



A closer look into the transformer architecture

April 28, 2025

Did you complete the reading assignments?

Estimating (GPU) RAM usage of a LLM

- I have a 7 billion parameter model. How much RAM do I need to host it?
- RAM Usage = 7 billion params x 2 bytes per parameter = 14 billion bytes*
 - **Kilo** bytes = $2^{10} = 1,024$ \approx 1 thousand bytes
 - **Mega** bytes = $2^{20} = 1,048,576$ \approx 1 million bytes
 - **Giga** bytes = $2^{30} = 1,073,741,824$ \approx 1 billion bytes
 - **Tera** bytes = $2^{40} = 1,099,511,627,776$ \approx 1 trillion bytes
- Unlike regular ML models where a parameter is 4 bytes (32 bit float), LLMs are usually 2 bytes

* Actually we need a bit more for the intermediate computations

Recap: How a transformer works

She heals patients daily.



Large Language Model



Describe a typical doctor.

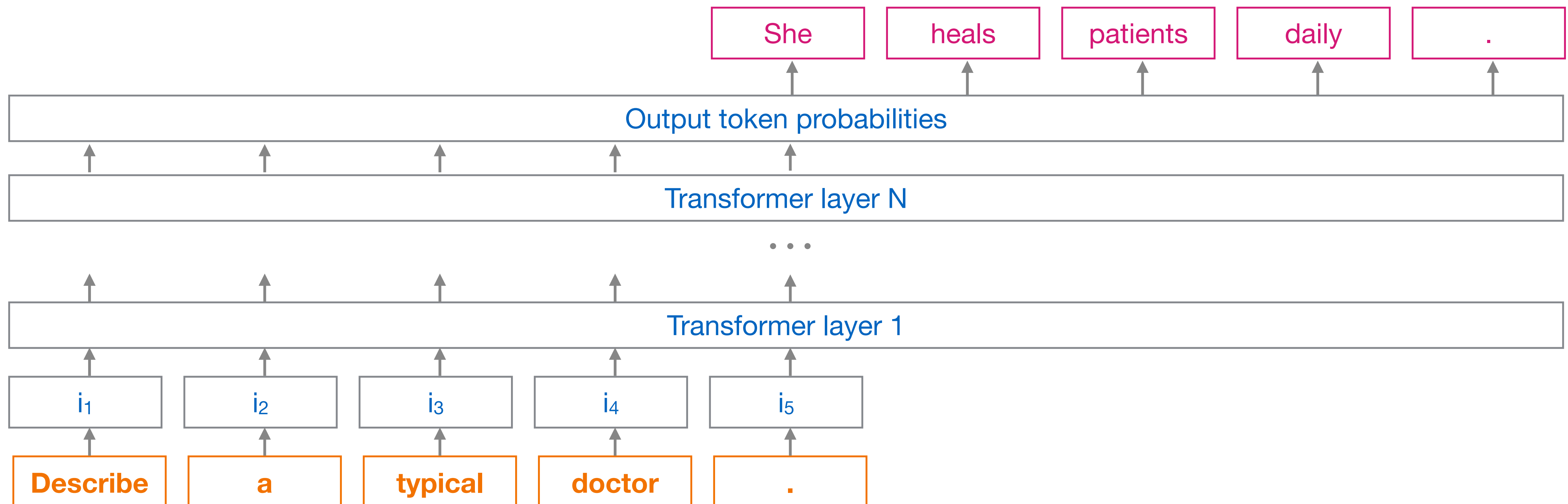
Recap: LLMs and causal language modeling



Input
aka prompt



Output we
usually observe



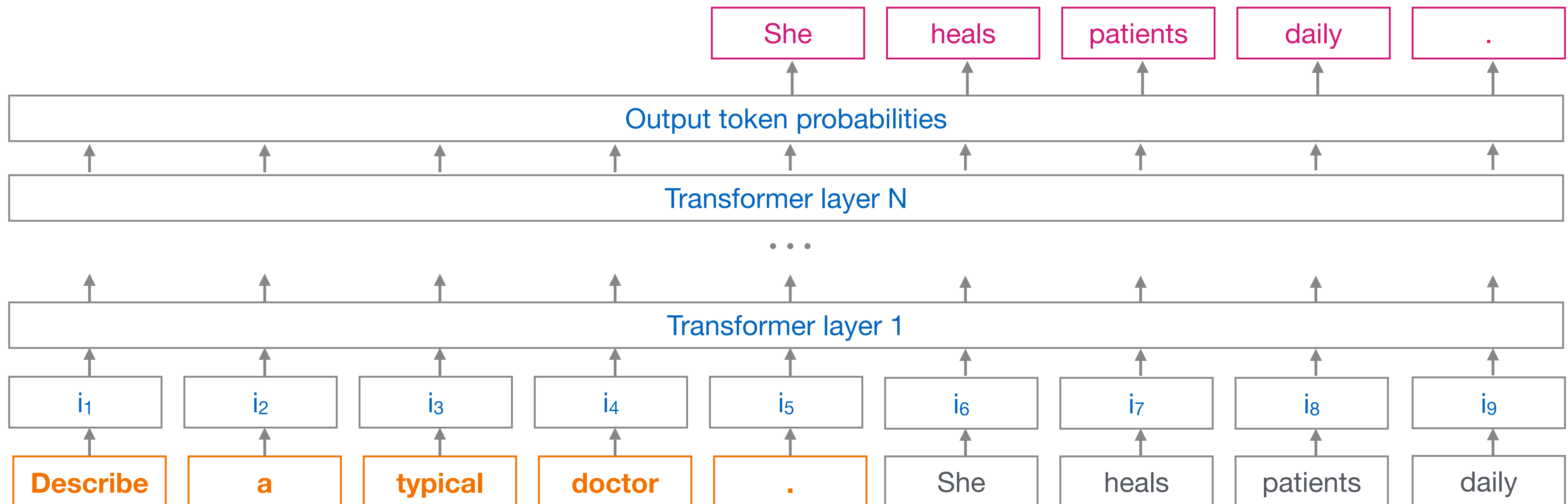
Recap: LLMs and causal language modeling



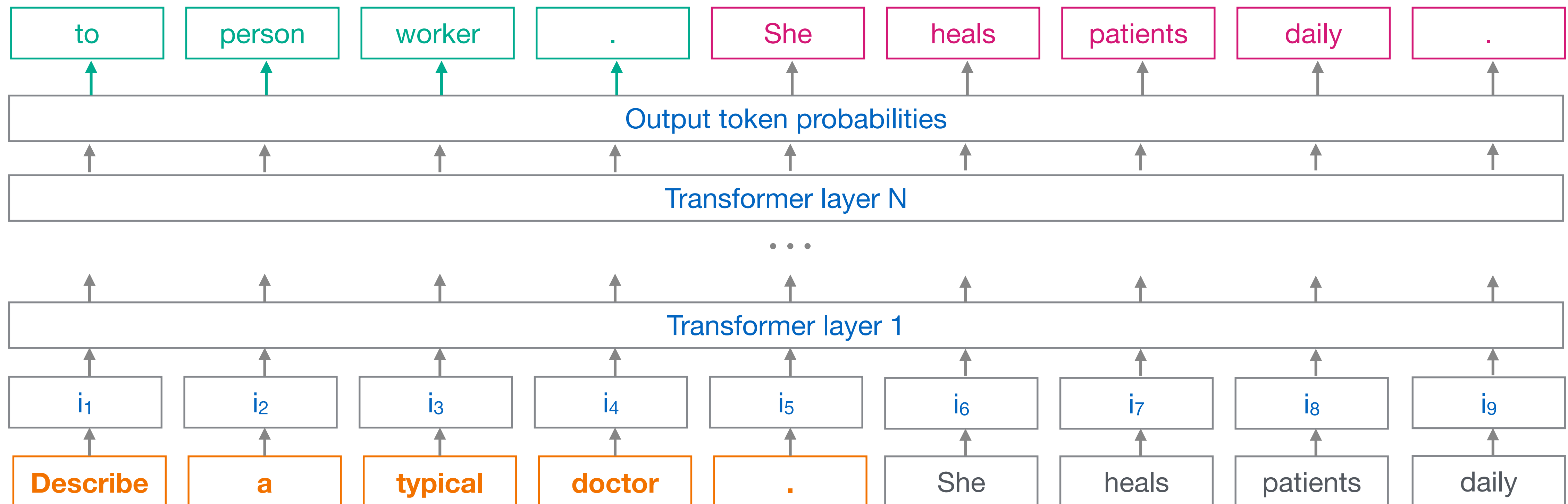
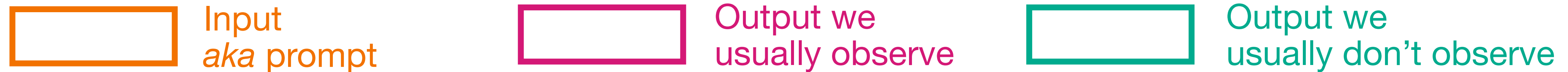
Input
aka prompt



