

Transformers

CS 4804: Introduction to AI

Fall 2025

<https://tuvllms.github.io/ai-fall-2025/>

Tu Vu



Logistics

- Homework 0 (due September 16th)
 - accuracy
- Quiz 0 (due tomorrow)
 - genuine attempt
- Final project group information (tomorrow)

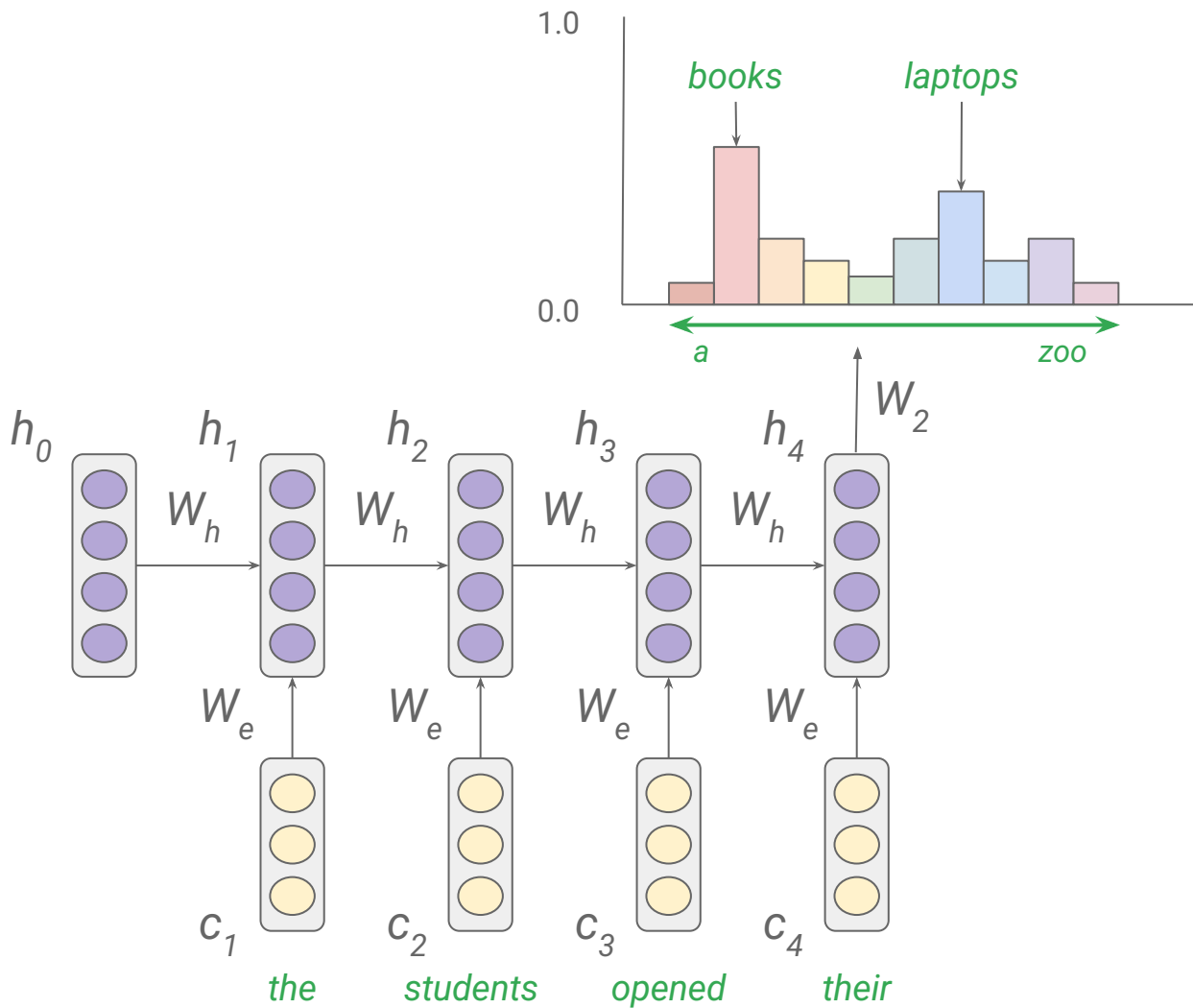
Recurrent neural networks (RNNs)

hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c^t)$$

output distribution

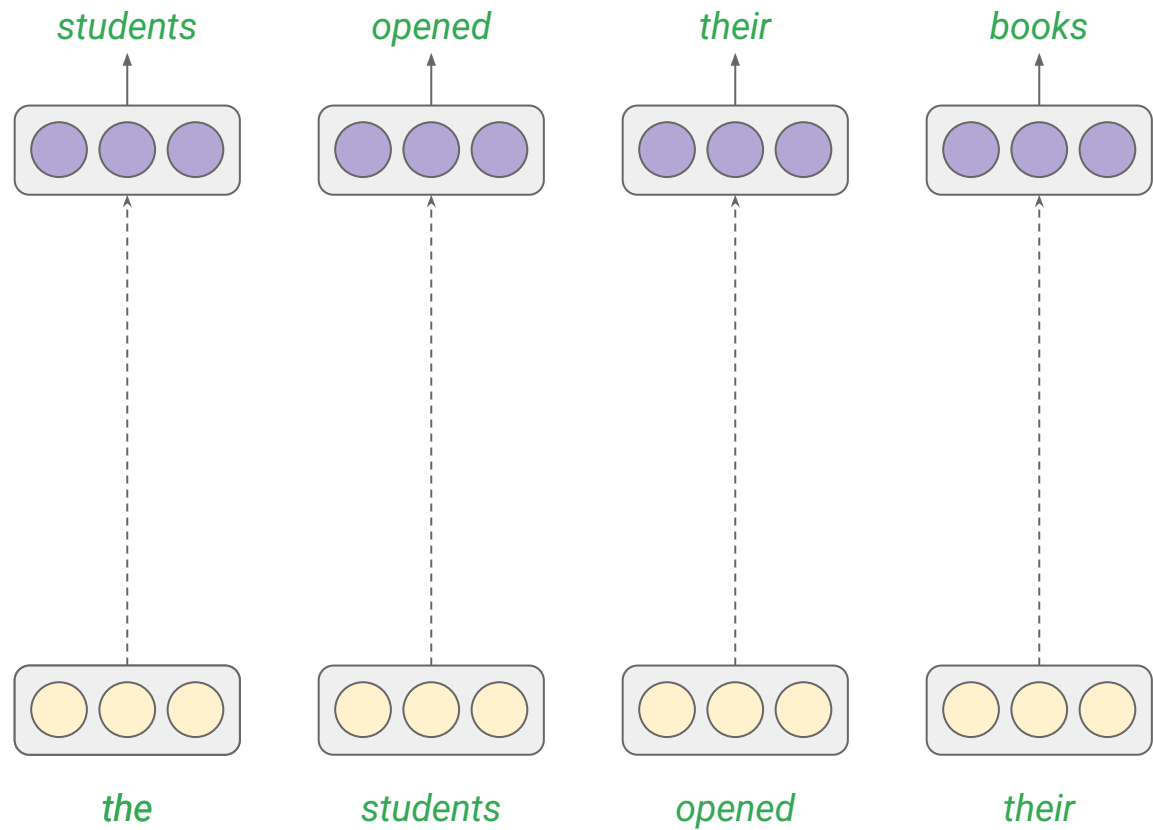
$$\hat{y} = \text{softmax}(W_2 h^{(n-1)})$$



Problems with RNNs

- Bottleneck representation issue
- Lack of parallelism

Seq2Seq



Transformers

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

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

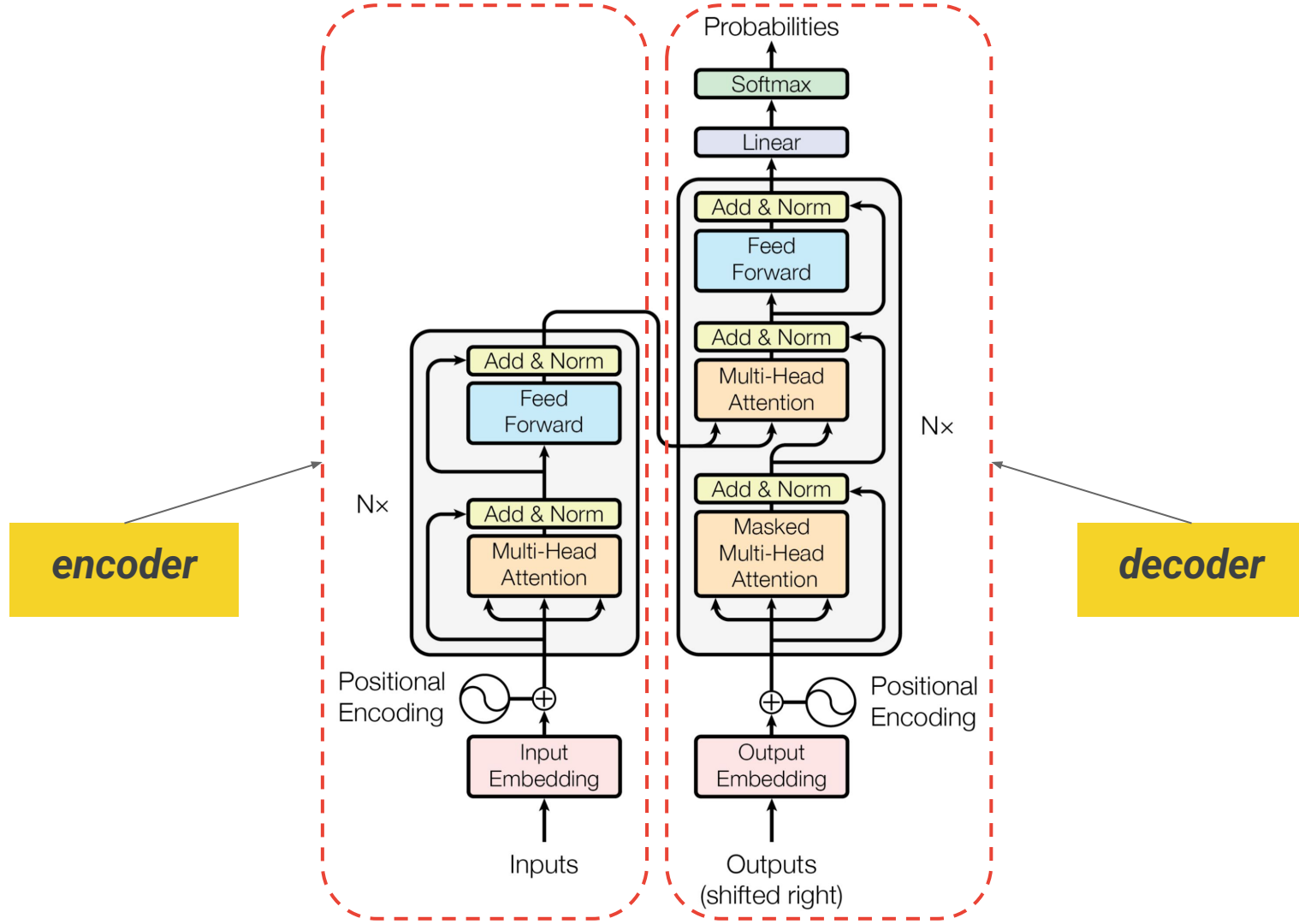
Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* †
illia.polosukhin@gmail.com

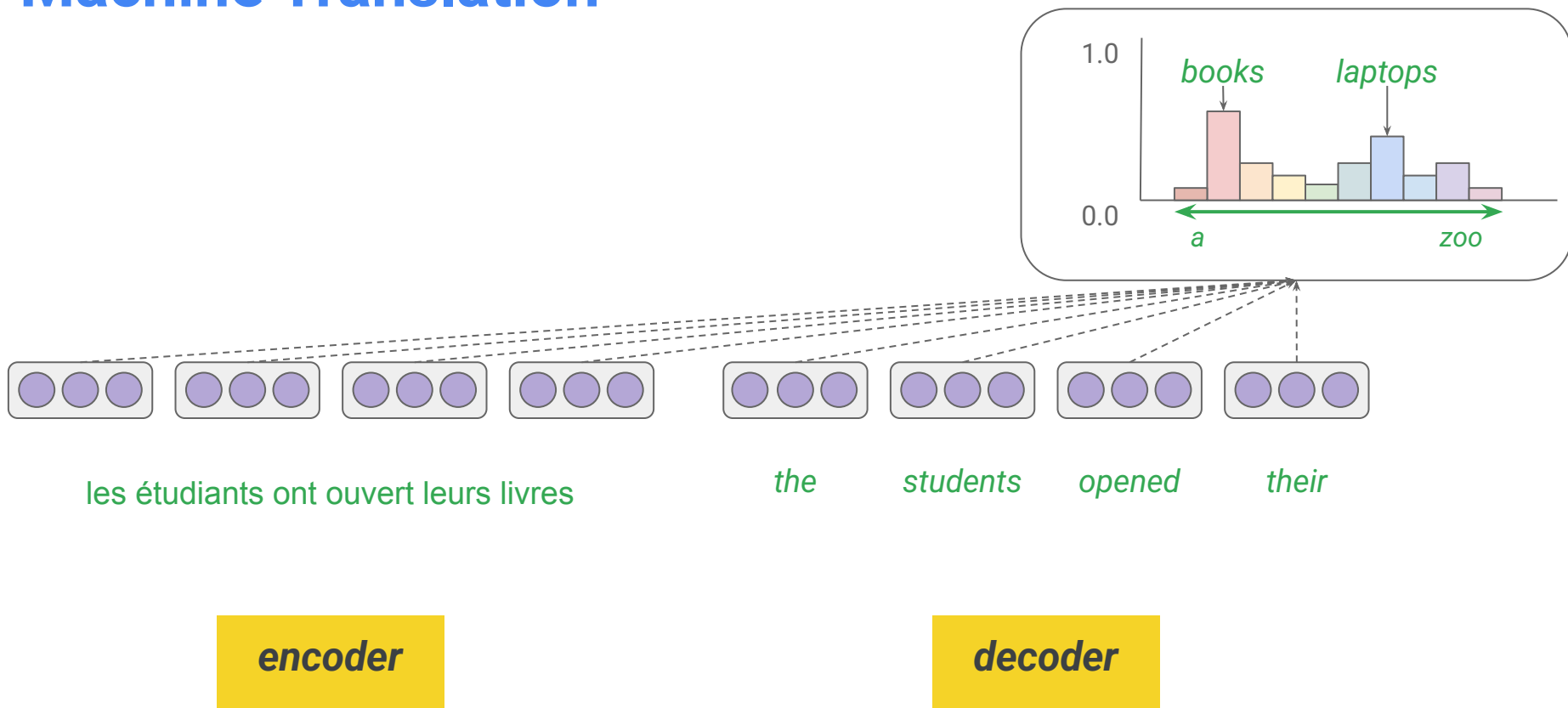


Transformers

- Before 2017
 - Recurrent neural networks (RNNs)
 - LSTM (Long Short-Term Memory)
 - Convolutional neural networks (CNNs)
- These days
 - Transformers



