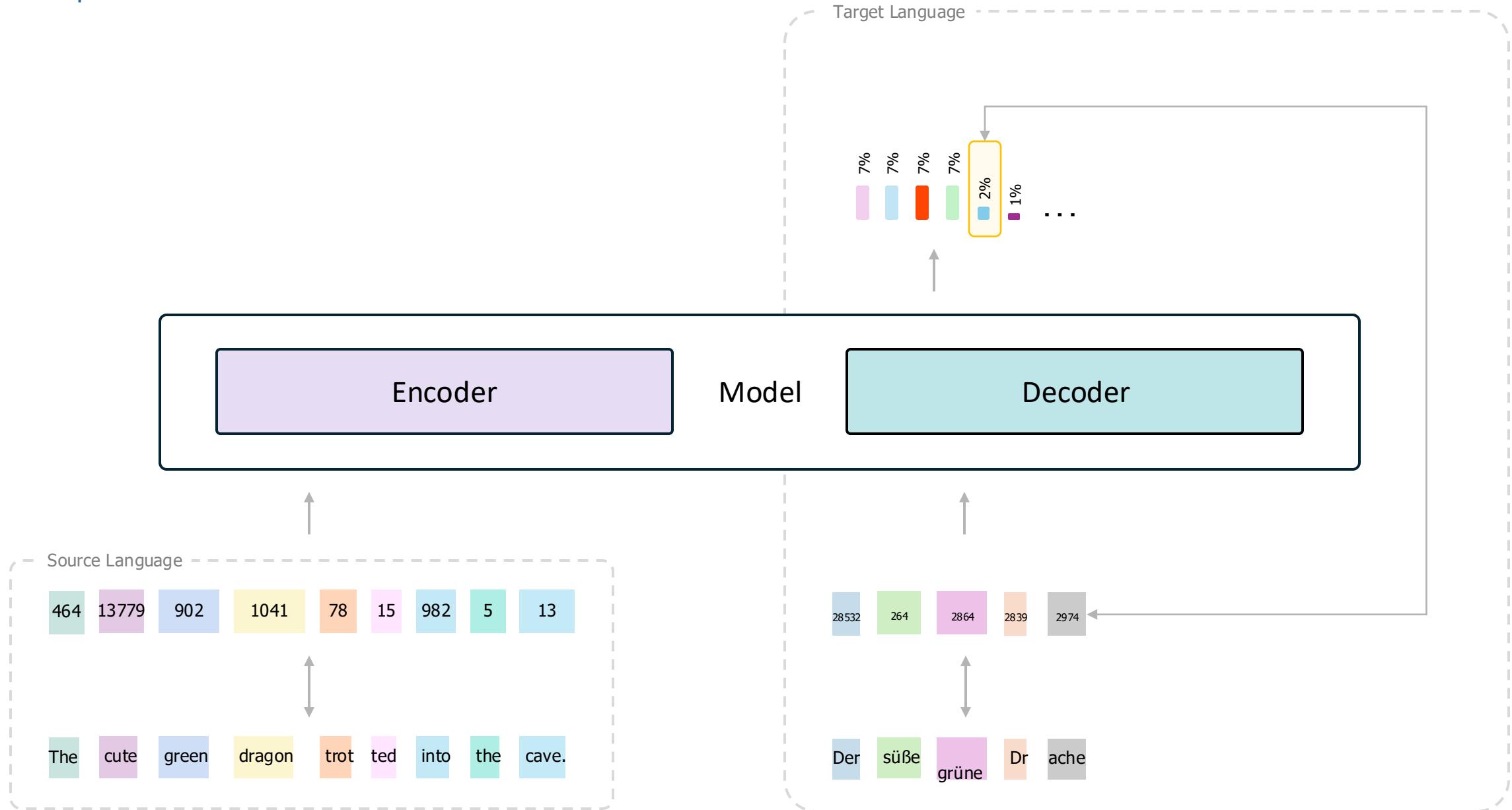
## Recap: Transformer Architecture for Machine Translation



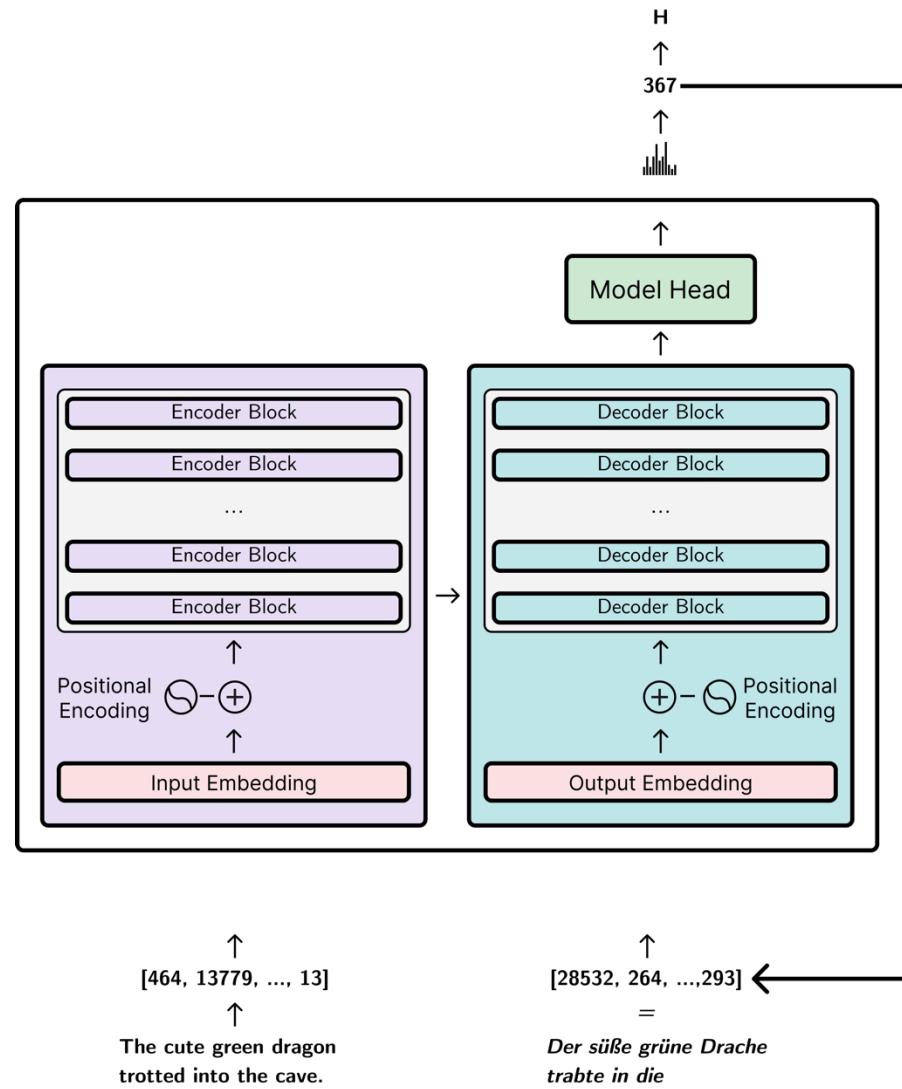
## Recap: Transformer Architecture for Machine Translation



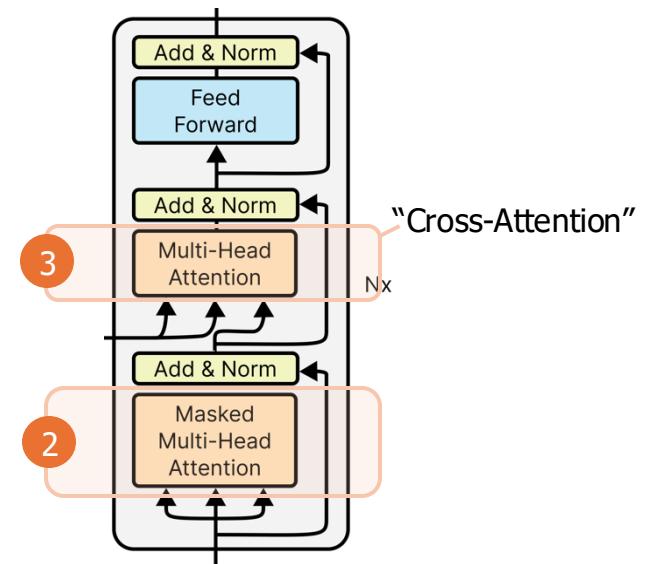
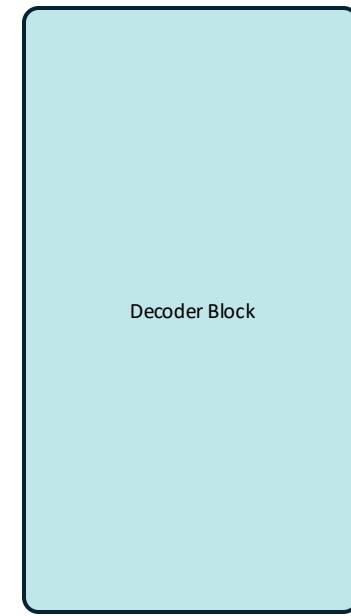
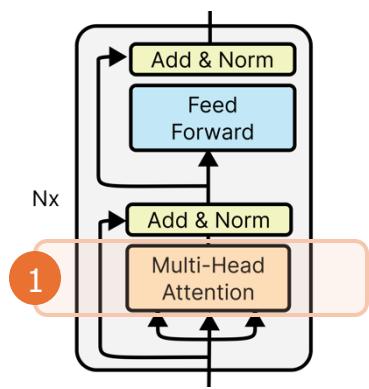
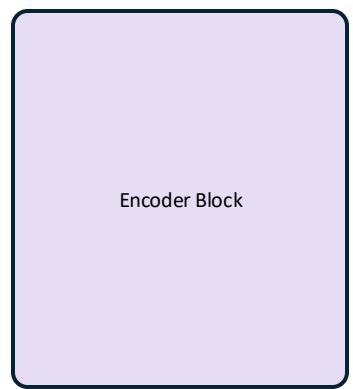
## Recap: Transformer Architecture for Machine Translation



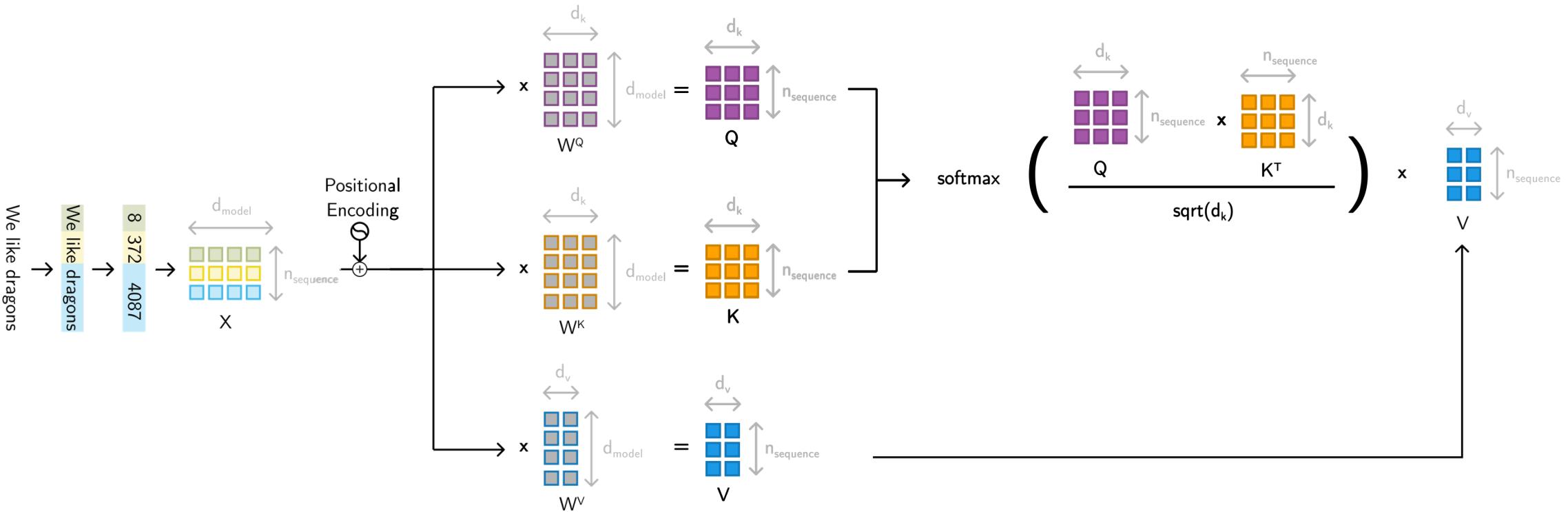
## Recap: Transformer Architecture for Machine Translation



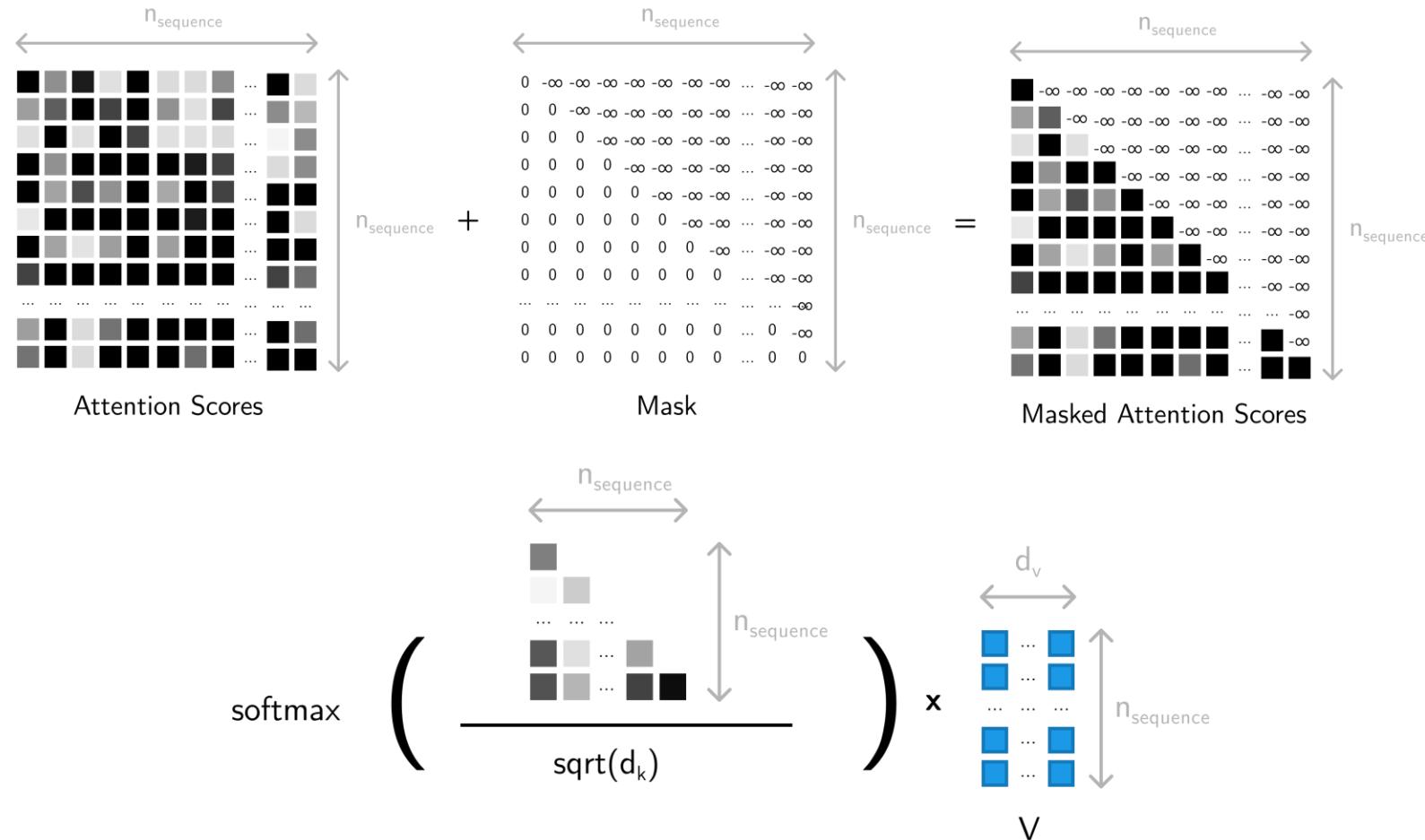
## Recap: Encoder and Decoder Blocks



## Recap: Attention Mechanism

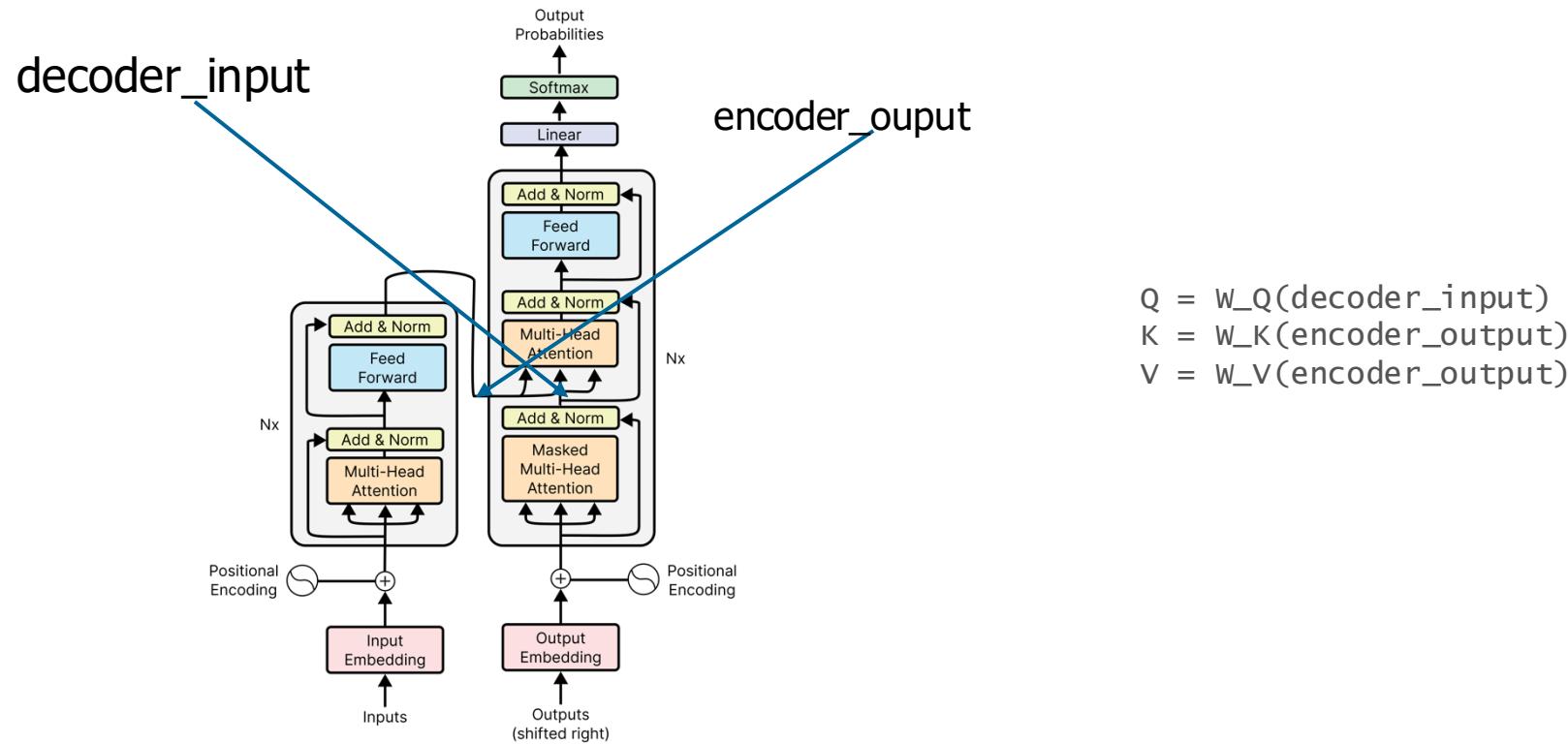


## Recap: Masked Attention Mechanism

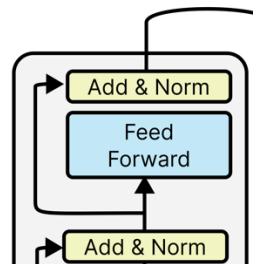


## Recap: Cross-Attention

 The query (Q) comes from the decoder layer. The key (K) and value (V) come from the encoder's output.