Output we
usually observe

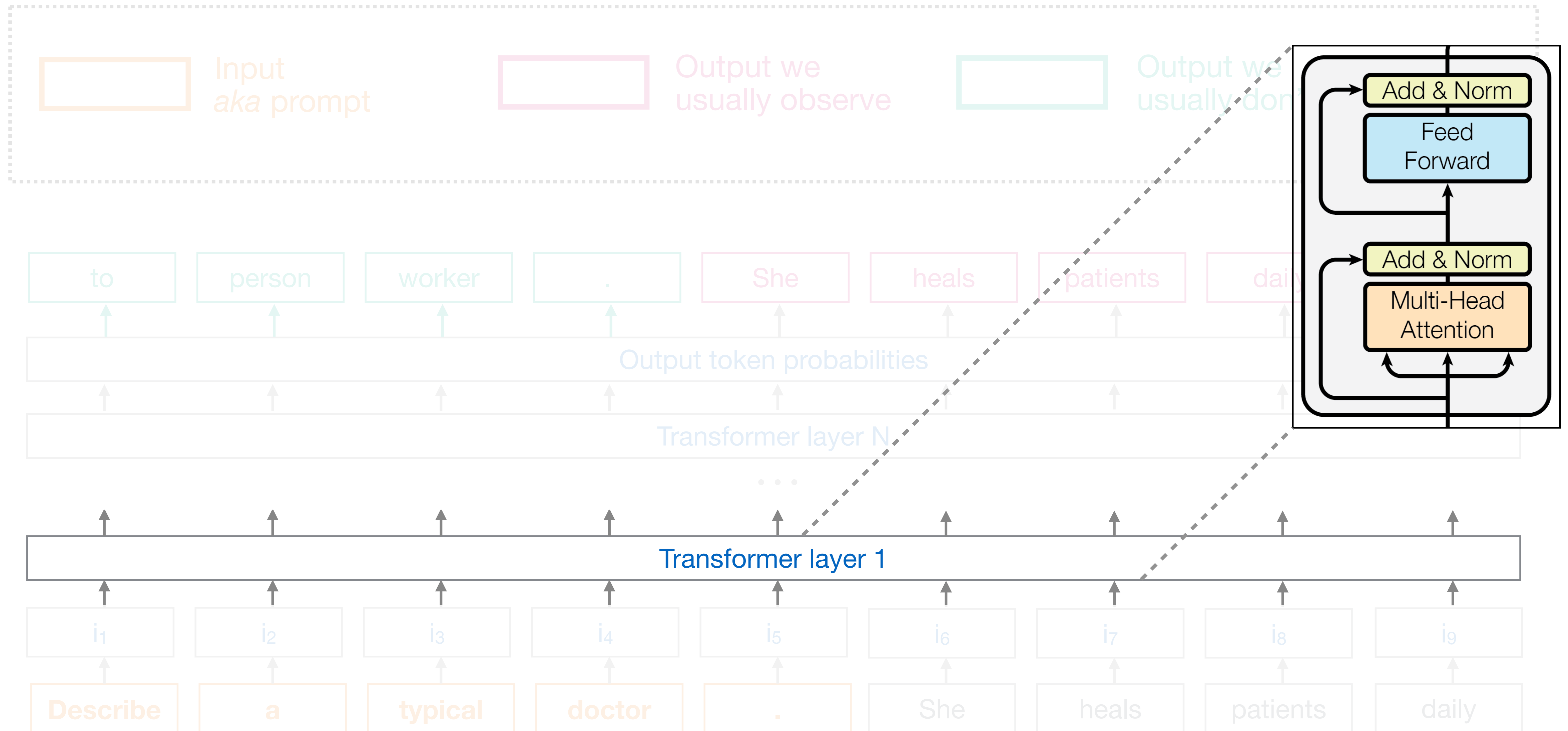


Recap: LLMs and causal language modeling

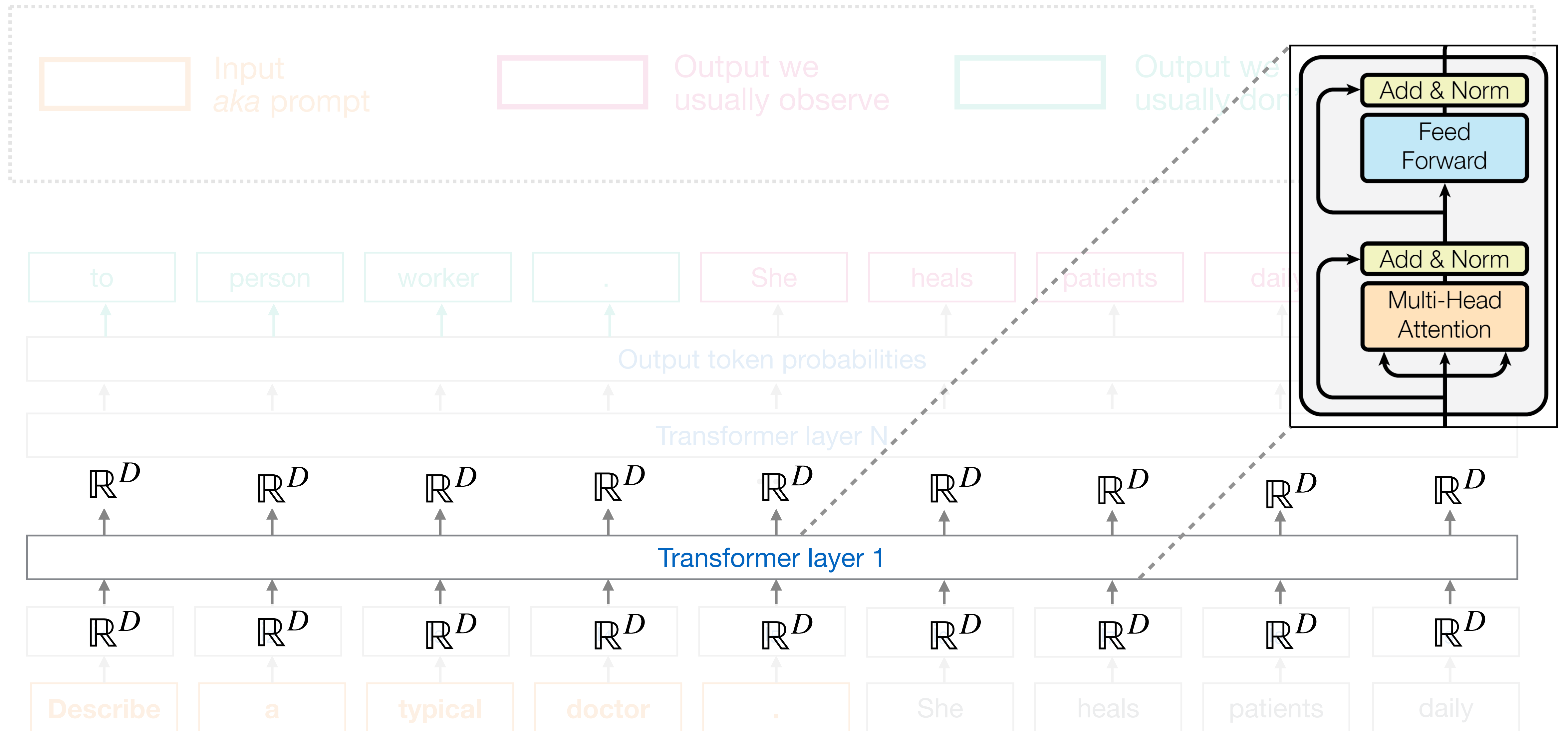


Exercise

Recap: LLMs and causal language modeling

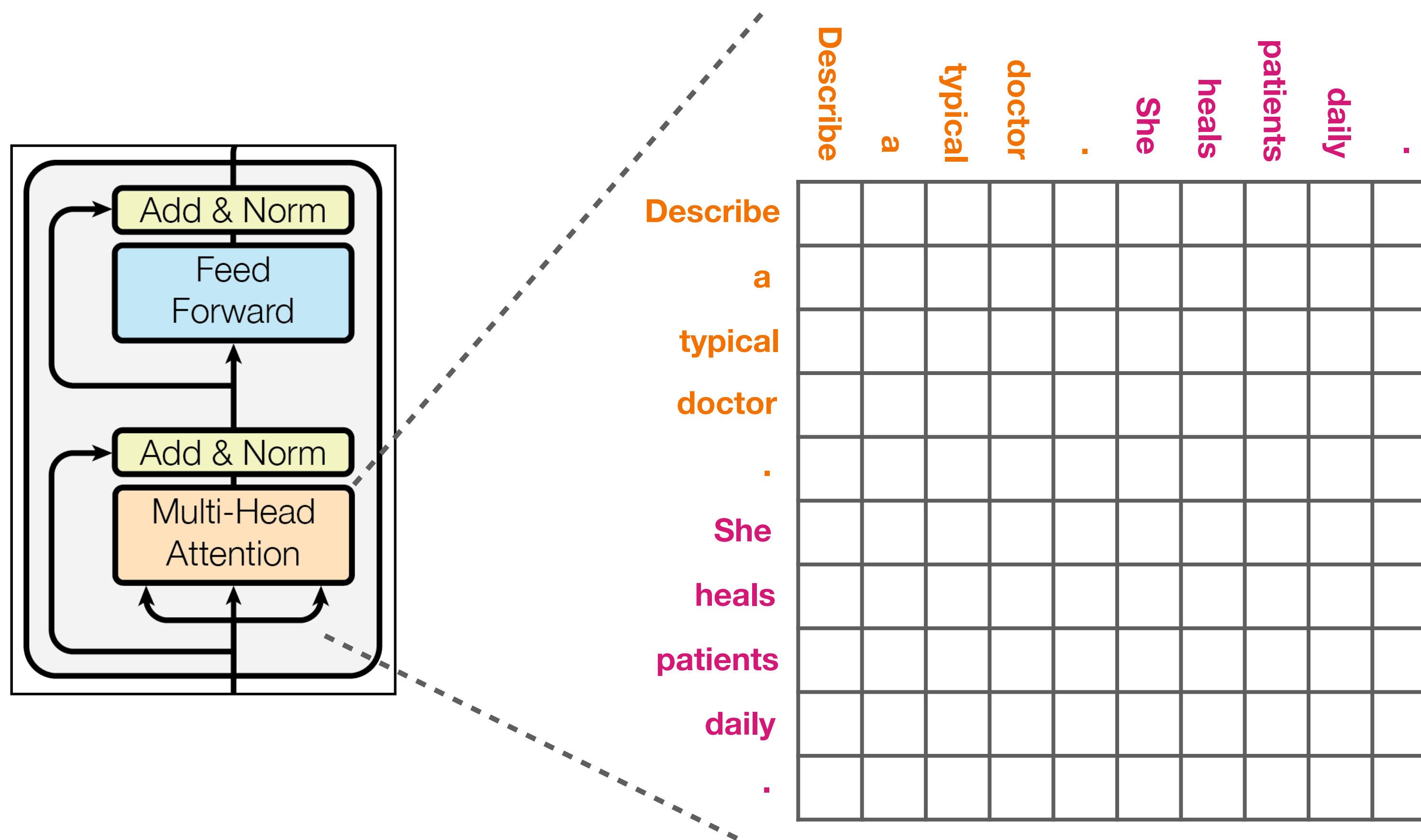


Recap: LLMs and causal language modeling



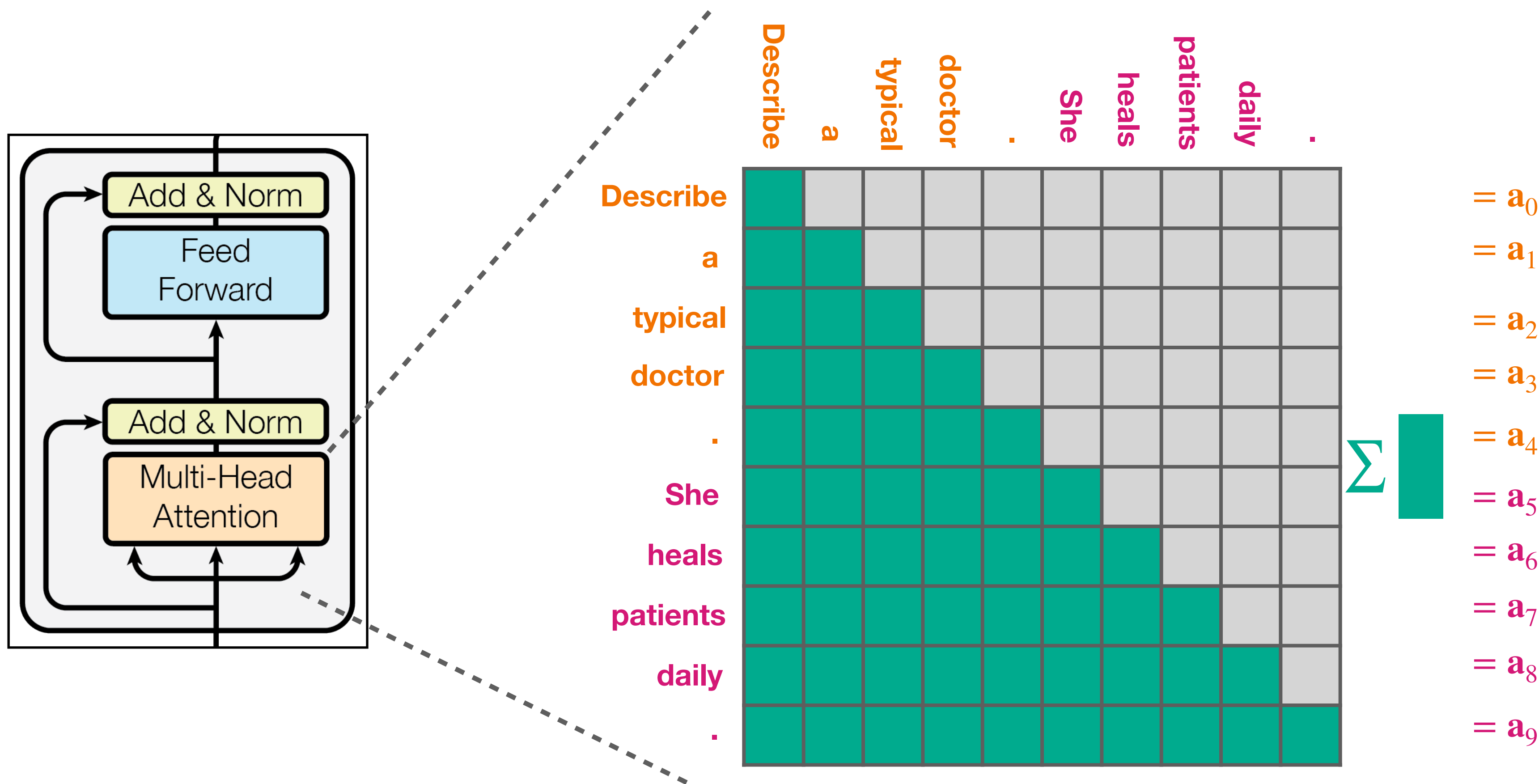
Self attention between all token pairs

When processing token i , how much each token j should contribute?



Self attention between all token pairs

When processing token i , how much each token j should contribute?



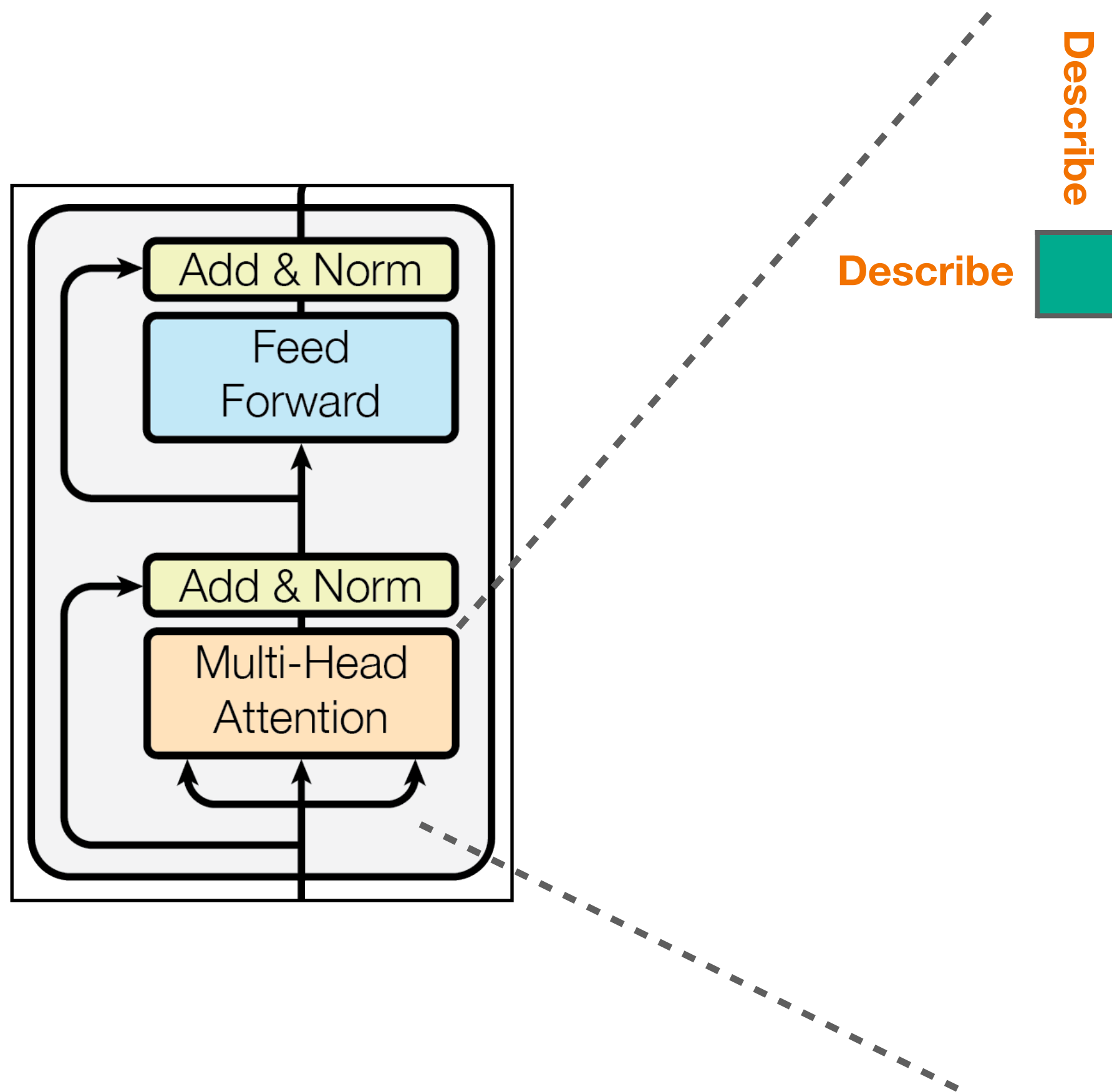
Causal LLMs

Attention of token i depends only on current or previous tokens

$$\mathbf{a}_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing token i , how much each token j should contribute?



Σ

$= \mathbf{a}_0$

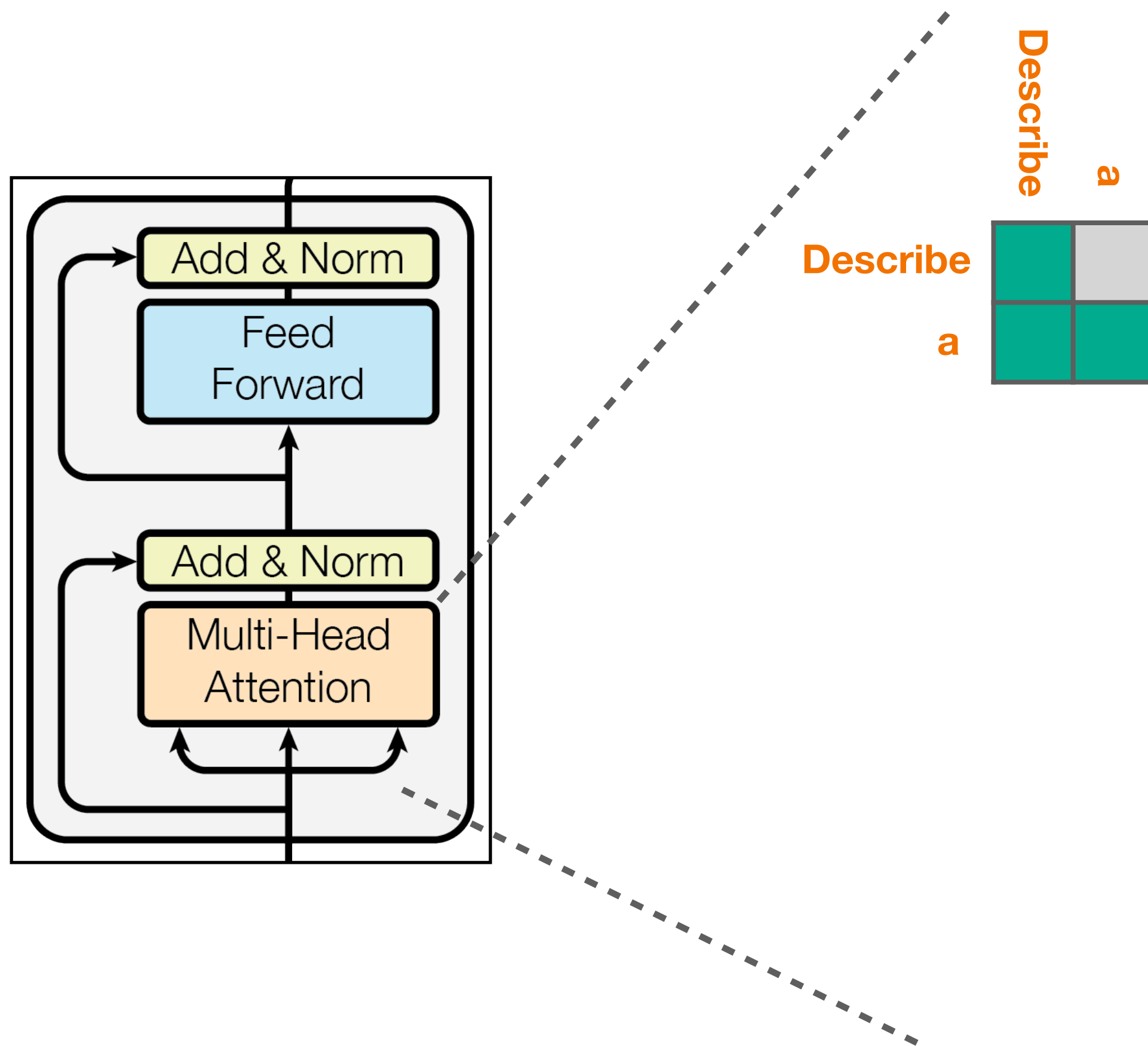
Causal LLMs

Attention of token i depends only on current or previous tokens

$$\mathbf{a}_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing token i , how much each token j should contribute?