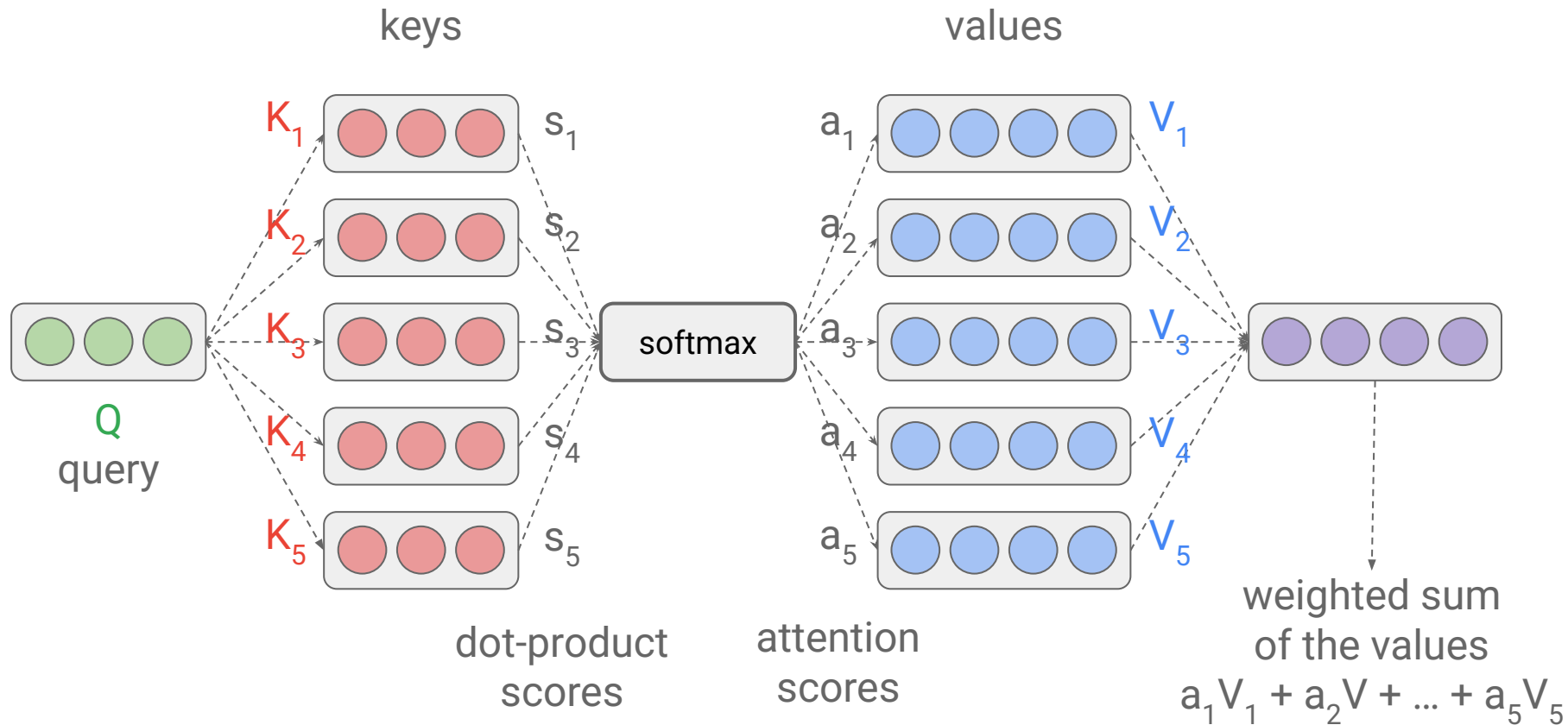
Machine Translation



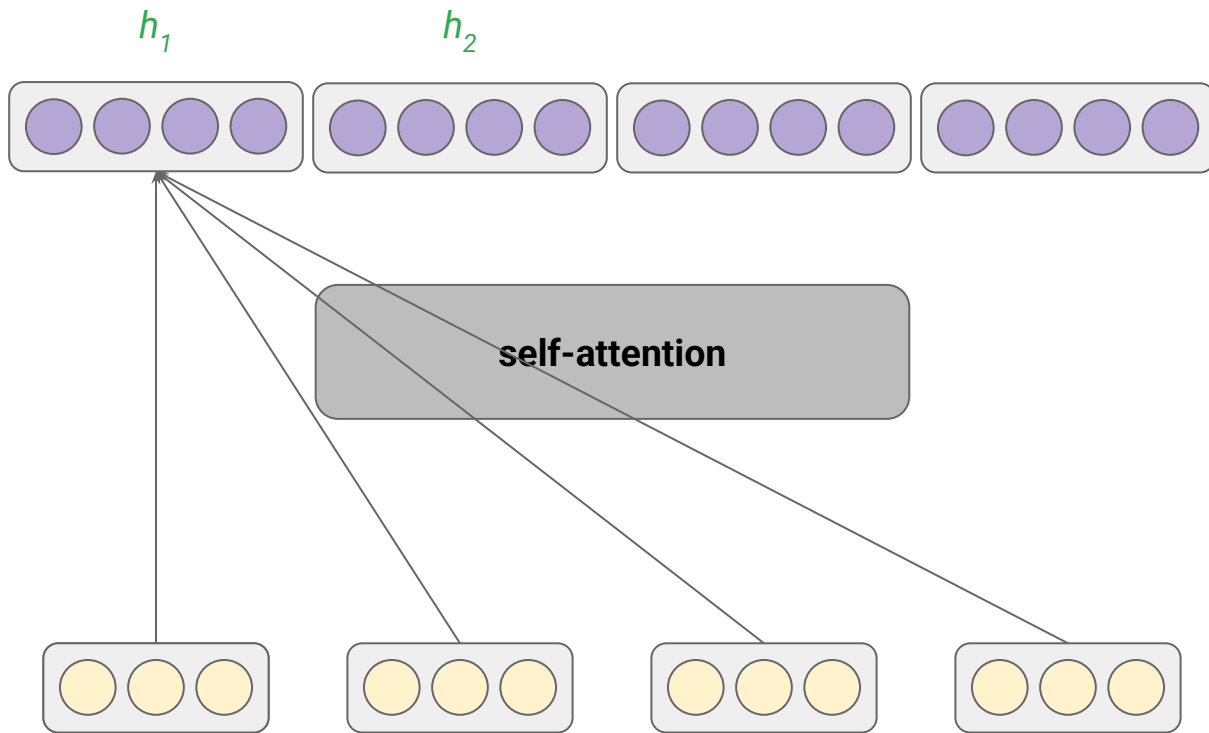
Different model architectures

- Encoder-only
 - BERT
- Encoder-decoder
 - T5
- Decoder-only
 - GPT

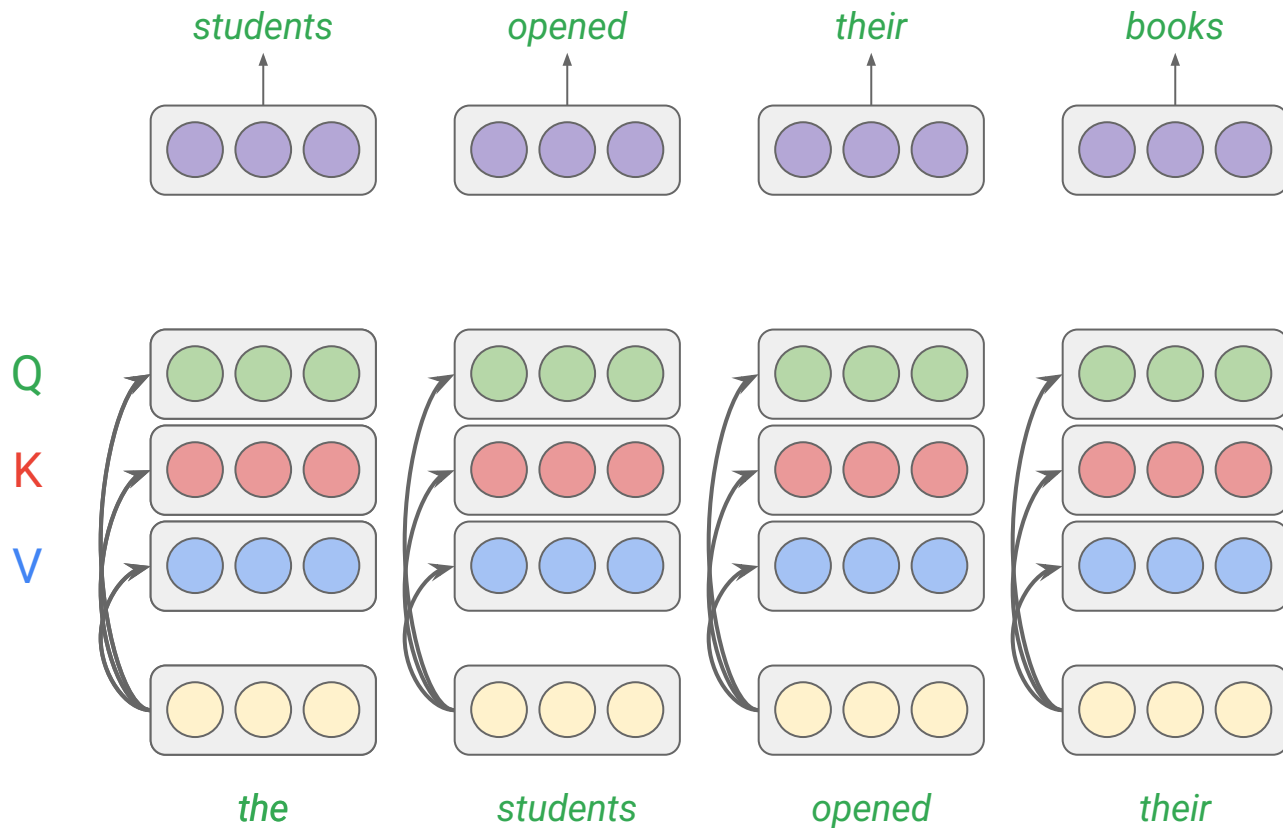
Attention mechanism



Self-attention



Self-attention (cont'd)