## Recap: Feed Forward Networks, Residual Connections and Layer Norms



```
import torch
from torch import nn

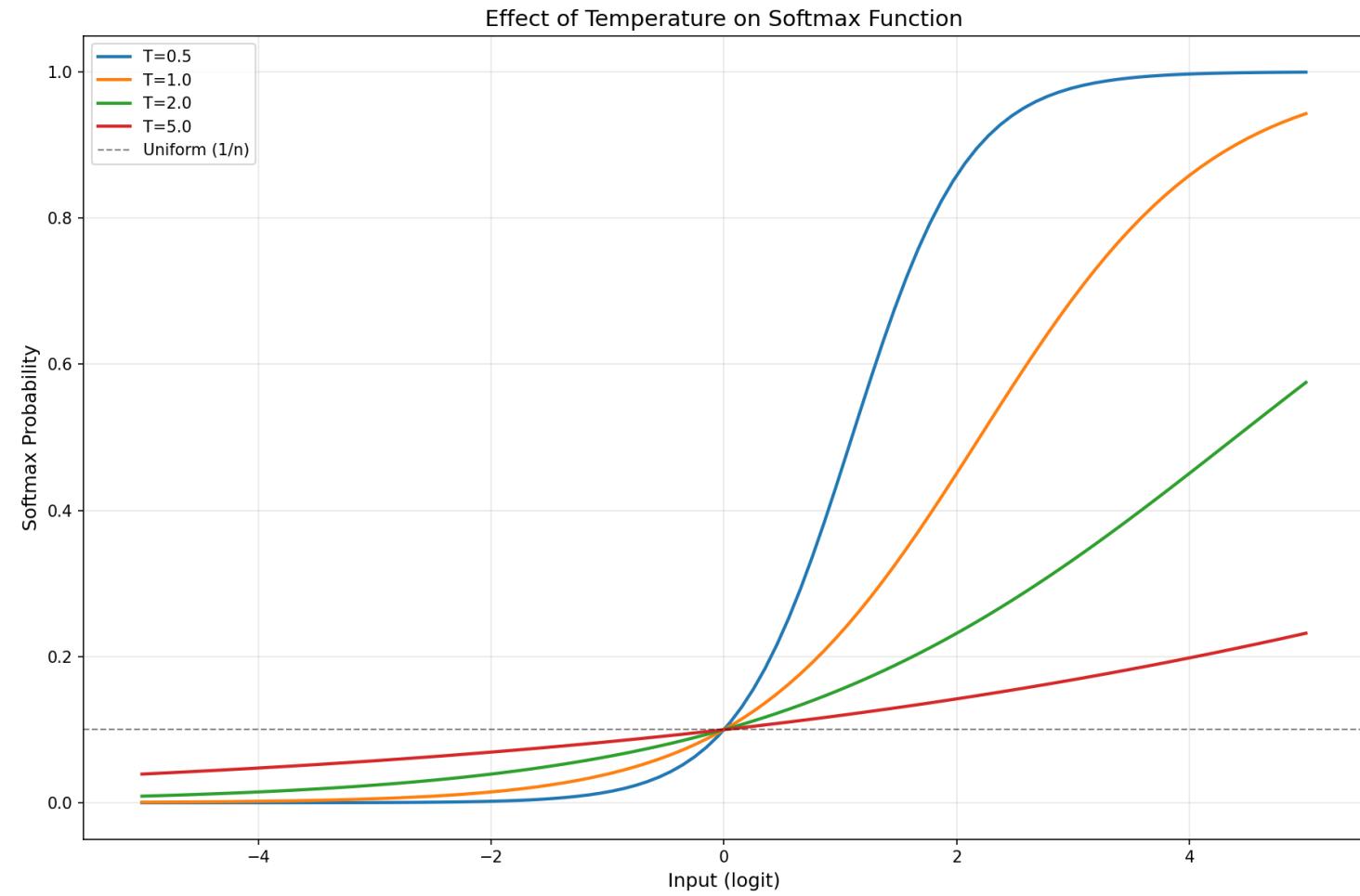
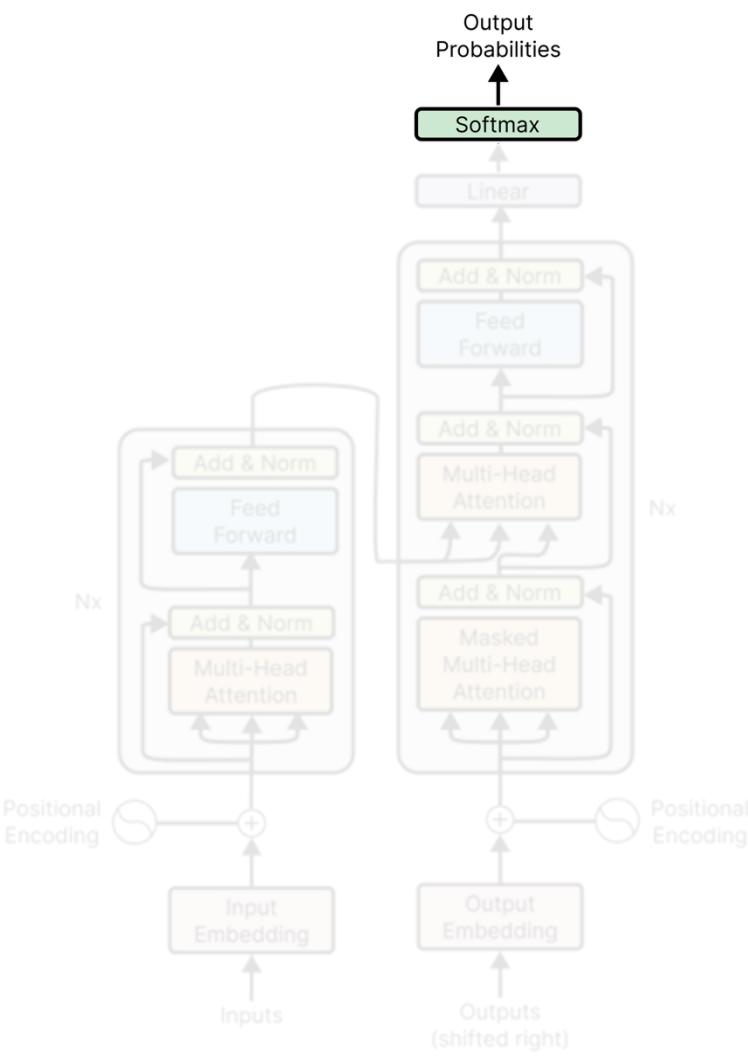
d_model = 128
d_ff = 4*d_model
n_sequence = 200 # sample sequence length
x = torch.randn(n_sequence, d_model) # [n_sequence, d_model]

norm = nn.LayerNorm(d_model)
feed_forward = nn.Sequential(
    nn.Linear(d_model, dim_ff),
    nn.ReLU(),
    nn.Linear(dim_ff, d_model)
)

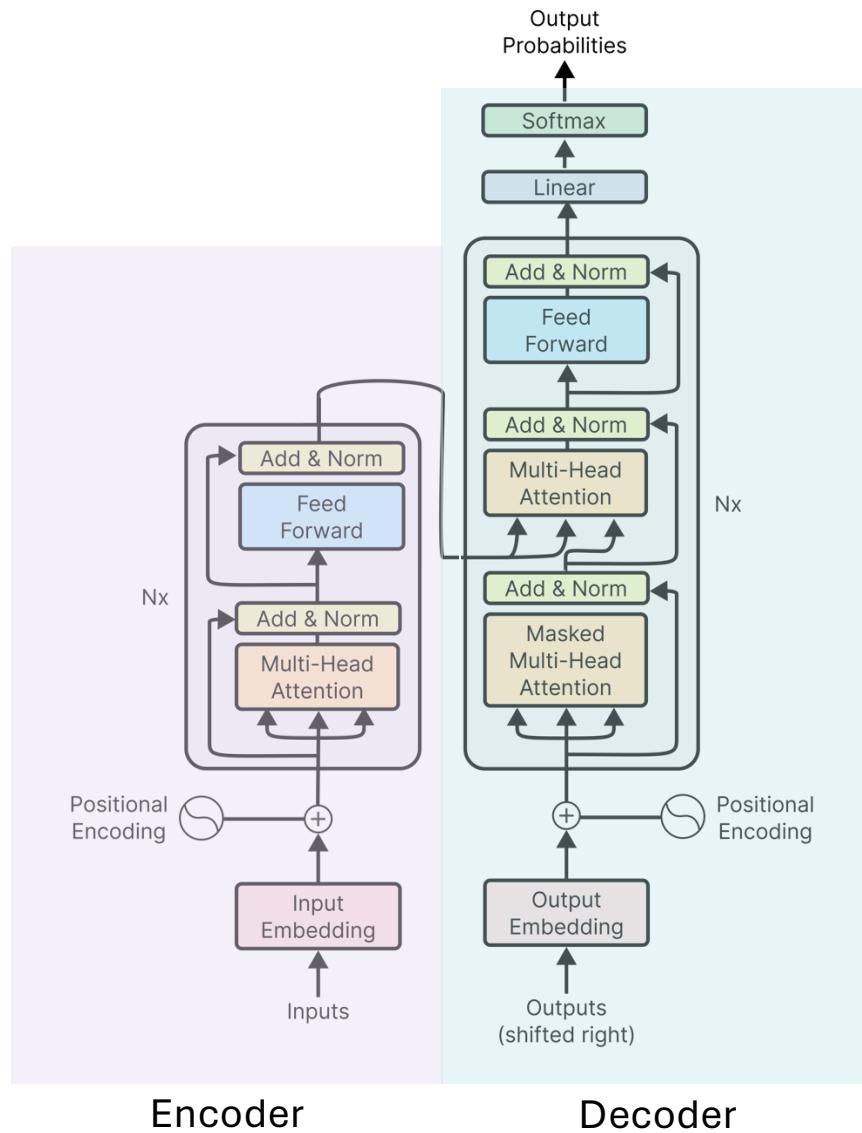
output = norm(x + feed_forward(x))
```

1. For details of the layer norm take a look at <https://docs.pytorch.org/docs/stable/generated/torch.nn.LayerNorm.html>

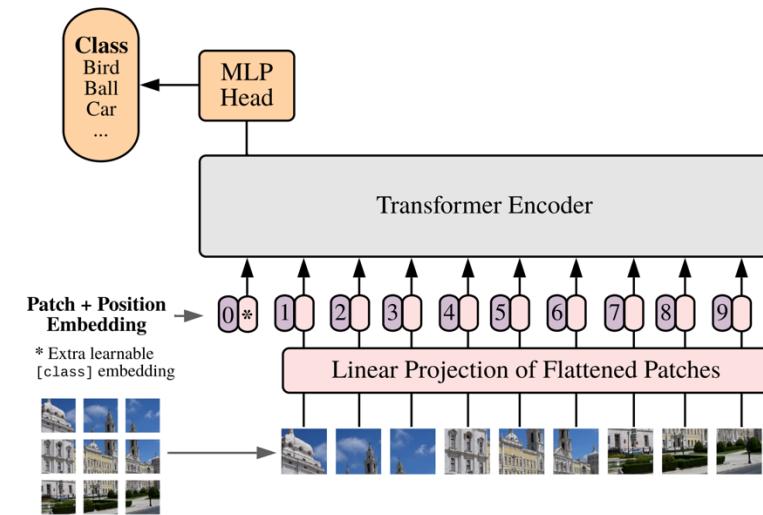
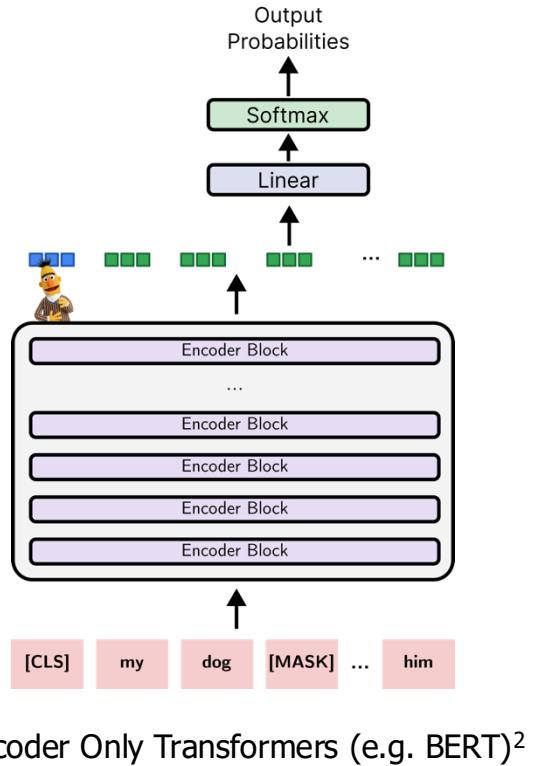
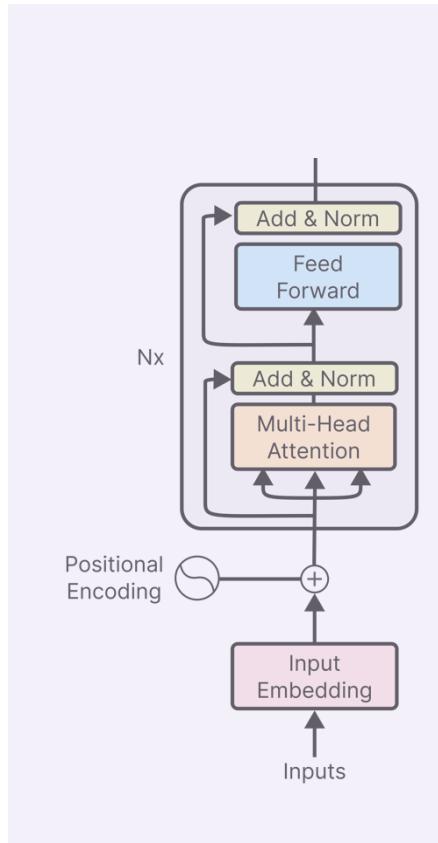
## Recap: Final Softmax with Temperature Parameter



