

Attention mechanisms & Transformers

CS 5624: Natural Language Processing

Spring 2025

<https://tuvllms.github.io/nlp-spring-2025>

Tu Vu



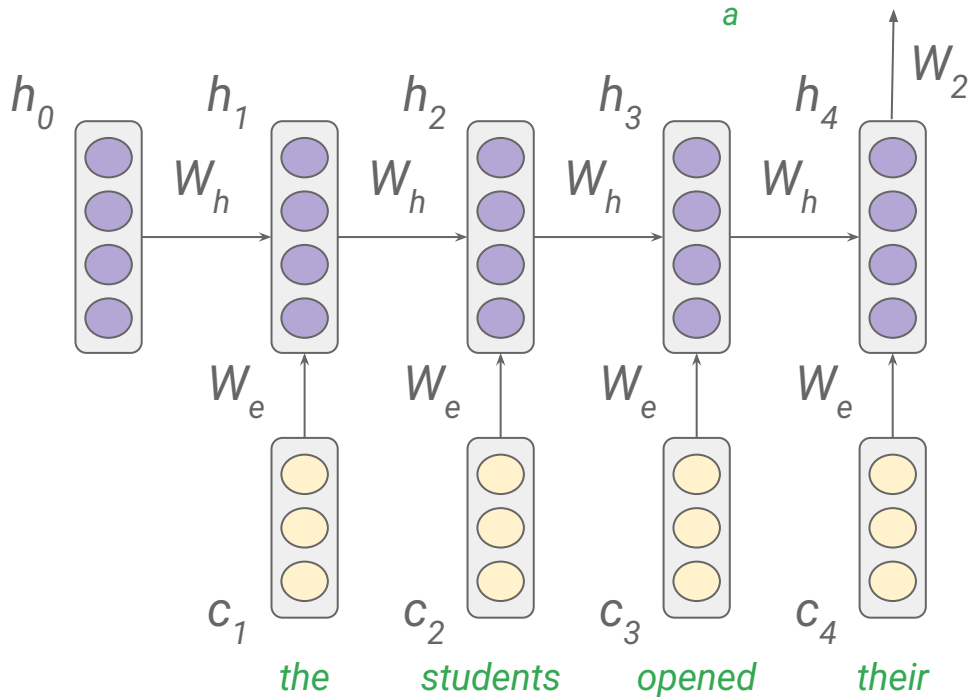
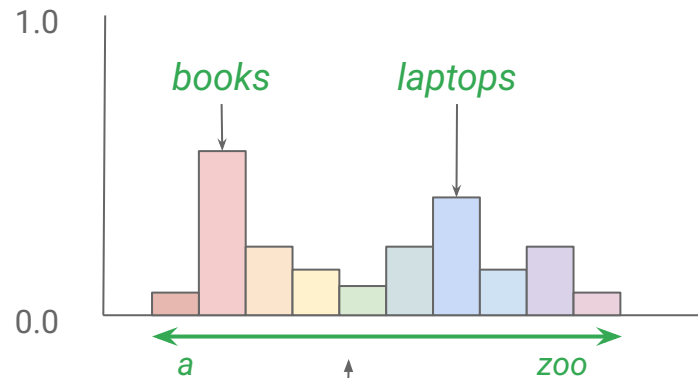
Logistics

- Homework 1 & Quiz 1 are on their way
- Final project proposal due on February 28

Recurrent neural networks (RNNs)

hidden states

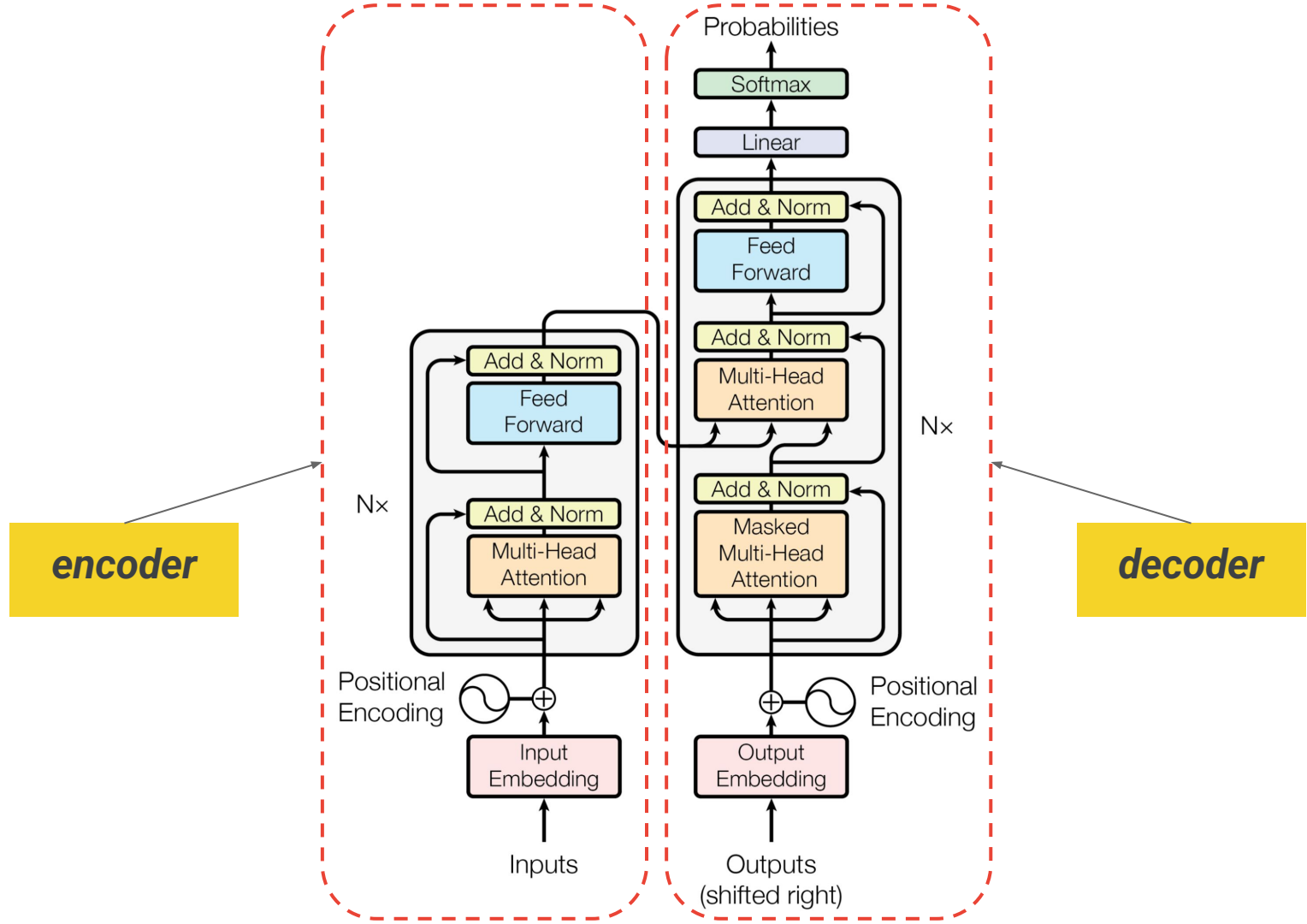
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c^t)$$