$$Q = X \cdot W_Q$$

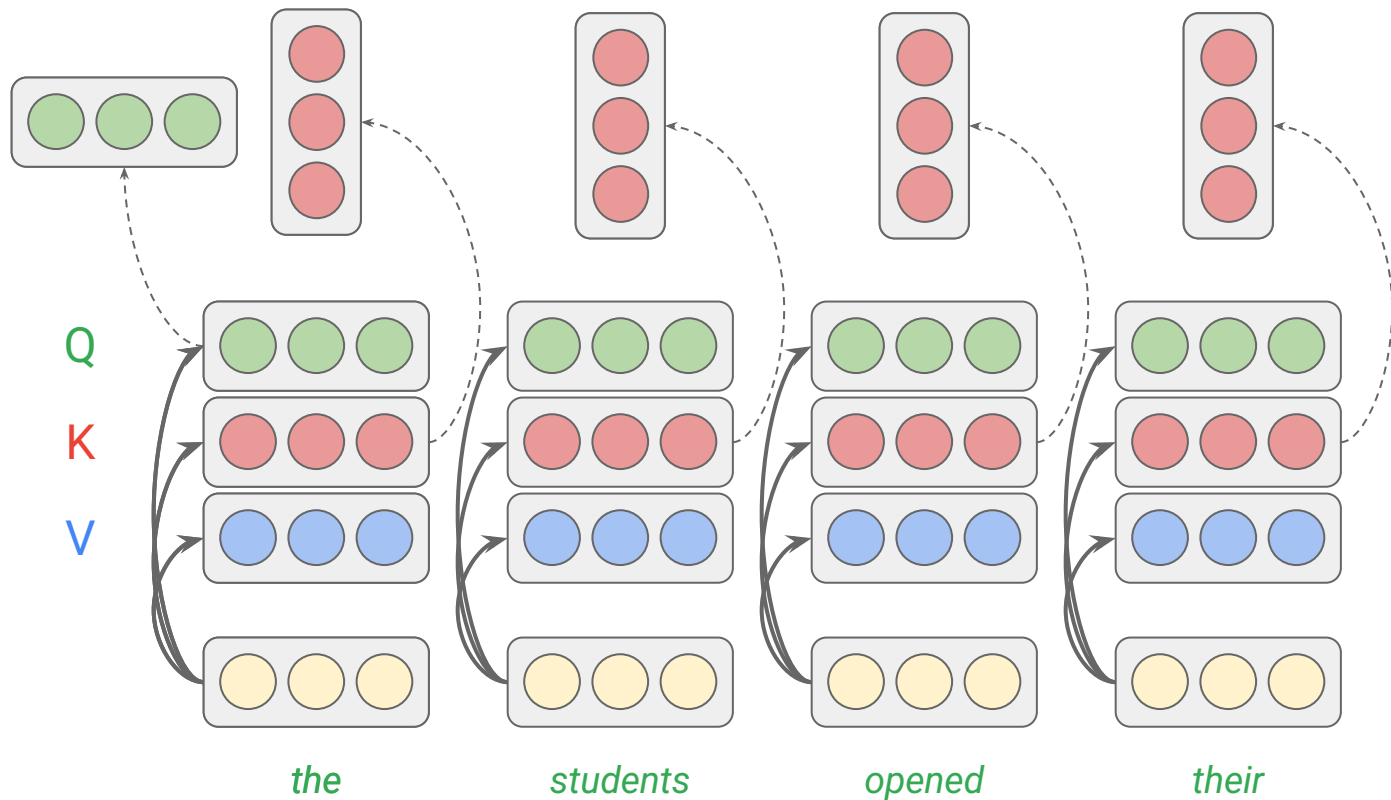
$$K = X \cdot W_K$$

$$V = X \cdot W_V$$

**linear
projections**

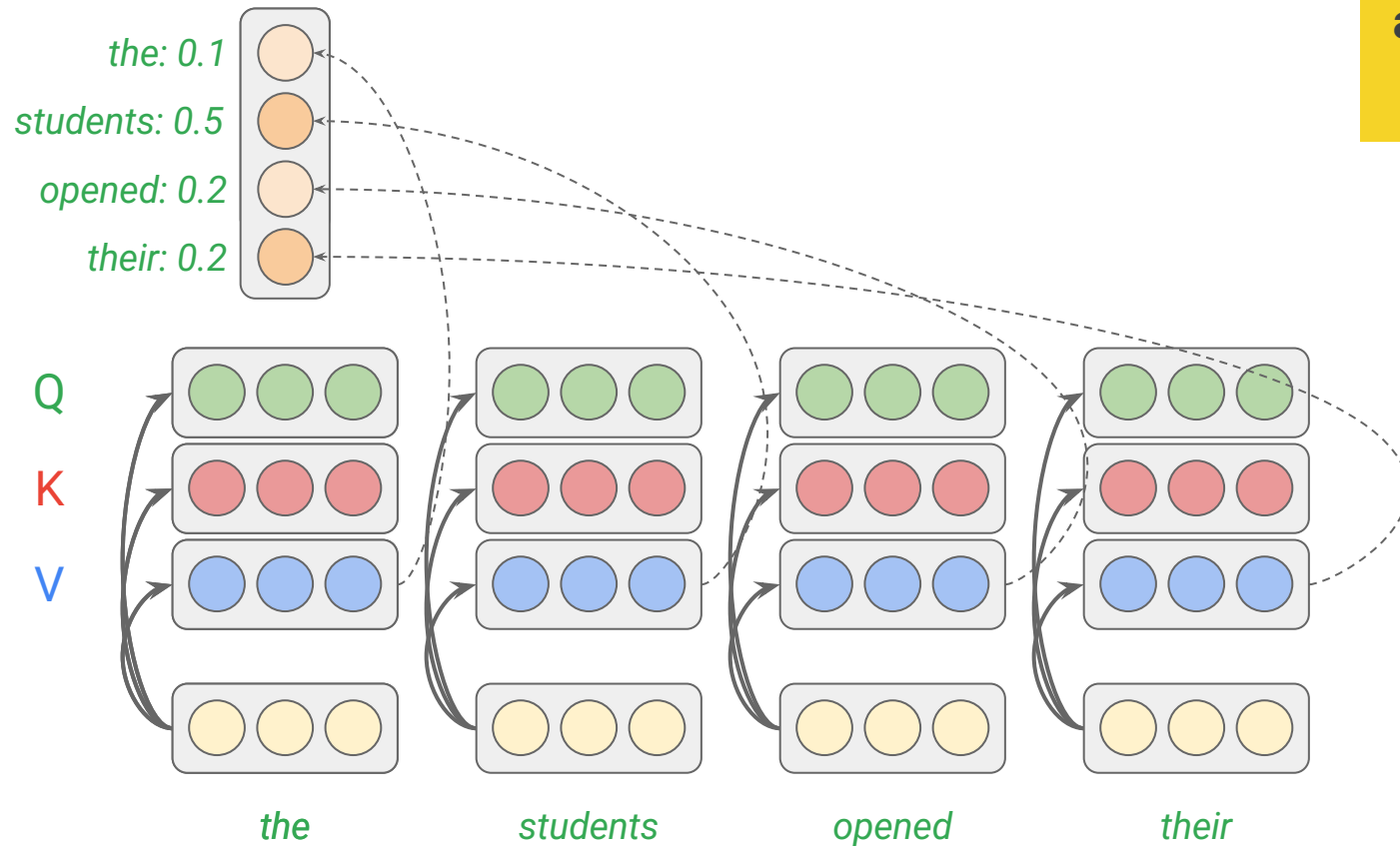
Self-attention (cont'd)

*all computations
are parallelized*



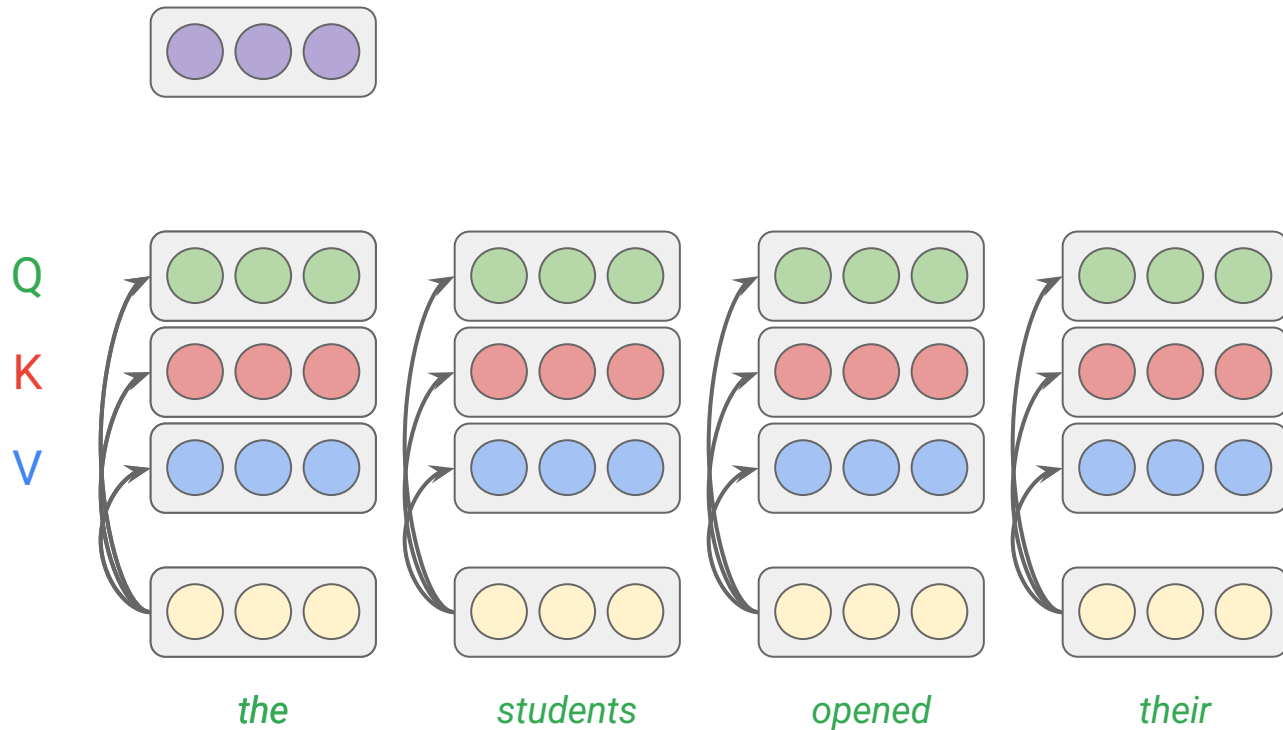
Self-attention (cont'd)

**all computations
are parallelized**



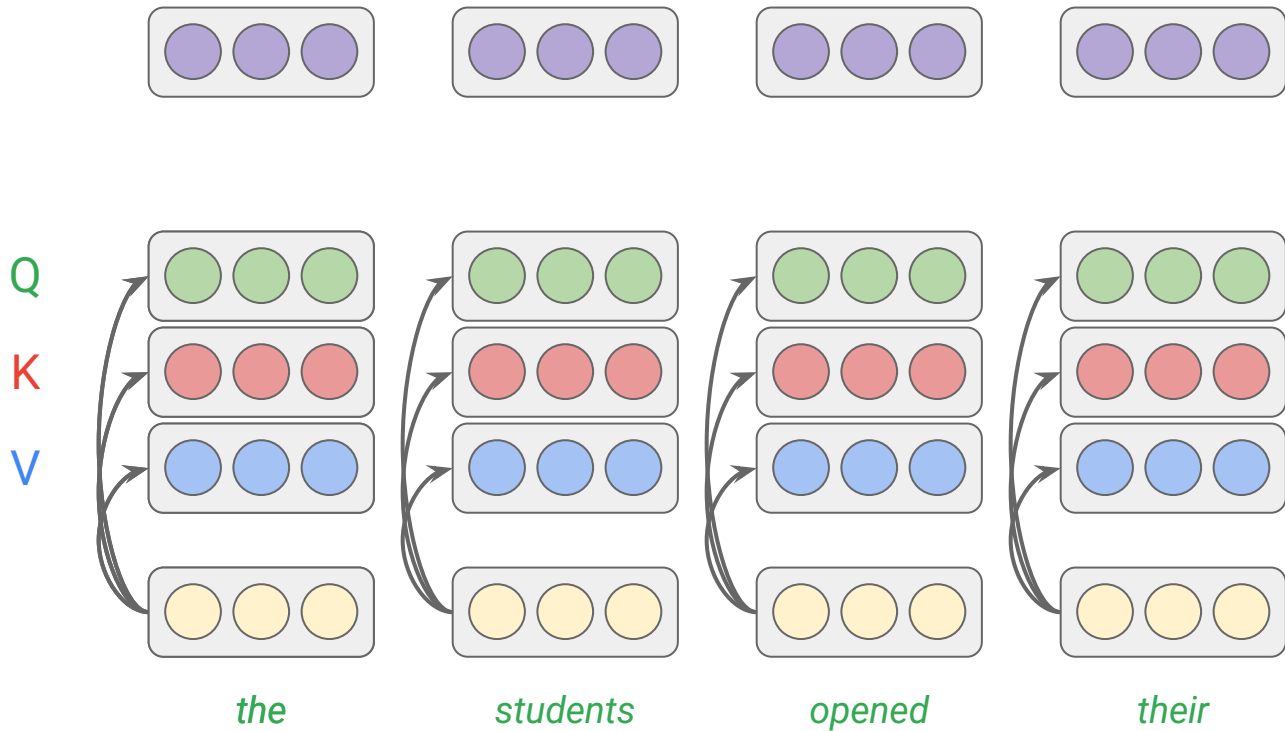
Self-attention (cont'd)

*all computations
are parallelized*



Self-attention (cont'd)