$$\text{softmax}(x_i) = \frac{e^{\frac{x_i}{T}}}{\sum_j e^{\frac{x_j}{T}}}$$



## Encoder Only Models

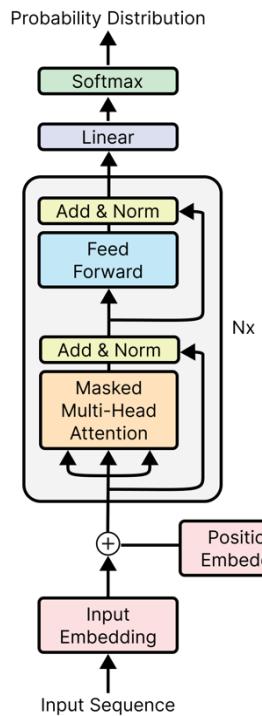


Encoder

Encoder Only Transformers (e.g. BERT)<sup>2</sup>

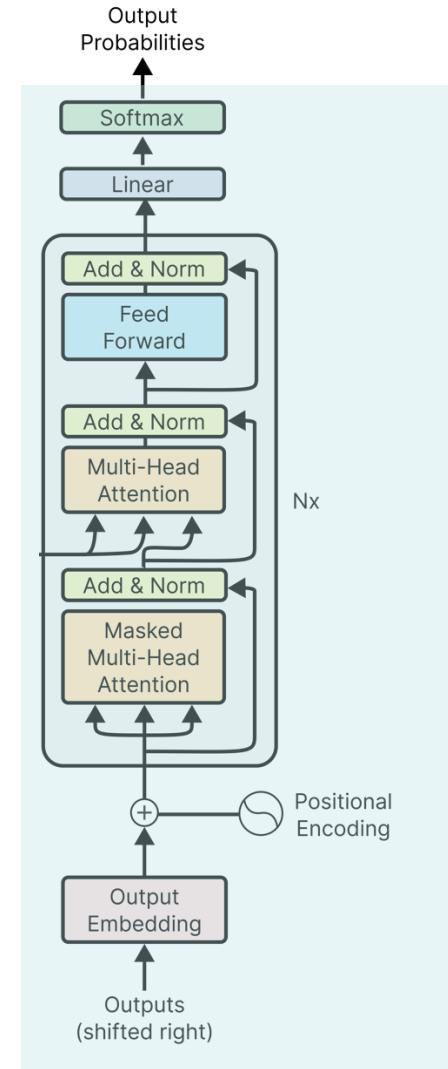
Encoder Only Vision Transformer

## Decoder Only Models



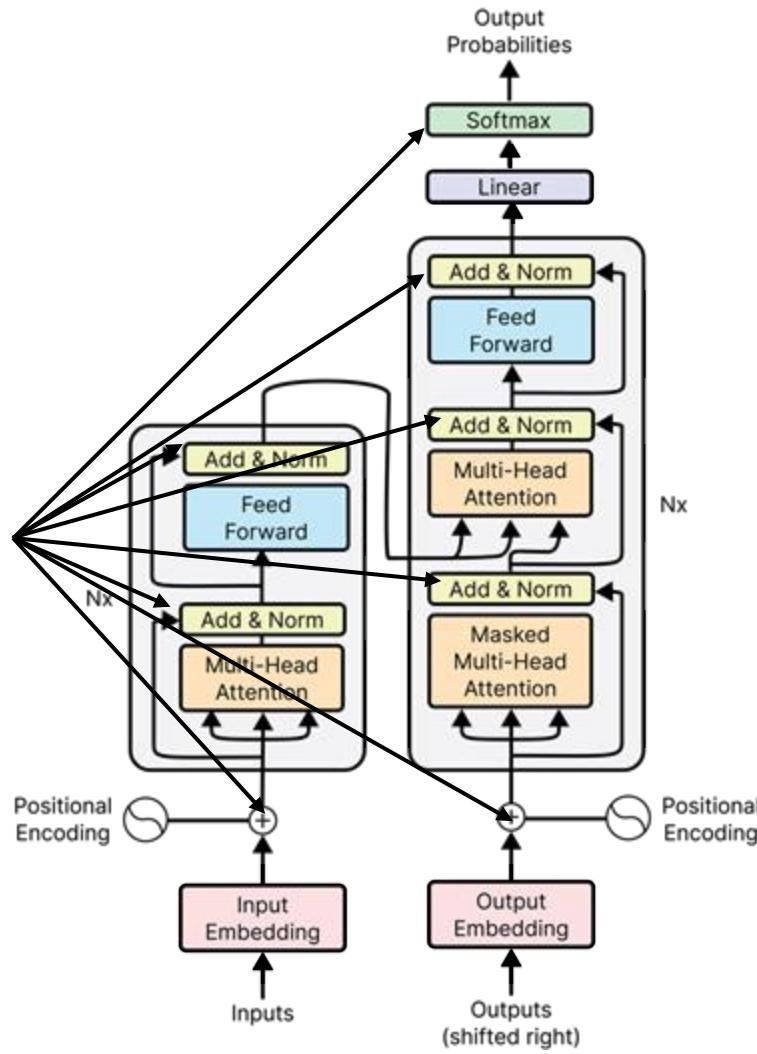
Decoder Only Transformers (e.g. GPT)<sup>3</sup>

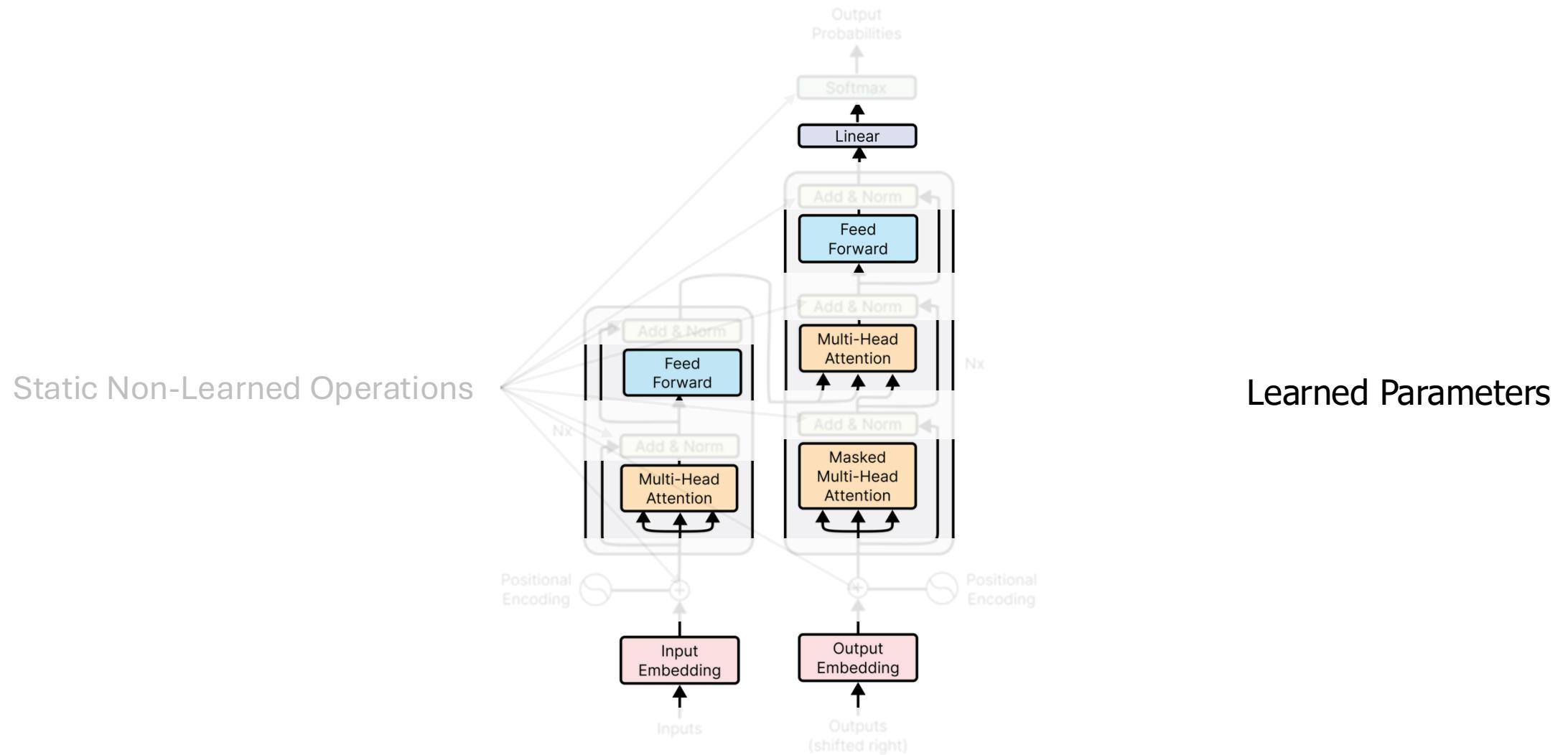
[...] Pretty much any LLM you can think of [...]



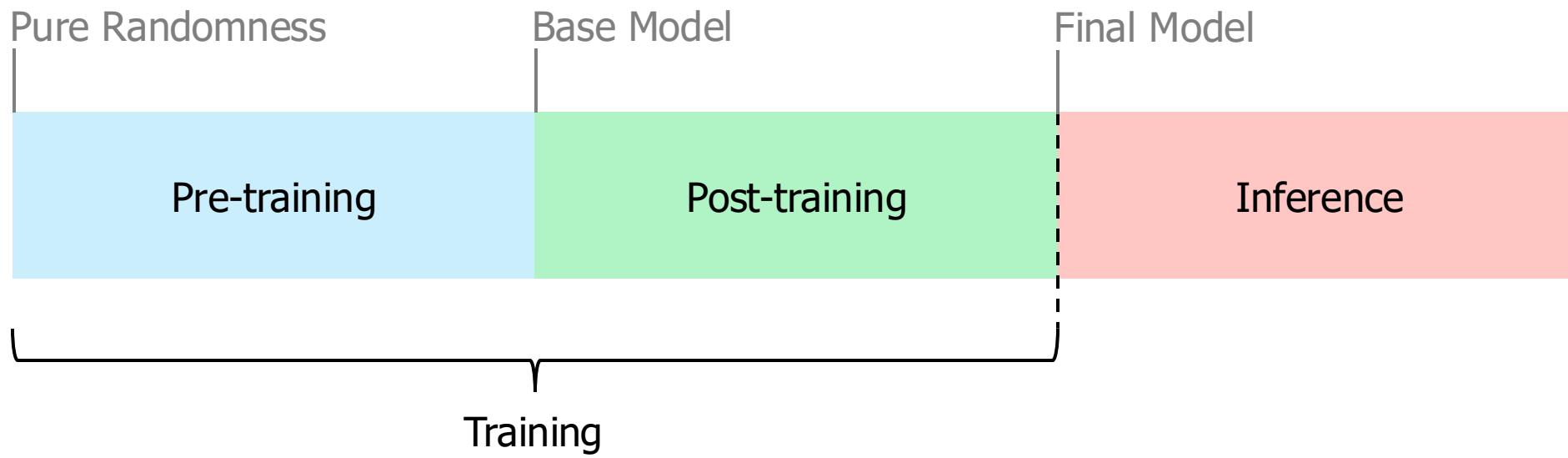
Decoder

## Non-Learned Operations

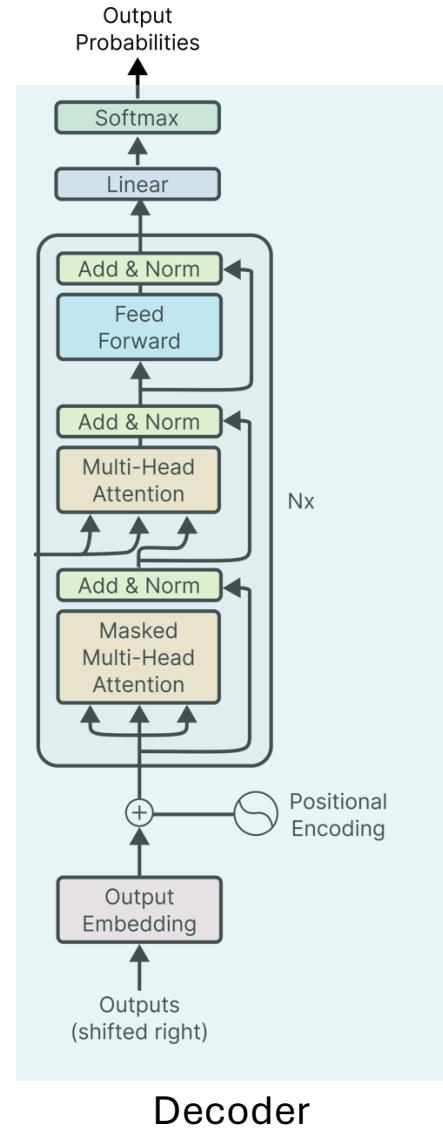




# How Computers Learn To Read and Write Training a Large Language Model



# Training Of Decoder Only Models



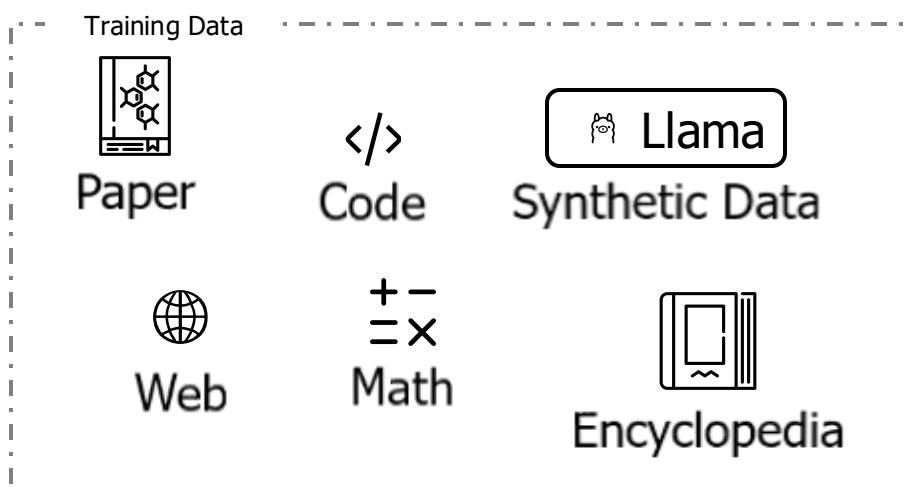
 START

*model.py* ~500 lines of code defining the models architecture with randomly initialized weights

```
# Random weight initialization example
layer = nn.Linear(512, 256)
nn.init.normal_(layer.weight, mean=0.0, std=0.02)
```

 METHODOLOGY

### *Self-Supervised Training*

 END RESULT

A *Base Model (Pre-trained Model)* that can recall factual knowledge and produce syntactically and semantically coherent text. "*Text completion on steroids*"

*but*

- can't follow instructions
- is aware of external alignment objectives

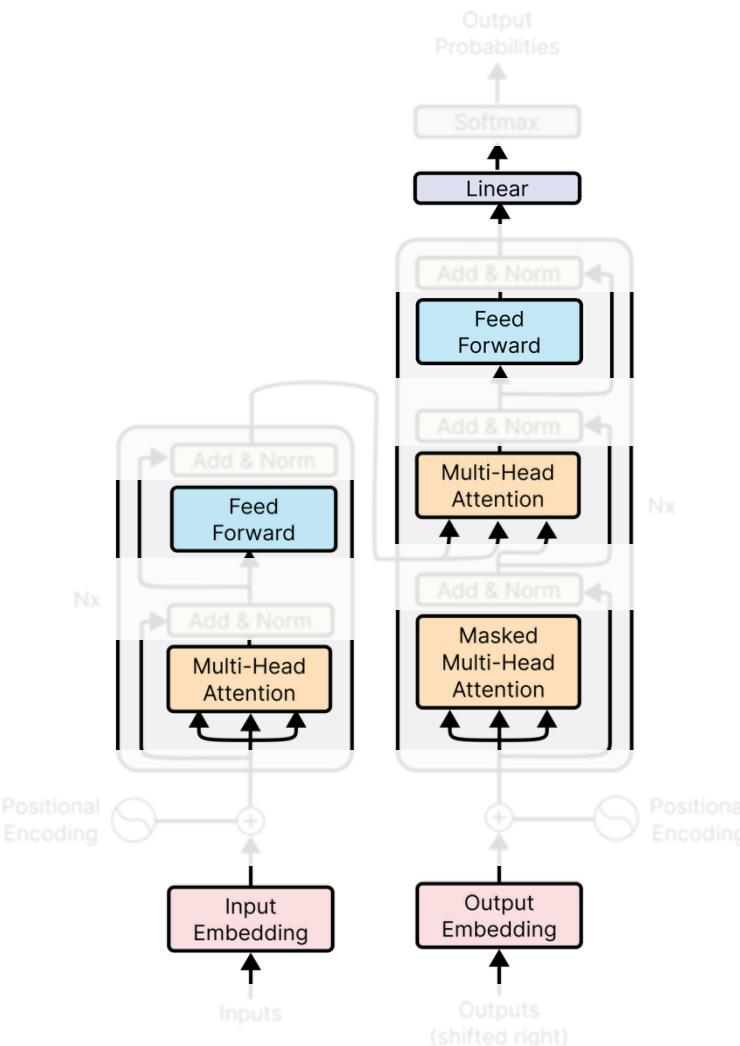
⌚ START

Base Model / Pre-trained Model

💡 METHODOLOGY & 🏆 END RESULT(S)

- Supervised Finetuning (SFT) / Instruction Tuning (IT)
- Preference Finetuning (PrefT)
- Reinforcement Learning with Verifiable Results (RLVR)

→ *Post-training strategies will be covered by Prof. Heitzinger*



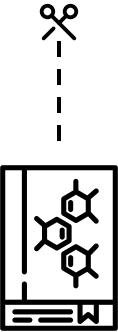
*model.py* ~500 lines of code defining the models architecture with randomly initialized weights

```
# Random weight initialization example
layer = nn.Linear(512, 256)
nn.init.normal_(layer.weight, mean=0.0, std=0.02)
```

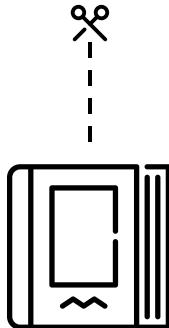
# Pre-training

Post-training    Inference

Training Data



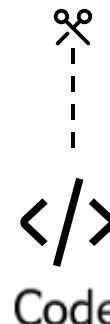
Paper



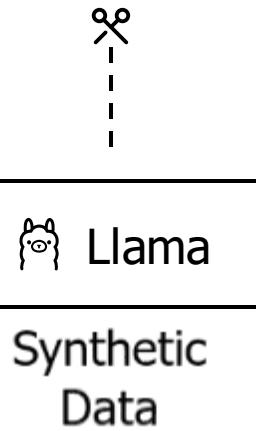
Encyclopedia



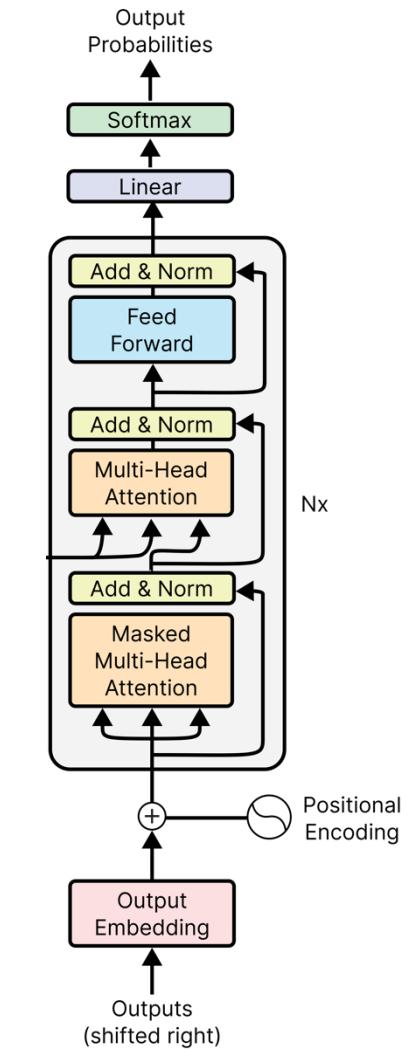
Web



Code



$15 * 10^{12}$   
 $\sim 15\,000\,000\,000\,000$





Web



Training Sample

The cute green dragon trotted into the cave.