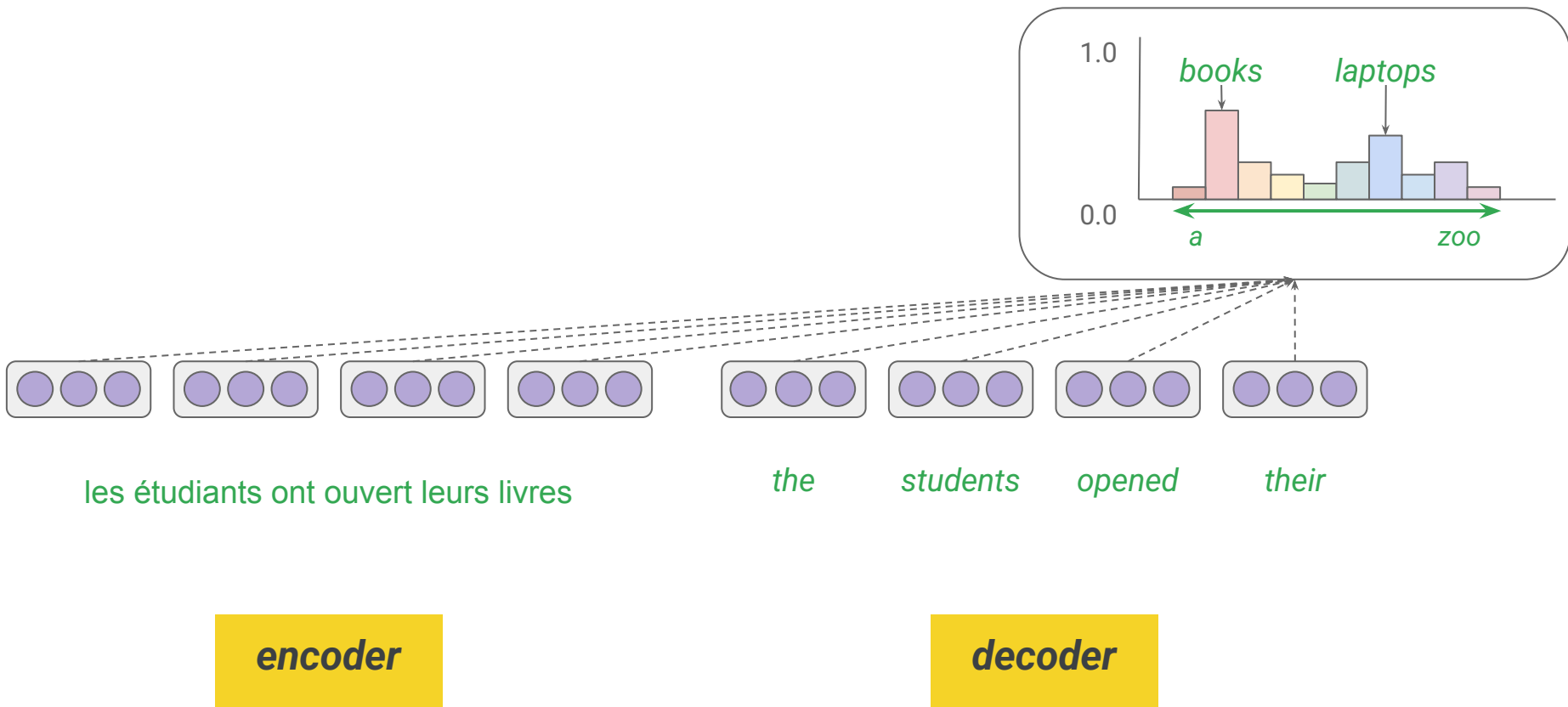


output distribution

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

Encoder-decoder architecture

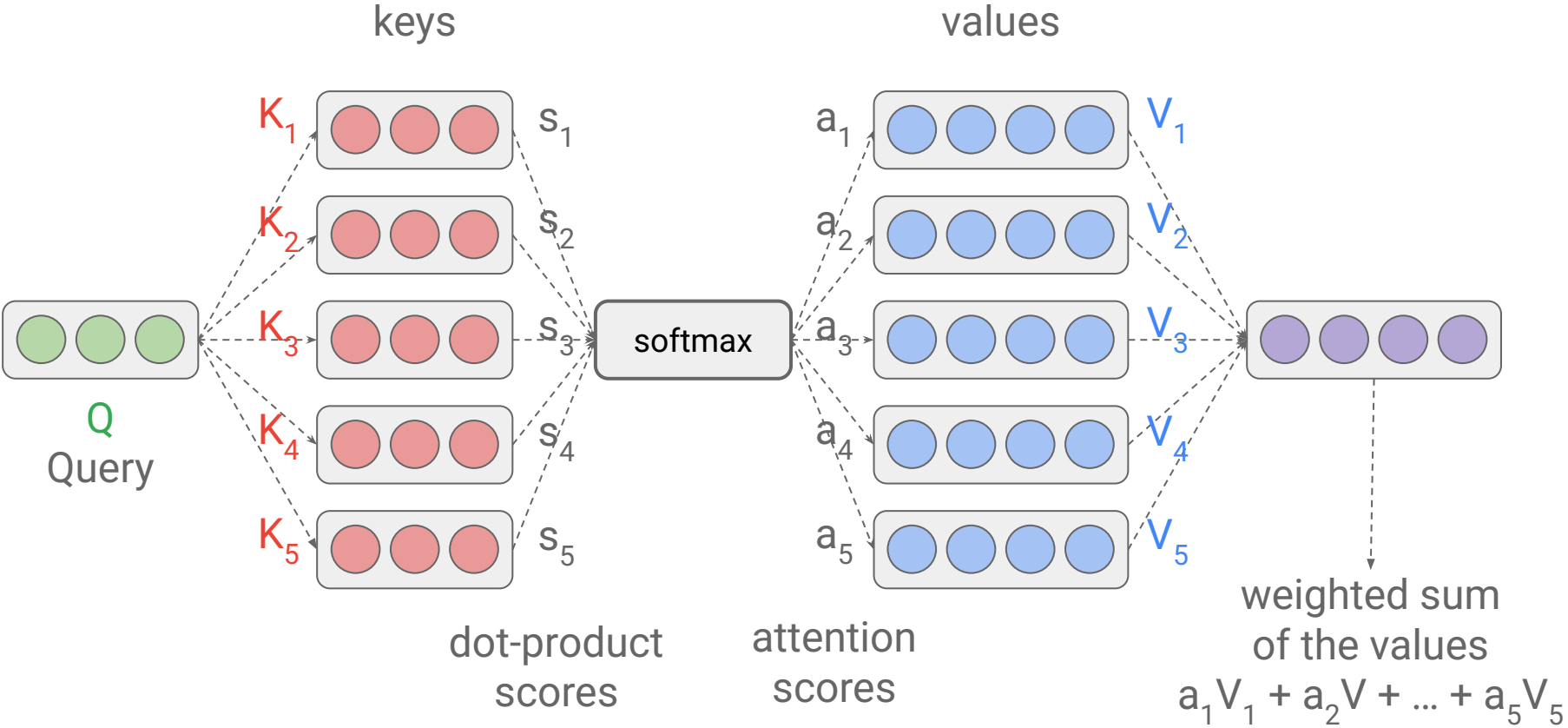




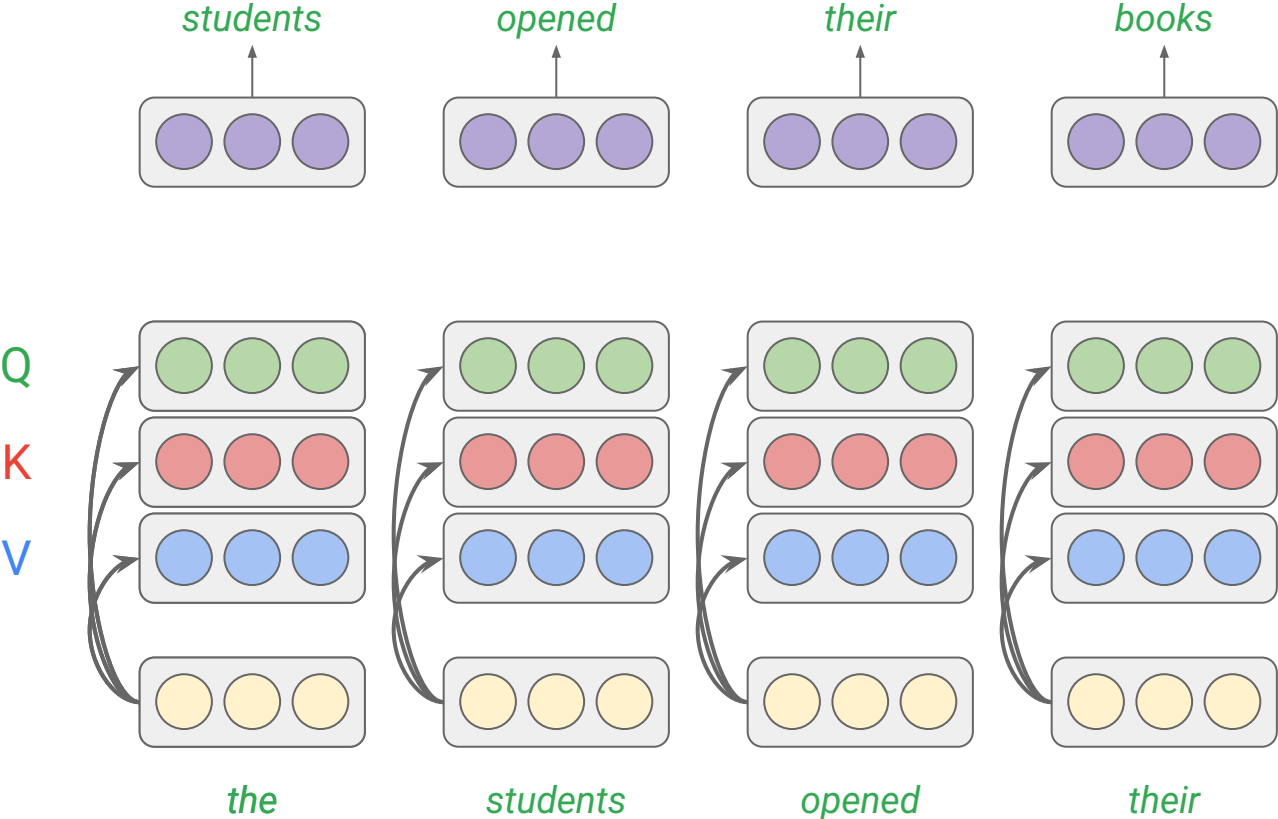
Different model architectures