Σ

$= a_0$
 $= a_1$

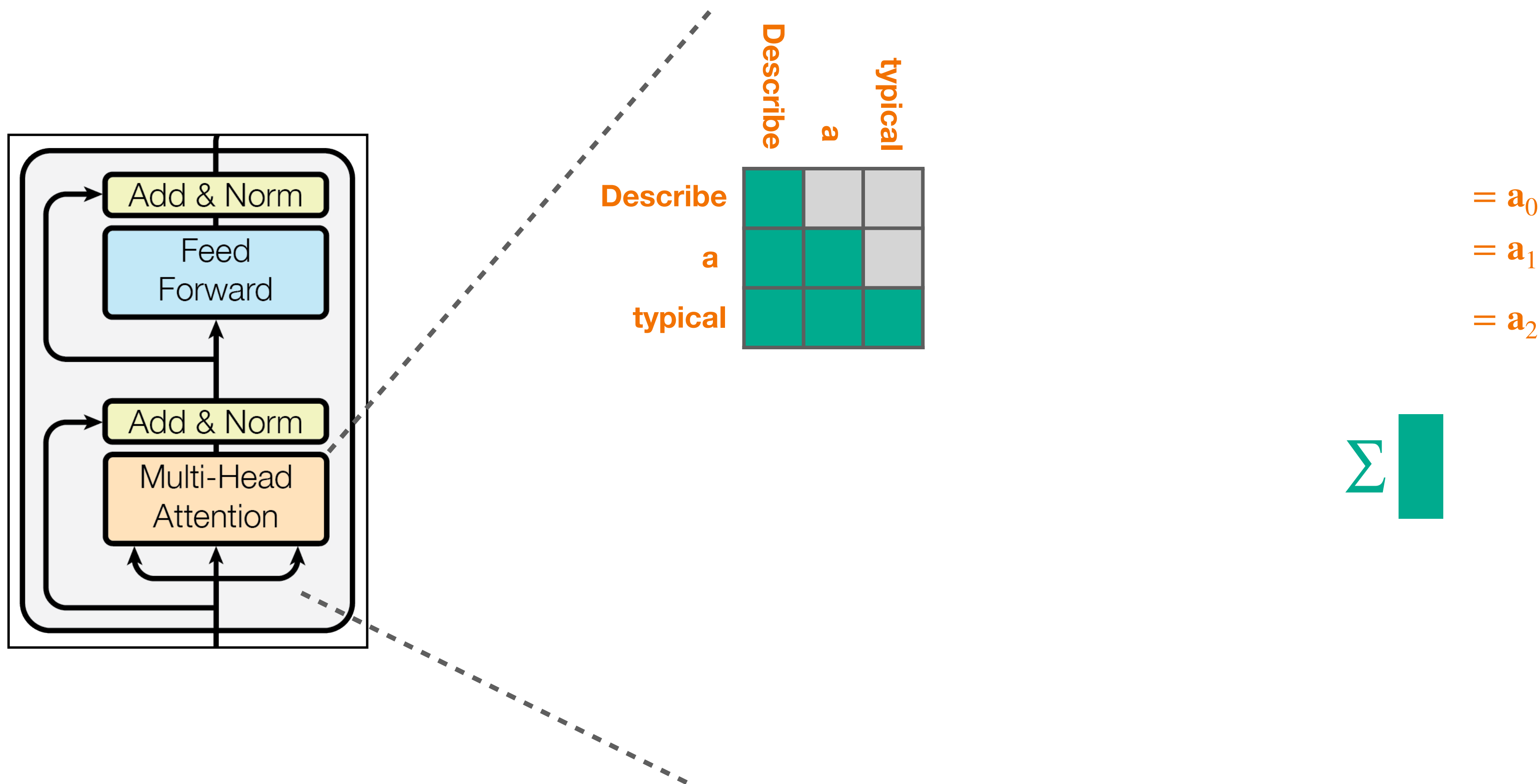
Causal LLMs

Attention of token i depends only on current or previous tokens

$$a_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing token i , how much each token j should contribute?



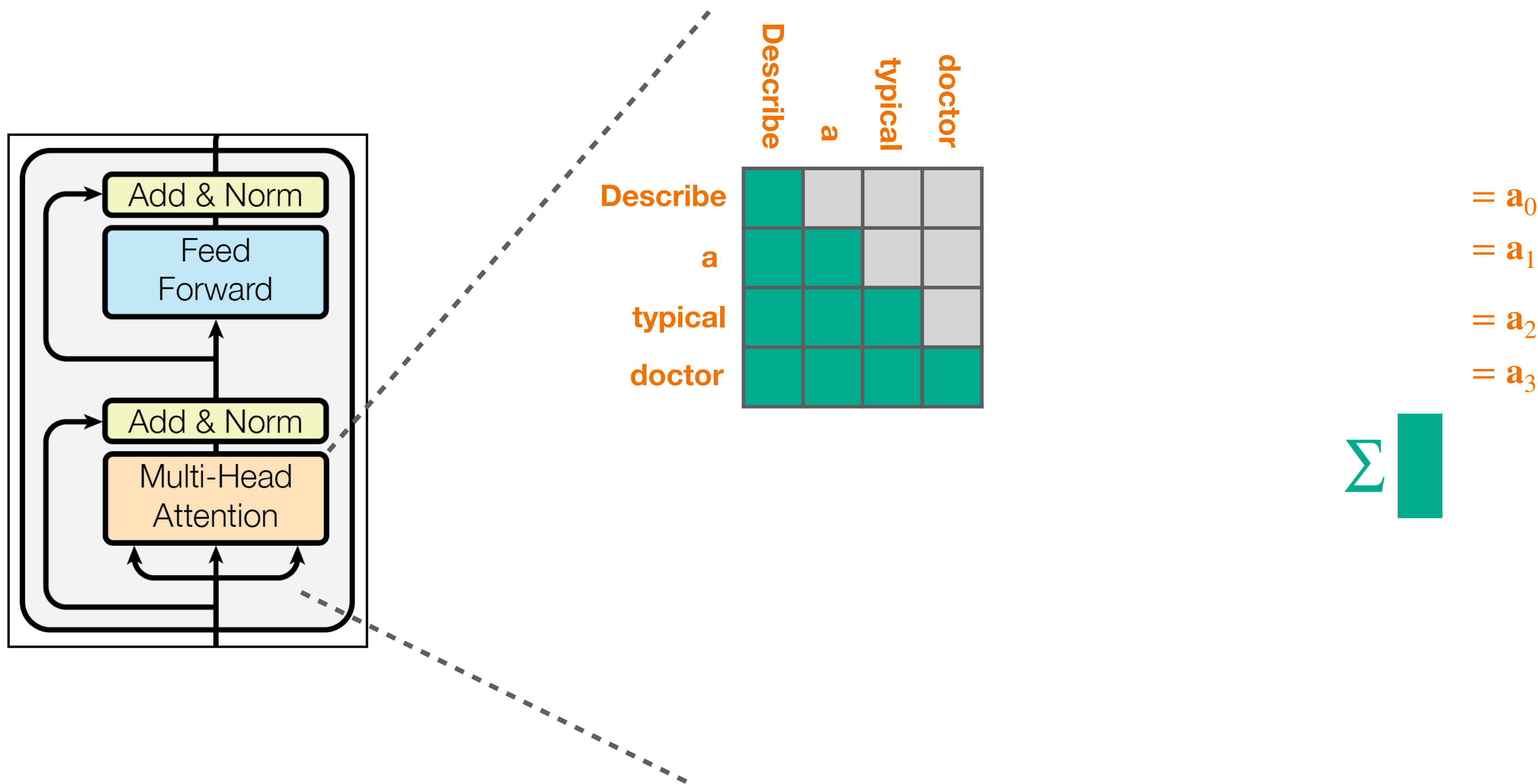
Causal LLMs

Attention of token i depends only on current or previous tokens

$$a_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing token i , how much each token j should contribute?



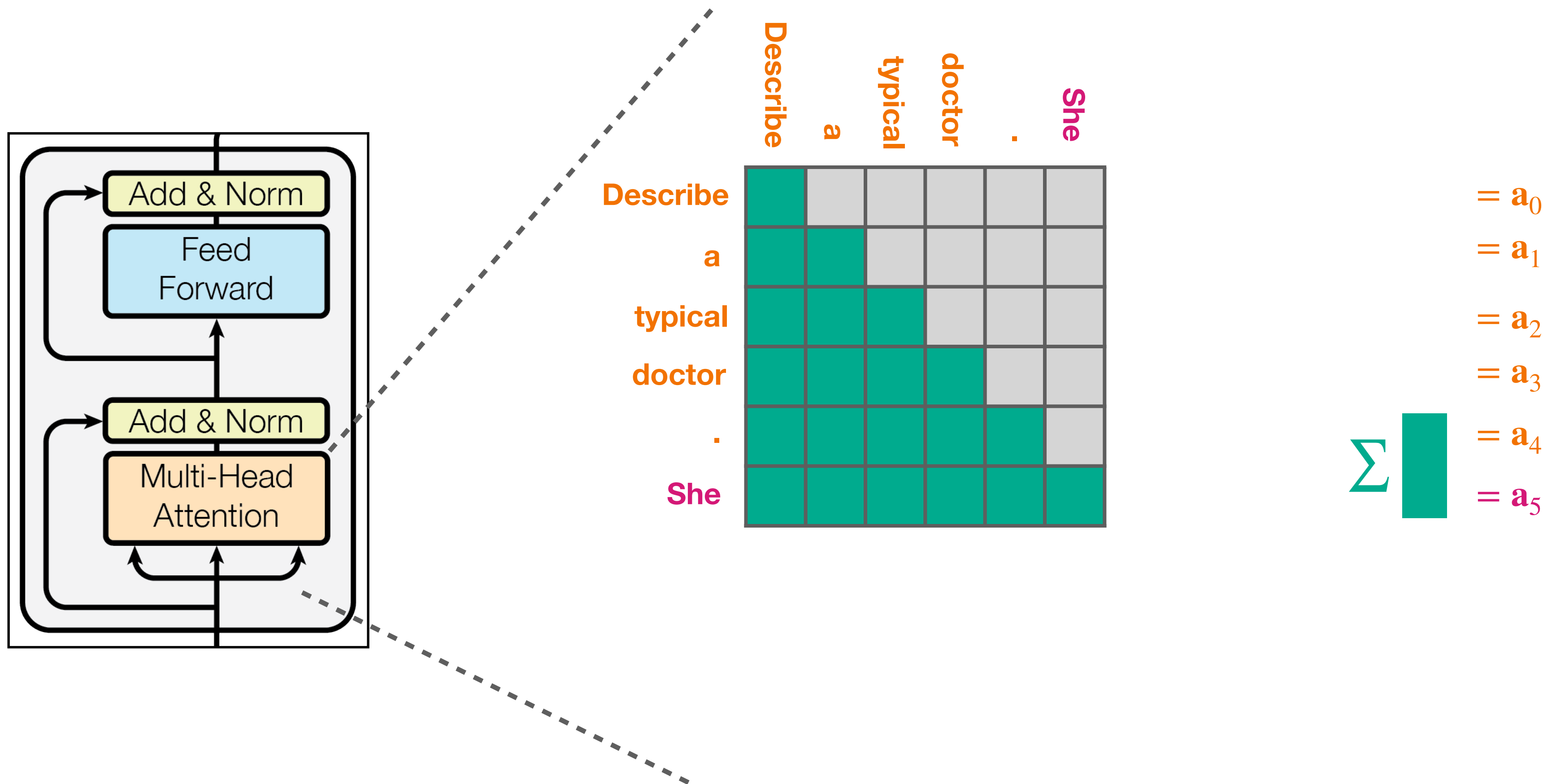
Causal LLMs

Attention of token i depends only on current or previous tokens

$$\mathbf{a}_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing token i , how much each token j should contribute?



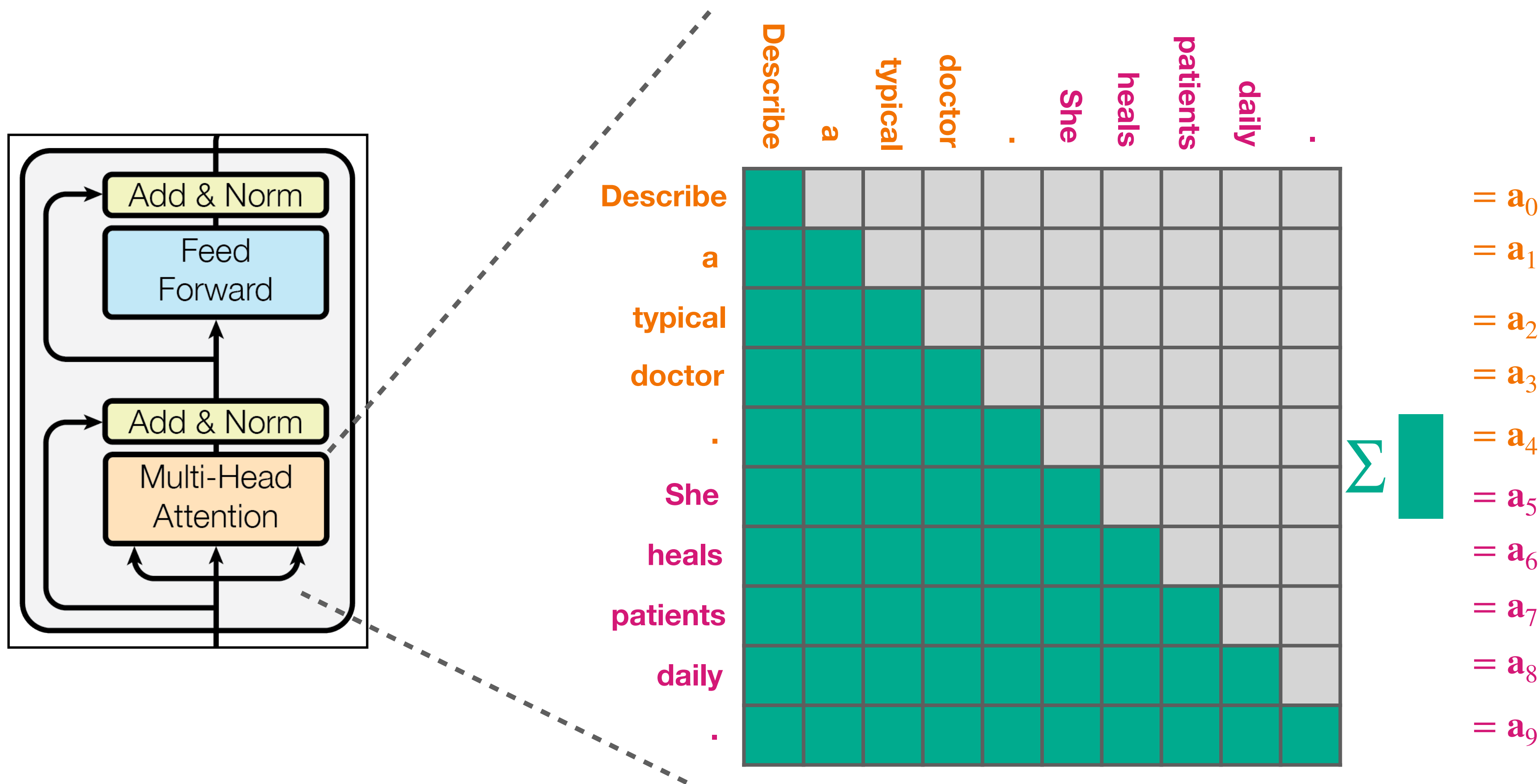
Causal LLMs

Attention of token i depends only on current or previous tokens

$$\mathbf{a}_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing *token i*, how much each *token j* should contribute?



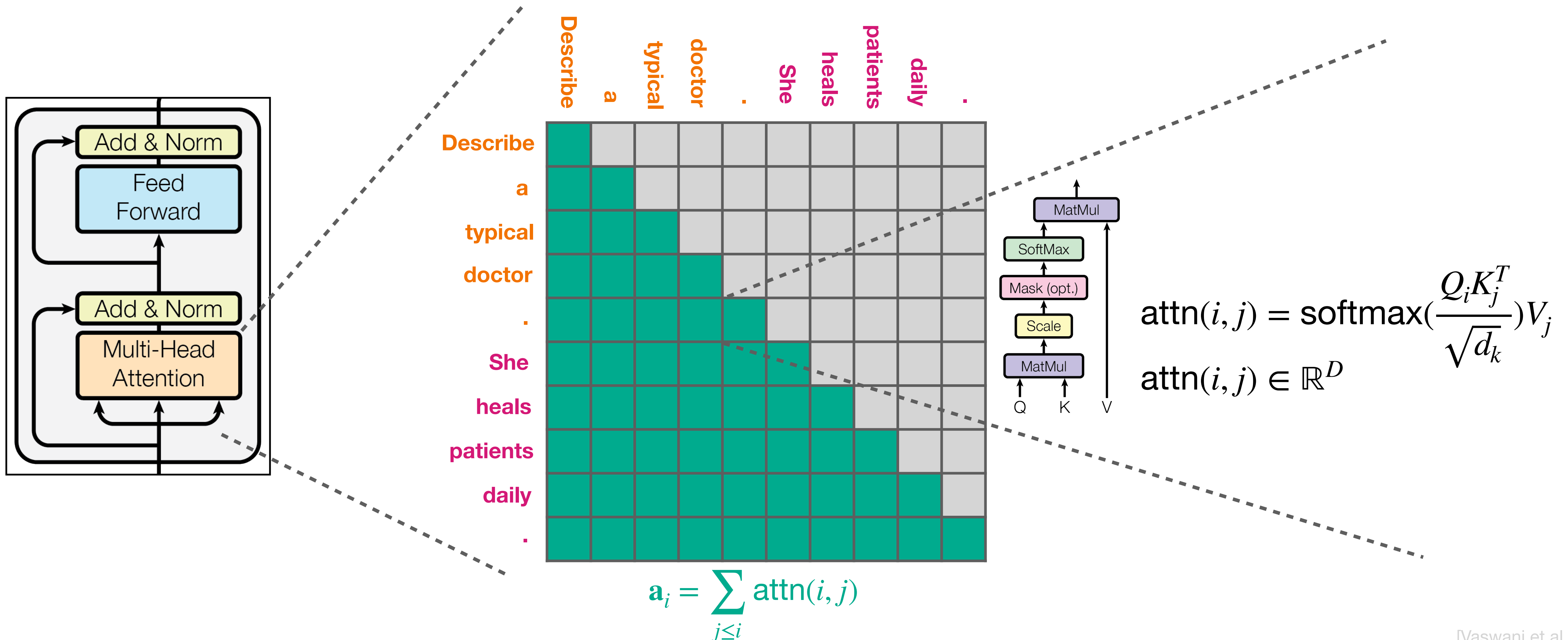
Causal LLMs

Attention of *token i* depends only on current or previous tokens

$$\mathbf{a}_i = \sum_{j \leq i} \text{attn}(i, j)$$

Self attention between all token pairs

When processing token i , how much each token j should contribute?