Input Data

The cute green dragon trotted into the cave.

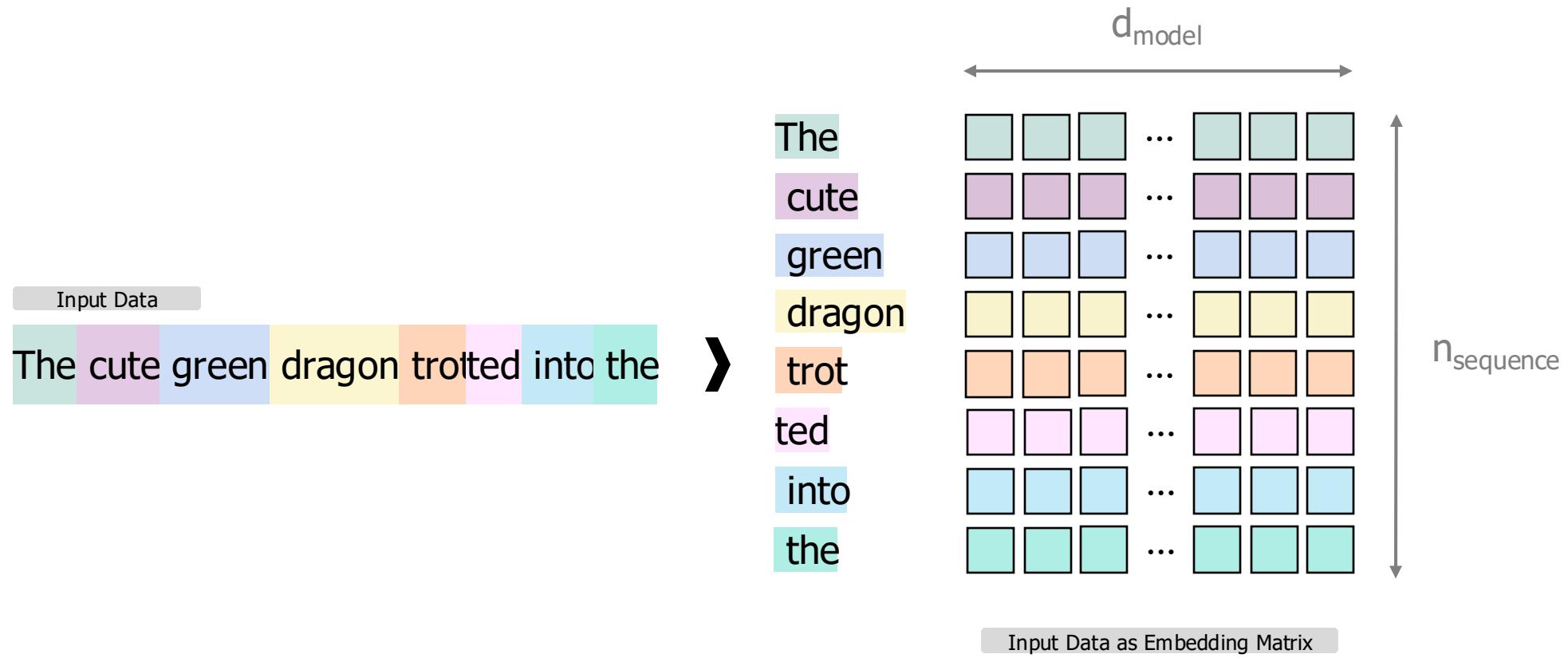
Sequence without last token



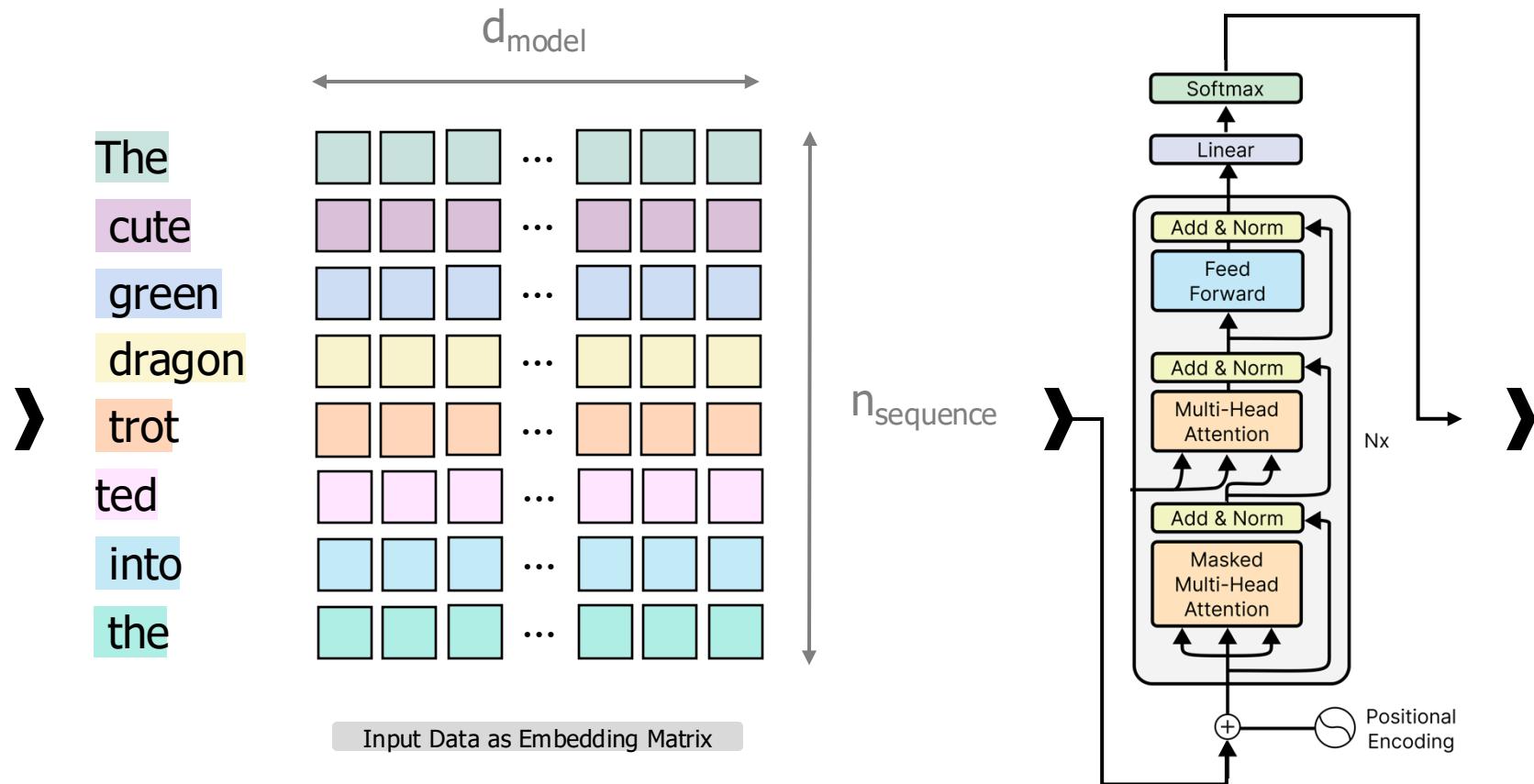
Target Data

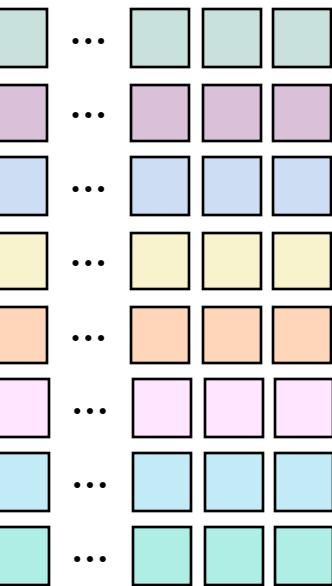
The cute green dragon trotted into the cave.