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

Attention mechanism

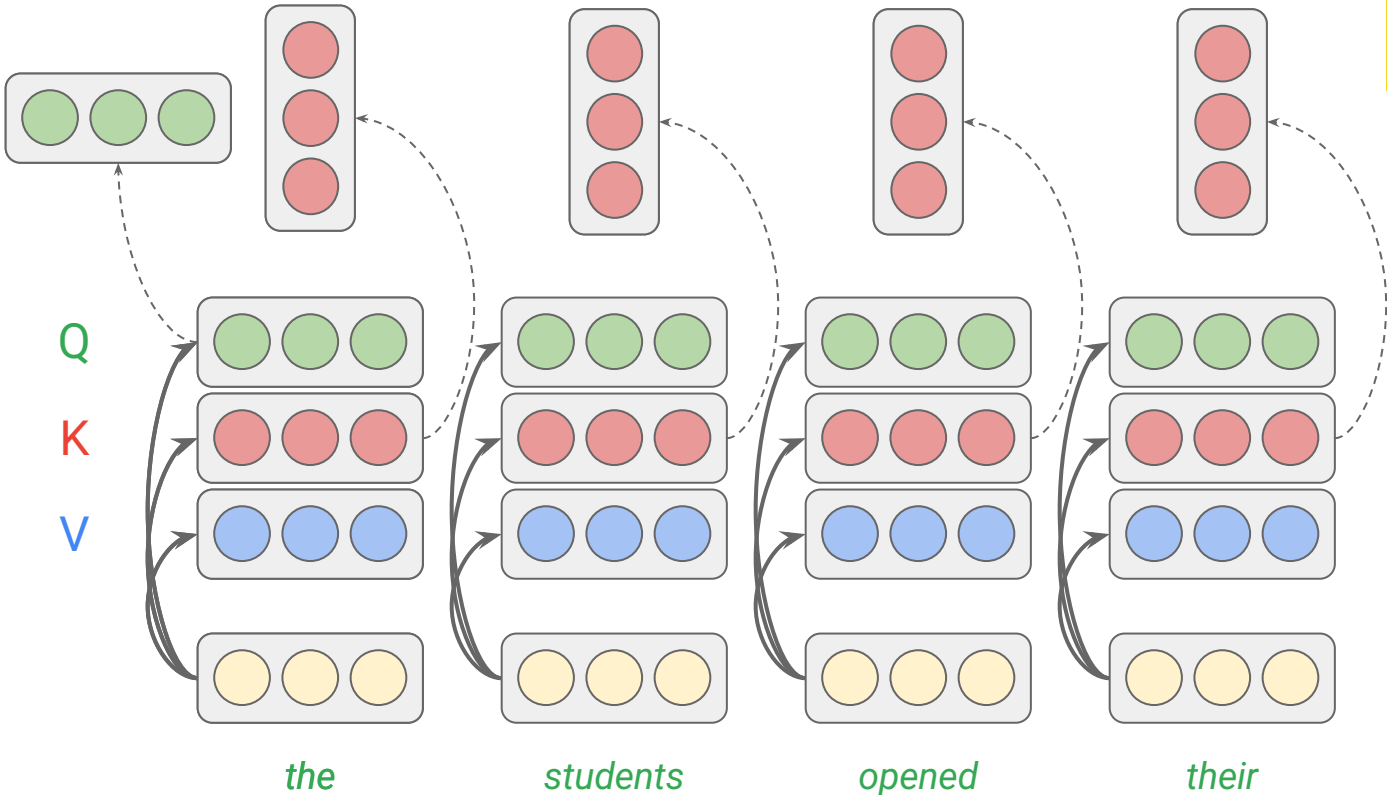


Self-attention



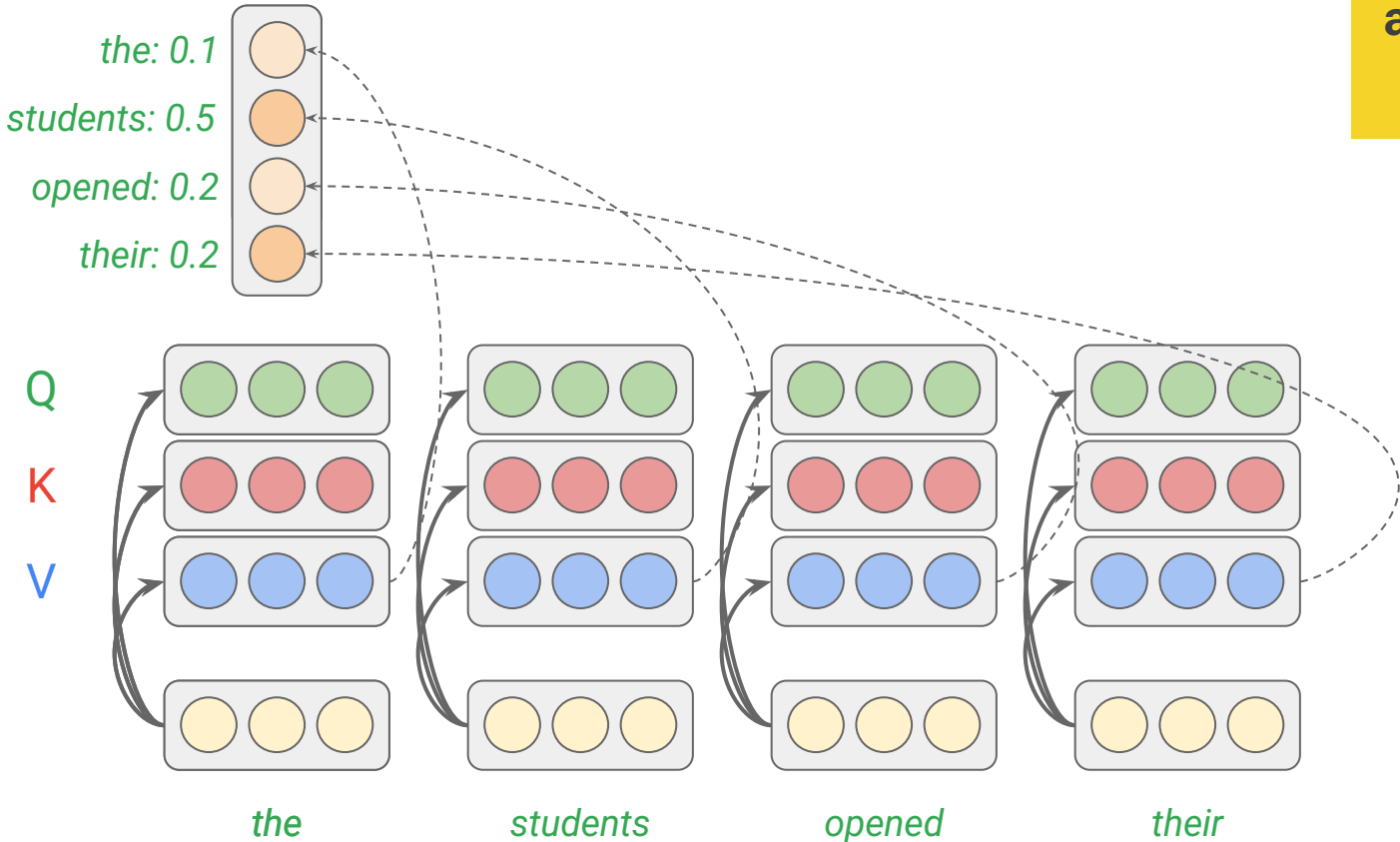
Self-attention (cont'd)

all computations are parallelized



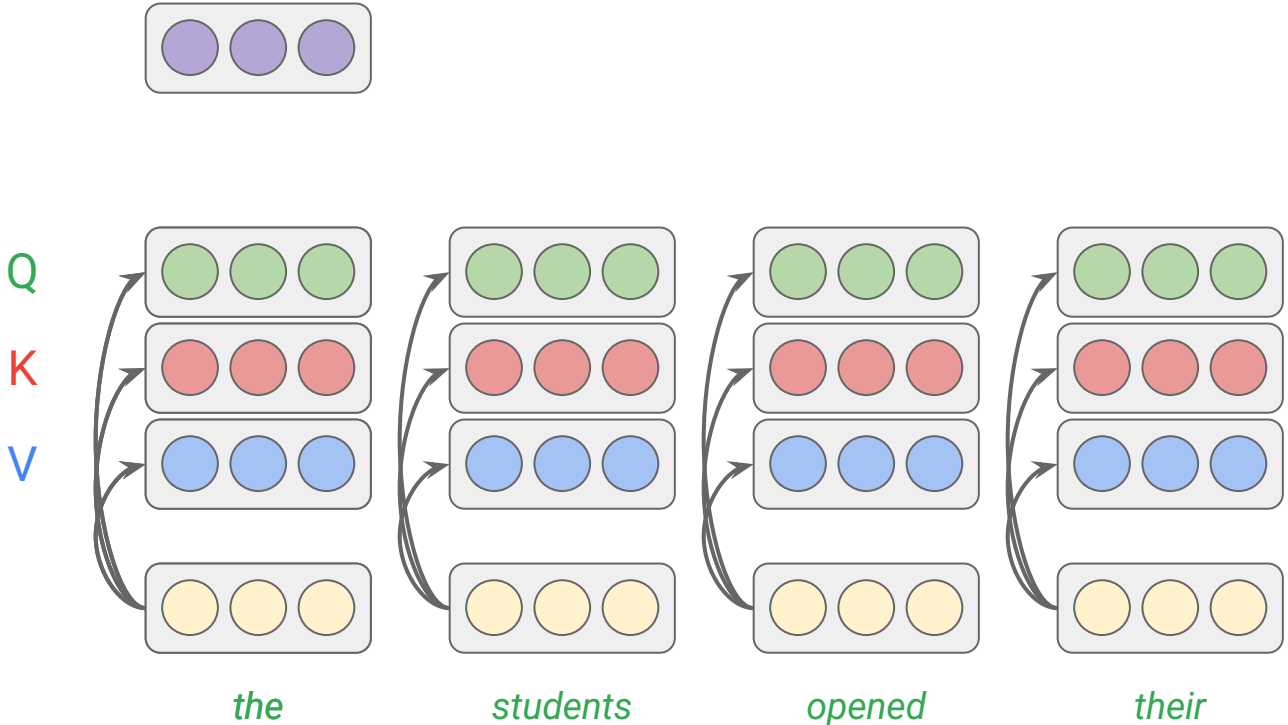
Self-attention (cont'd)