The nitty gritty

$$Q_i \in \mathbb{R}^D$$



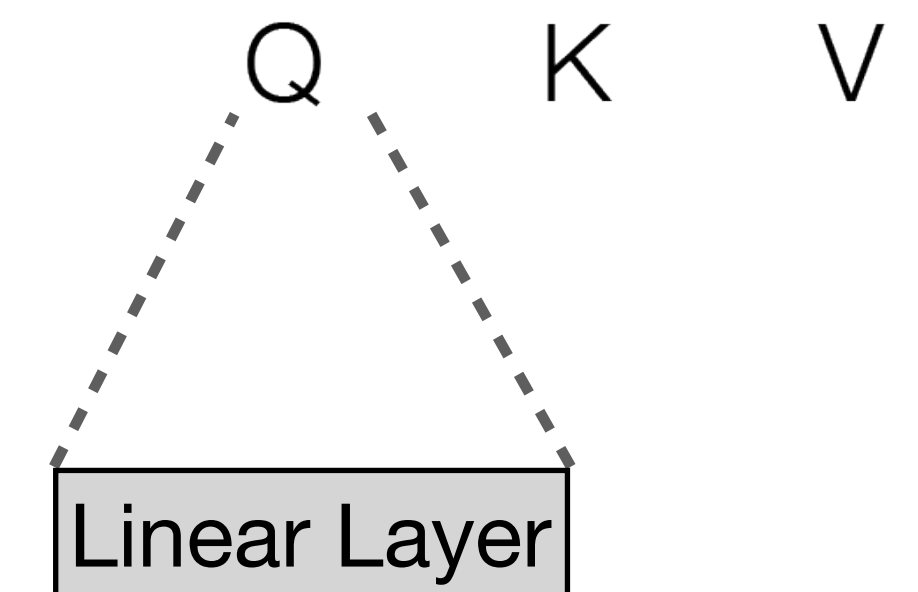
$$K_i \in \mathbb{R}^D$$



$$V_i \in \mathbb{R}^D$$



$$T_i \in \mathbb{R}^D$$



The nitty gritty

$\text{attn}(i, j) \in \mathbb{R}^D$ 

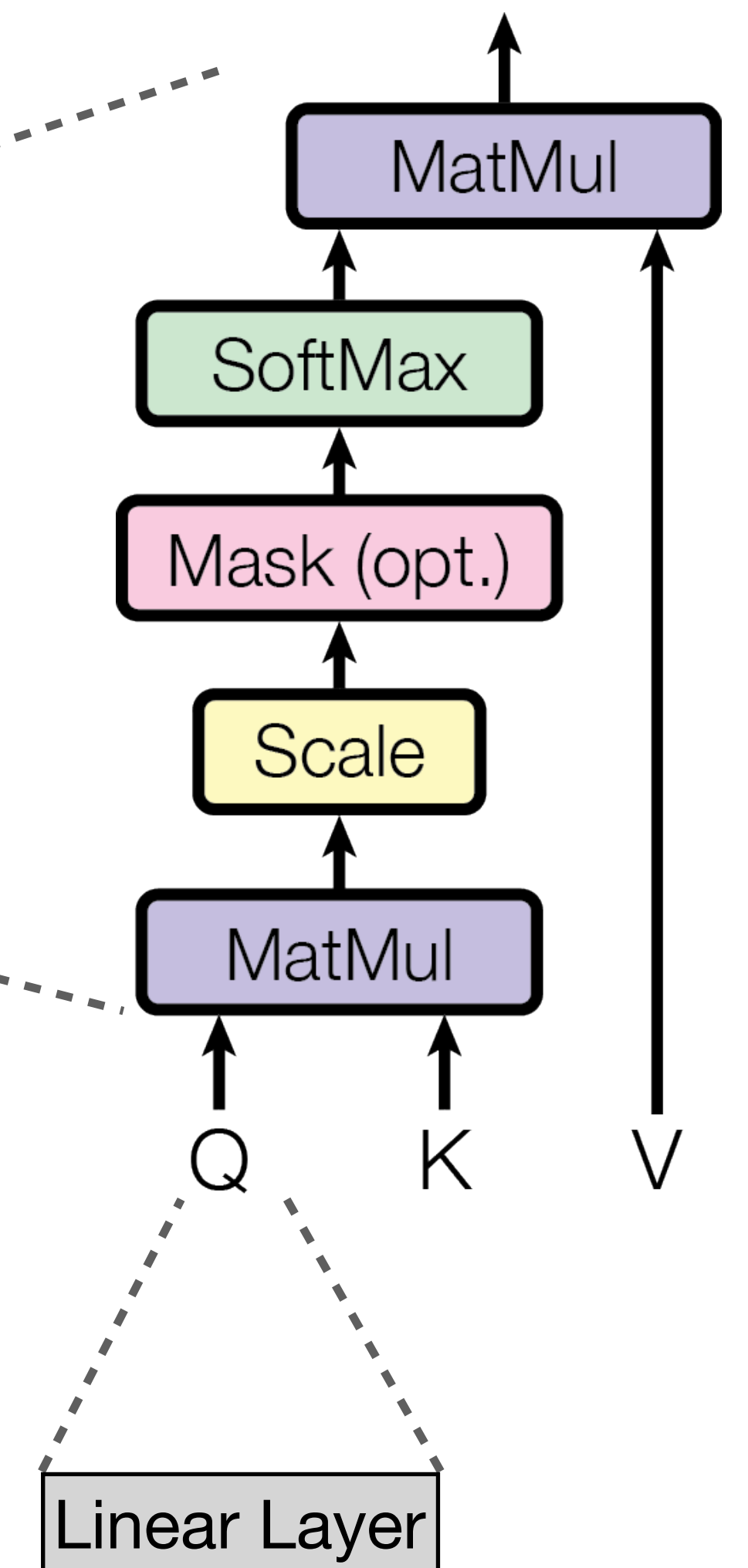
$$\text{attn}(i, j) = \text{softmax}\left(\frac{Q_i K_j^T}{\sqrt{d_k}}\right) V_j$$

$Q_i \in \mathbb{R}^D$ 

$K_i \in \mathbb{R}^D$ 

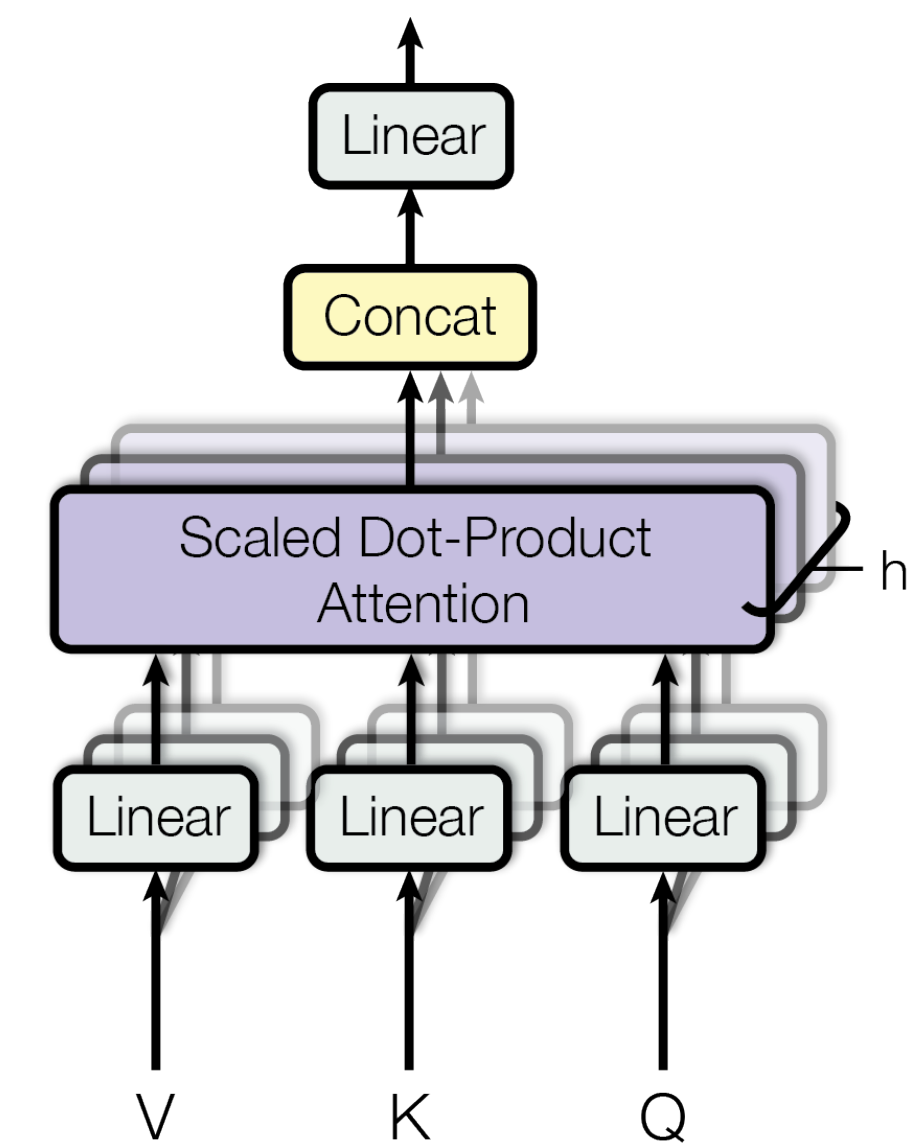
$V_i \in \mathbb{R}^D$ 

$T_i \in \mathbb{R}^D$ 

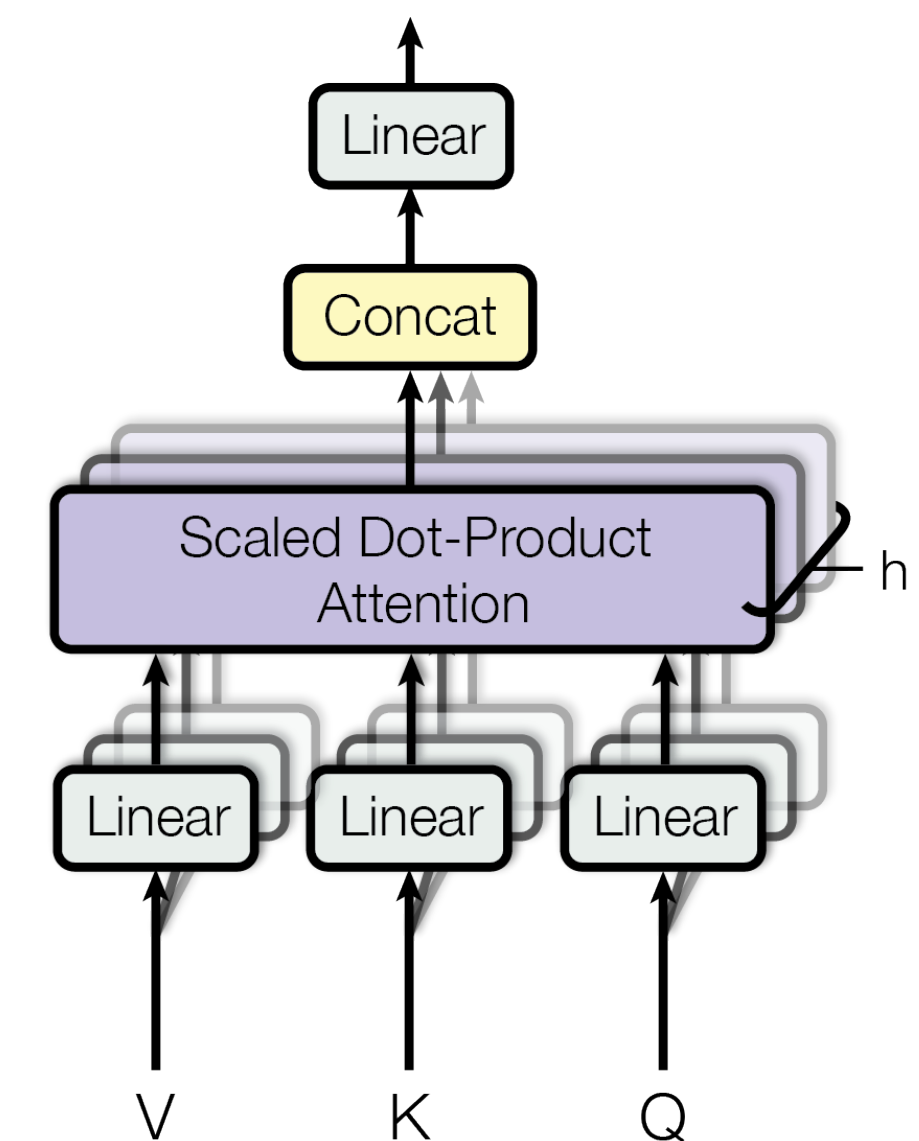
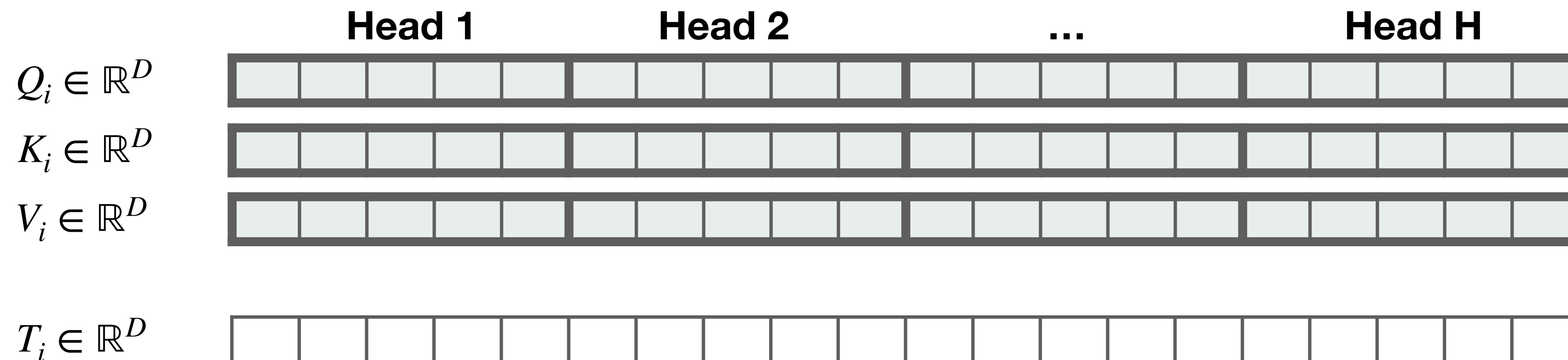