Sequence shifted by one token

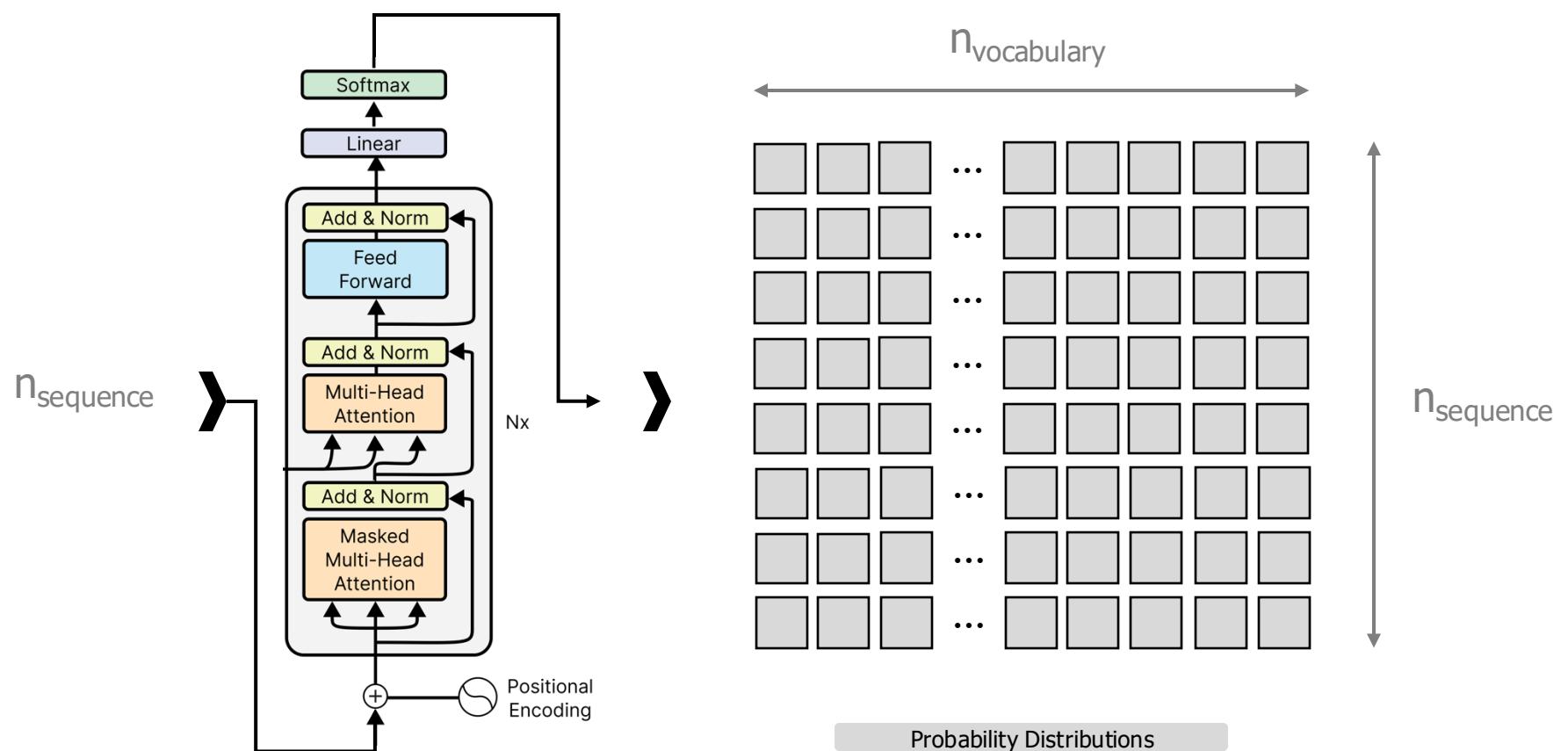


the

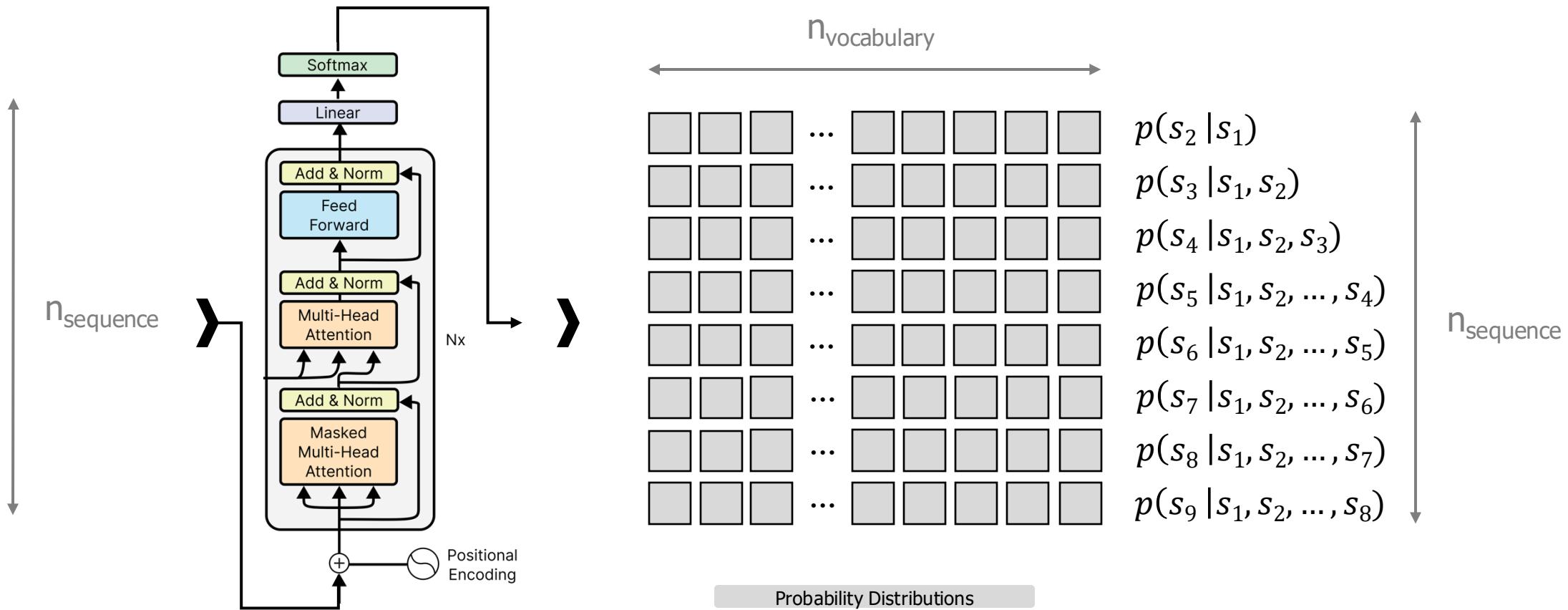


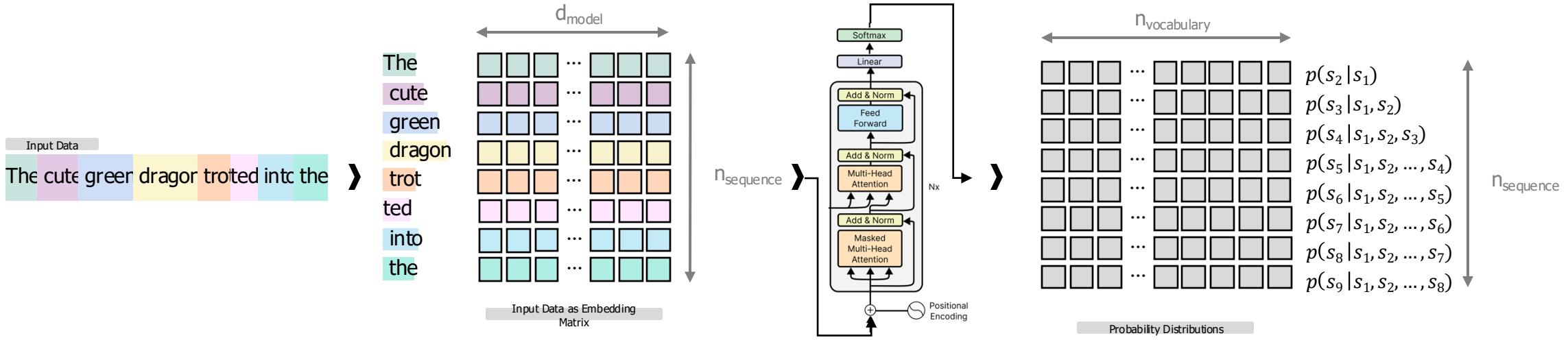
$d_{\text{model}}$ 

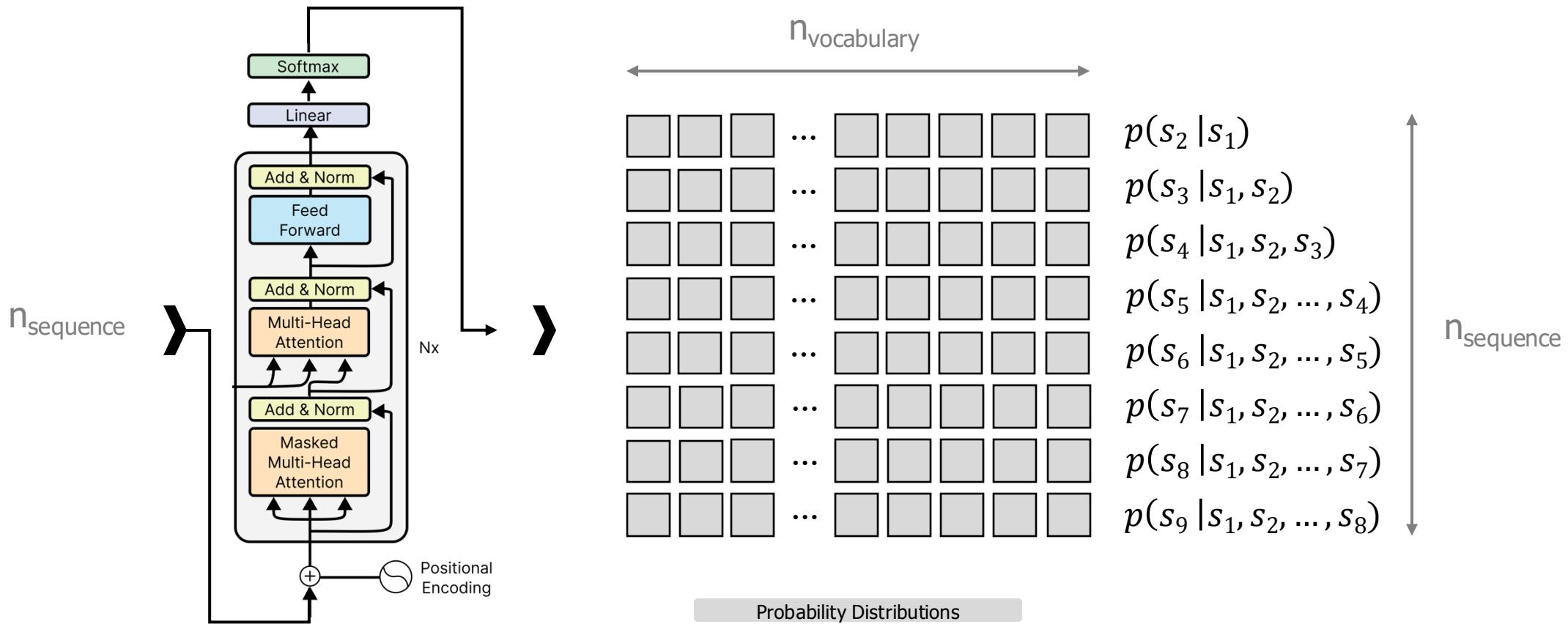
Data as Embedding Matrix



Probability Distributions







 The goal of generative language modelling is to capture the “true distribution of language”:  $p(X)$

### A few definitions:

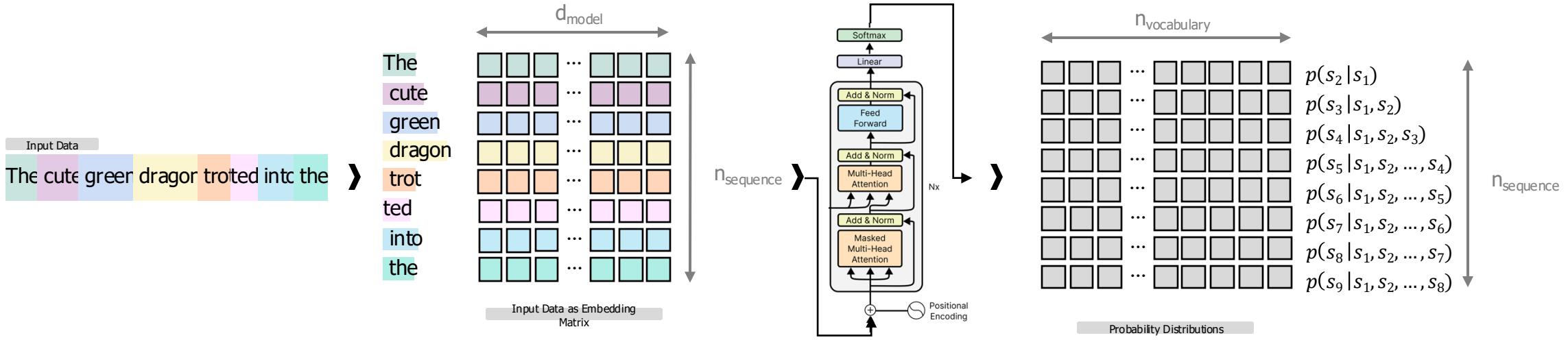
- $V$  ... vocabulary of unique tokens  $s_i \in V$
- $n$  ... maximum sequence length
- $X$  ... space of all possible token sequences

$X = \{(s_1, s_2, \dots, s_k) | s_i \in V, 0 \leq k \leq n, k \in \mathbb{Z}\}$  where  $k$  is the sequence length

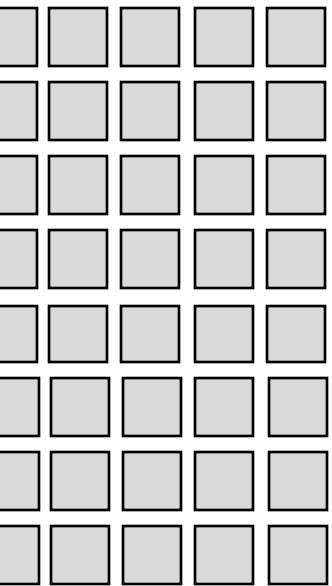
- $x \in X$  ... a particular sequence of tokens
- The probability of a given sequence  $x$

$$\begin{aligned} p(x) &= p(s_1) \cdot p(s_2 | s_1) \cdot p(s_3 | s_1, s_2) \cdot \dots \cdot p(s_k | s_1, \dots, s_{k-1}) \\ &= p(s_1) \cdot \prod_{i=2}^k p(s_i | s_1, \dots, s_{i-1}) \end{aligned}$$

 The probability of a sequence is the product of the probability of the first token and the conditional probabilities of each subsequent token given all previous tokens.



vocabulary

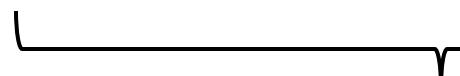


$$\begin{aligned} p(s_2 | s_1) \\ p(s_3 | s_1, s_2) \\ p(s_4 | s_1, s_2, s_3) \\ p(s_5 | s_1, s_2, \dots, s_4) \\ p(s_6 | s_1, s_2, \dots, s_5) \\ p(s_7 | s_1, s_2, \dots, s_6) \\ p(s_8 | s_1, s_2, \dots, s_7) \\ p(s_9 | s_1, s_2, \dots, s_8) \end{aligned}$$

$n_{\text{sequence}}$

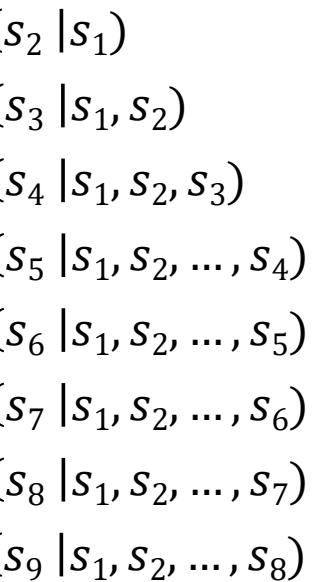
vs

Target Data  
The cute green dragon trotted into the cave.

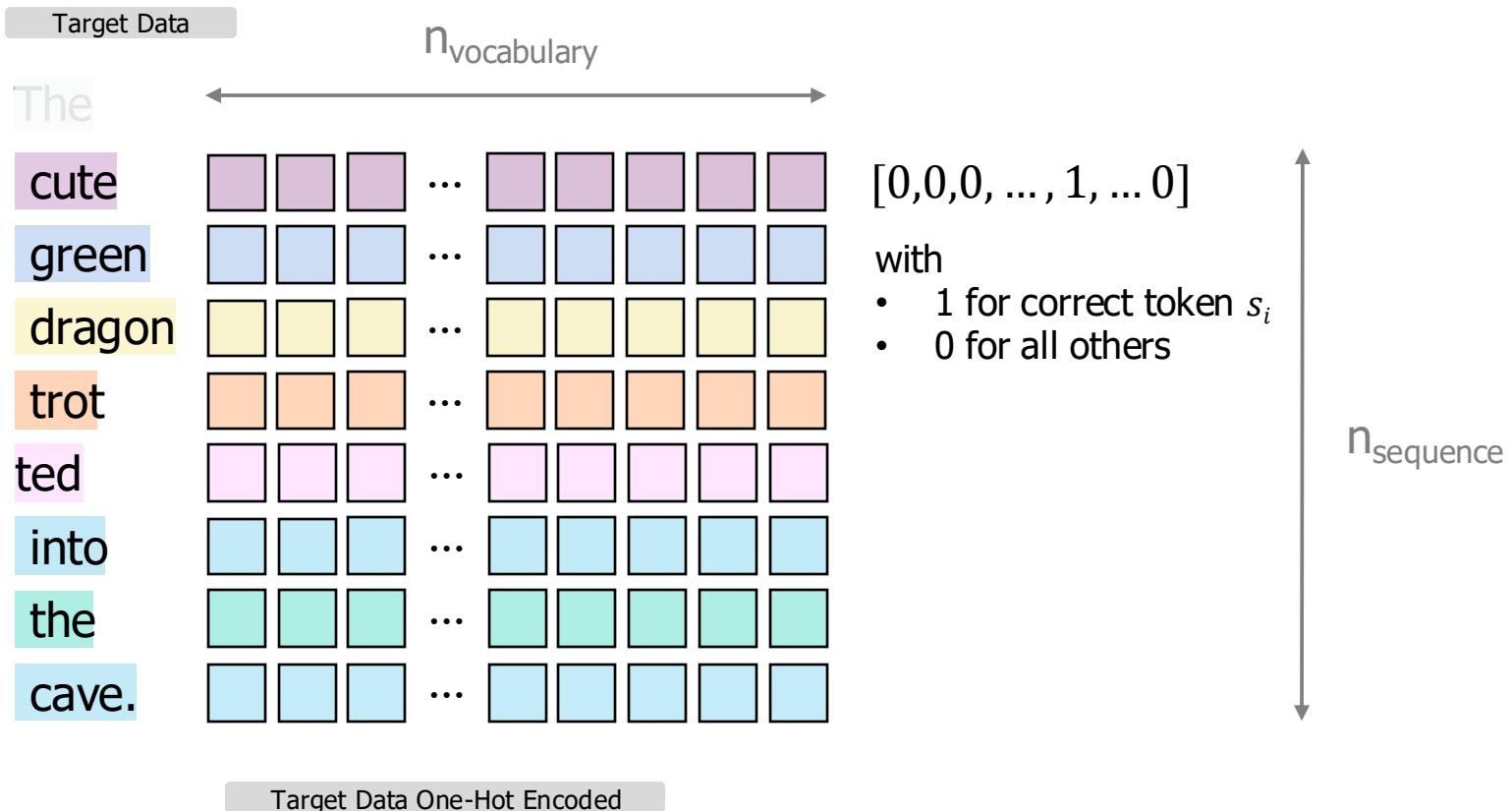


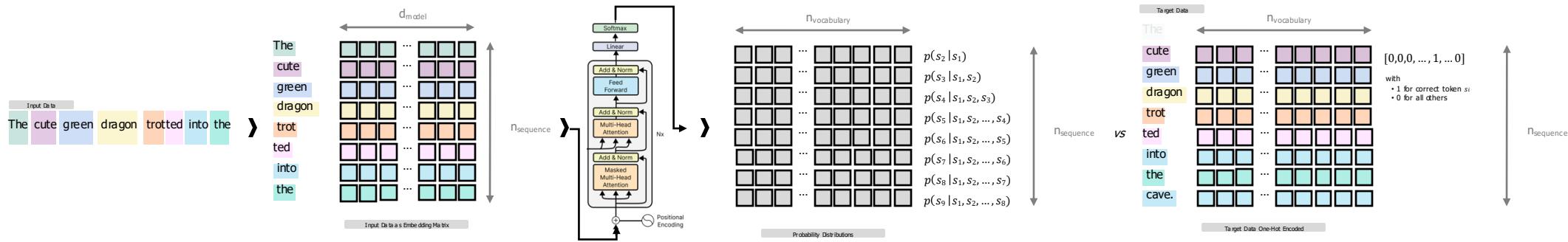
Sequence shifted by one token

Distributions



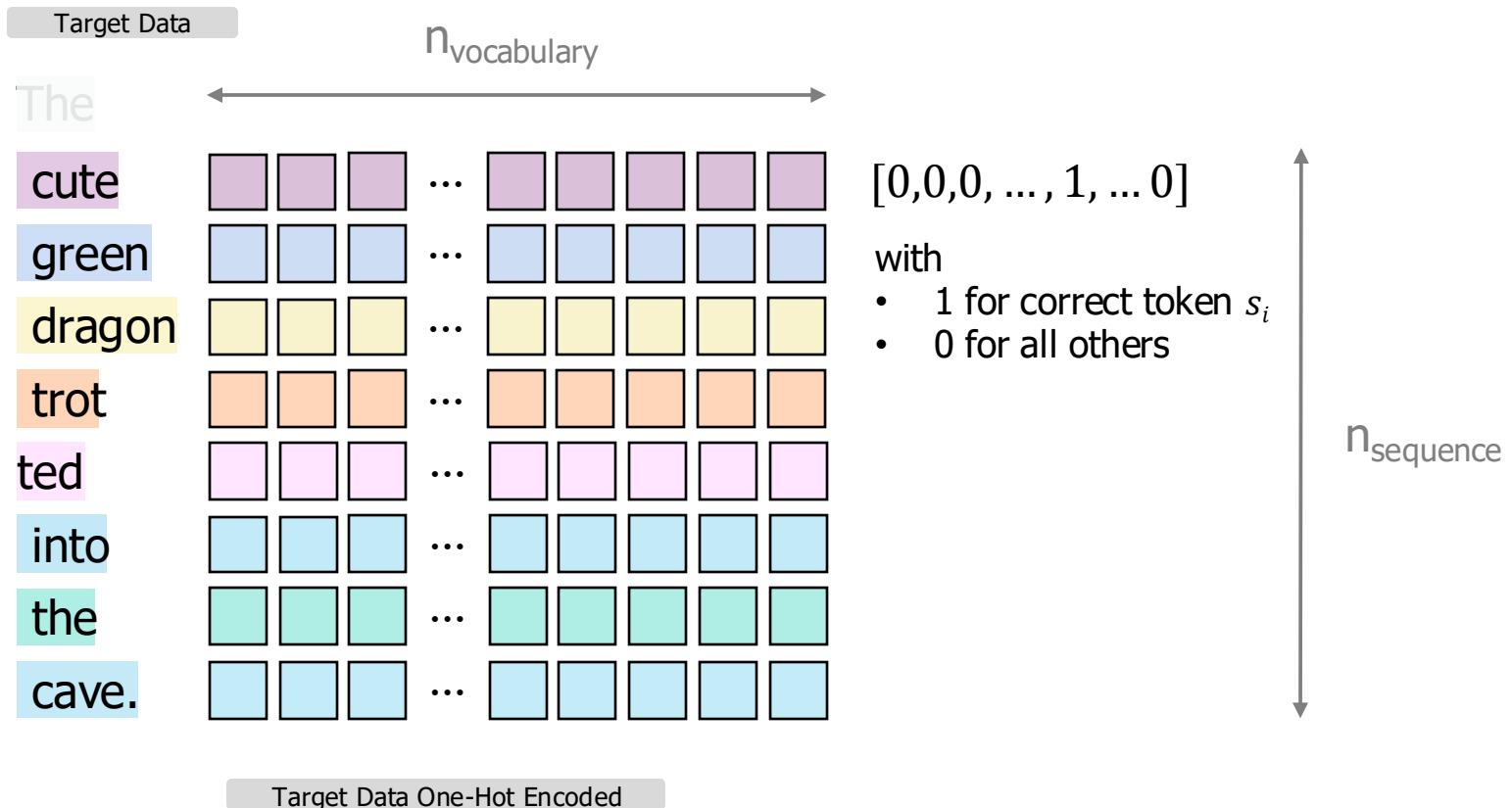
VS





$(s_2 | s_1)$   
 $(s_3 | s_1, s_2)$   
 $(s_4 | s_1, s_2, s_3)$   
 $(s_5 | s_1, s_2, \dots, s_4)$   
 $(s_6 | s_1, s_2, \dots, s_5)$   
 $(s_7 | s_1, s_2, \dots, s_6)$   
 $(s_8 | s_1, s_2, \dots, s_7)$   
 $(s_9 | s_1, s_2, \dots, s_8)$

$n_{sequence}$  vs



 The probability of a sequence is the product of the probability of the first token and the conditional probabilities of each subsequent token given all previous tokens.

### A few more definitions:

- $V_i$  ... the random variable for a token at position  $i$ , and  $s_i$  is the actual token observed at that position.
- $p(V_i | s_1, \dots, s_{i-1})$  ... the probability distribution over all possible tokens that will come next.
- $y_i$  ... a one-hot encoded vector of dimension equal to the vocabulary size holding 1 for the correct token at position  $i$  and 0 otherwise. This acts as the ground truth.

$$\text{🎯 } L = -\sum_i y_i \log(p(V_i | s_1, \dots, s_{i-1})) = -\sum_i \log(p(s_i | s_1, \dots, s_{i-1}))$$

 The loss function is defined as the negative log-likelihood of the correct tokens across all positions in the sequence, given the preceding tokens.