all computations are parallelized



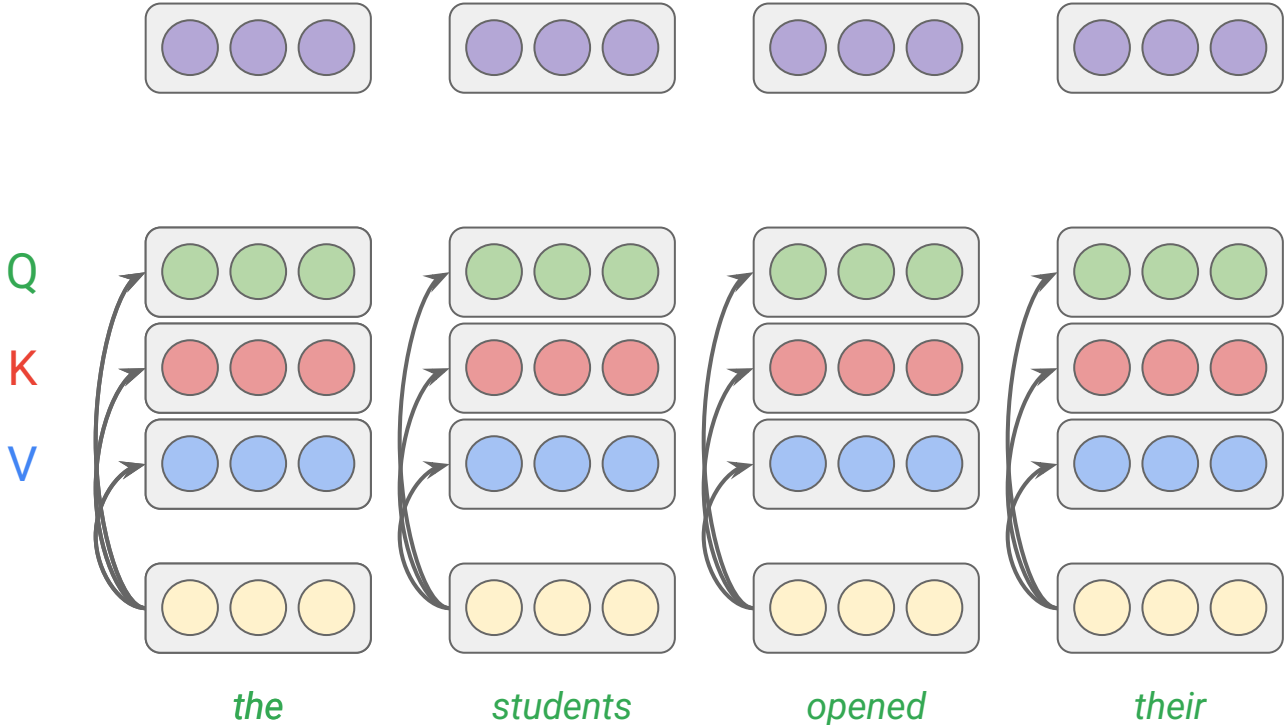
Self-attention (cont'd)

all computations are parallelized

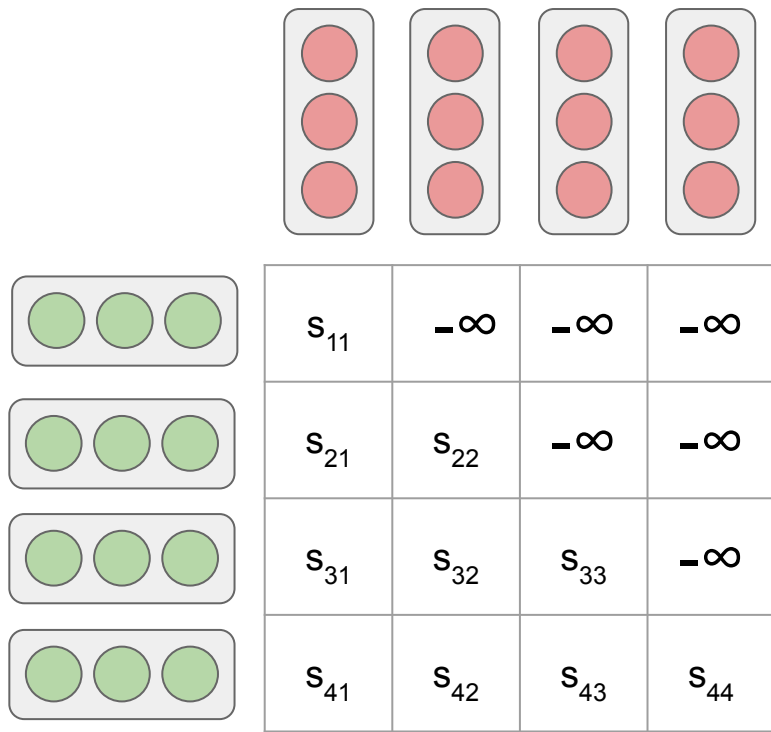


Self-attention (cont'd)

all computations are parallelized

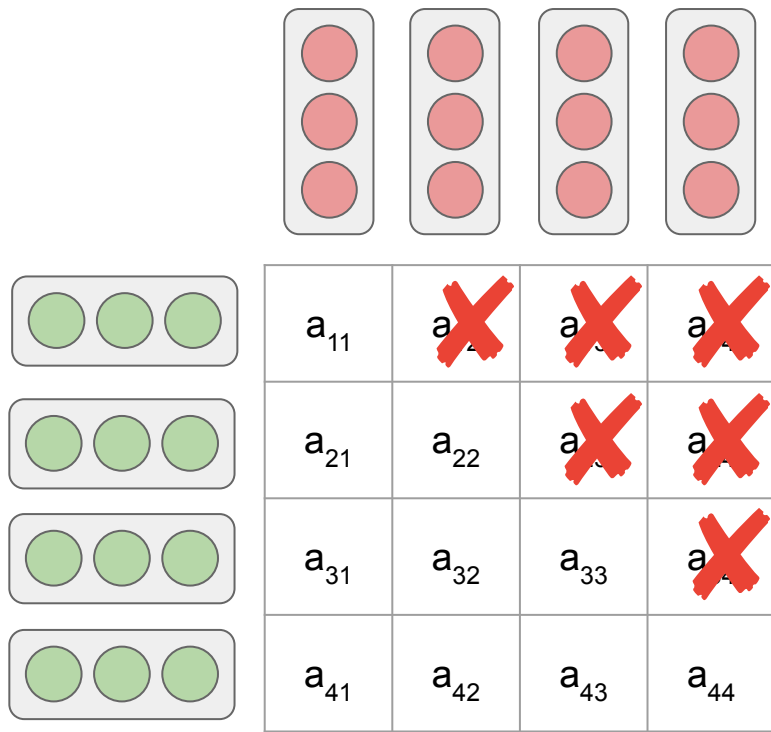


Self-attention in the decoder



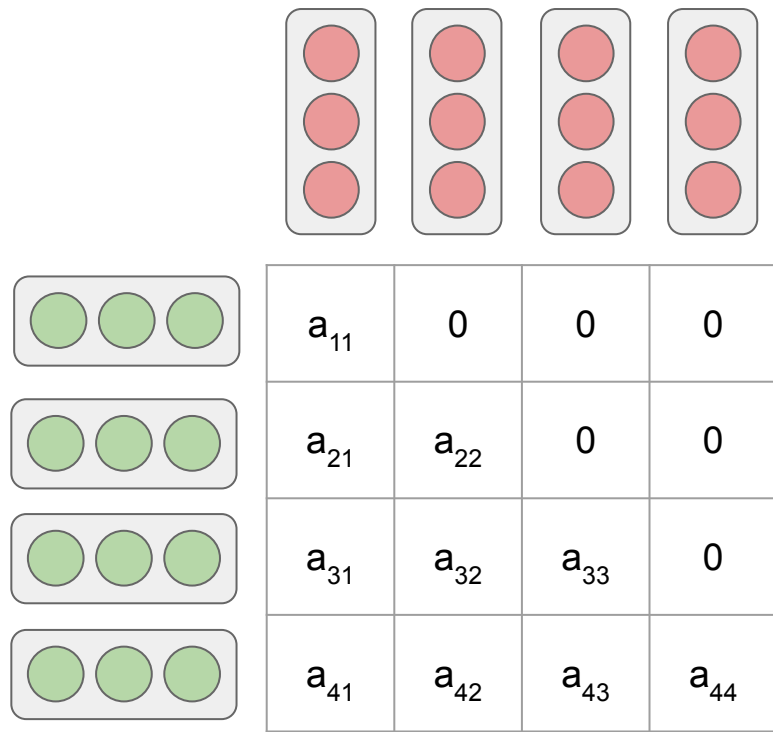
masking out (setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections

Self-attention in the decoder (cont'd)



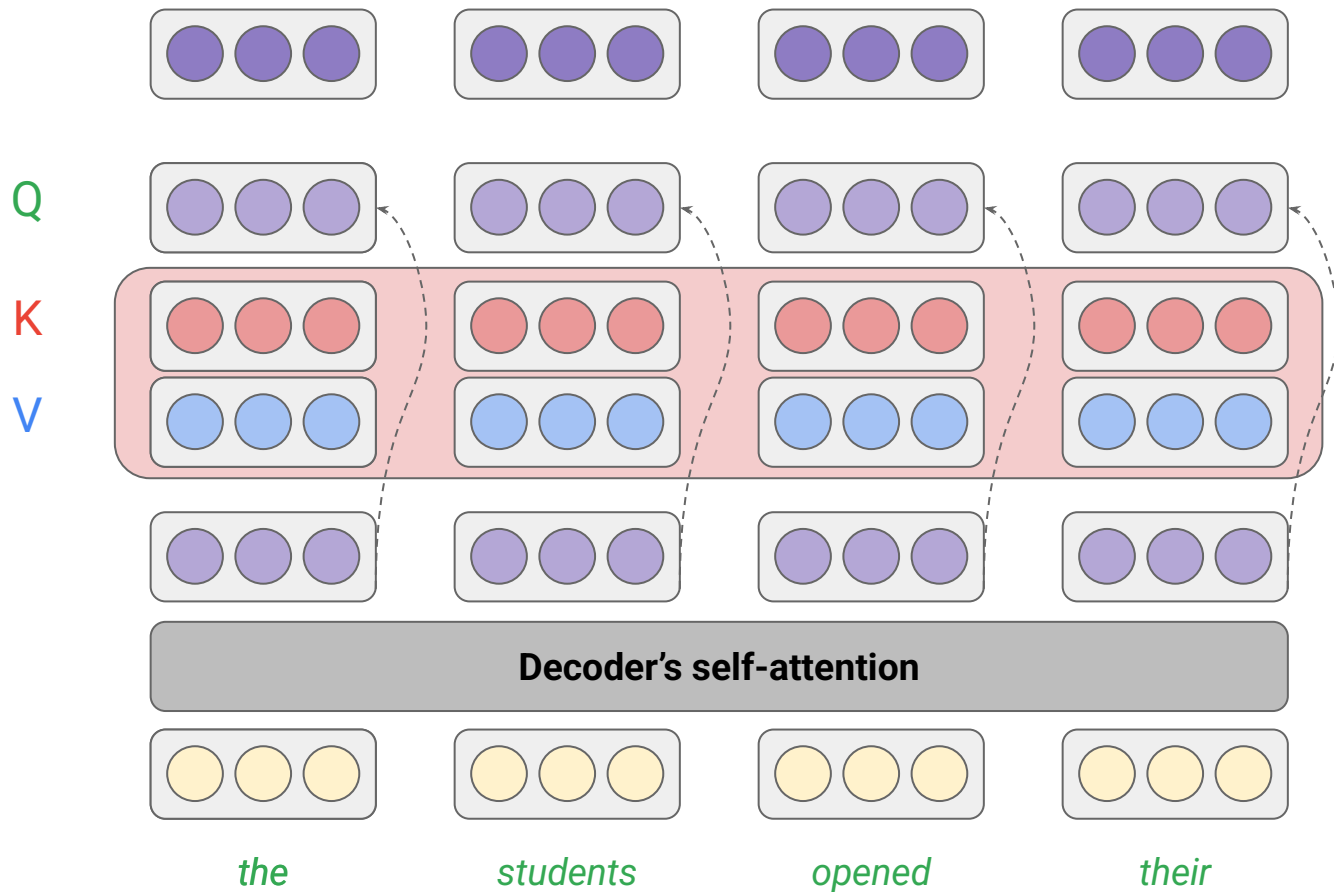
masking out all values in the input of the softmax which correspond to illegal connections

Self-attention in the decoder (cont'd)



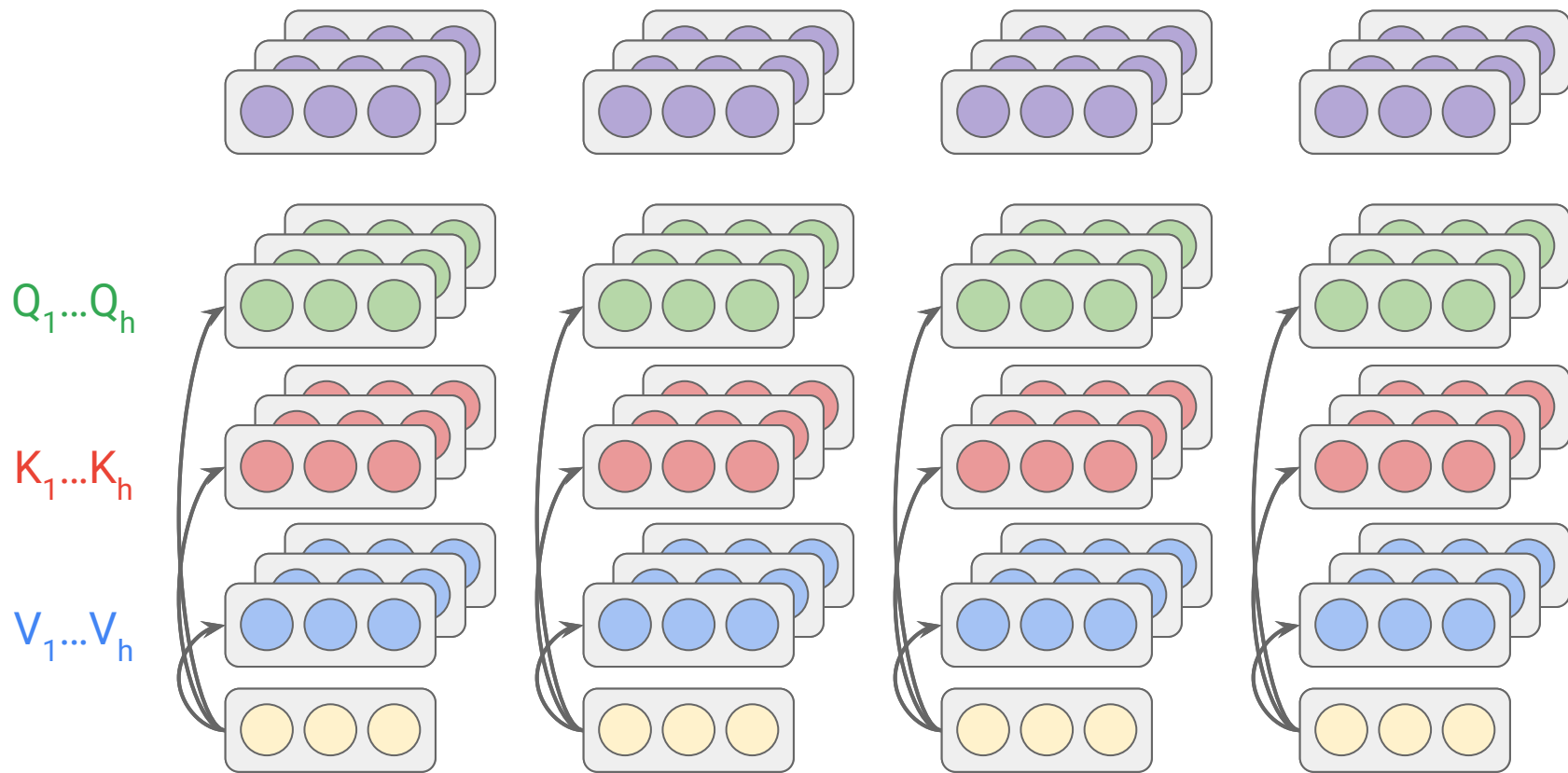
masking out all values in the input of the softmax which correspond to illegal connections