Attention is multi-headed

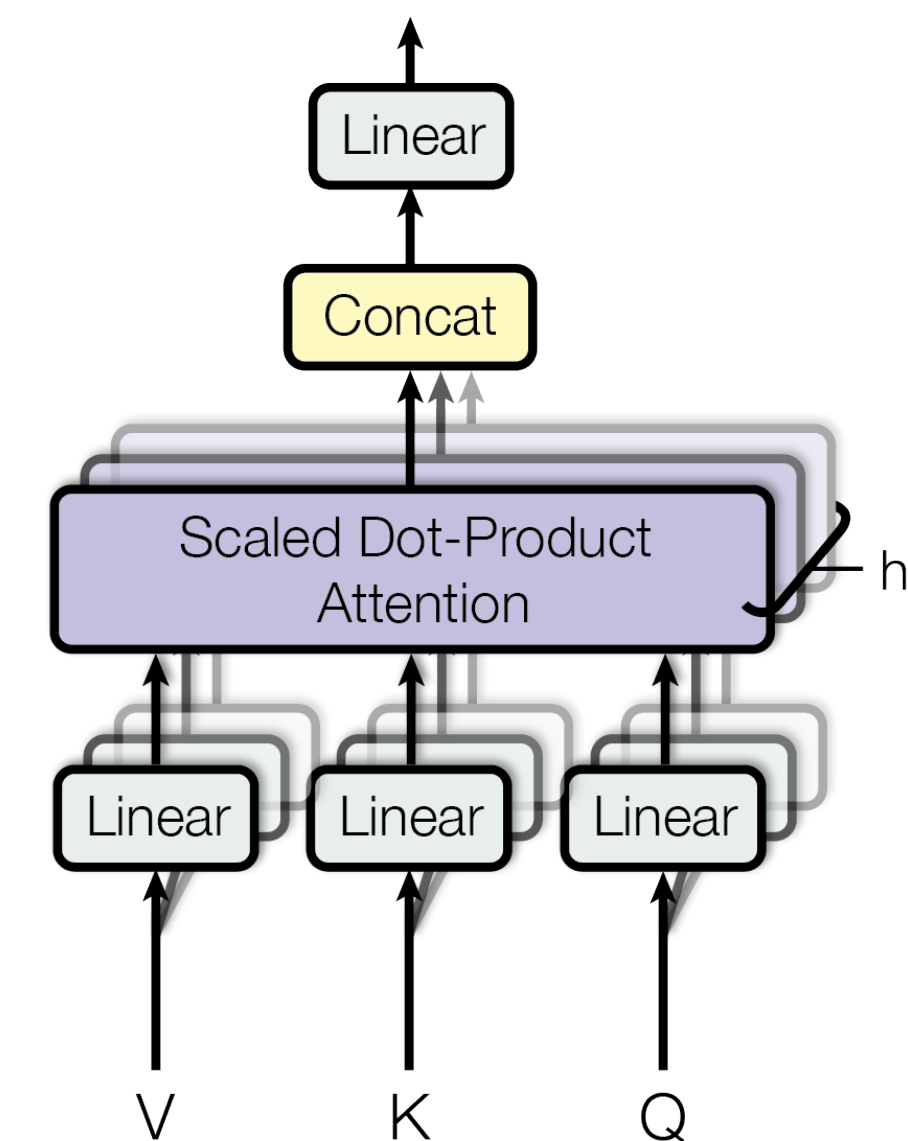
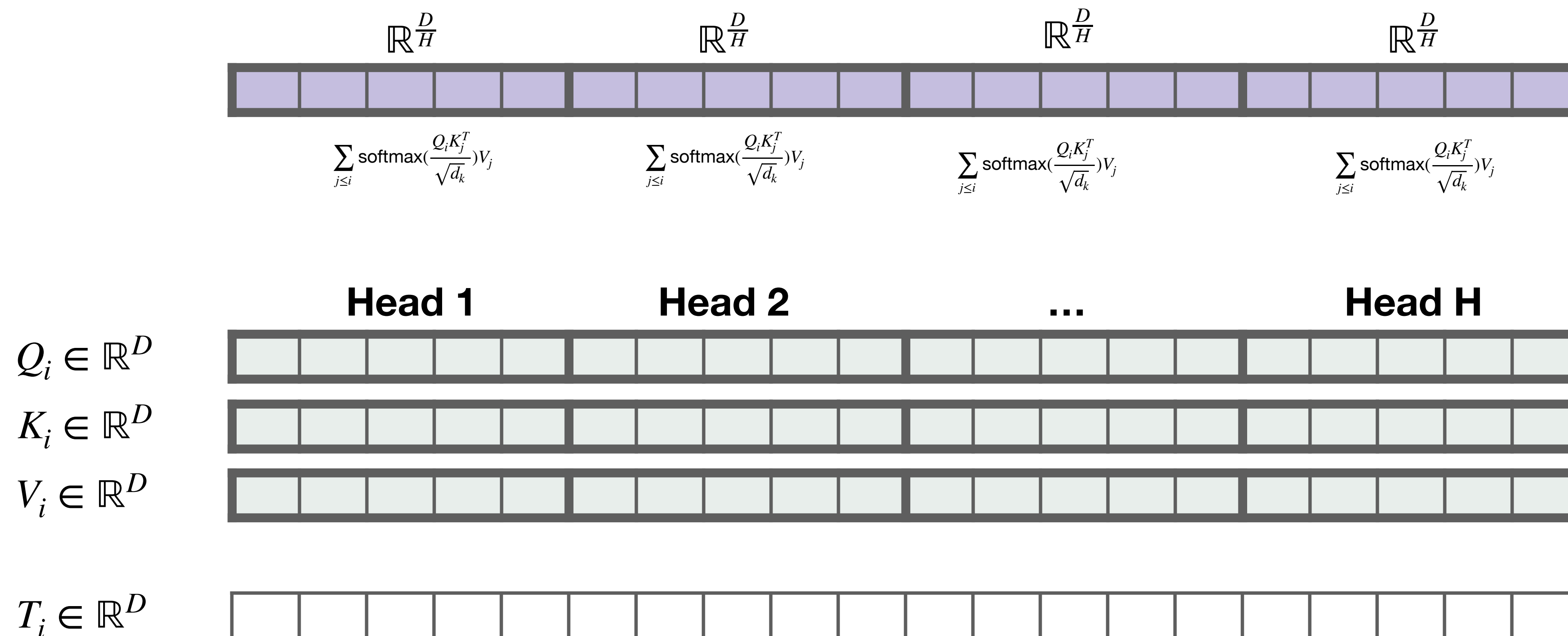
$$T_i \in \mathbb{R}^D$$



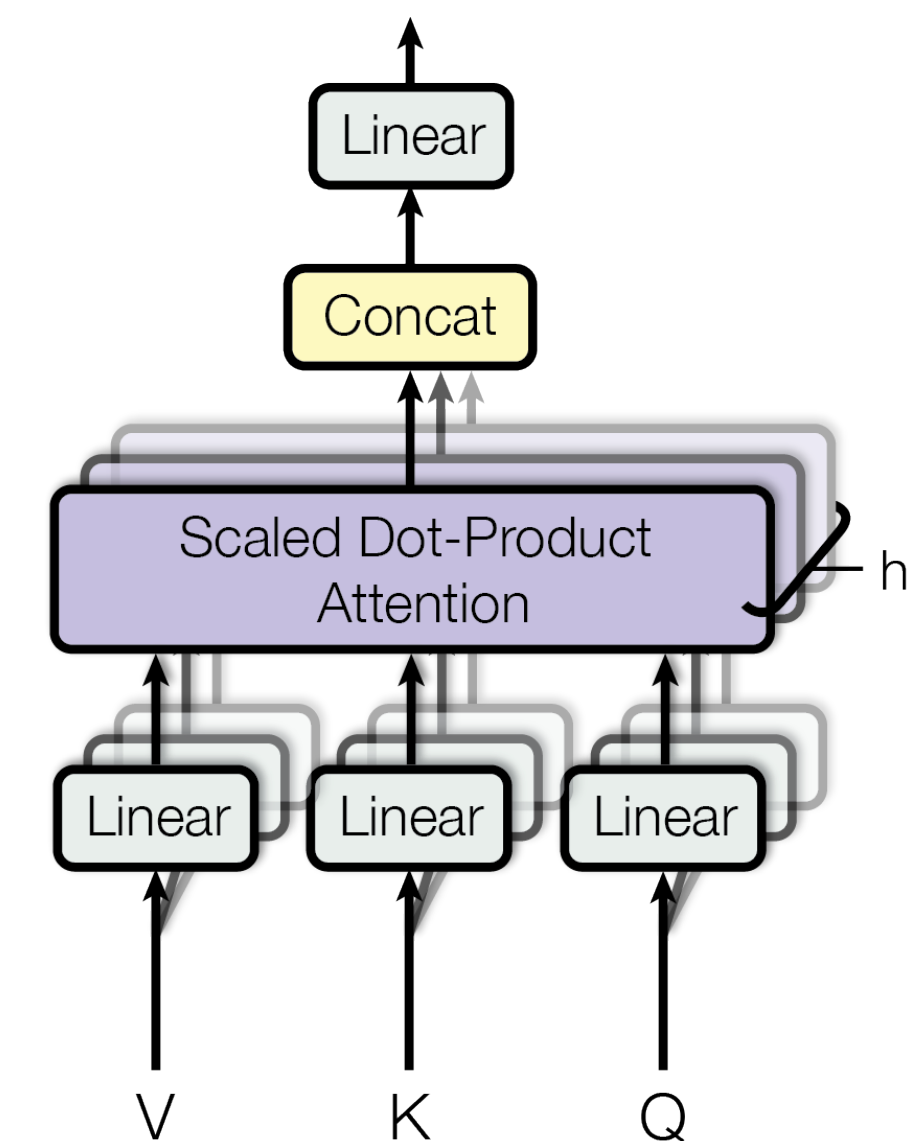
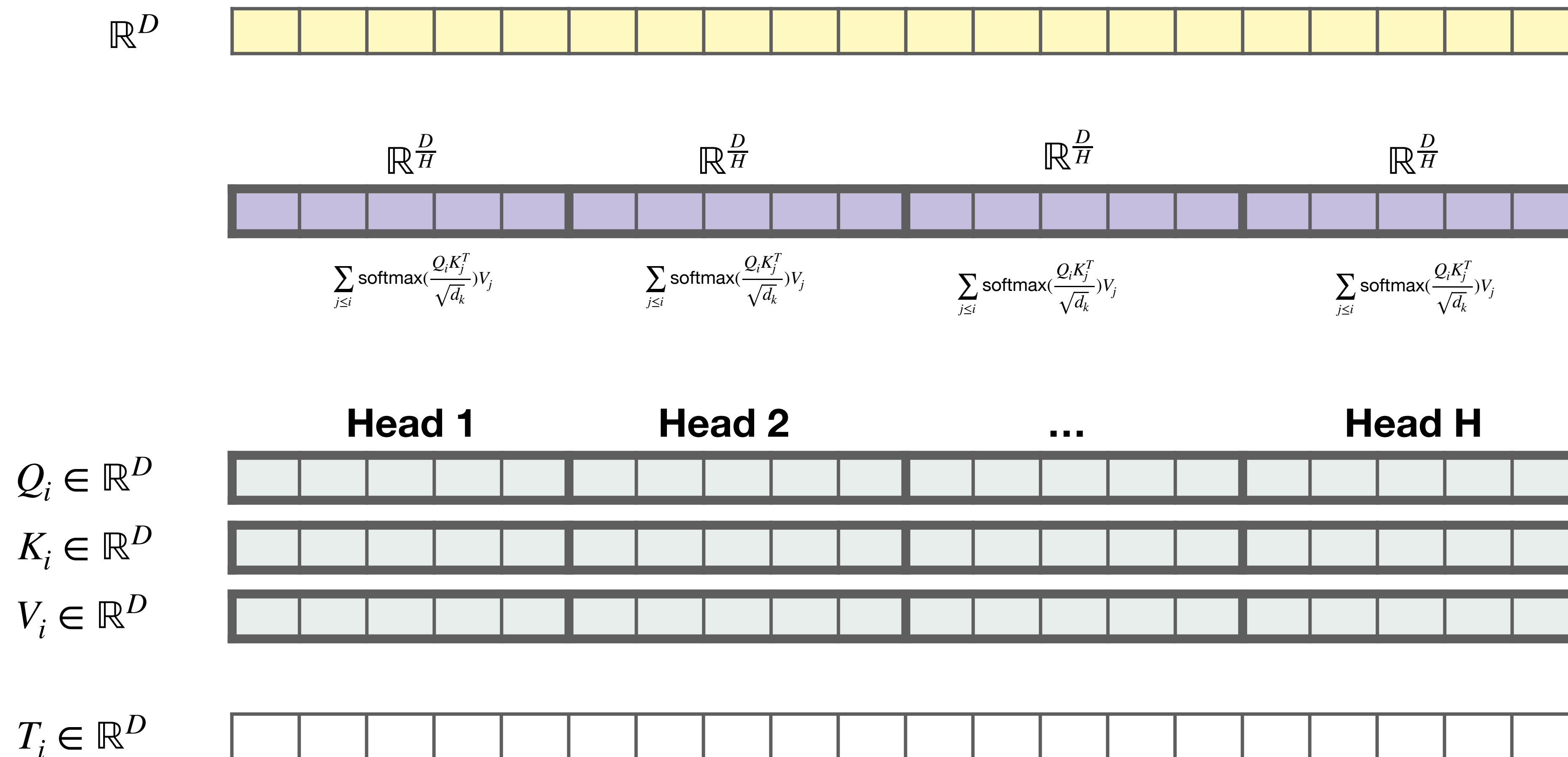
Attention is multi-headed



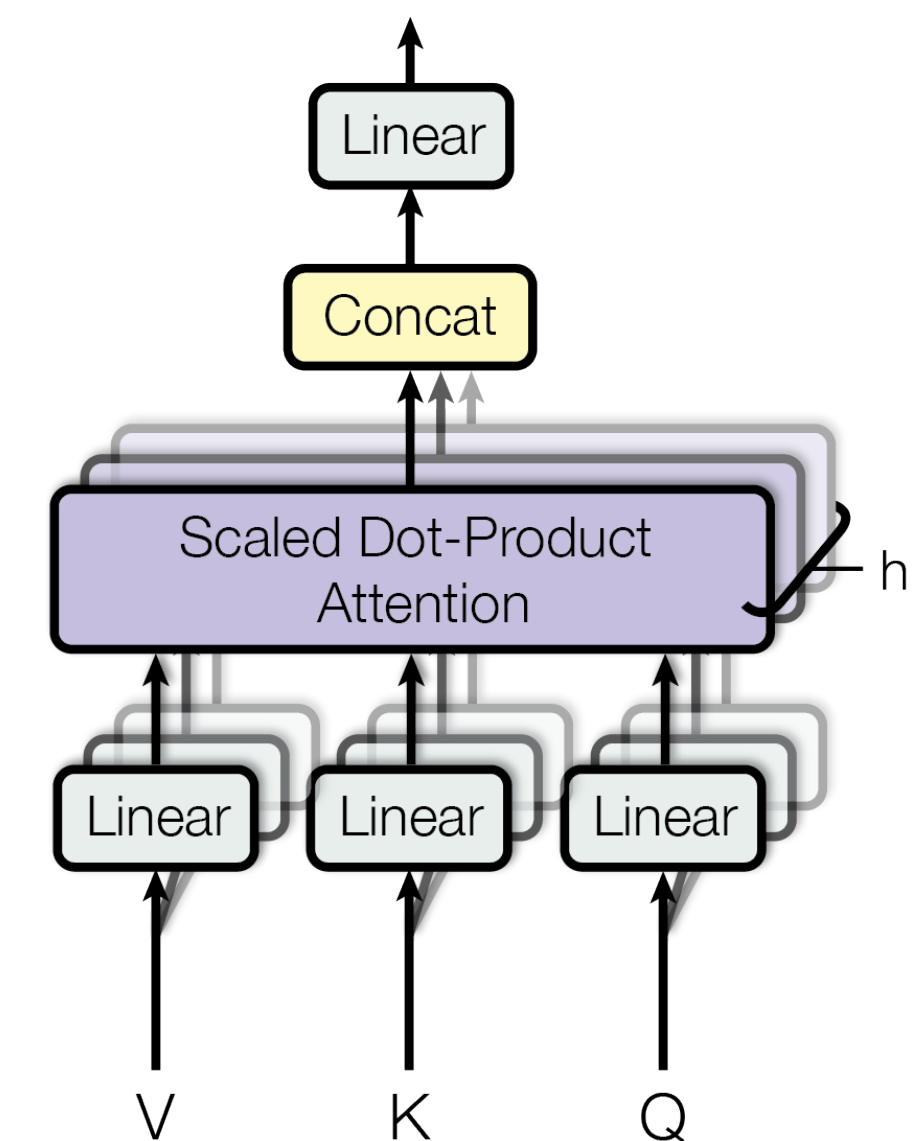
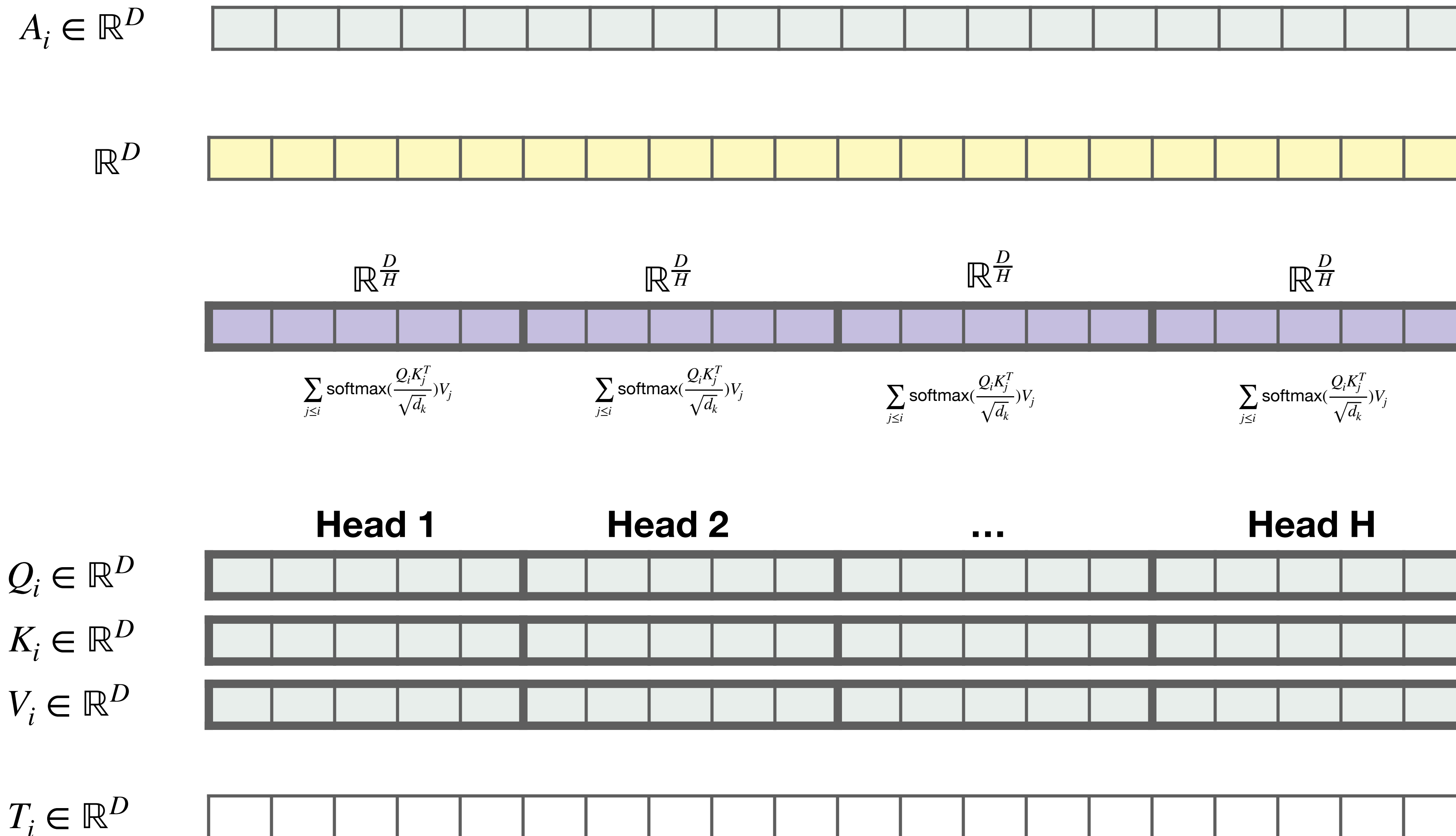
Attention is multi-headed