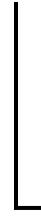


$$L = -\sum_i y_i \log(p(V_i|s_1, \dots, s_{i-1})) = -\sum_i \log(p(s_i|s_1, \dots, s_{i-1}))$$

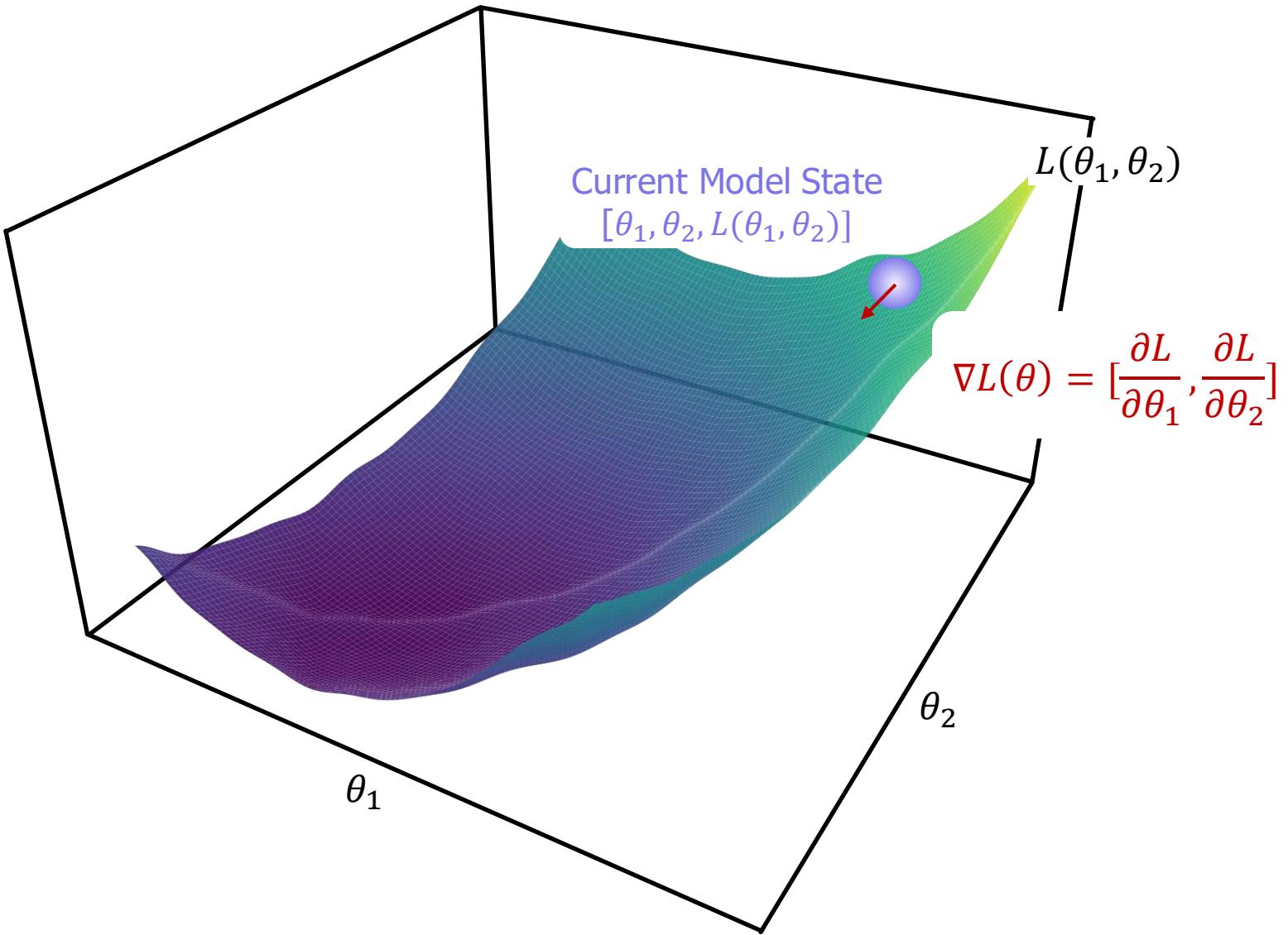


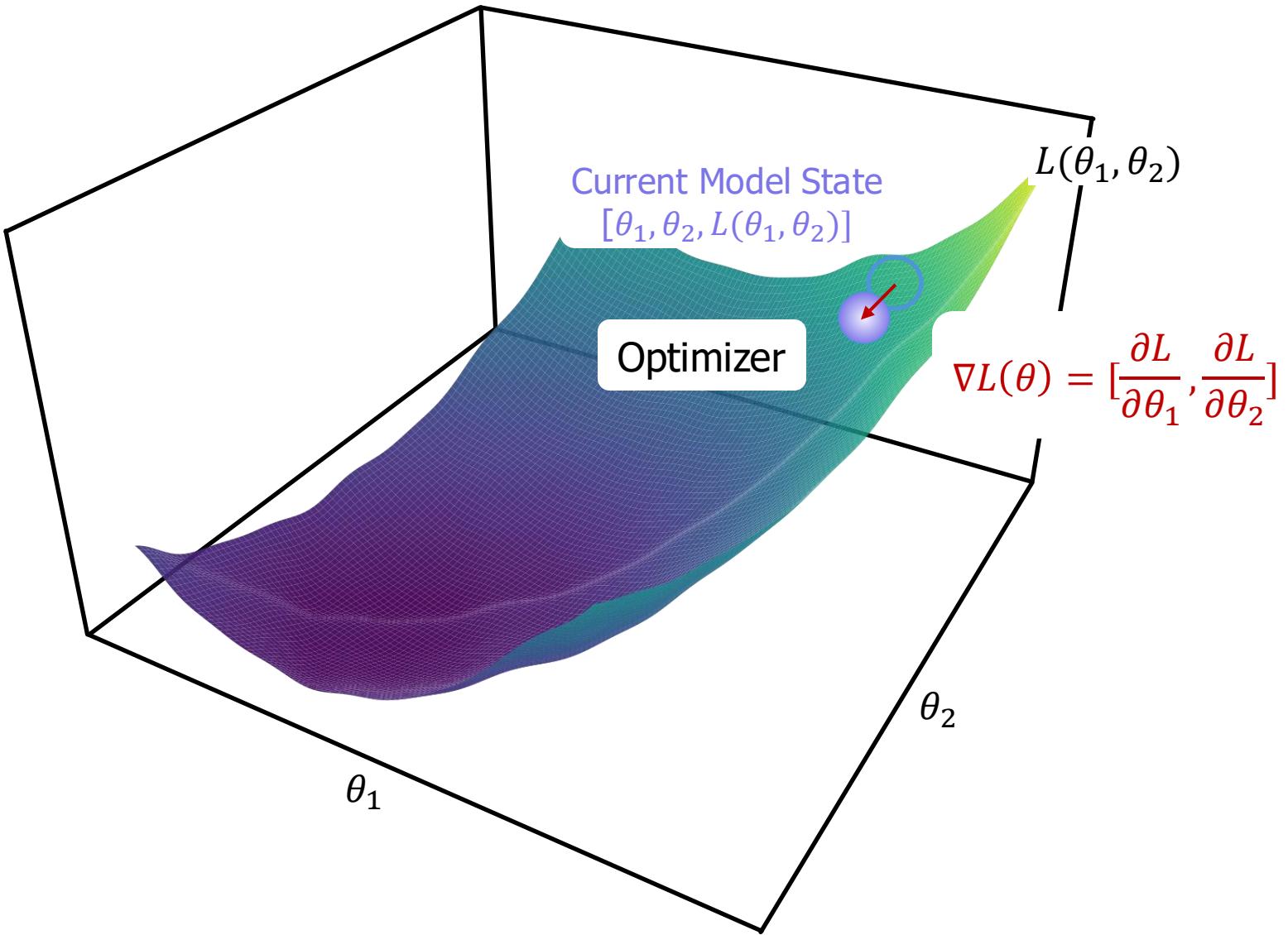
$$\nabla L(\theta) = \left[ \frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots, \frac{\partial L}{\partial \theta_n} \right]$$

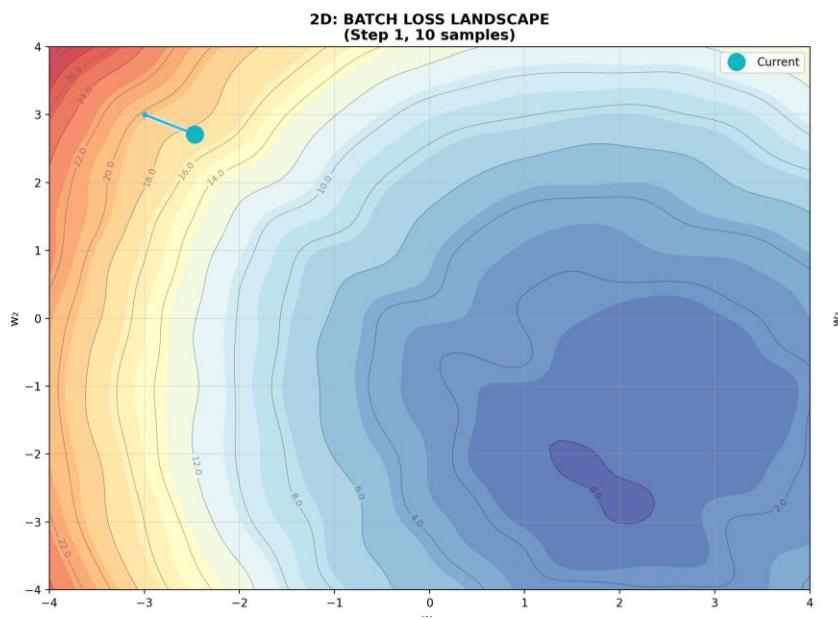
*"How does the loss change in respect to parameter  $\theta_1, \dots \theta_n$ "*

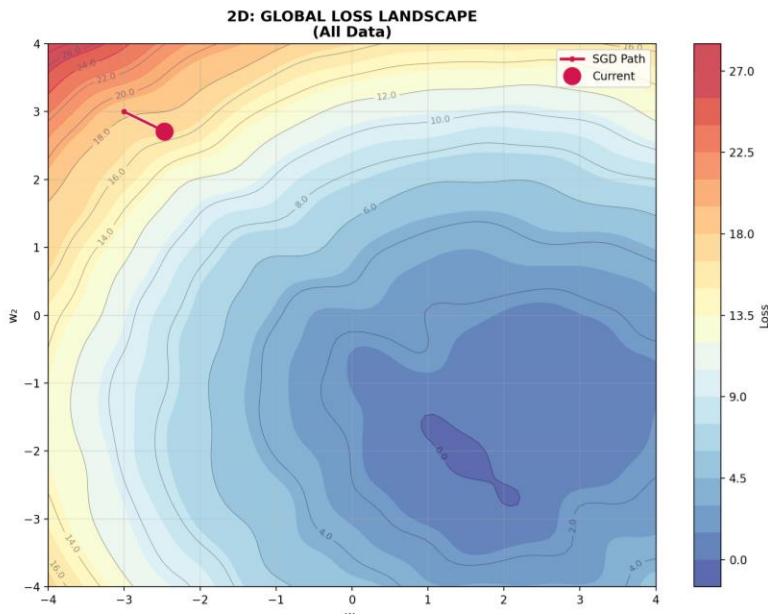
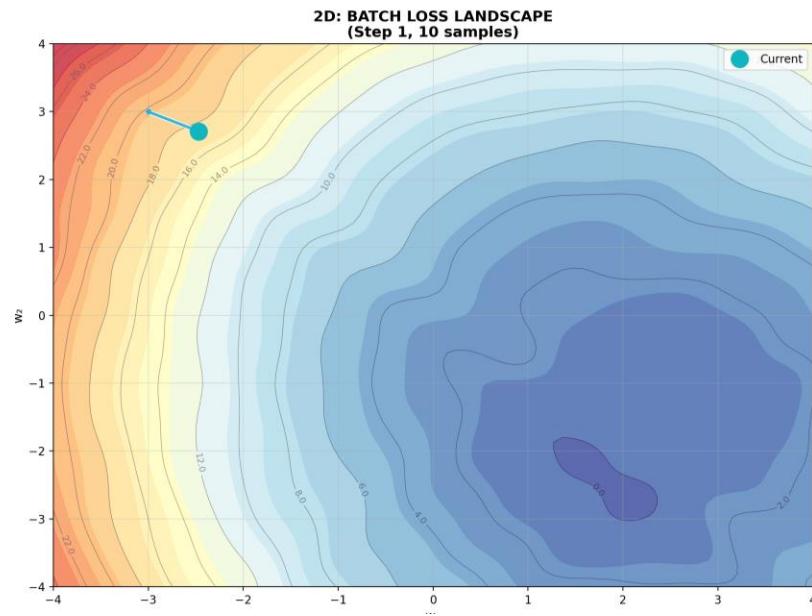


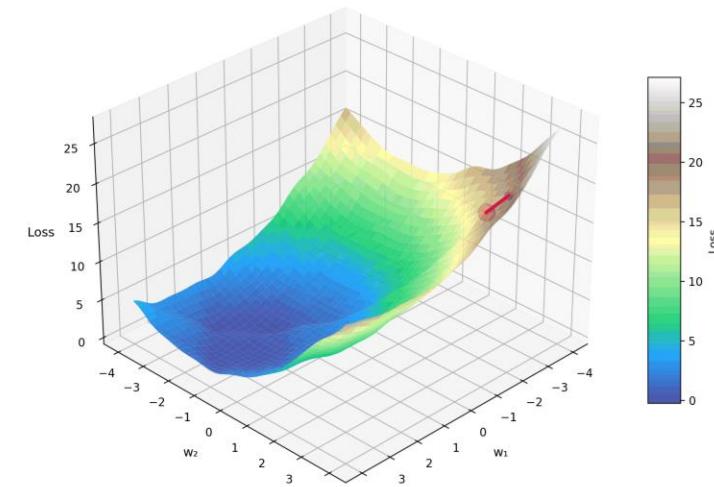
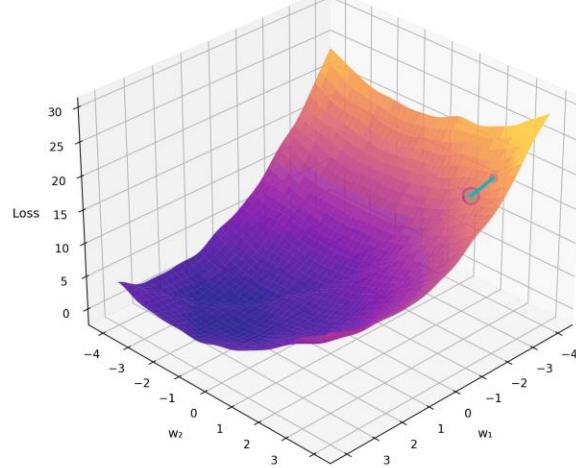
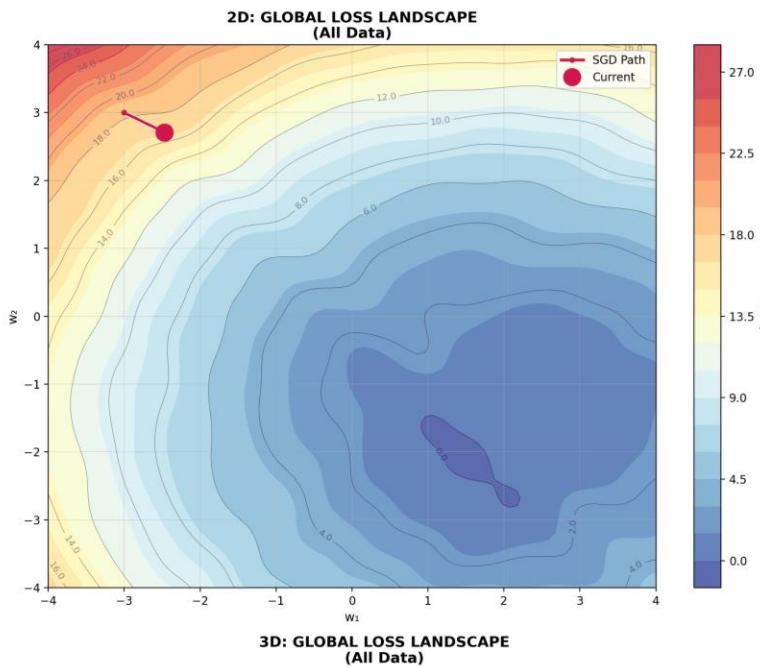
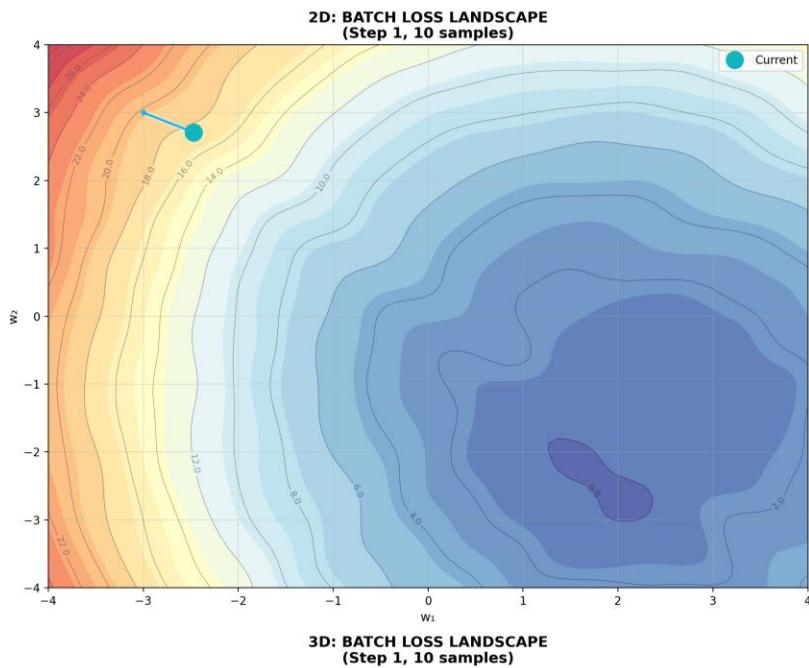
Optimizer (e.g. SGD) adapts weights based on gradient of loss