Cross-attention in the decoder

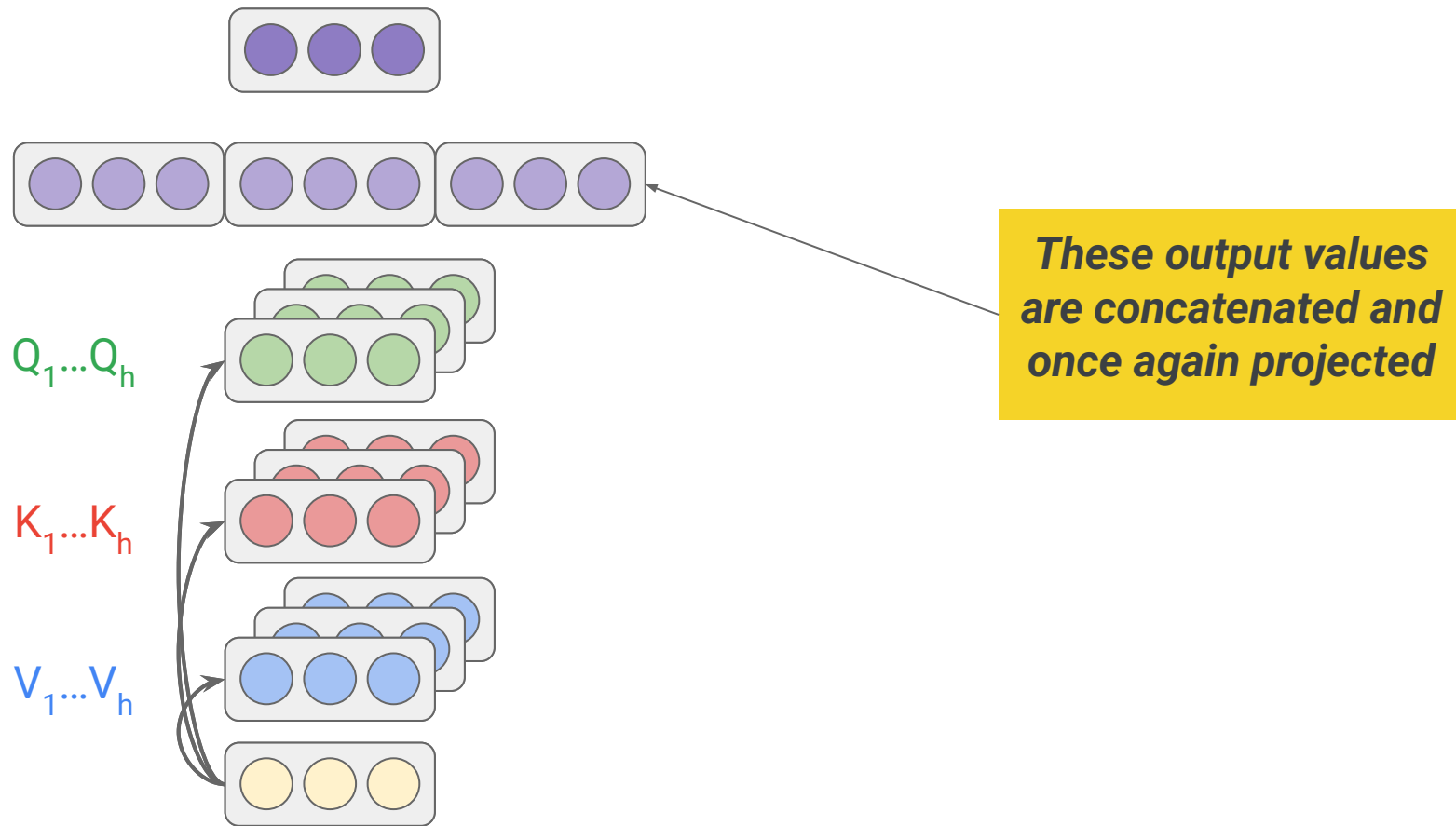


These K, V are the output of the encoder

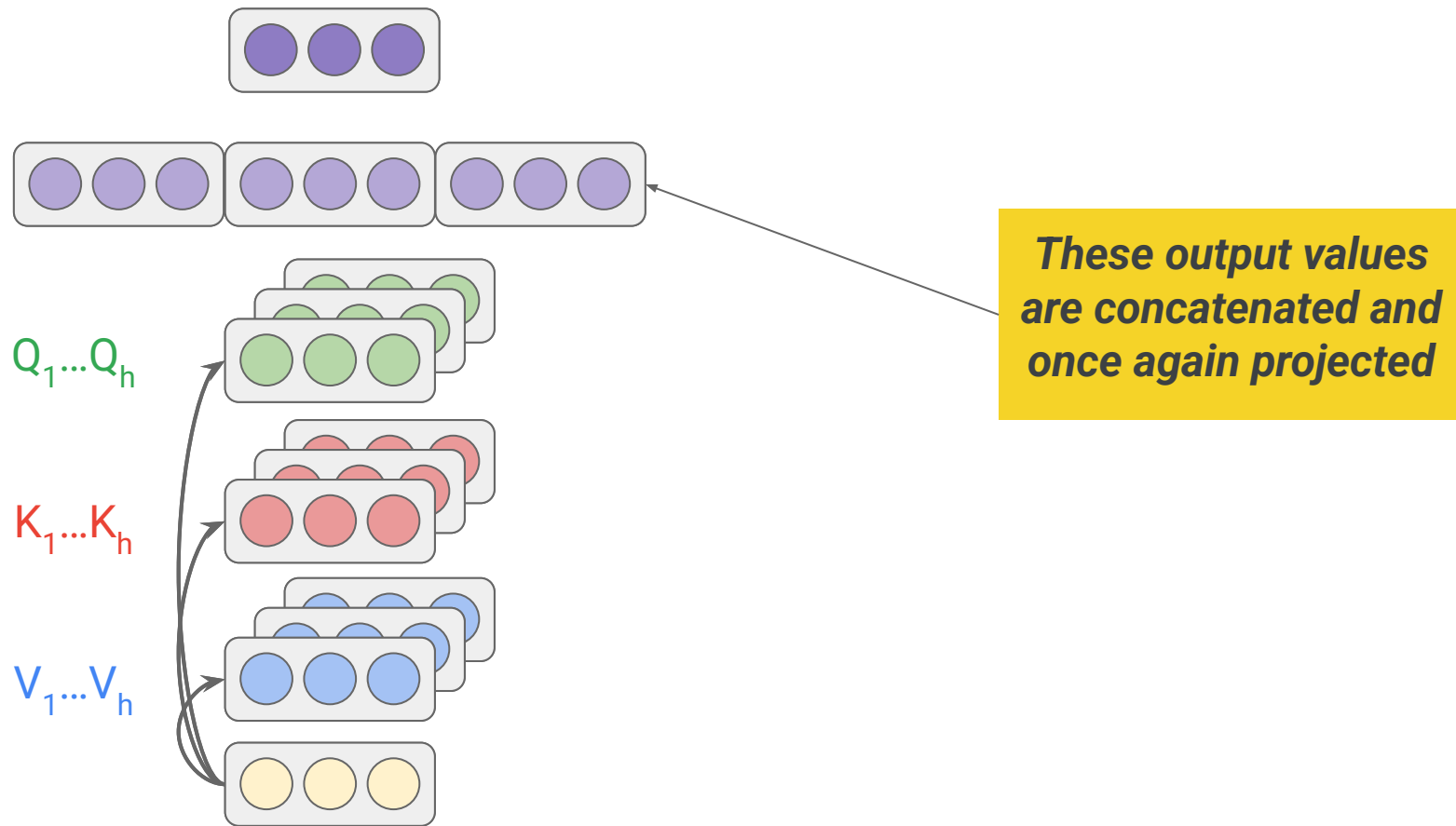
Multi-head attention



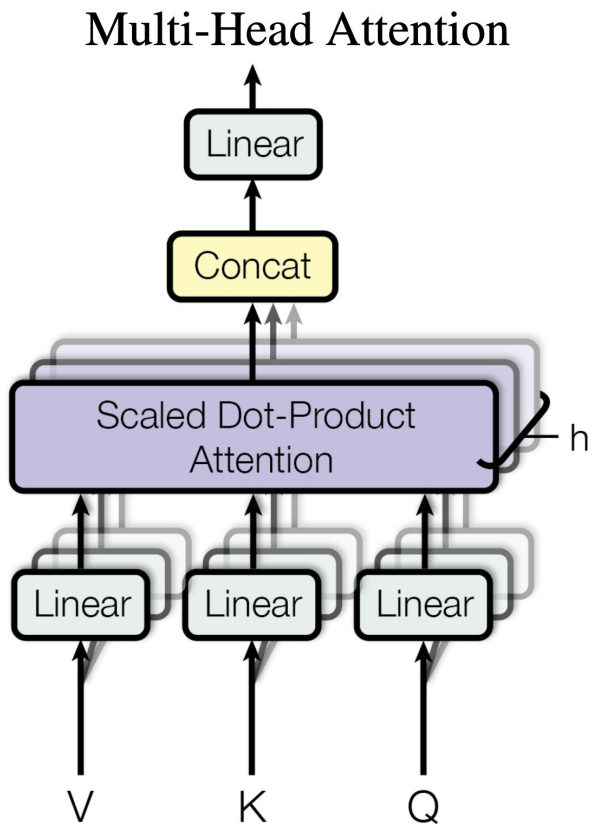
Multi-head attention (cont'd)