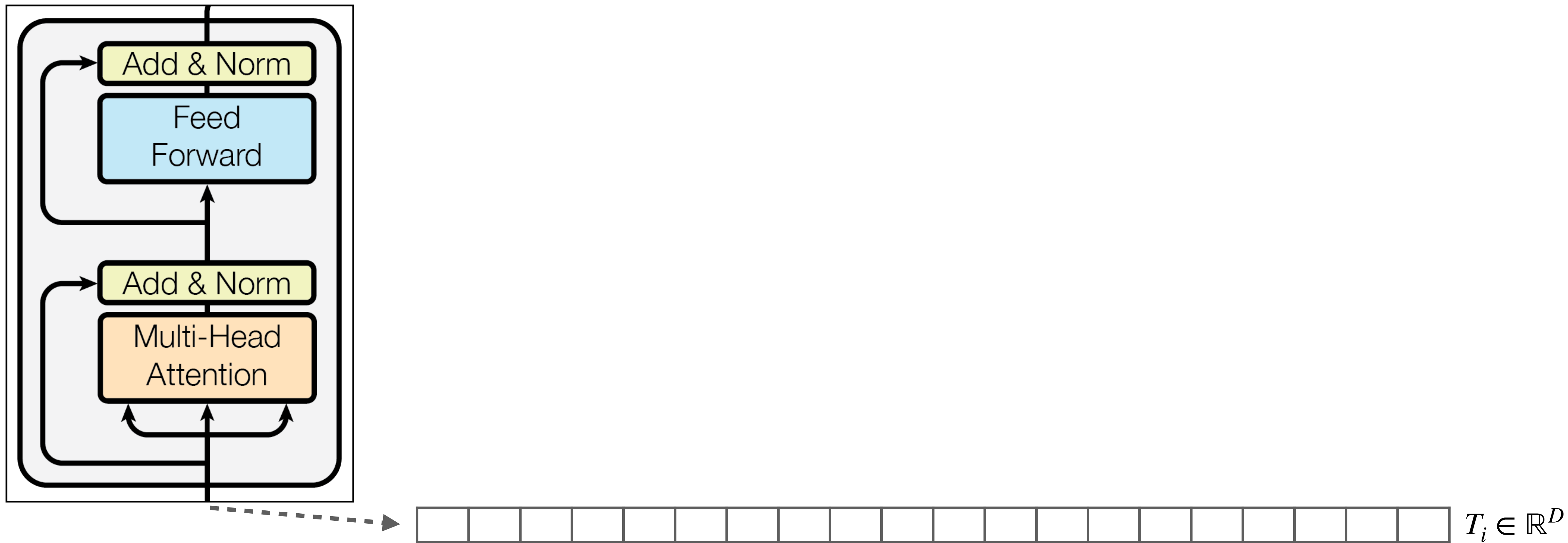
Attention is multi-headed



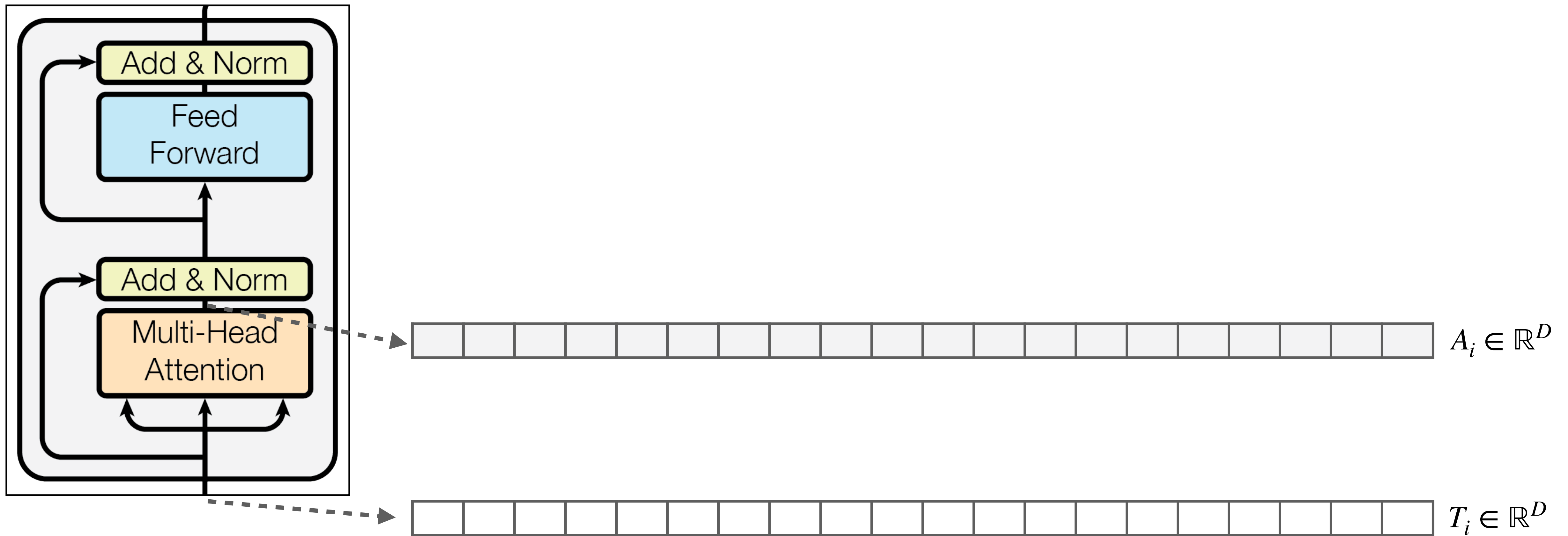
Attention is multi-headed



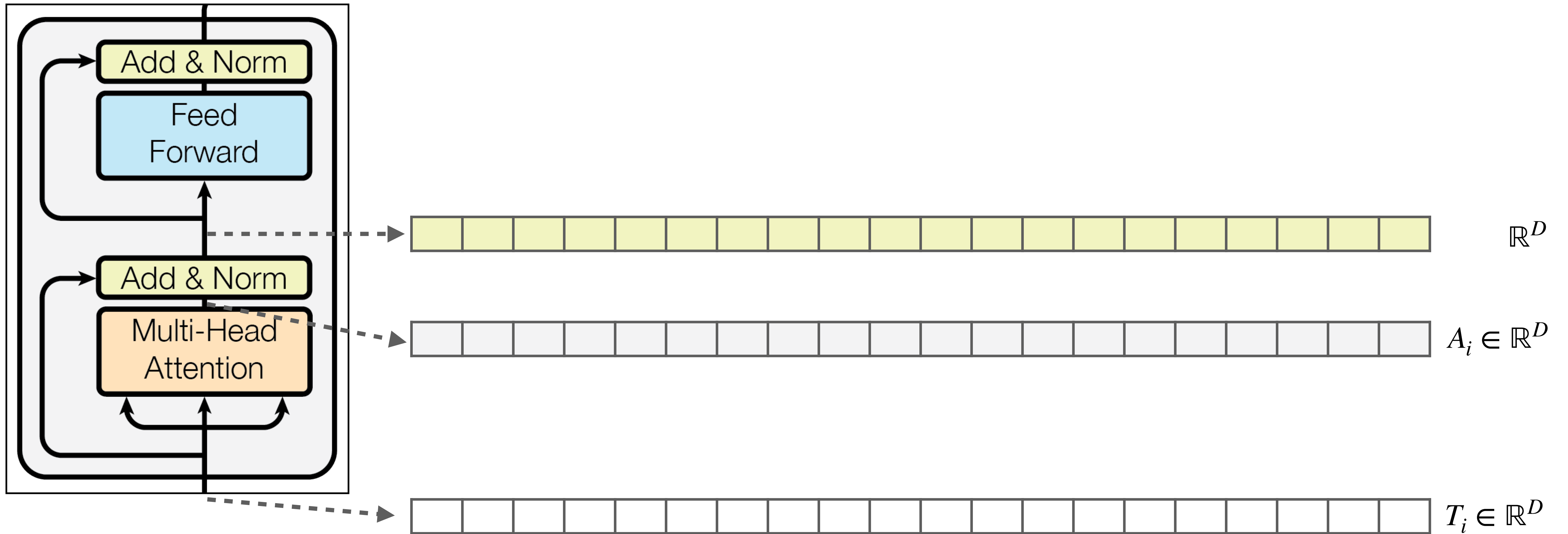
The full transformer layer



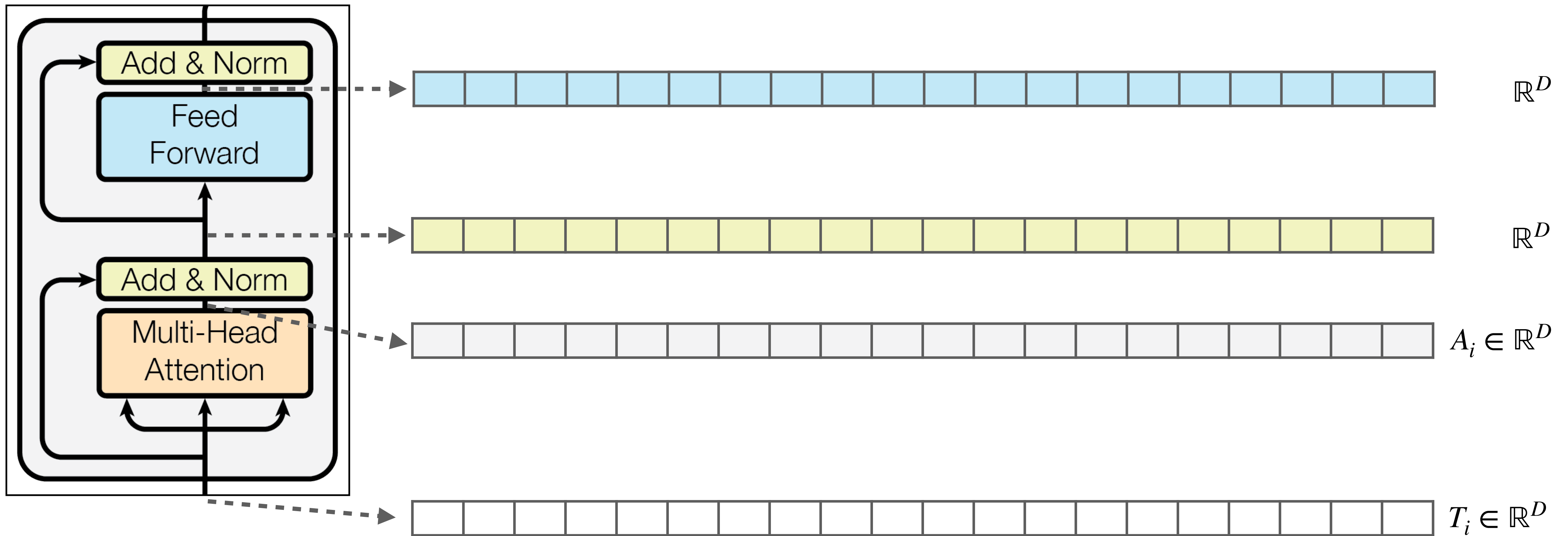
The full transformer layer



The full transformer layer



The full transformer layer



The full transformer layer



The full transformer layer

Both input and the output are \mathbb{R}^D



The full transformer layer

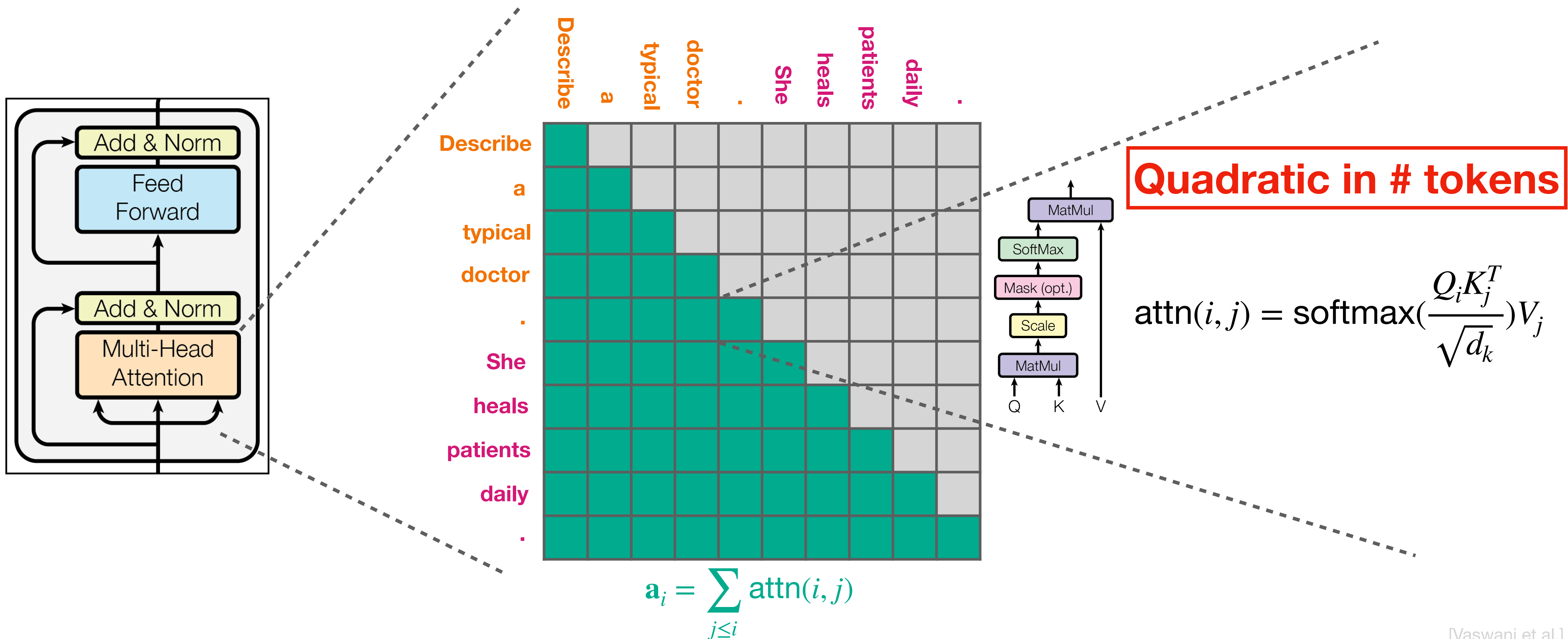
Both input and the output are \mathbb{R}^D