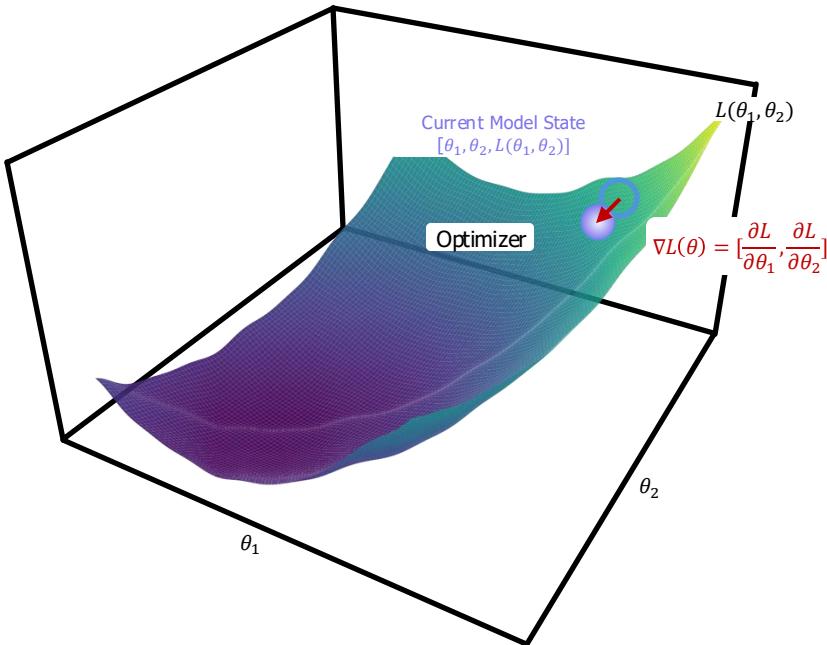




$$L = -\sum_i y_i \log(p(V_i|s_1, \dots, s_{i-1})) = -\sum_i \log(p(s_i|s_1, \dots, s_{i-1}))$$

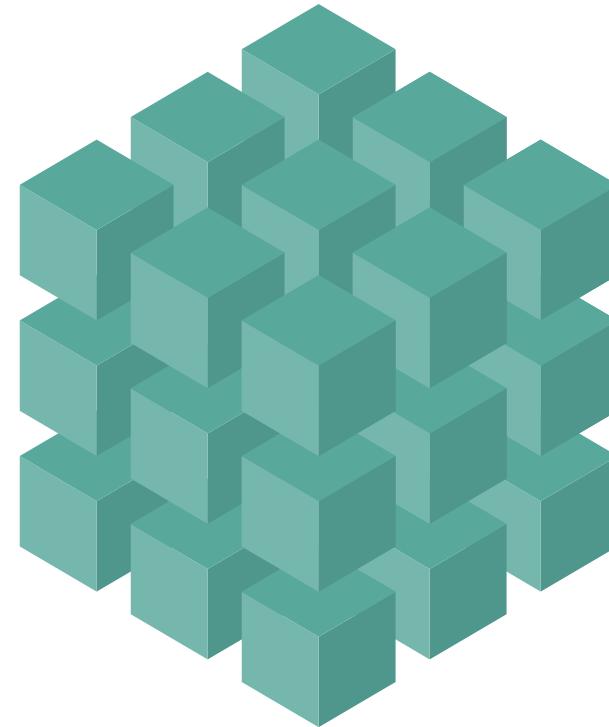
→  $\nabla L(\theta) = \left[ \frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots, \frac{\partial L}{\partial \theta_n} \right]$  "How does the loss change in respect to parameter  $\theta_1, \dots, \theta_n$ "

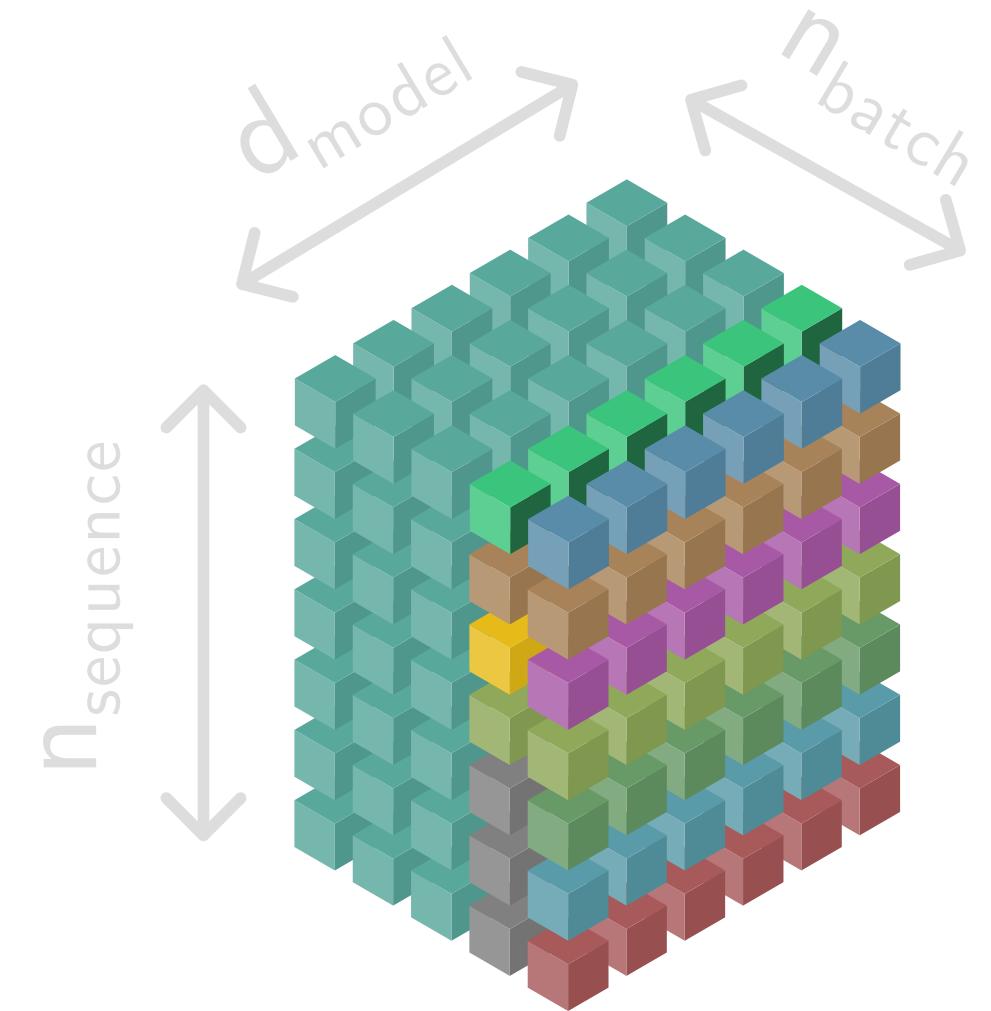
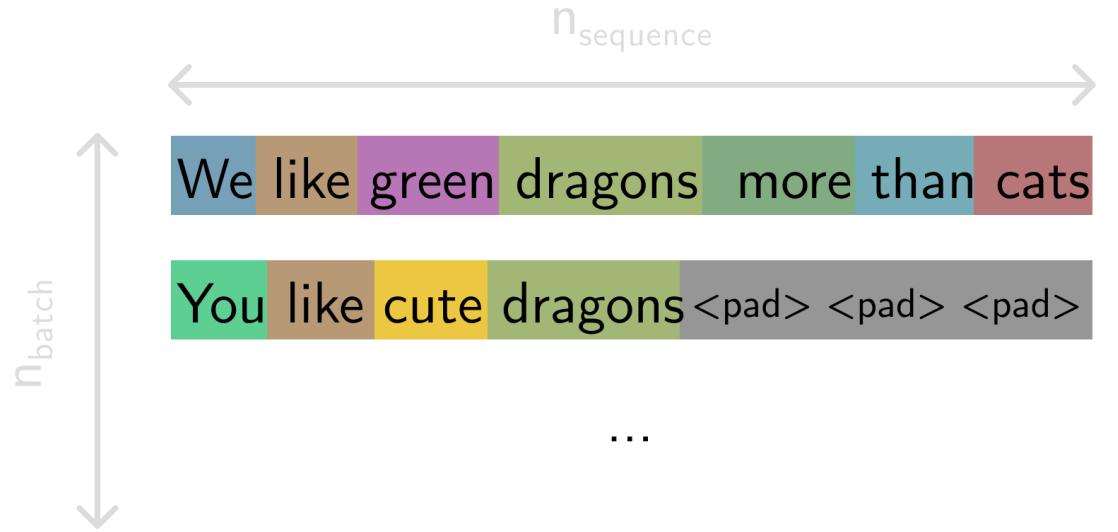
→ Optimizer (e.g. SGD, AdamW) adapts weights based on gradient of loss



## Side Note Batching

Simultaneously process multiple training samples

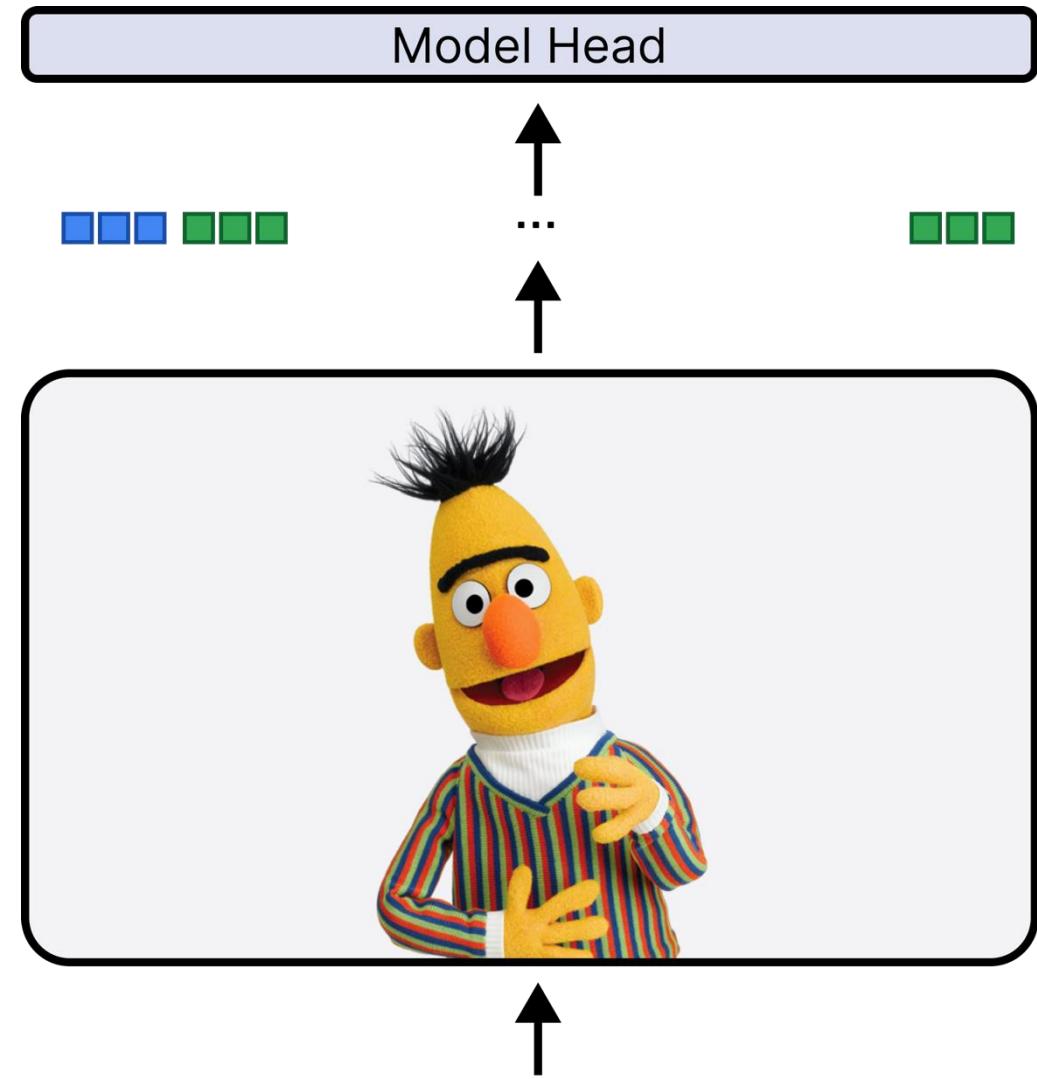




# Training Objective For Encoder Only Models

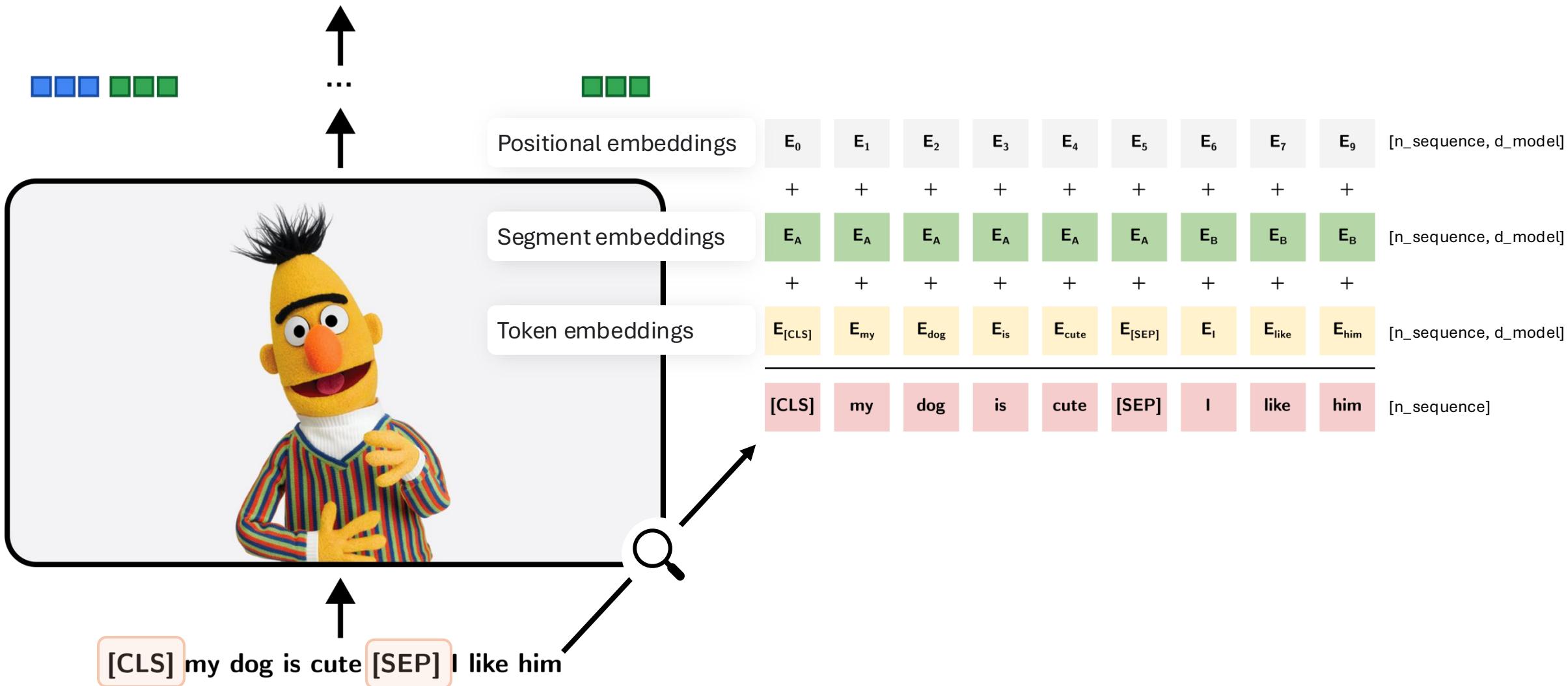
In reference to BERT - Pre-training of Deep Bidirectional Transformers for Language Understanding<sup>1</sup>