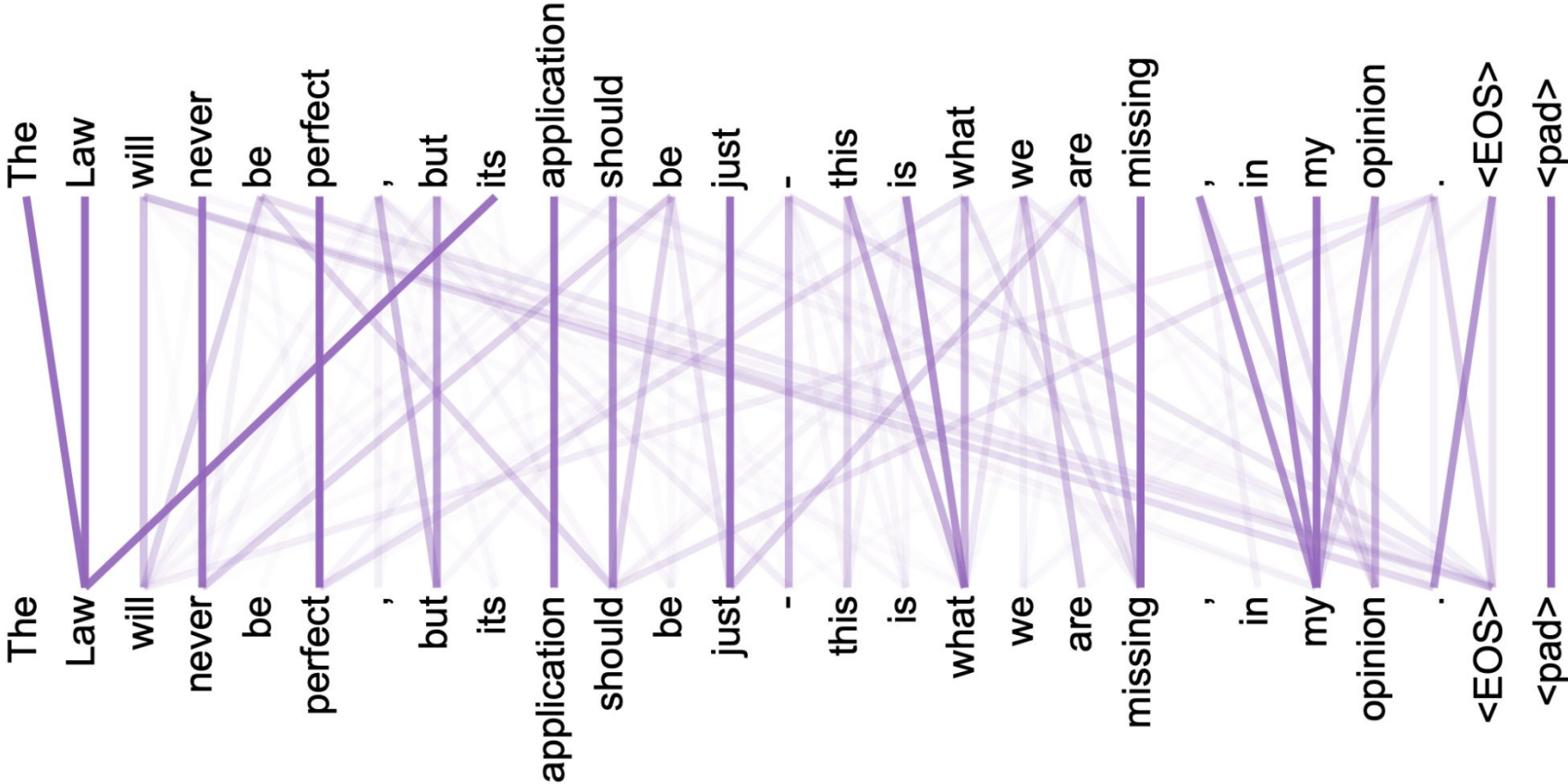
Multi-head attention (cont'd)



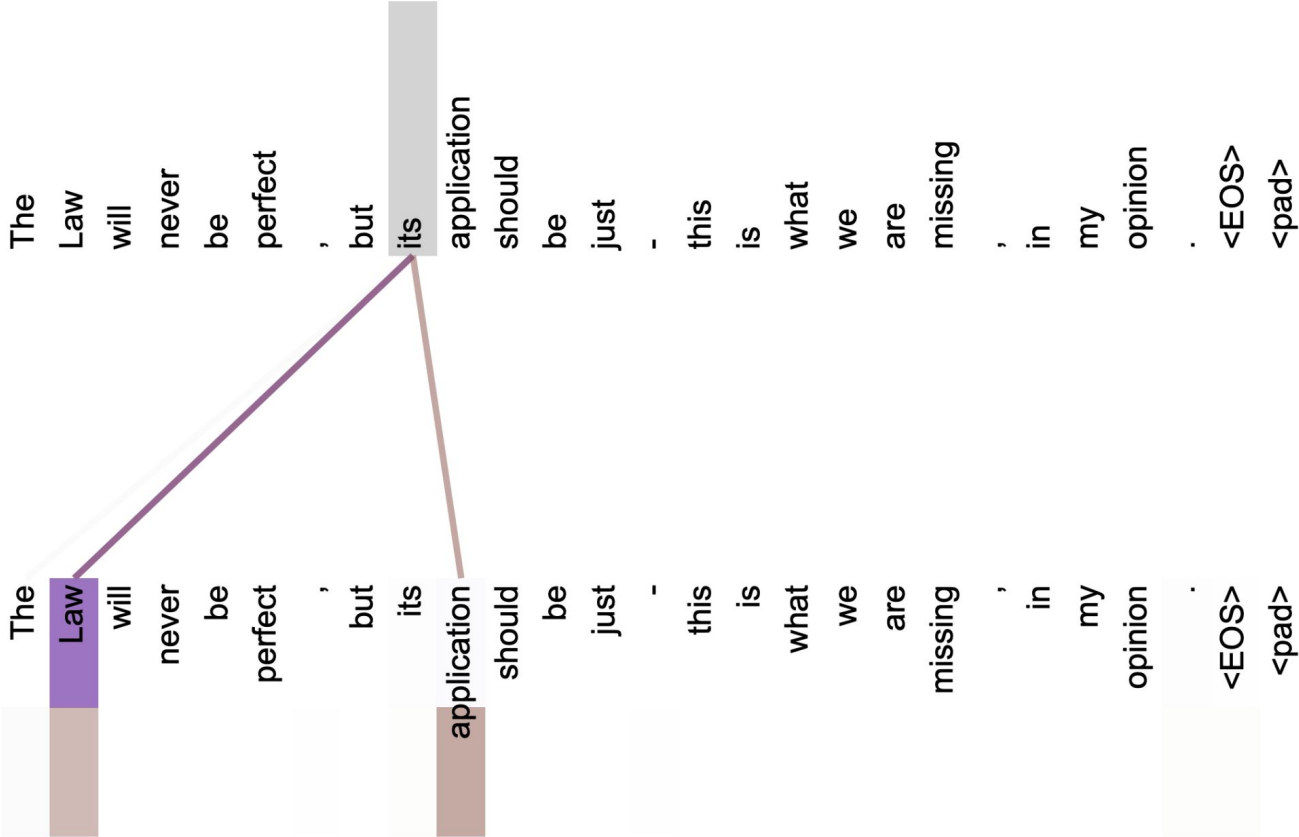
Multi-head attention (cont'd)



Attention visualizations



Attention visualizations (cont'd)