***all* computations are parallelized during training and sequential during inference**



All computations are parallelized

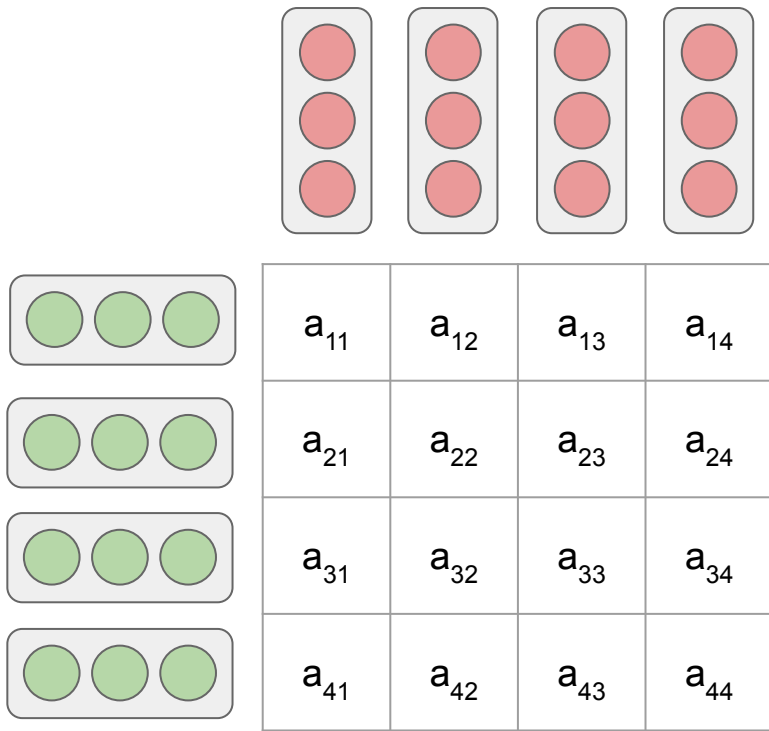
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



d_k : scaling factor

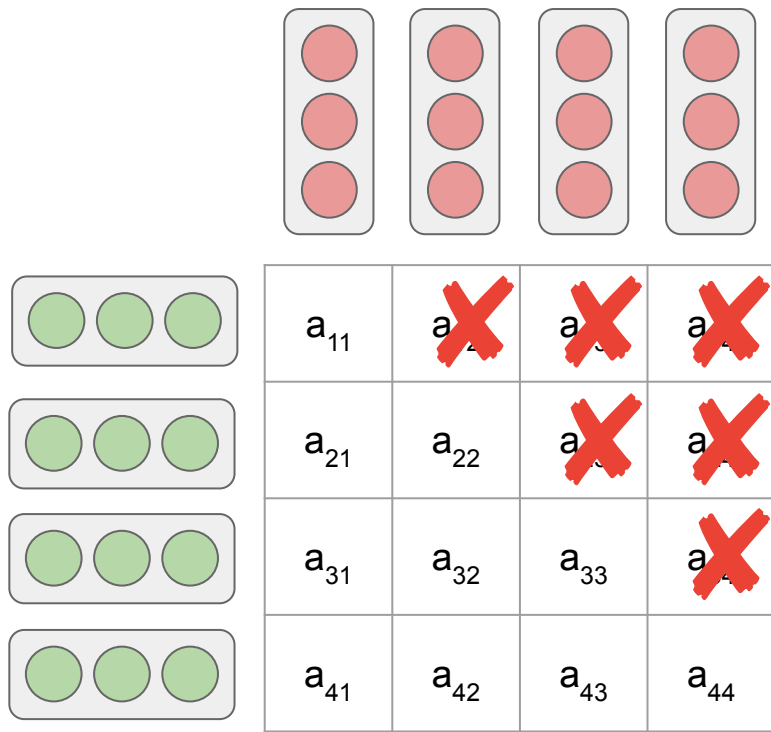
*large products push the softmax
function into regions where it
has extremely small gradients*

Quadratic complexity



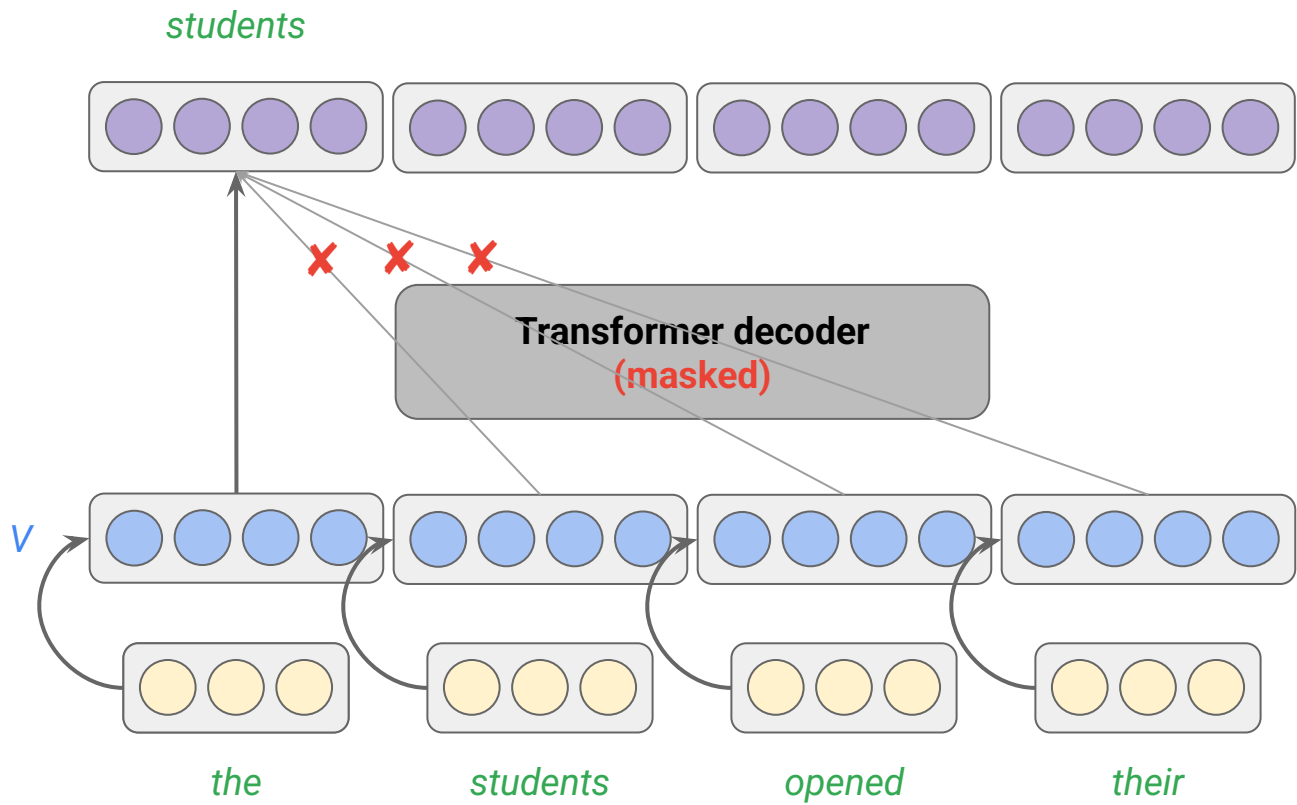
The time complexity of self-attention is quadratic in the input length $O(n^2)$

Self-attention in the decoder

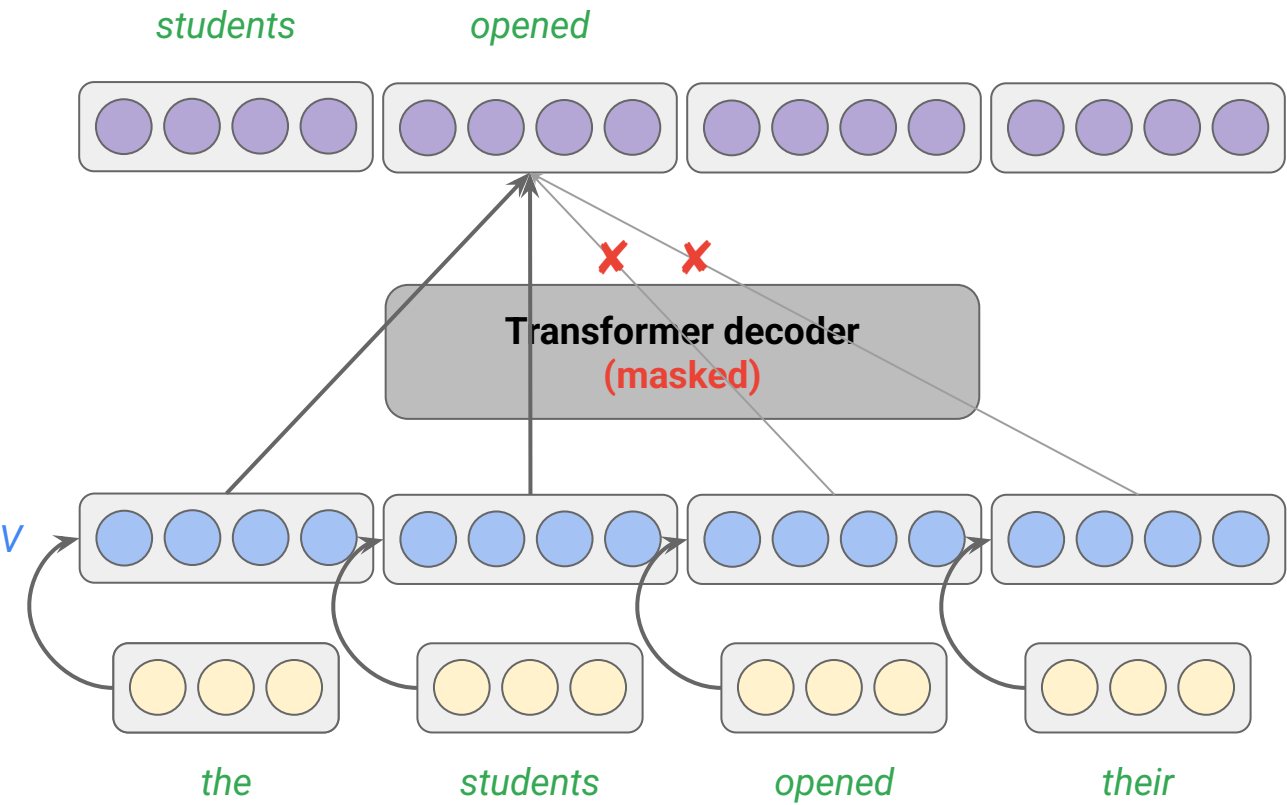


masking out all values in the input of the softmax which correspond to illegal connections

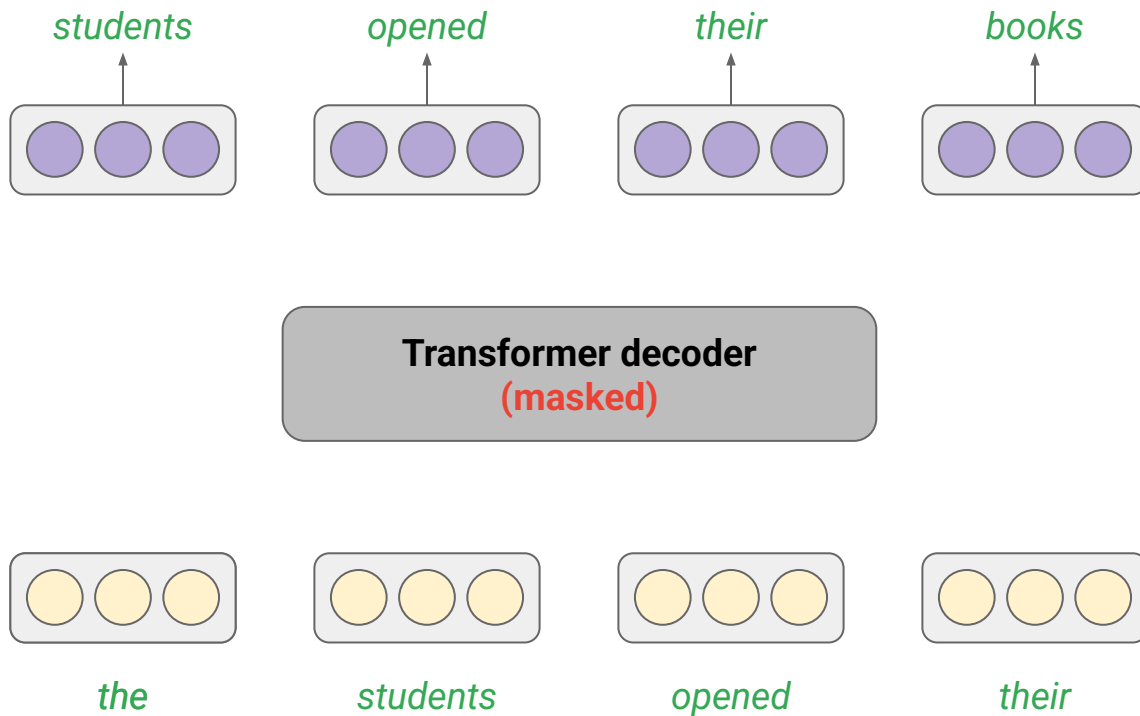
Transformer decoder



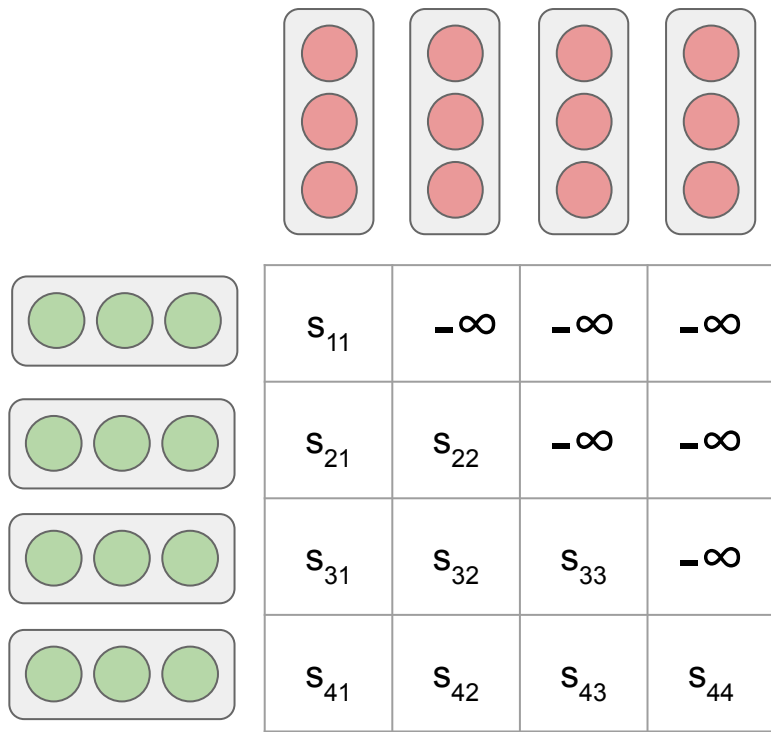
Transformer decoder (cont'd)