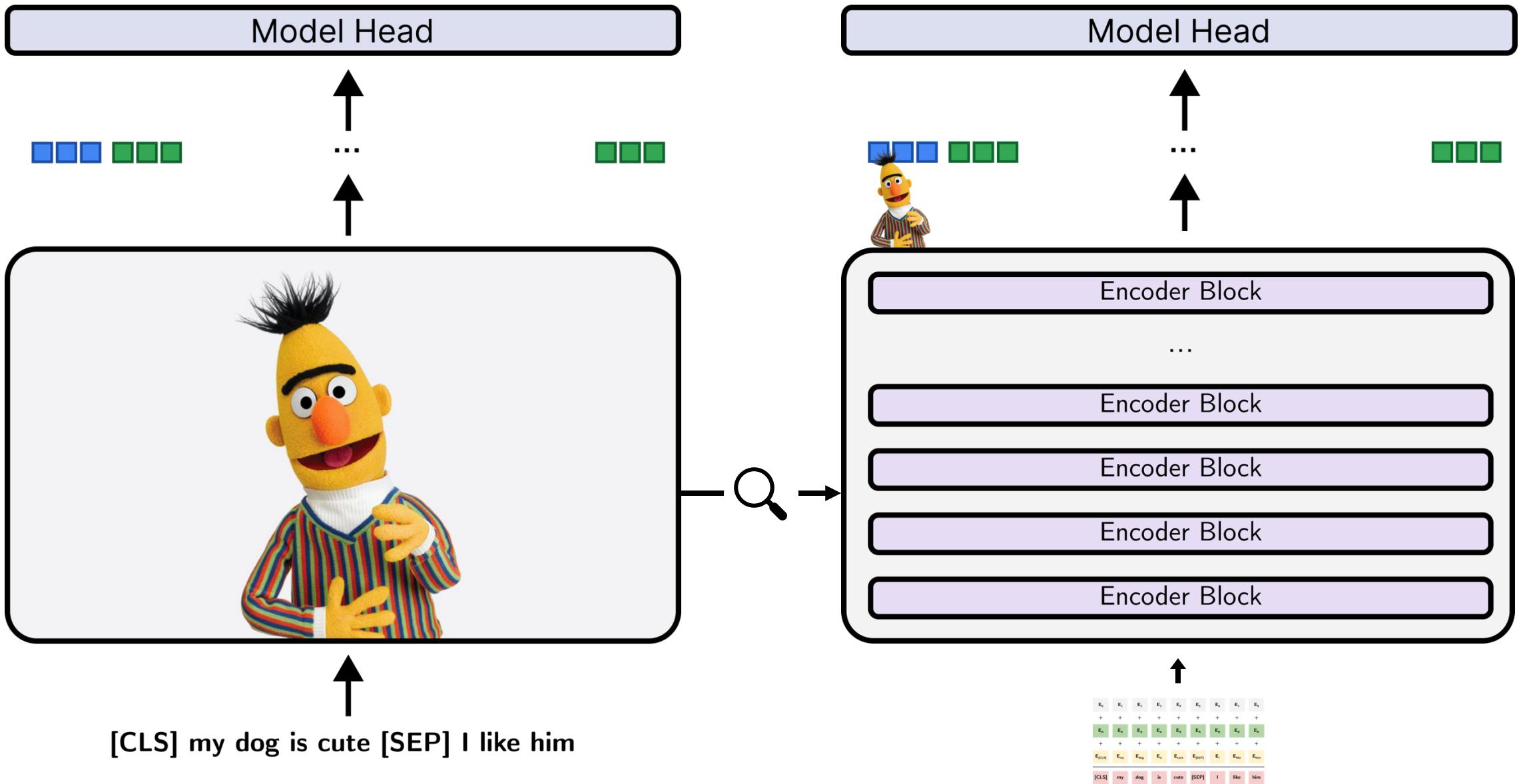
<sup>1</sup><https://arxiv.org/abs/1810.04805>



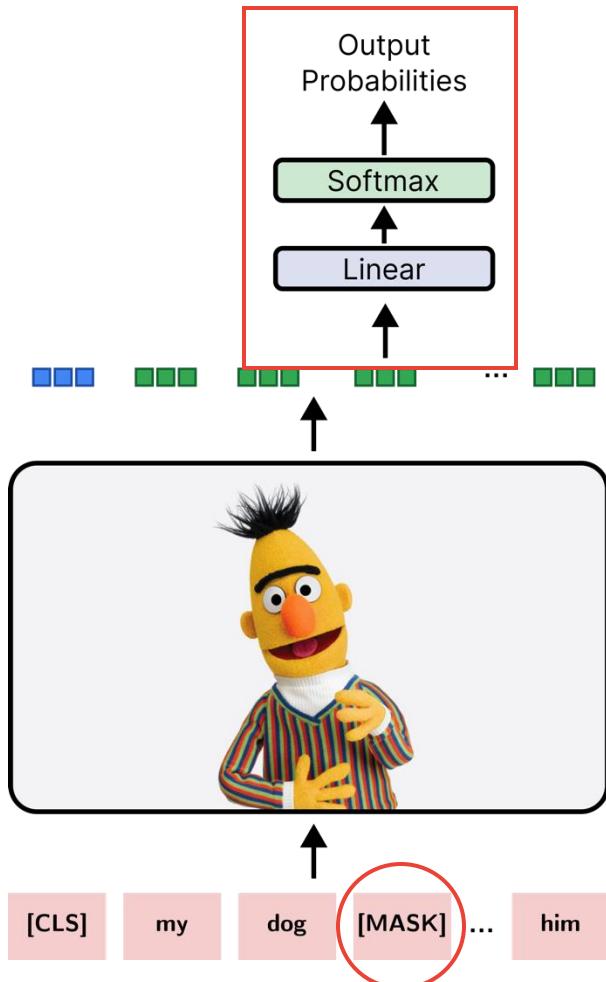
[CLS] my dog is cute [SEP] I like him

## Model Head



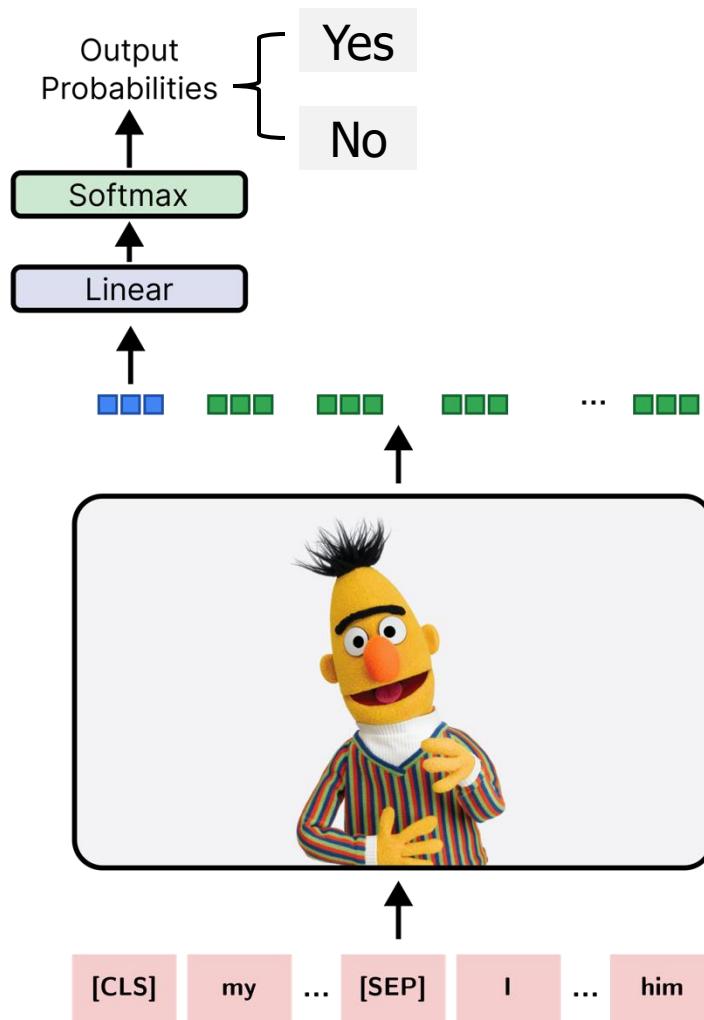


## Masked Token Prediction

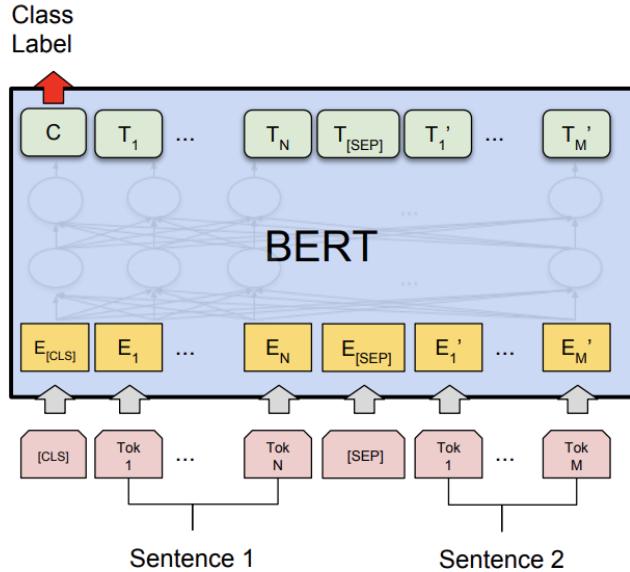


- › 80% of the time: Replace the word with the **[MASK]** token to learn how words relate to their surrounding context to make accurate predictions.
- › 10% of the time: Replace the token with a random token to prevent the model from simply learning that it should pay special attention to [MASK] tokens.
- › 10% of the time: Keep the word unchanged representation towards the actual observed word.

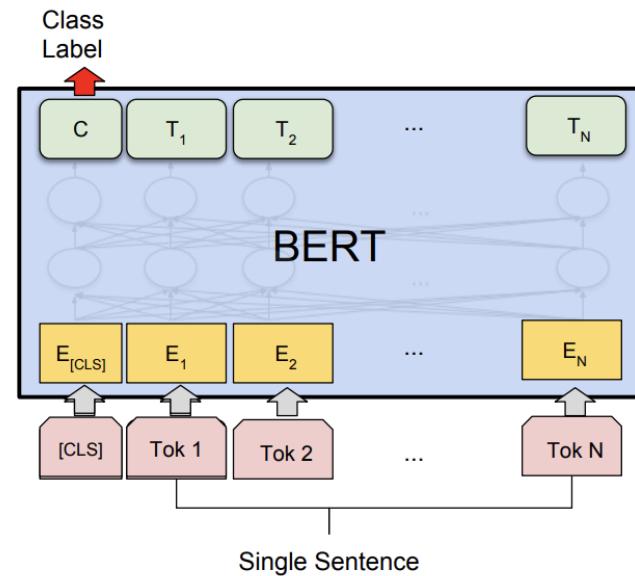
## Next Sentence Prediction



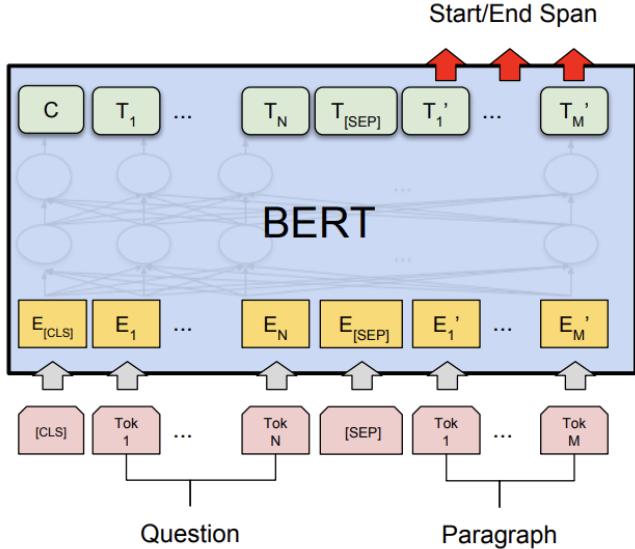
- › 50% of the time actual next sentence is used 50% of the time it is a random sentence from the corpus
- › [CLS] related output is used for binary classification



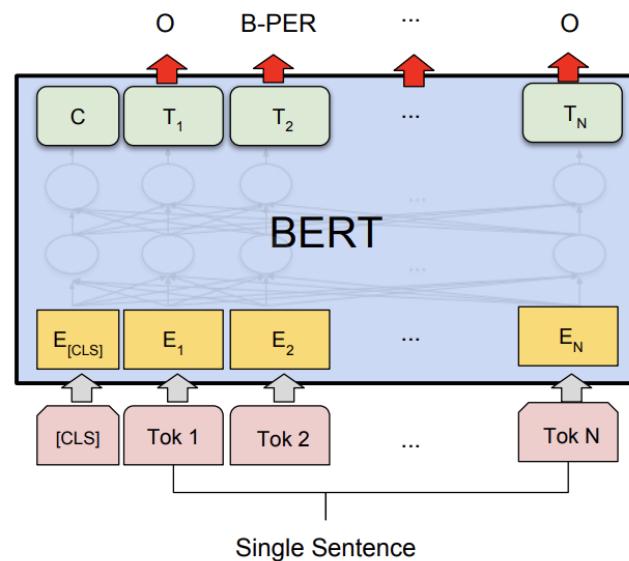
Sentence Pair  
Classification



Single Sentence  
Classification Tasks



Question Answering  
Tasks

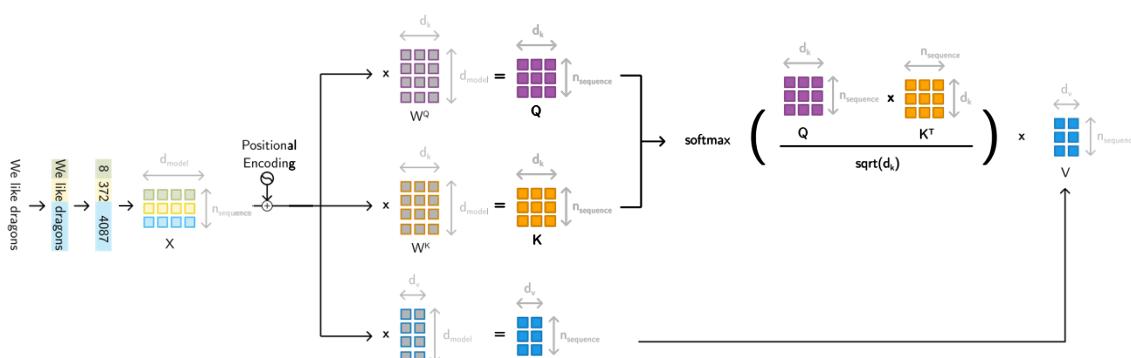


Single Token  
Tagging Tasks

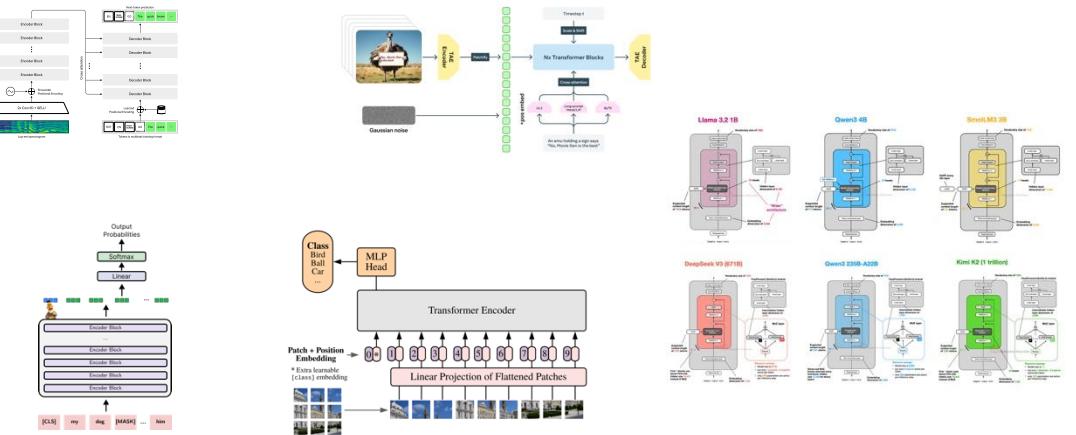
# What is so special about The Transformer Architecture?

1<https://arxiv.org/abs/1810.04805>

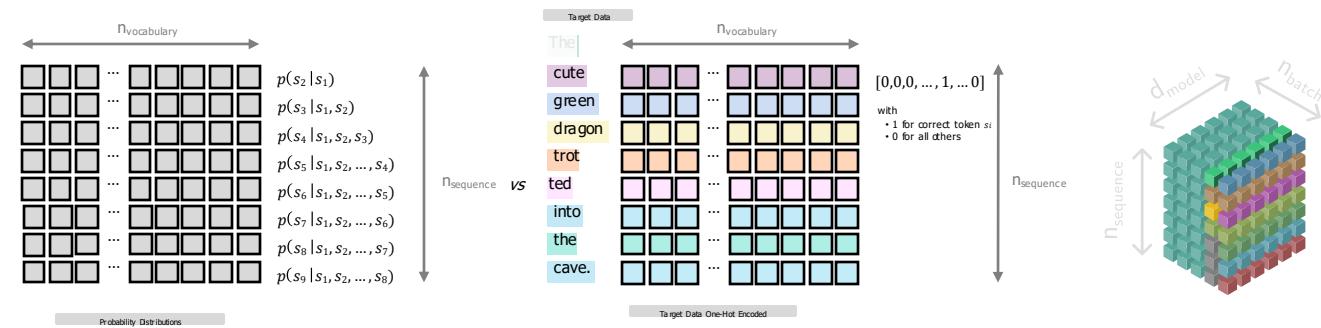
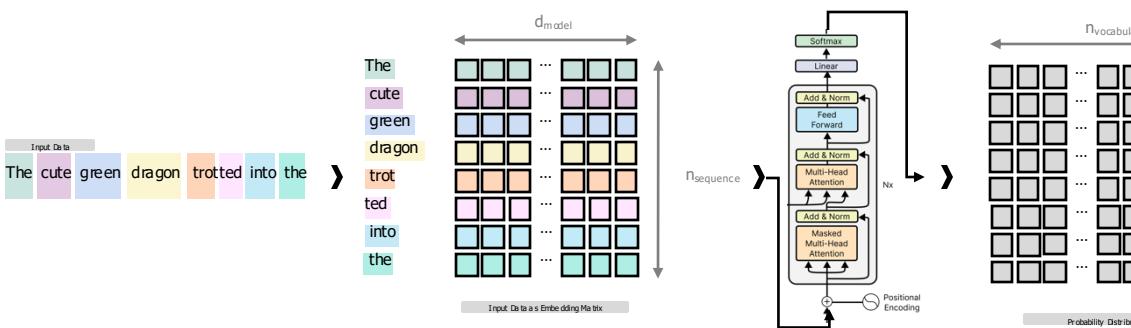
## 1. Long-range dependencies in sequences



## 3. Versatility



## 2. Fully parallelizable training paradigm



## Further Resource



Allen Institute for AI

Pre-training: OLMo, OLMo 2, OLMoE

Post-training: Tülu 3

<https://allenai.org/language-models>



HuggingFace – LLM Course

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



Annotated Transformer Implementation

<https://nlp.seas.harvard.edu/annotated-transformer/>



Jay Alammar – The Illustrated Transformer

<https://jalammar.github.io/illustrated-transformer/>



Sebastian Raschka – Build an LLM from Scratch

<https://www.youtube.com/@SebastianRaschka>