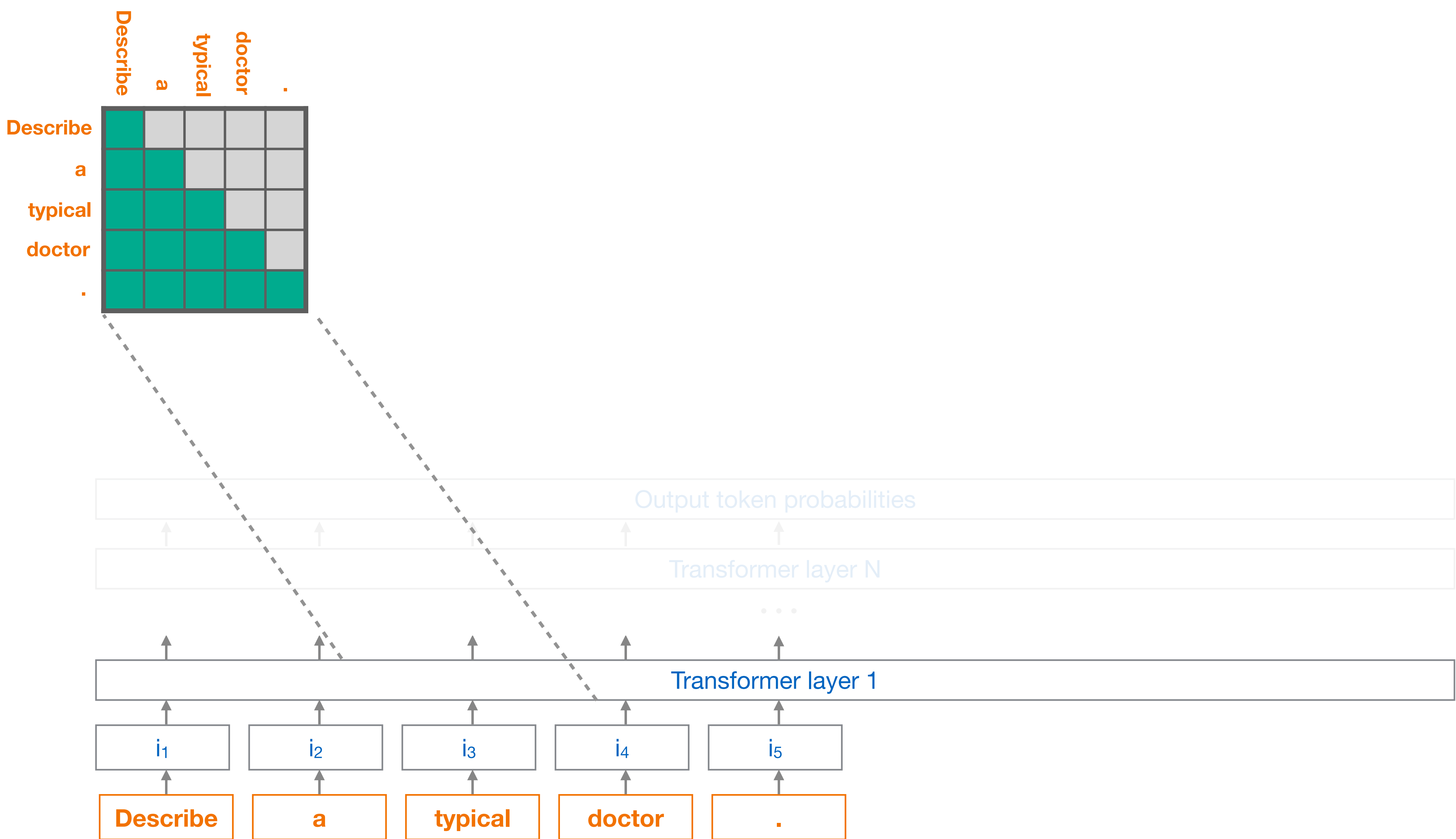


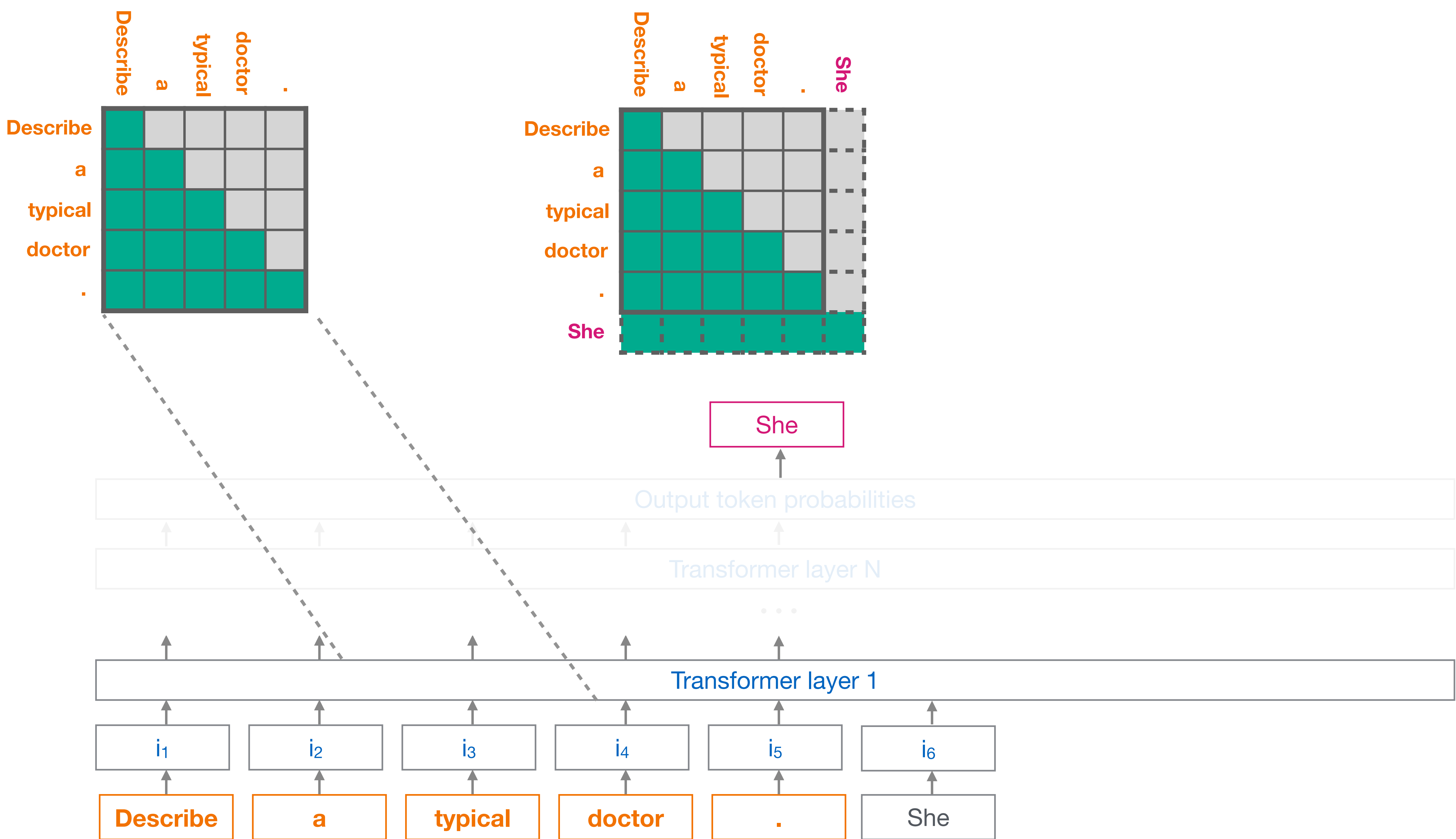
Which part is computationally most expensive?