Transformer decoder (cont'd)

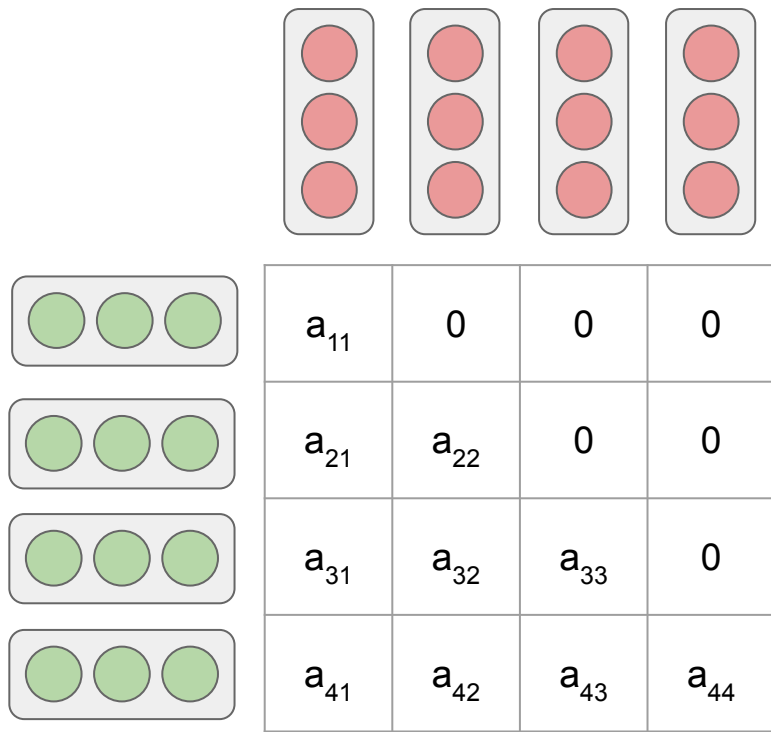


Self-attention in the decoder



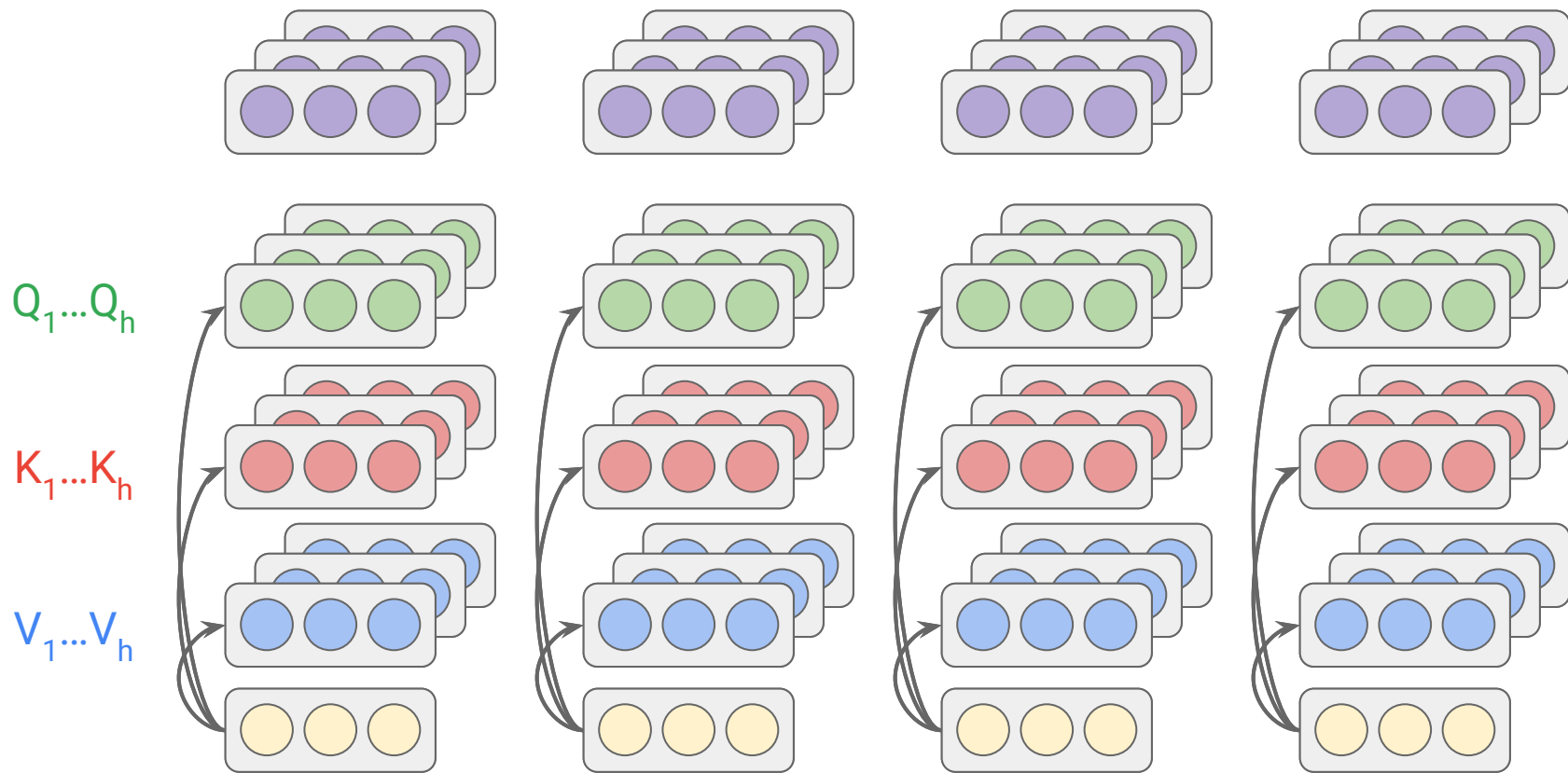
masking out (setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections

Self-attention in the decoder (cont'd)

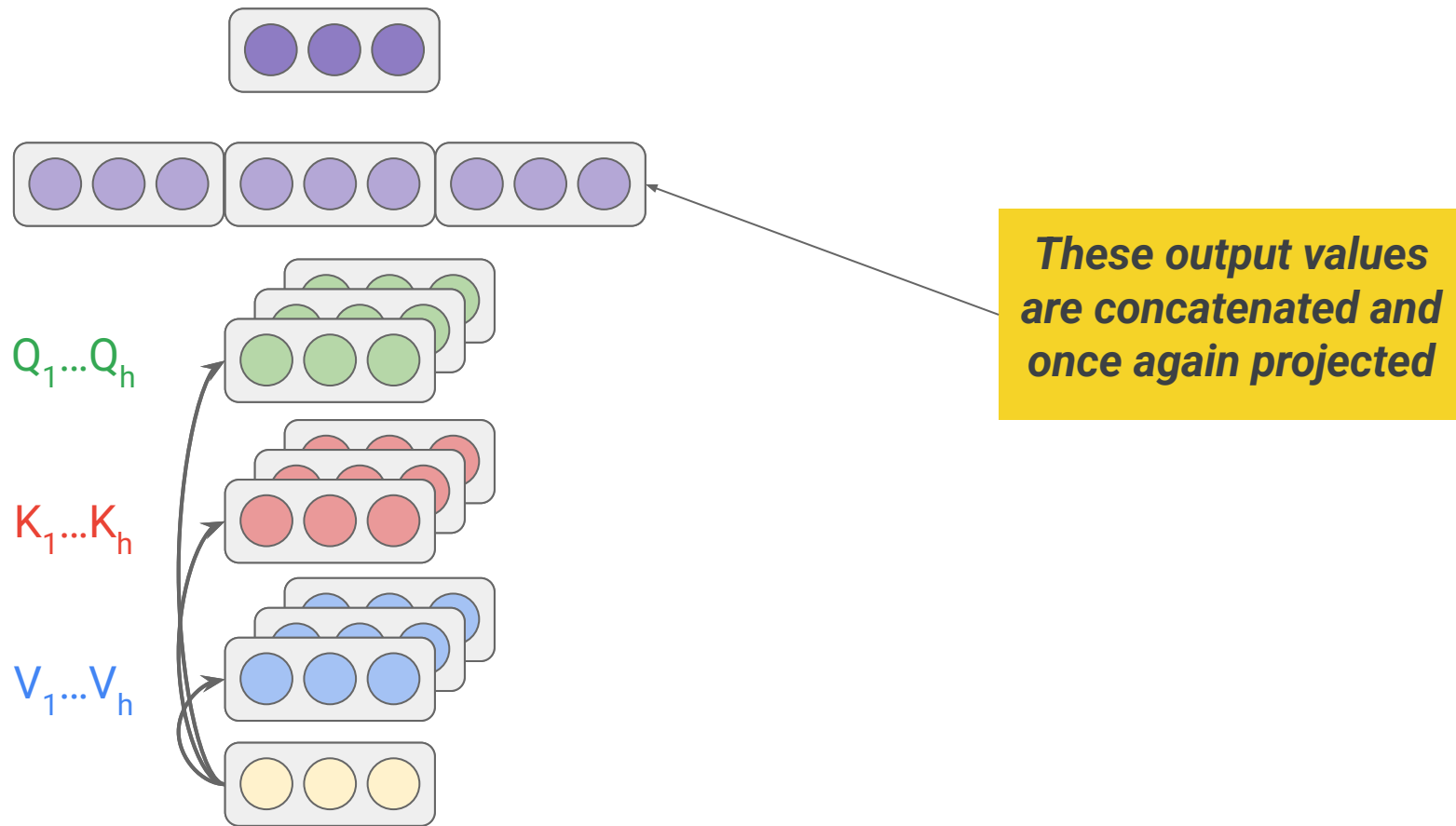


masking out all values in the input of the softmax which correspond to illegal connections