Position-wise Feed-Forward Networks

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

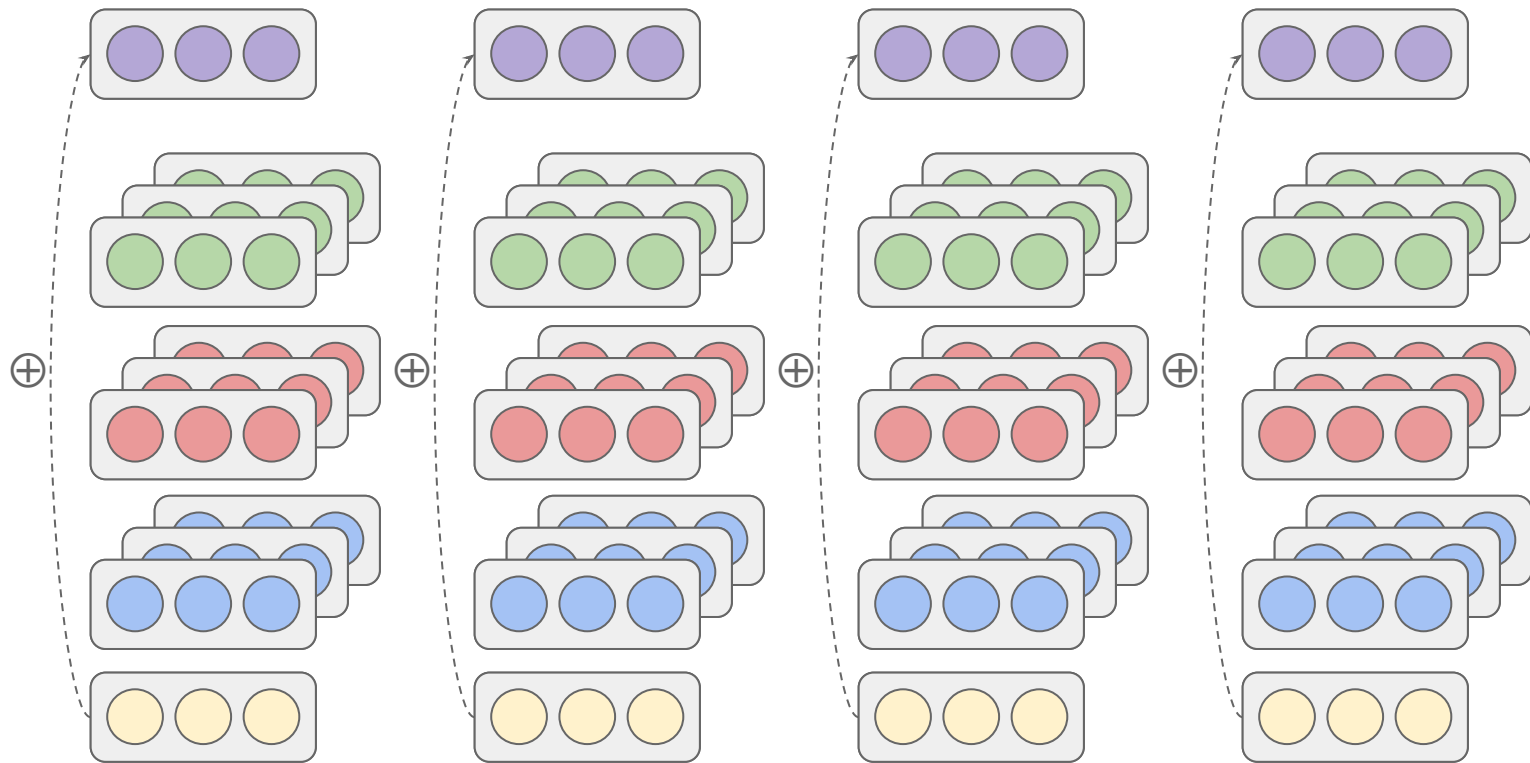


ReLU (Rectified
Linear Unit)

Residual connection and layer normalization

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

Residual connection

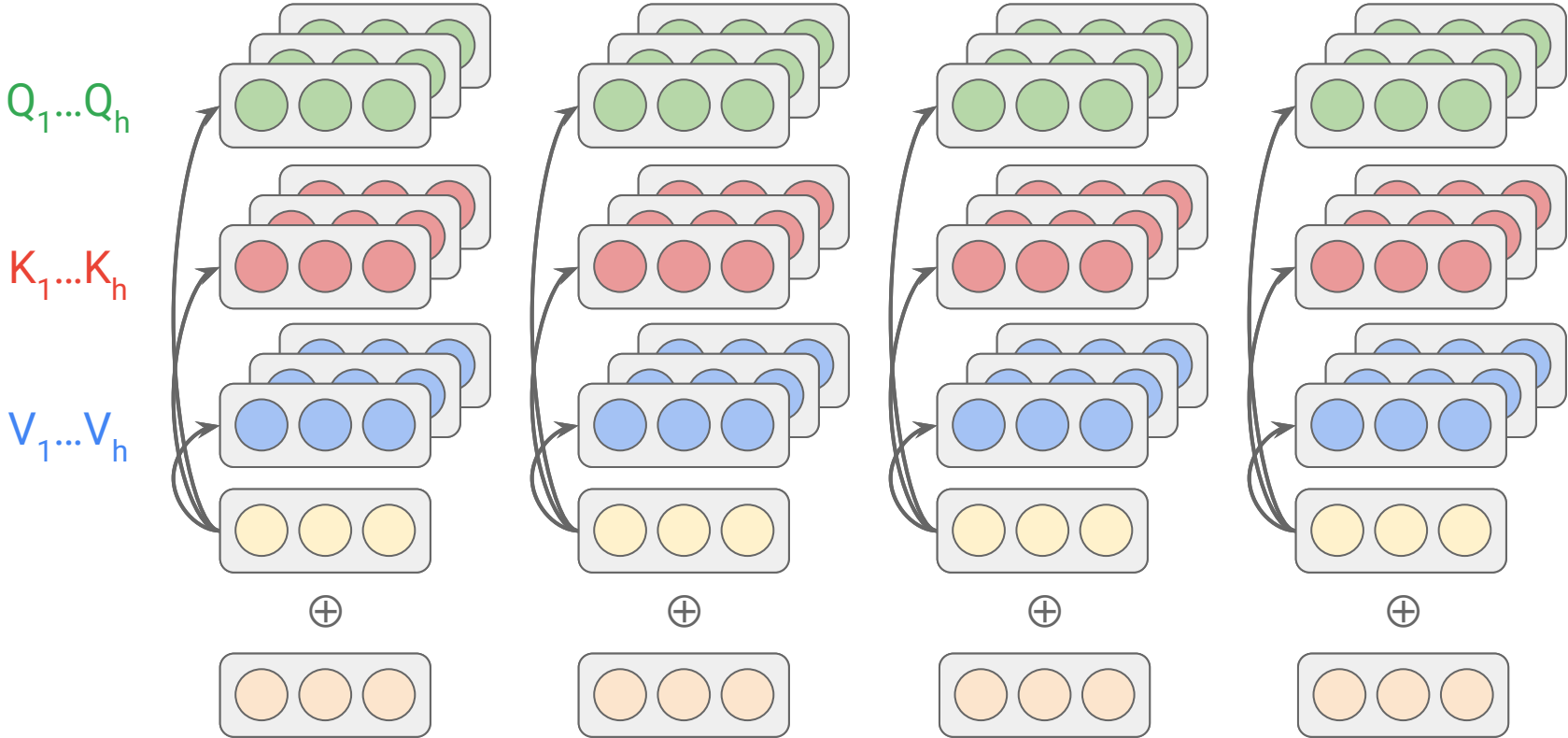


Positional Encoding

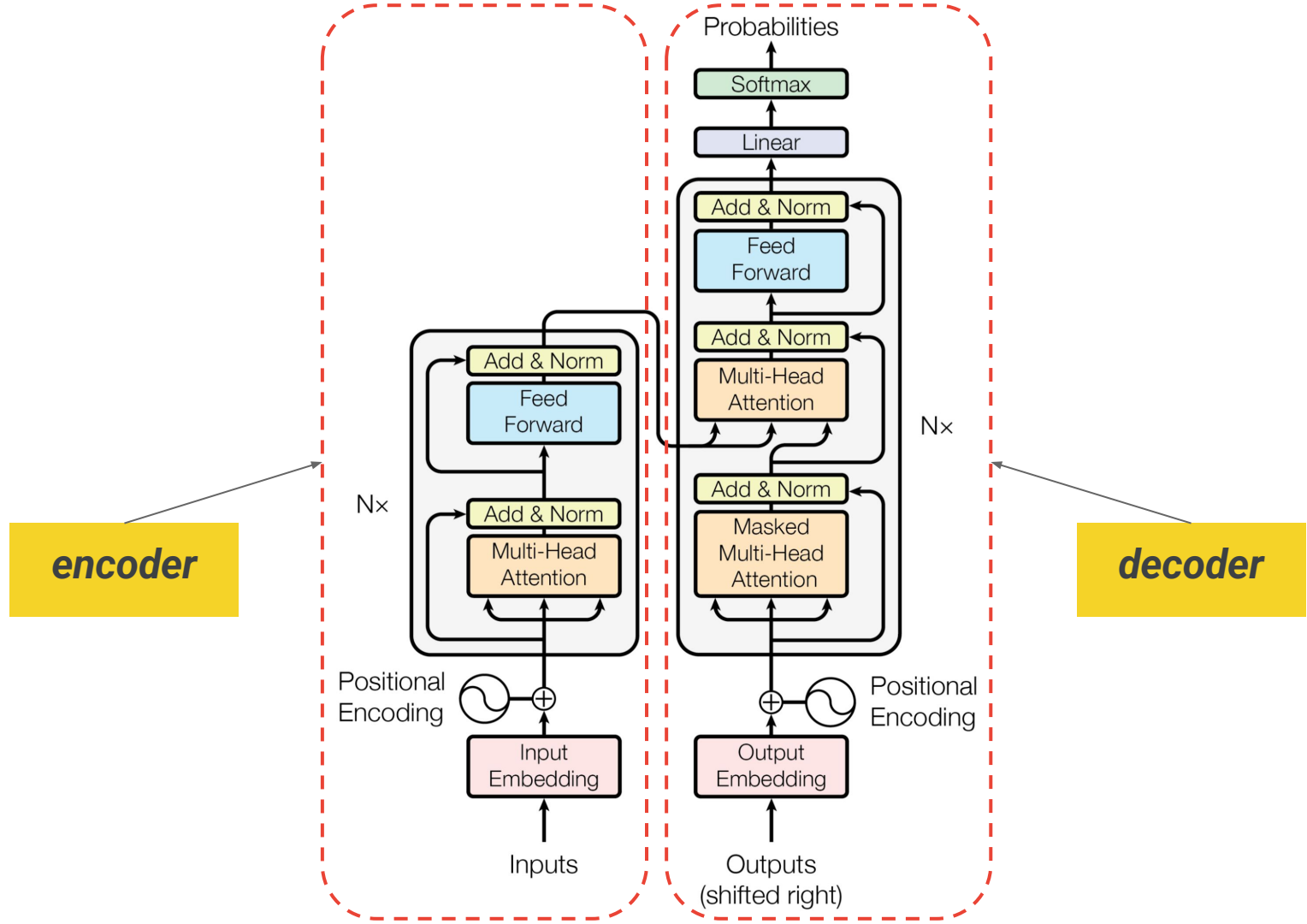
$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i / d_{\text{model}}})$$

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

Positional Encoding (cont'd)



Transformer block (putting it together)



Training and Test

Thank you!