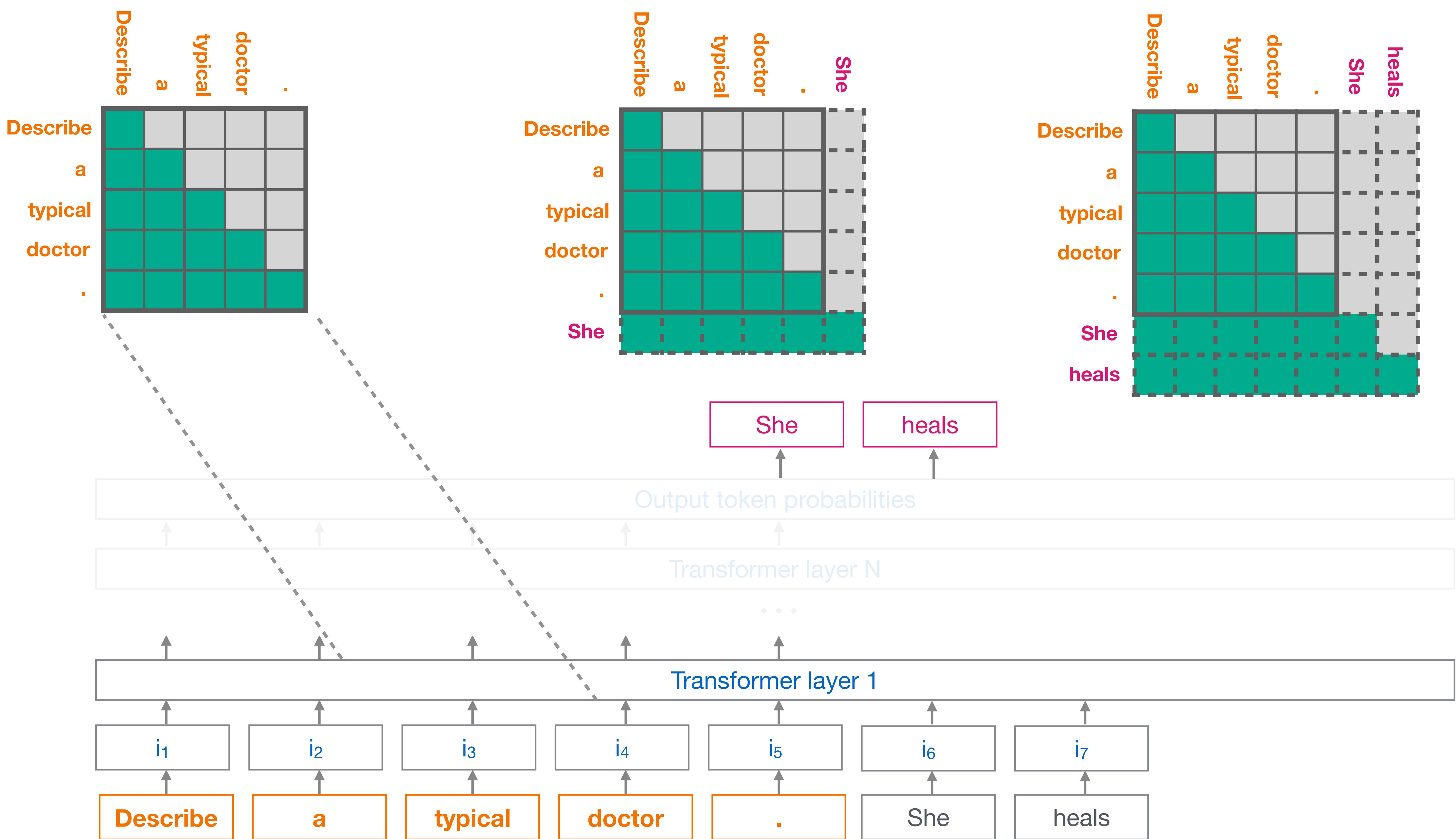
Self attention between all token pairs

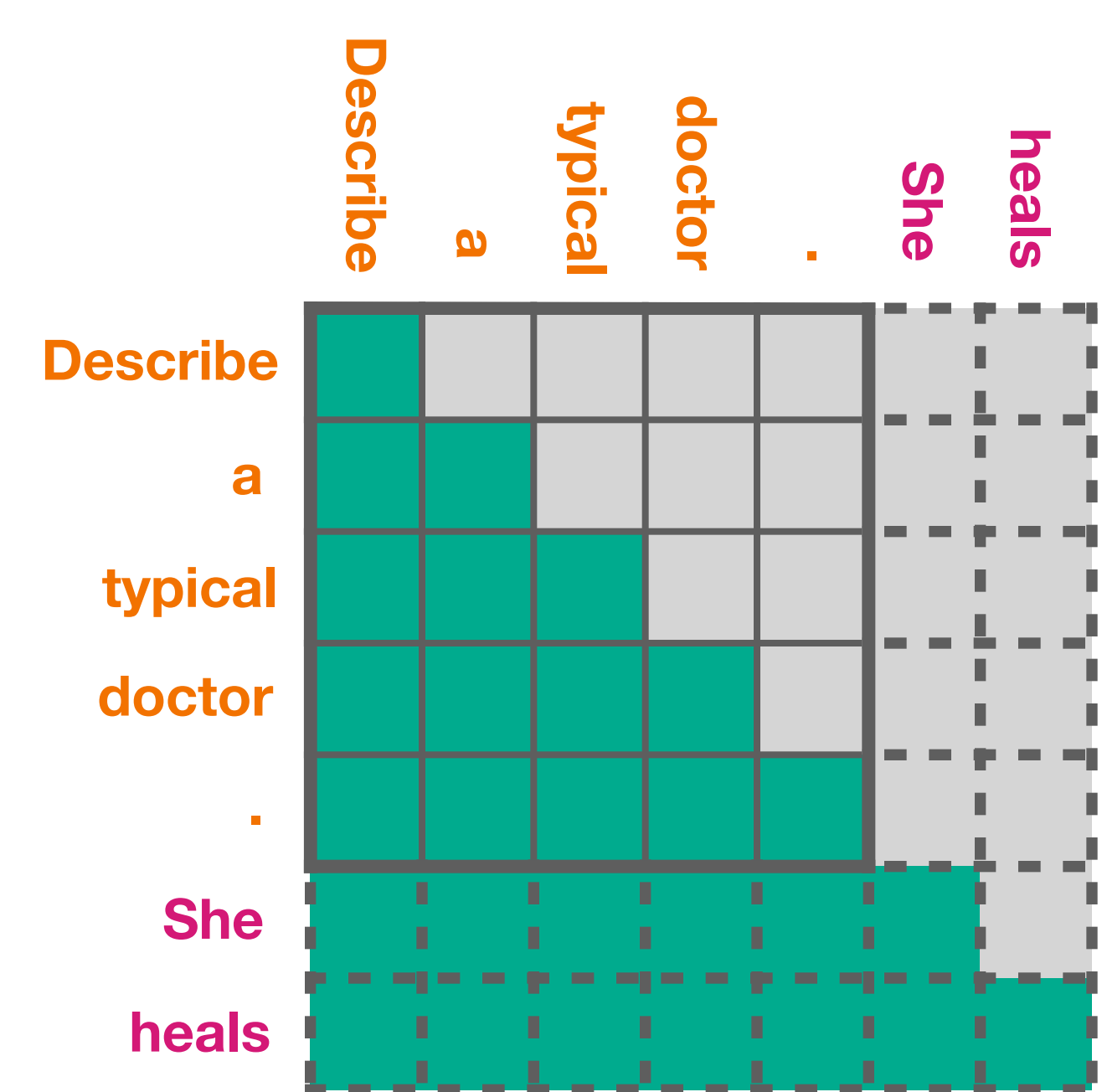
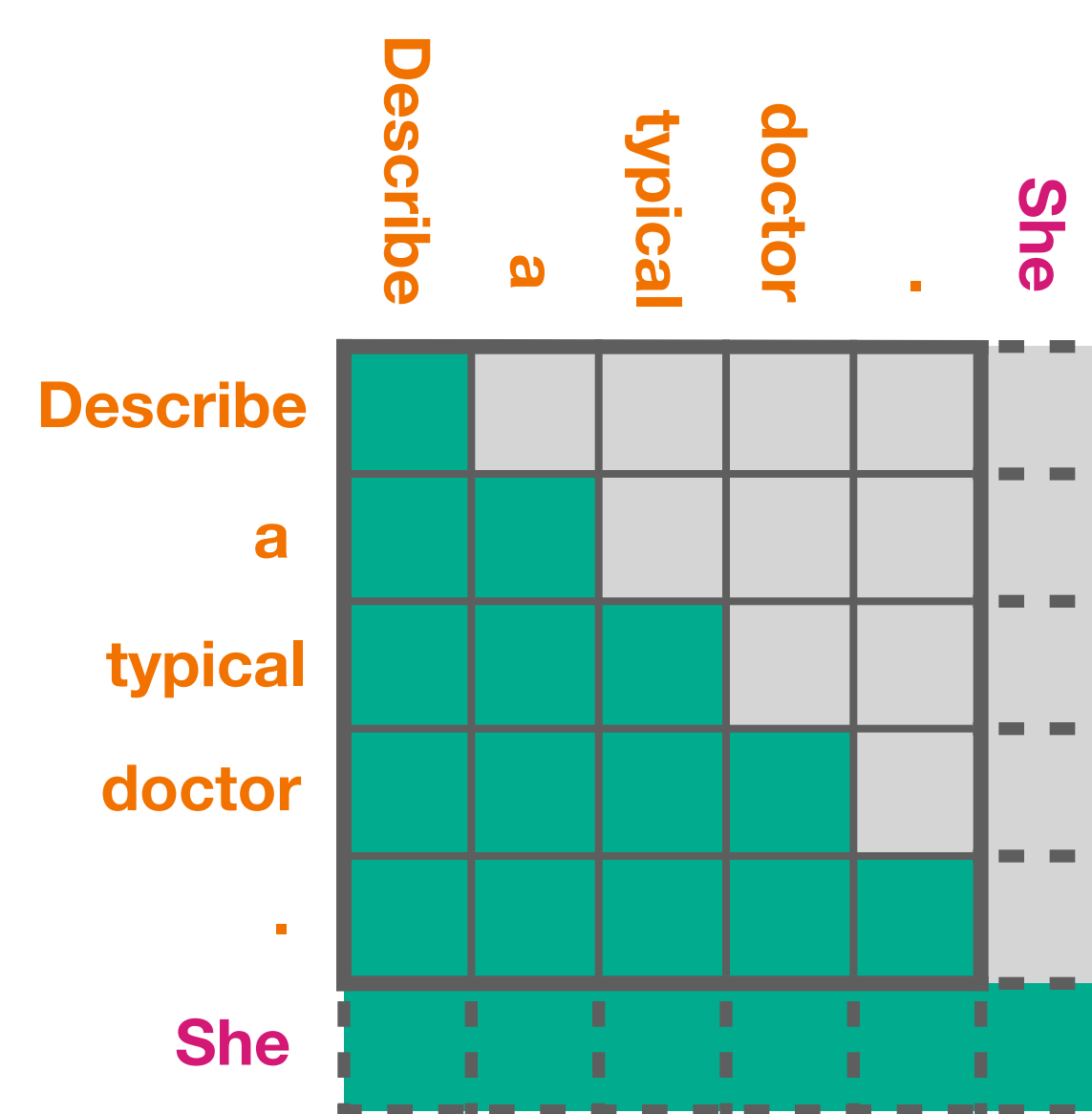
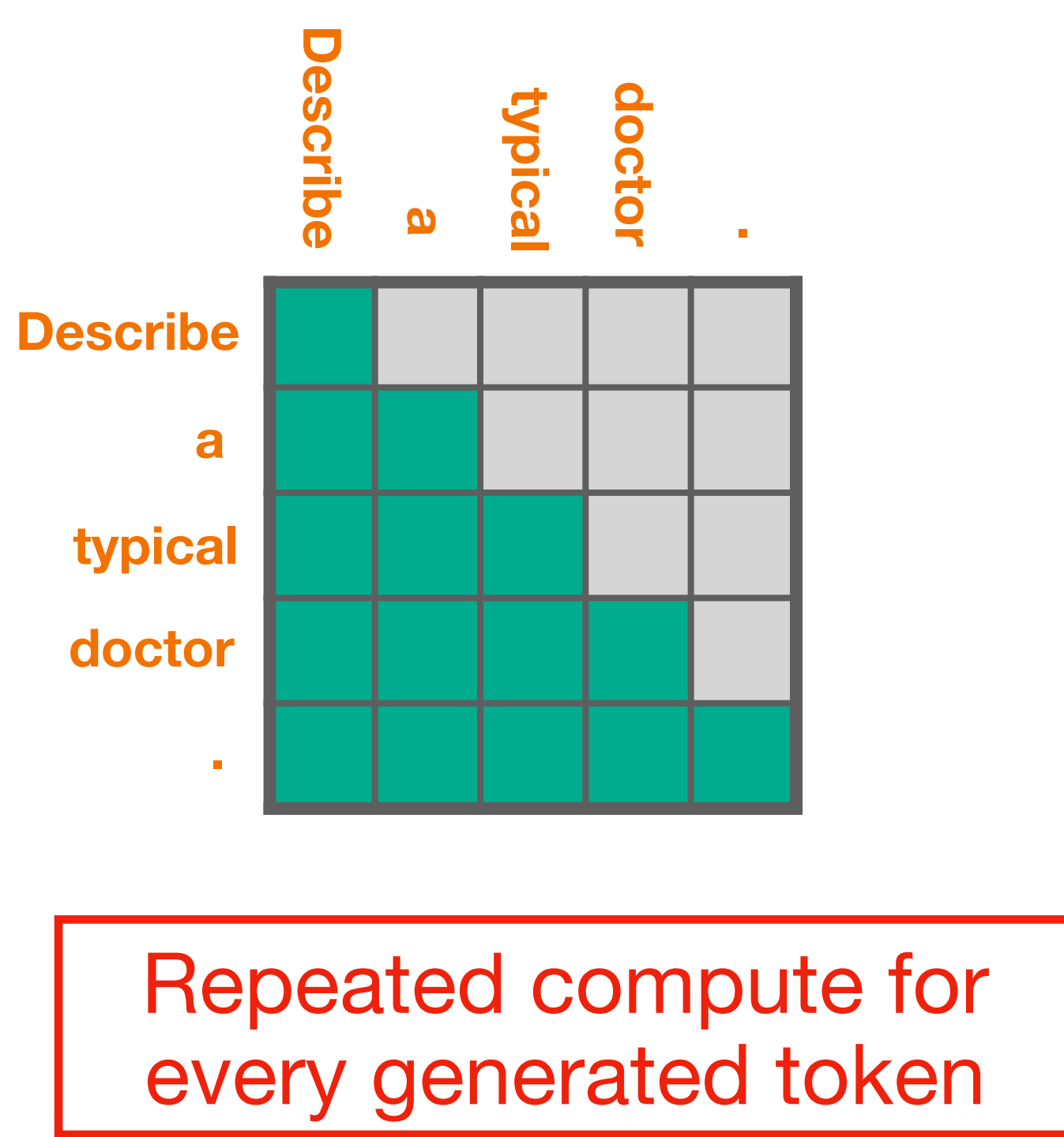
When processing token i , how much each token j should contribute?





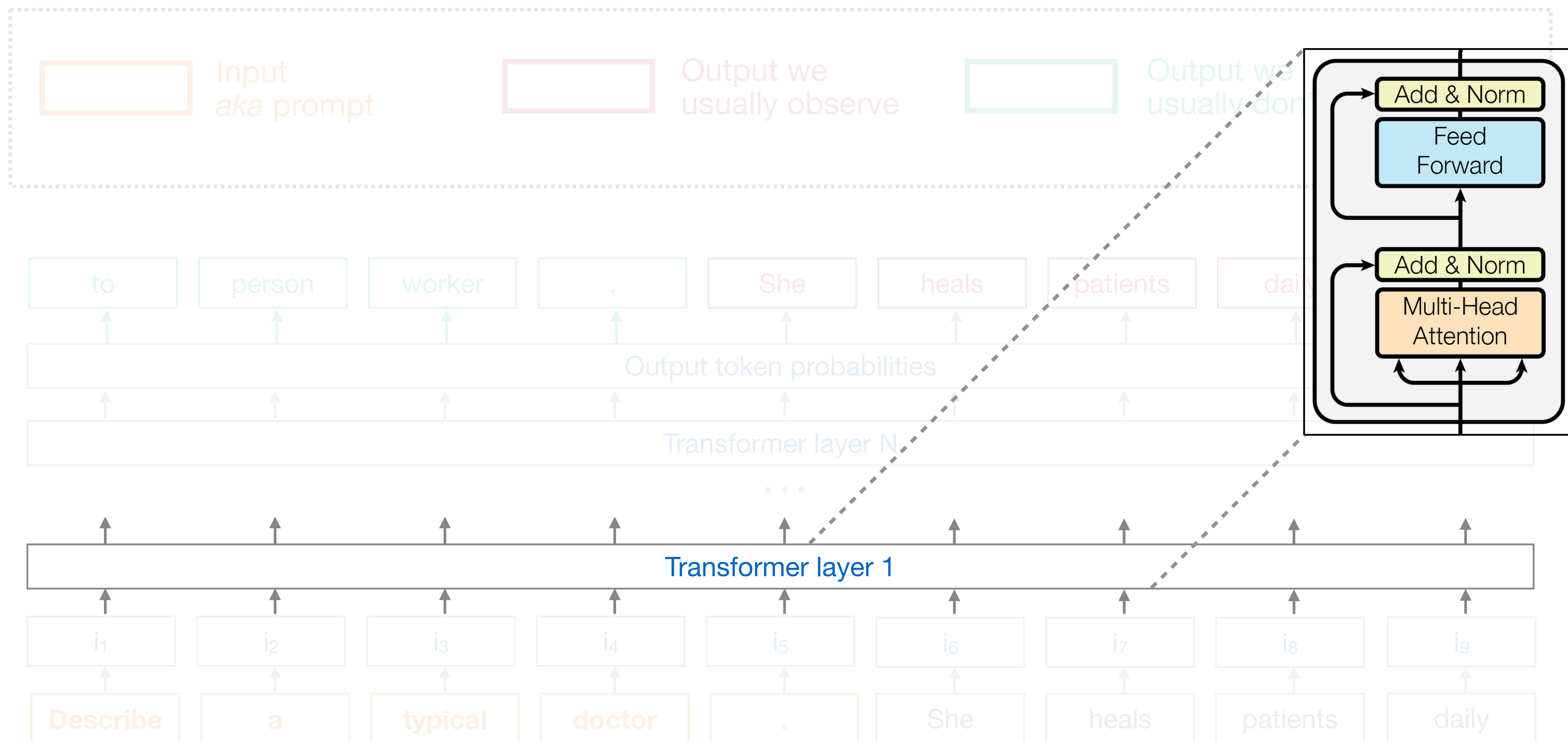




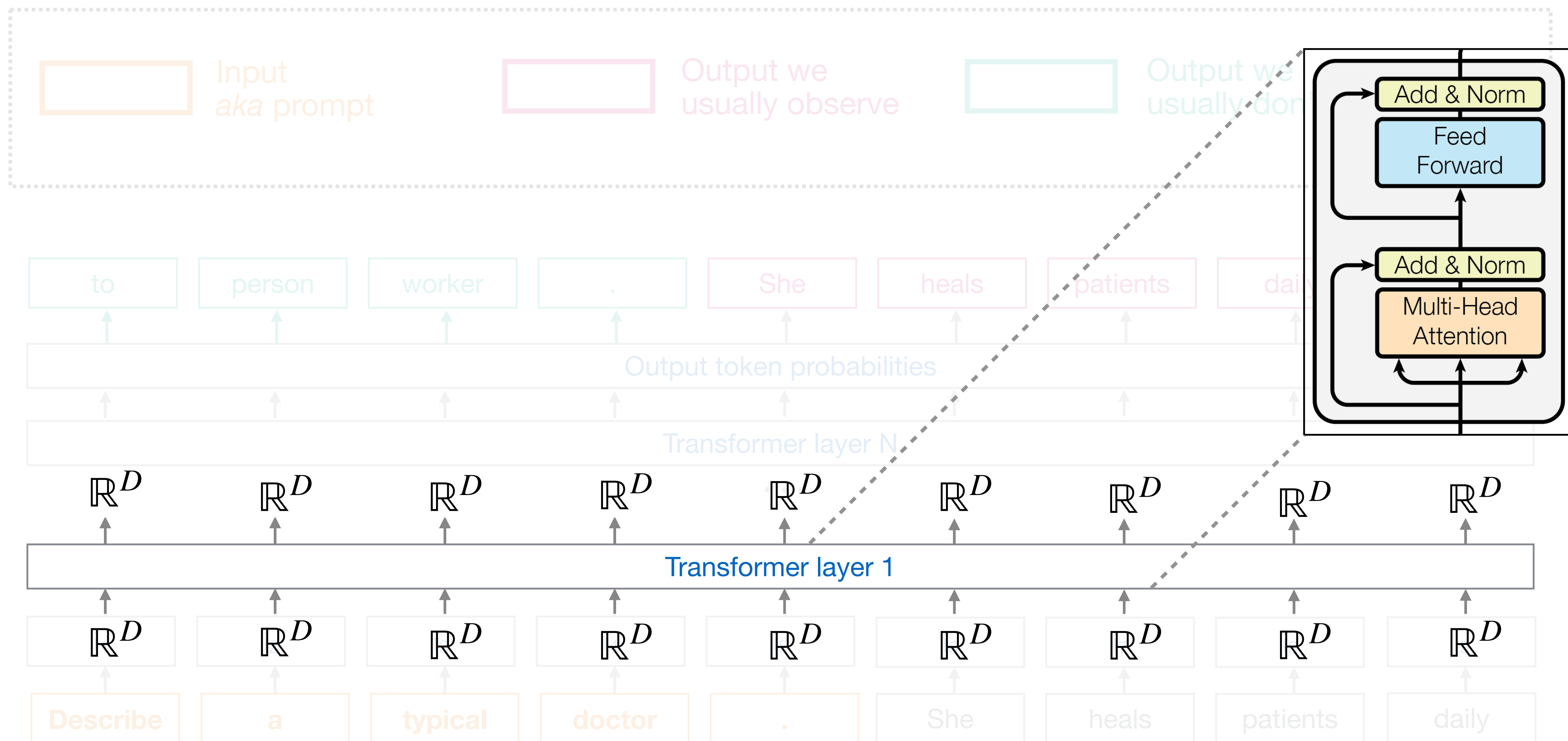


KV-Caching:
Cache the previous attention computation

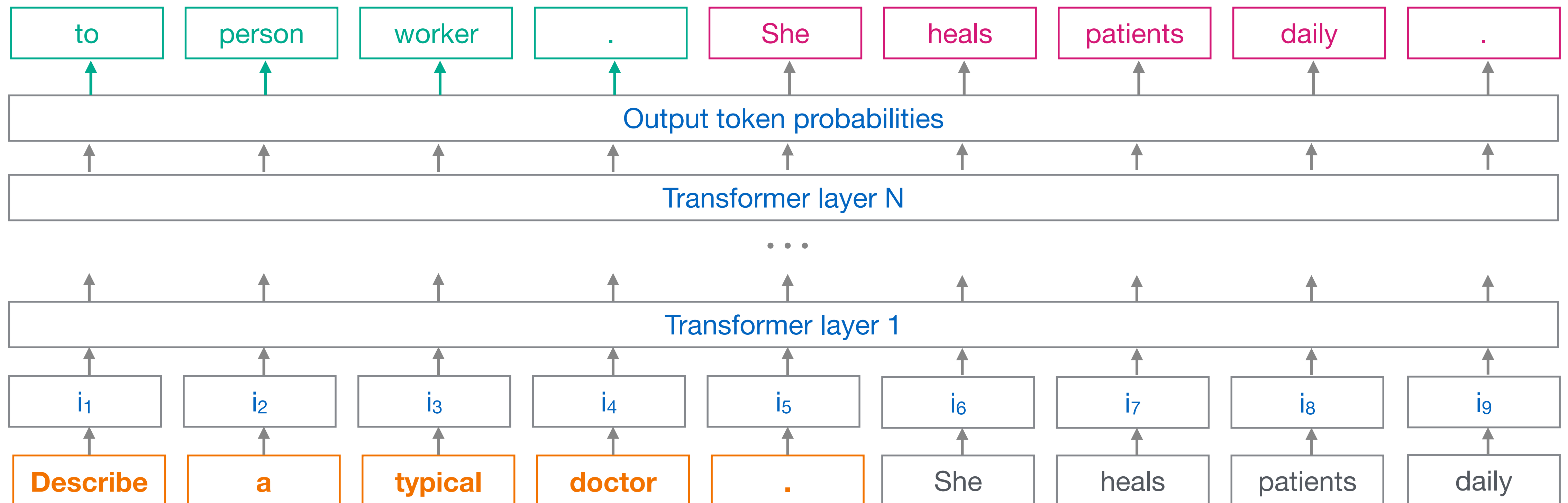
Recap: What is the output dimensionality of each transformer component?



Recap: What is the output dimensionality of each transformer component?



Recap: What is the output dimensionality of each layer?



Exercise

References

- [Attention is all you need](#)
- [Mastering tensor dimensions in transformers](#)
- [The illustrated transformer](#)