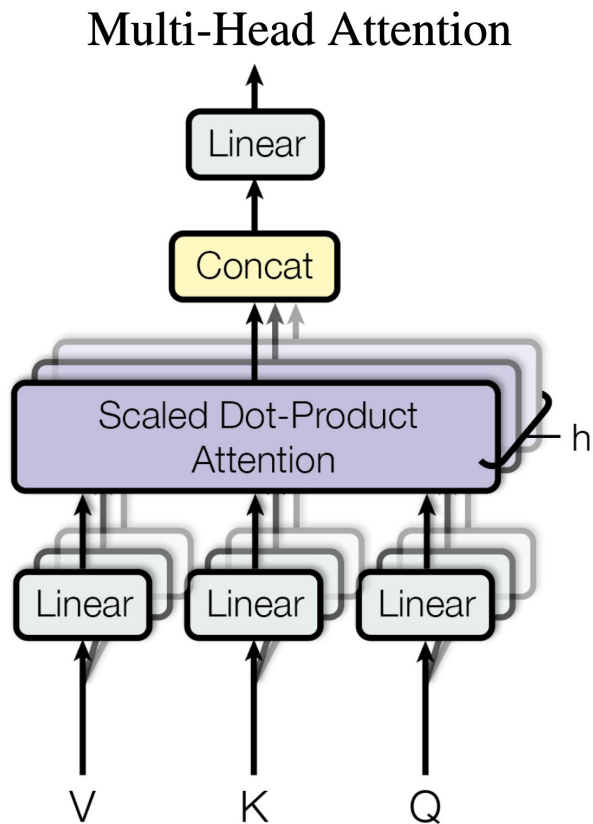
Multi-head attention



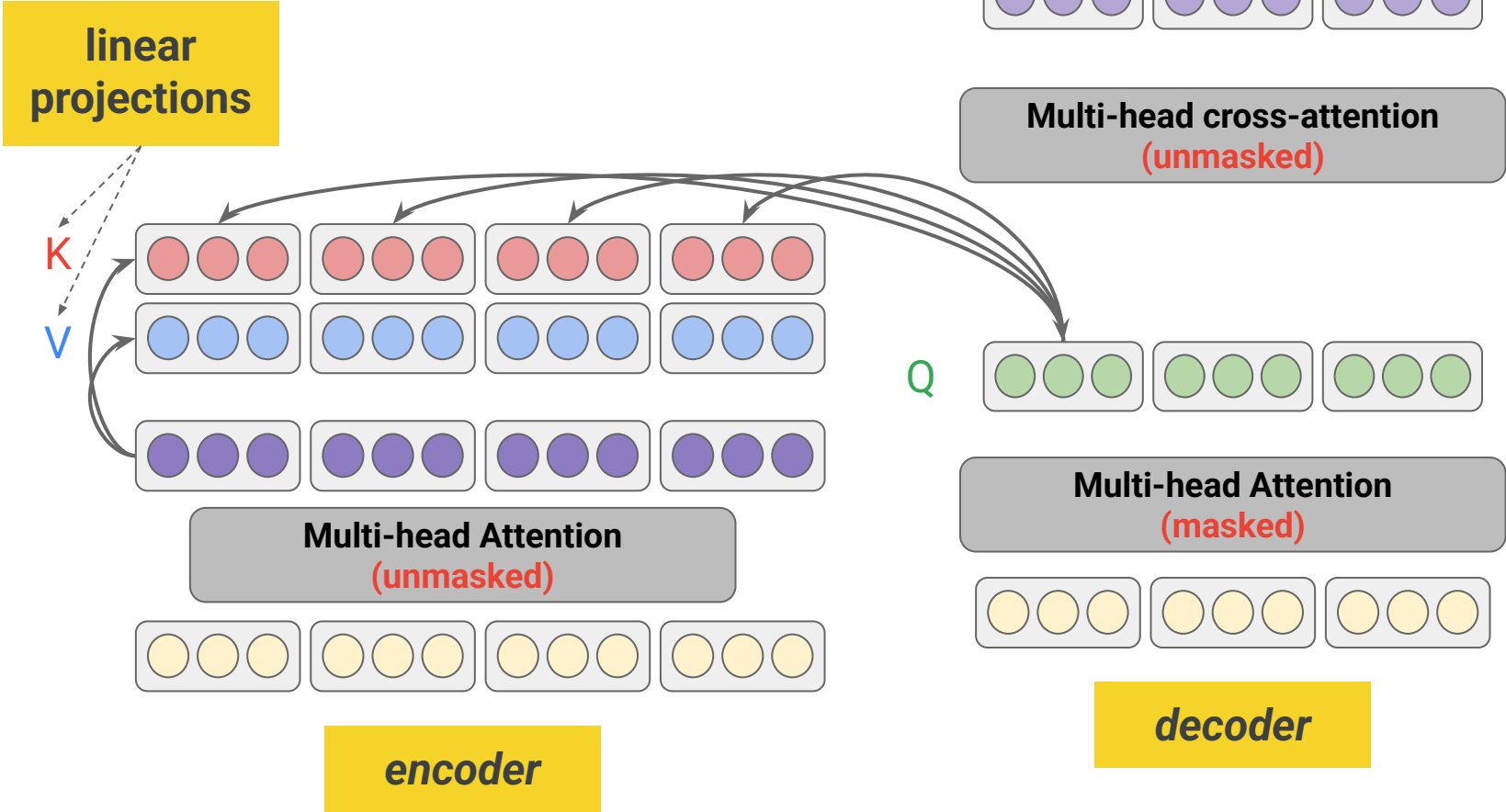
Multi-head attention (cont'd)



Multi-head attention (cont'd)



Cross-attention in the decoder



Cross-attention in the decoder (cont'd)

linear
projections

K

V

Multi-head Attention
(unmasked)

encoder

Multi-head cross-attention
(unmasked)

Q

Multi-head Attention
(masked)

decoder

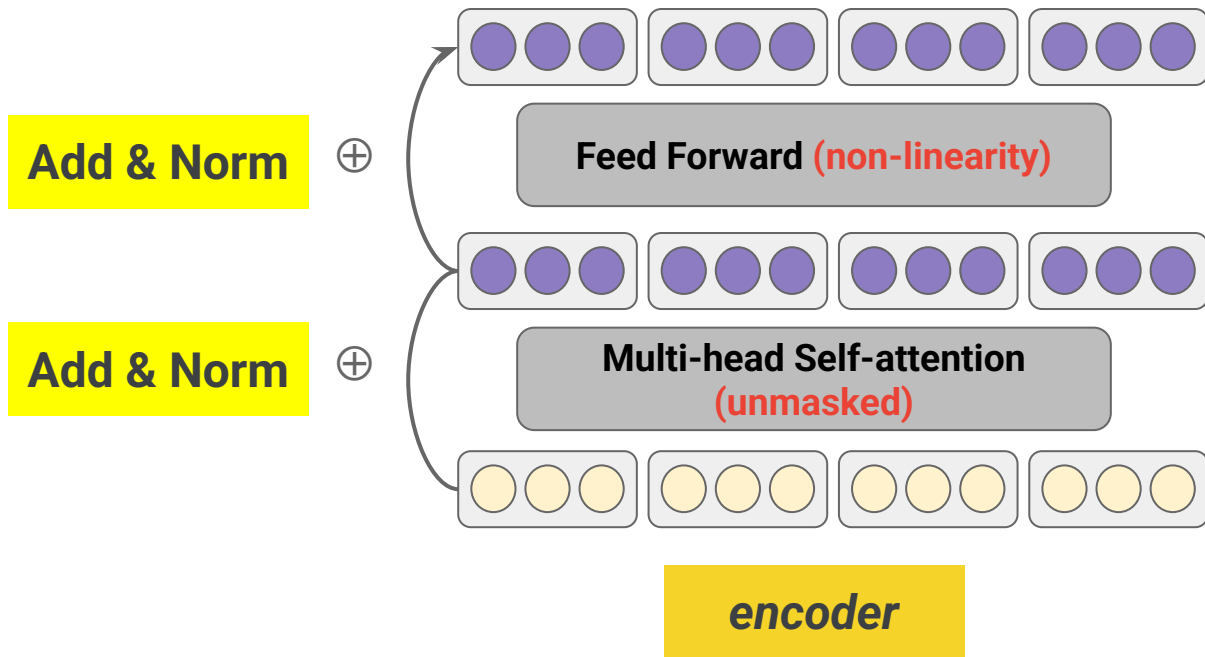
residual
connections

residual
connections

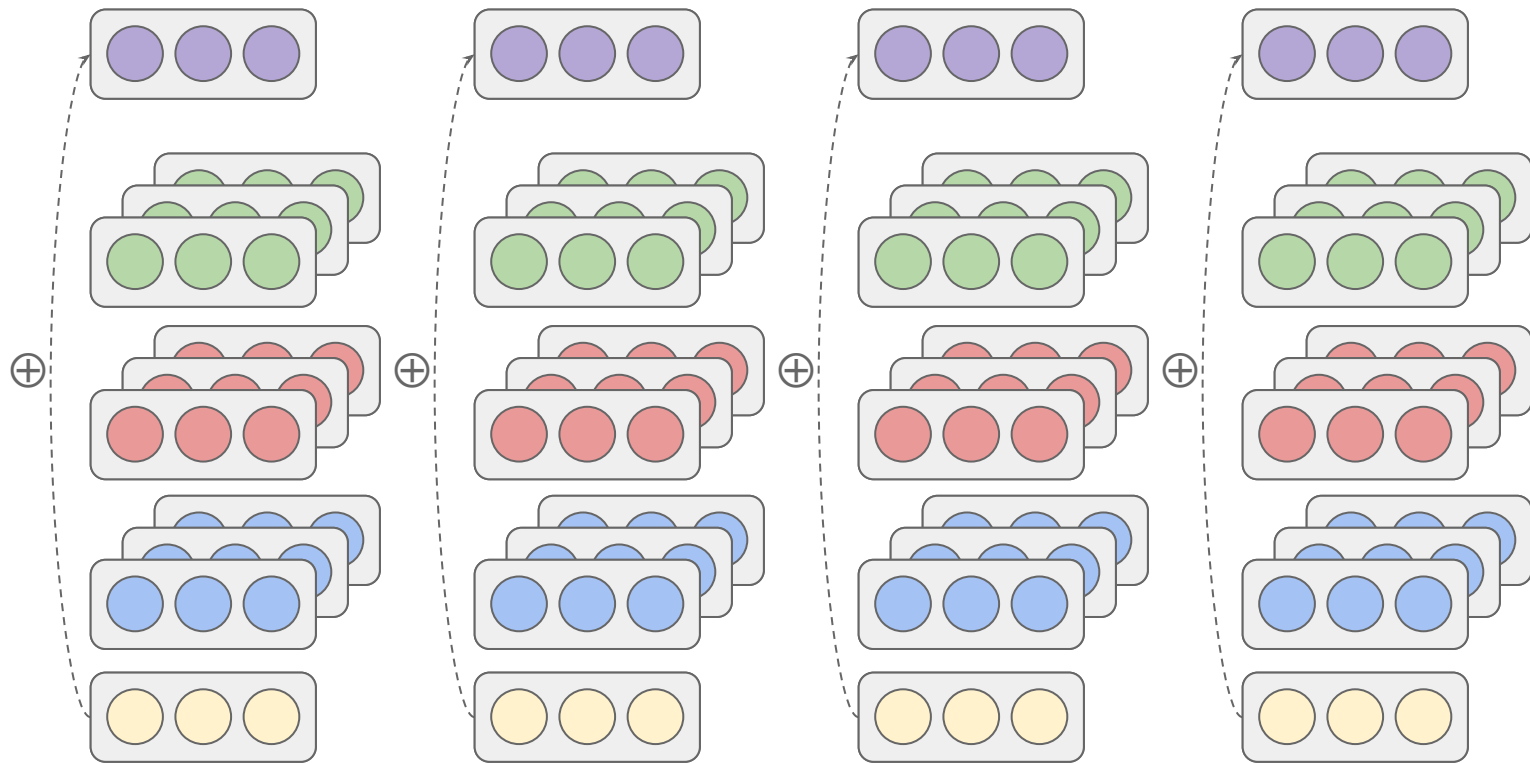
⊕

⊕

Encoder (one layer)



Residual connection



Residual connection

$$\text{output} = \text{sublayer}(x) + x$$

Layer normalization

$$\text{Norm}(z) = \frac{z - \mu}{\sigma} \cdot \gamma + \beta$$

where μ and σ are the mean and standard deviation of the activations, and γ, β are learnable parameters.

Each activation vector is normalized so that its components have mean 0 and variance 1. This prevents activations from becoming too large or too small as they propagate through the network.

Residual connection and layer normalization

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

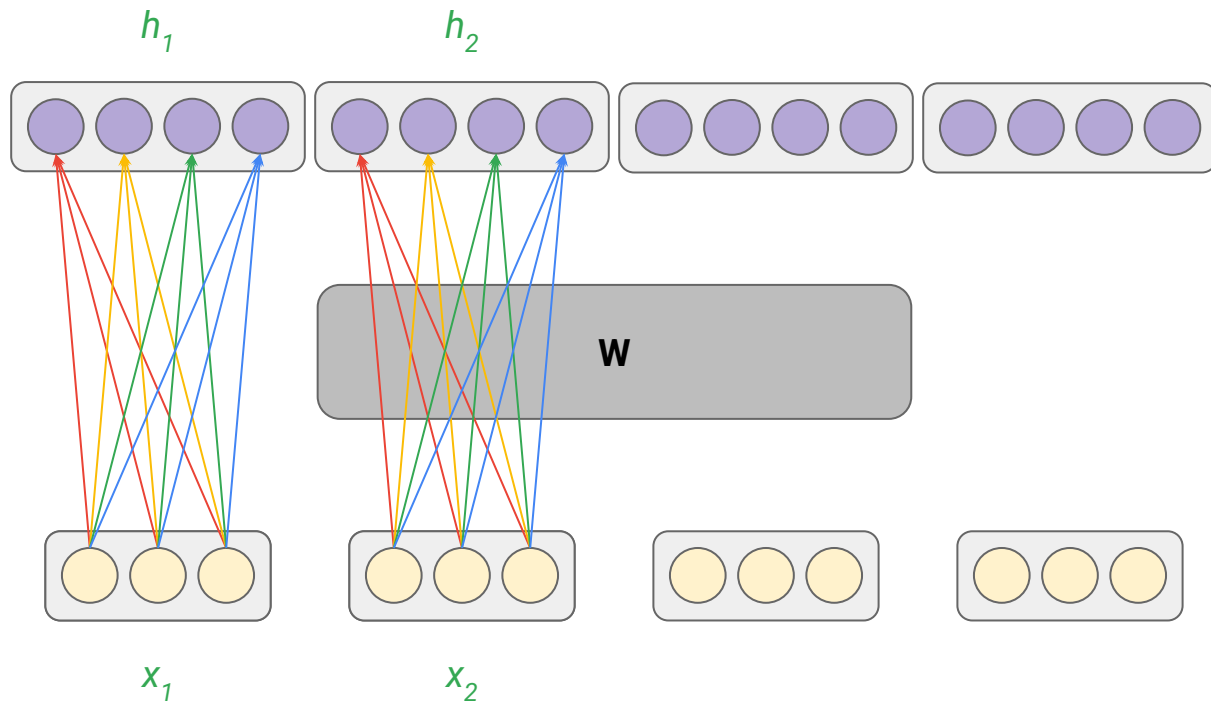
Position-wise Feed-Forward Networks

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

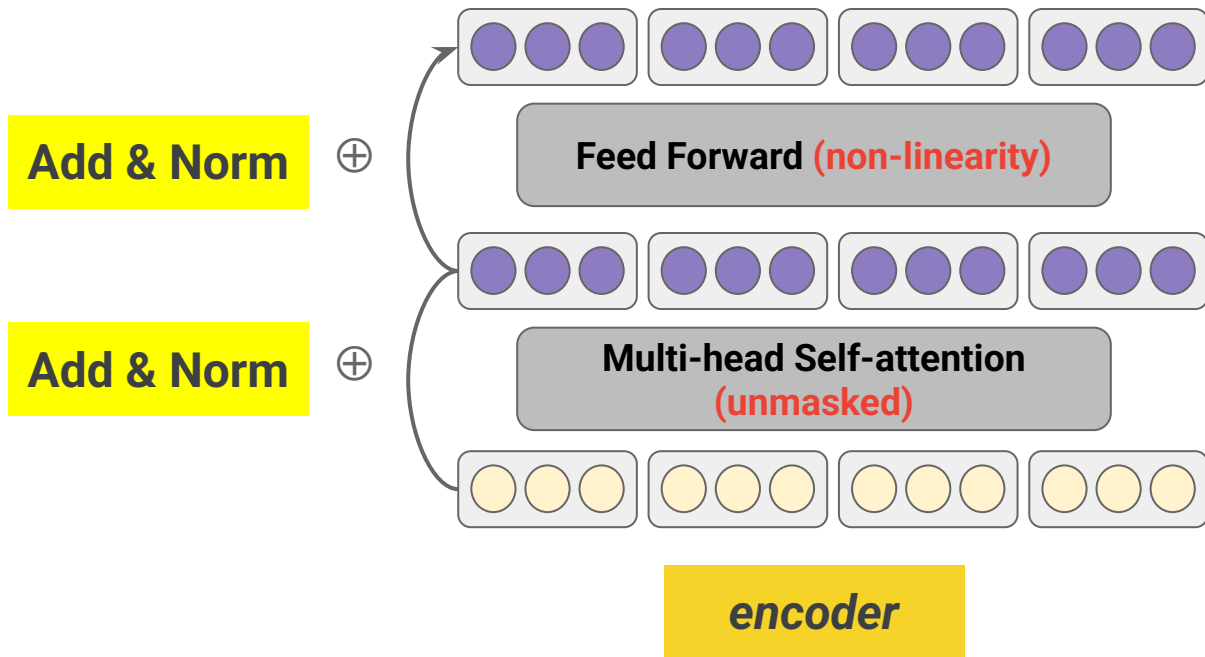


**ReLU (Rectified
Linear Unit)**

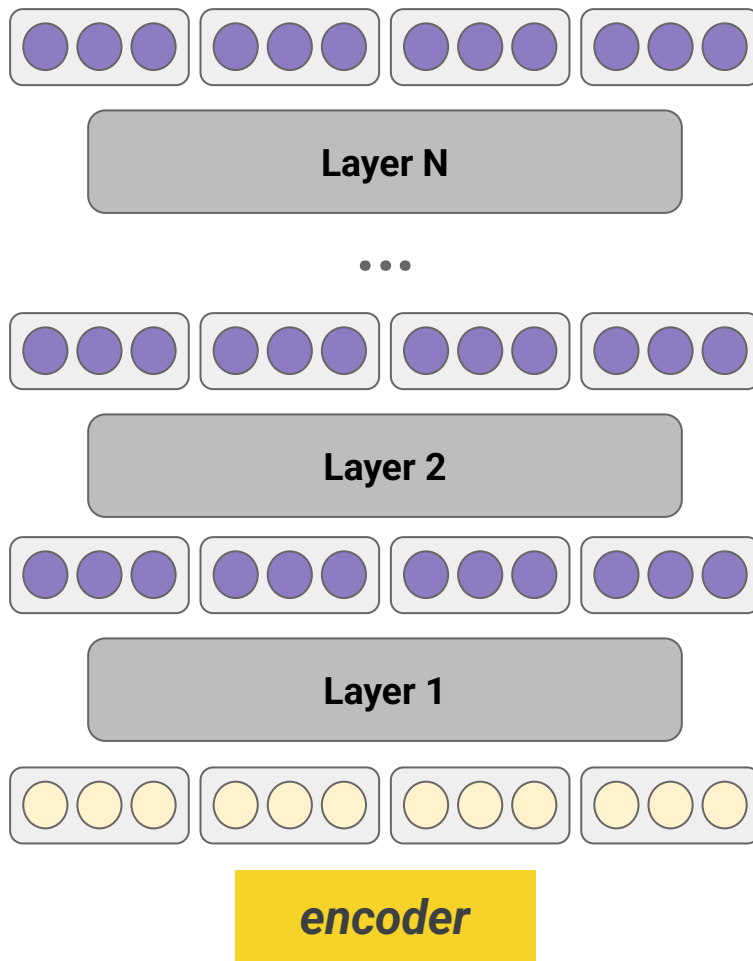
Position-wise Feed-Forward Networks (cont'd)



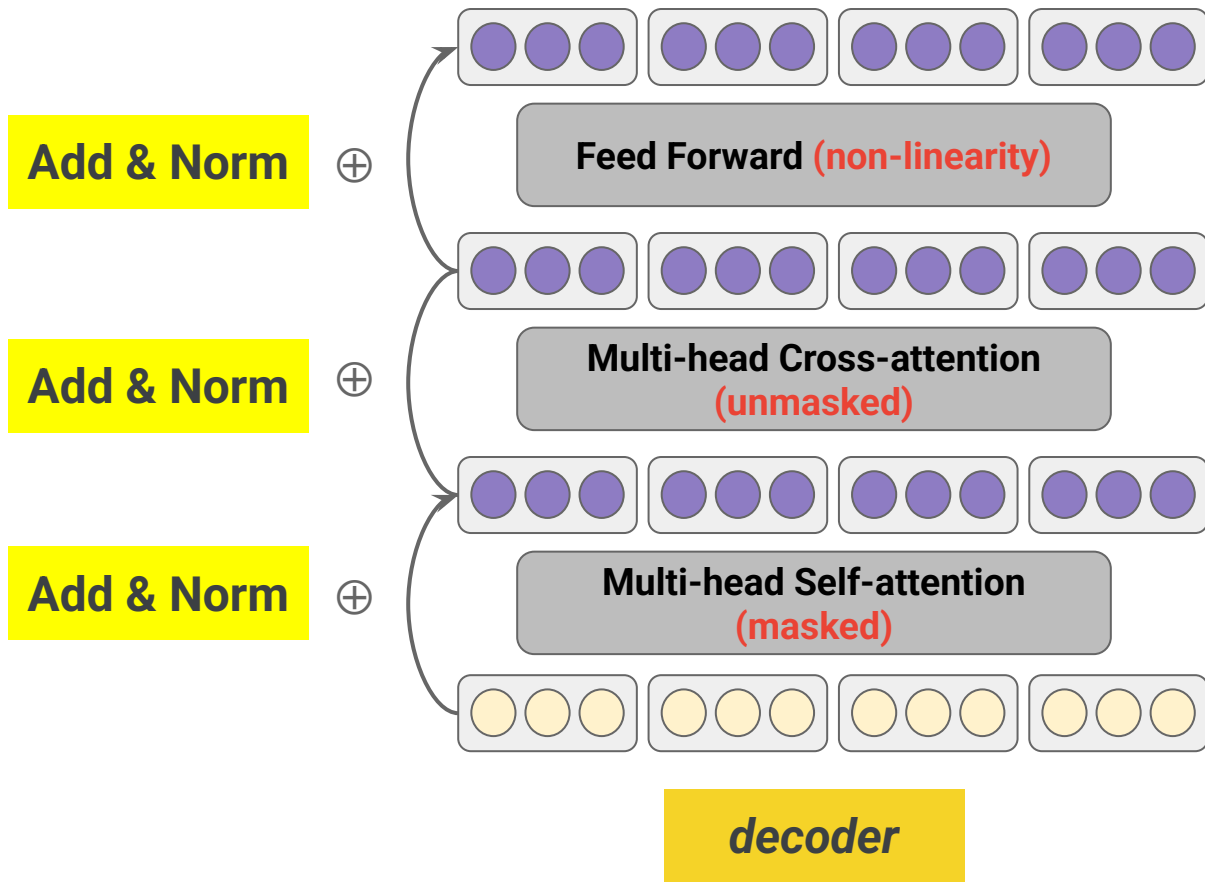
Encoder (one layer)



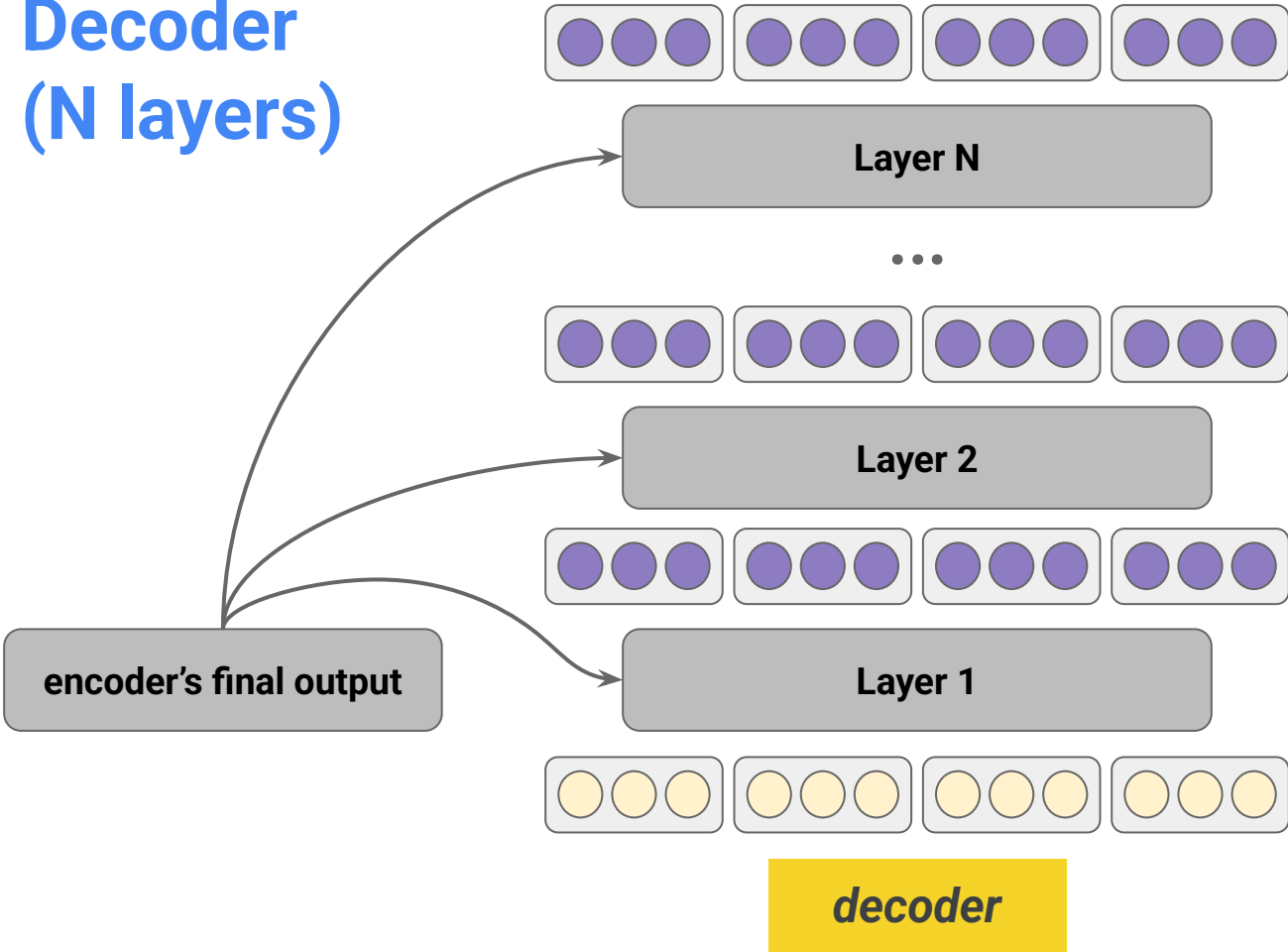
Encoder (N layers)



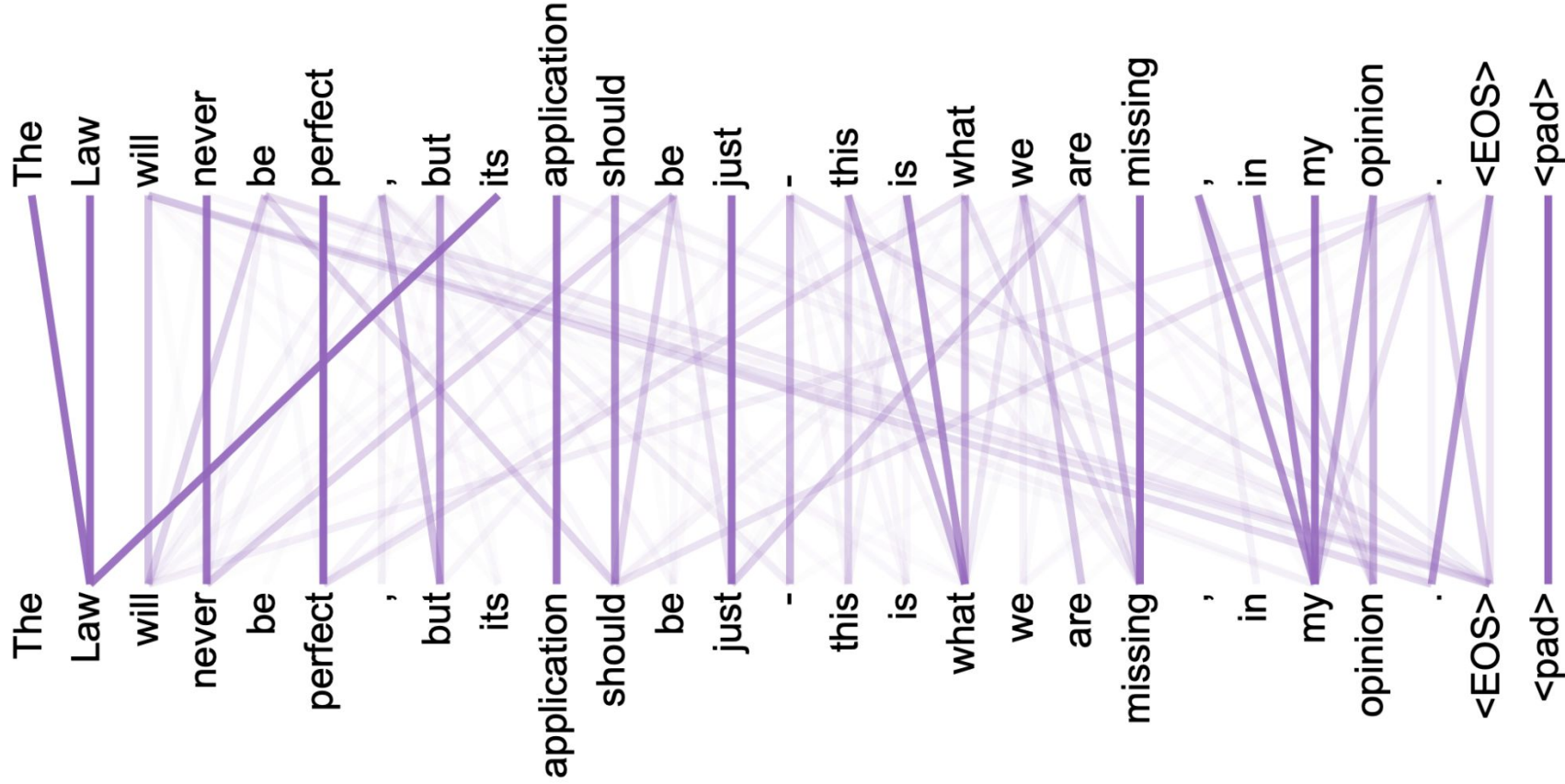
Decoder (one layer)



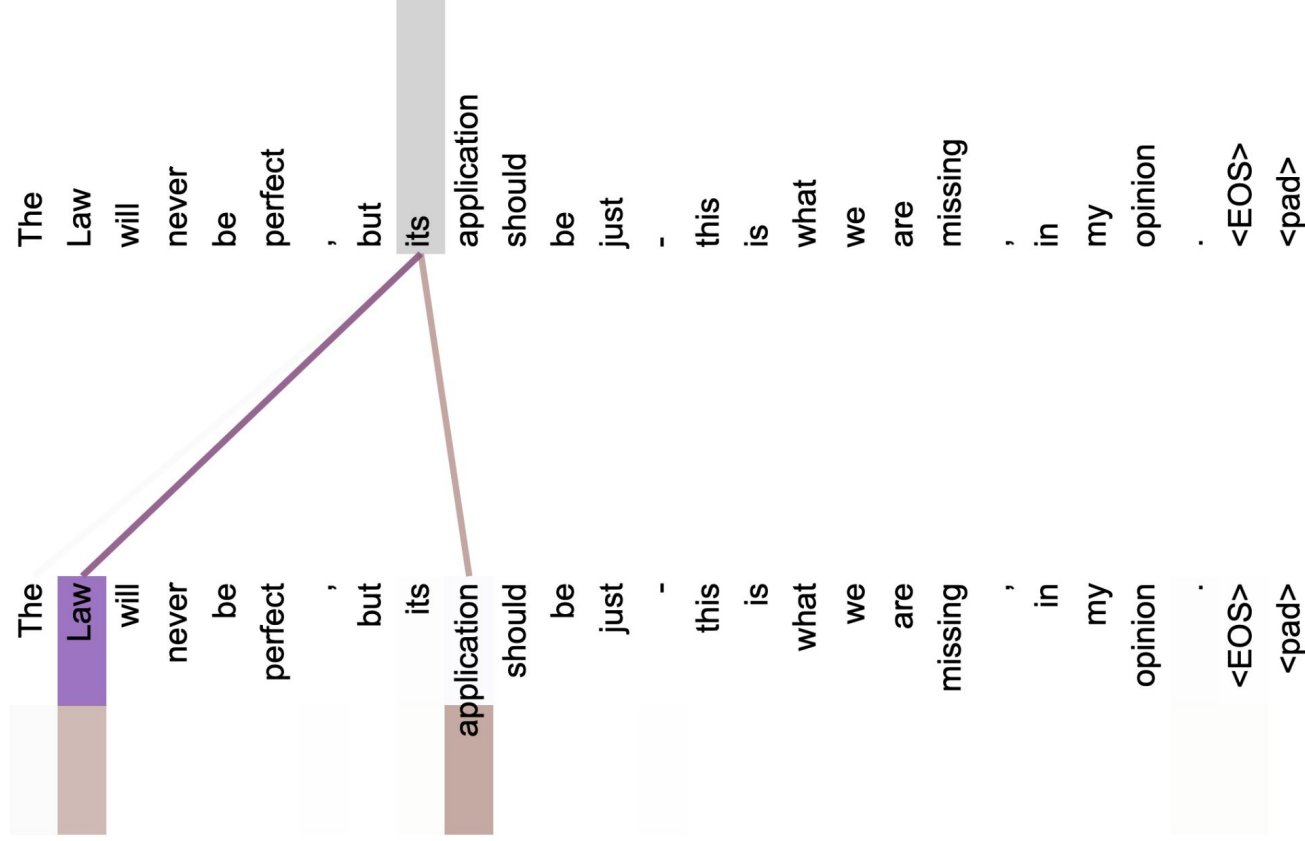
Decoder
(N layers)



Attention visualizations



Attention visualizations (cont'd)

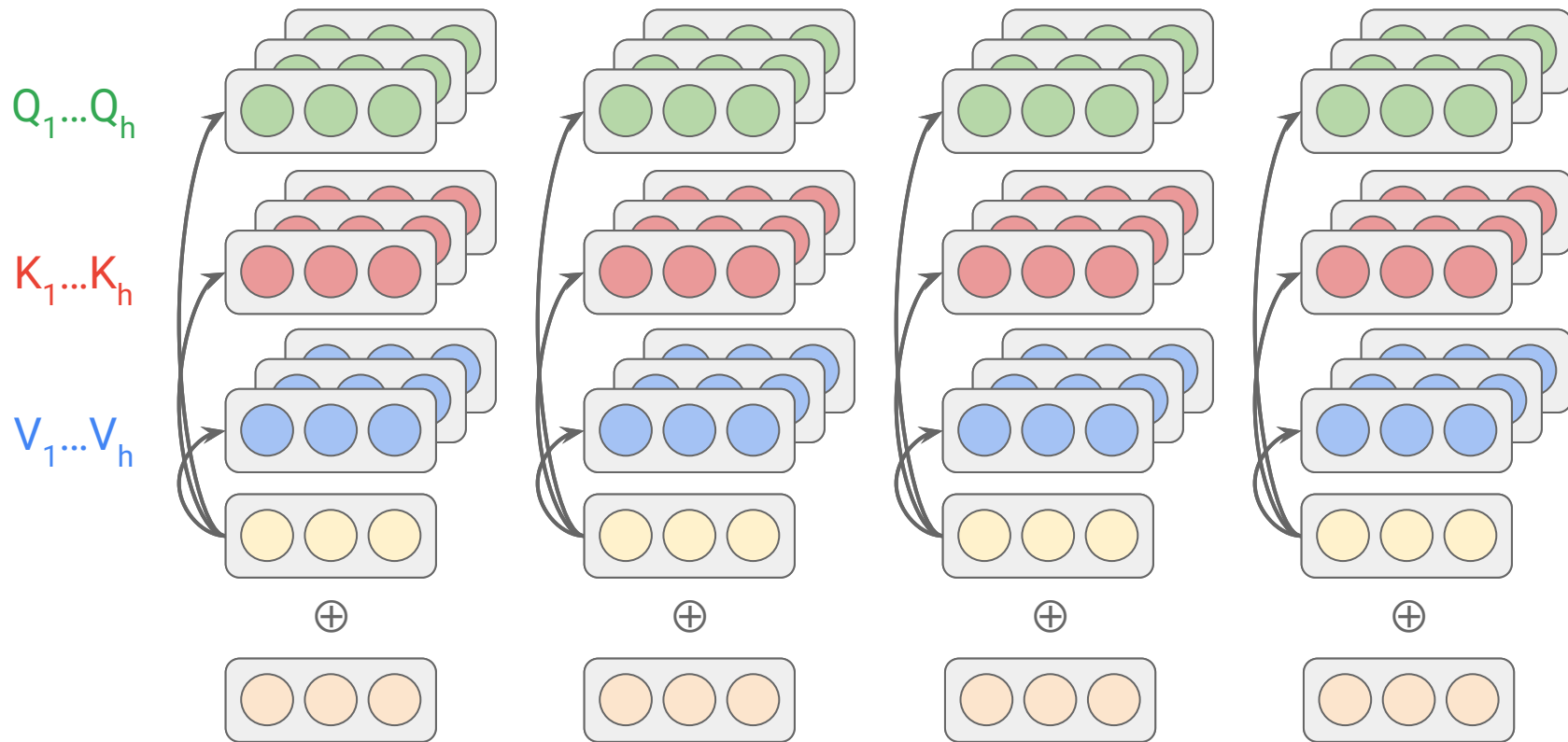


Sinusoidal positional encoding

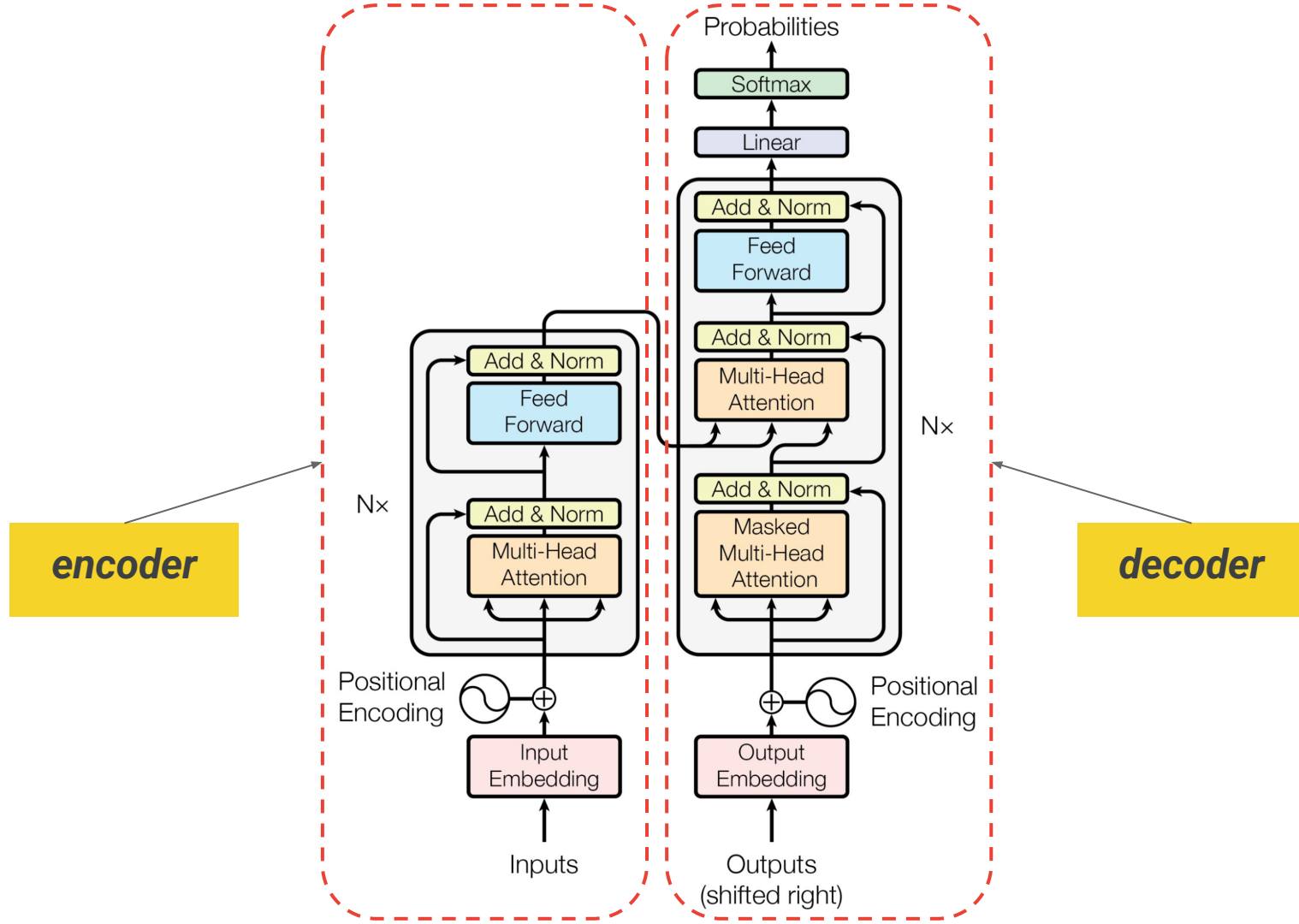
$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i / d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i / d_{\text{model}}})$$

Positional Encoding (cont'd)



Transformer block (putting it together)



Thank you!