

# Transformers (cont'd)

**CS 4804: Introduction to AI**

*Fall 2025*

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

**Tu Vu**

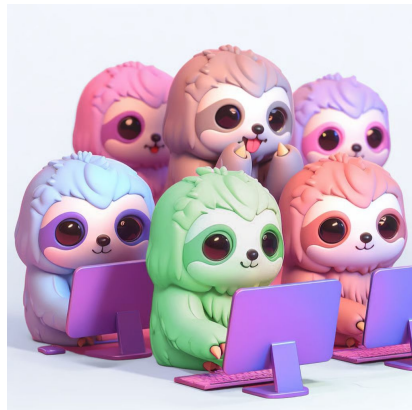


# Logistics

- Homework 0 (due today)
- Final project groups finalized
  - Project proposal template released today (due in 3 weeks)
  - Project research ideas
    - Check out [GPT-5 System Card](#) & [Gemini 2.5 Pro Technical Report](#)
    - New benchmarks
    - New prompting methods
    - Improving small models
    - ...

# Useful libraries

- Hugging Face (<https://huggingface.co>)
  - Transformers
  - Models / Datasets
- Unsloth (<https://unsloth.ai>)
  - Efficient fine-tuning
- vLLM (<https://github.com/vllm-project/vllm>)
  - Efficient inference
- [Together.AI](https://together.ai)
  - APIs for Open-weight models



together.ai

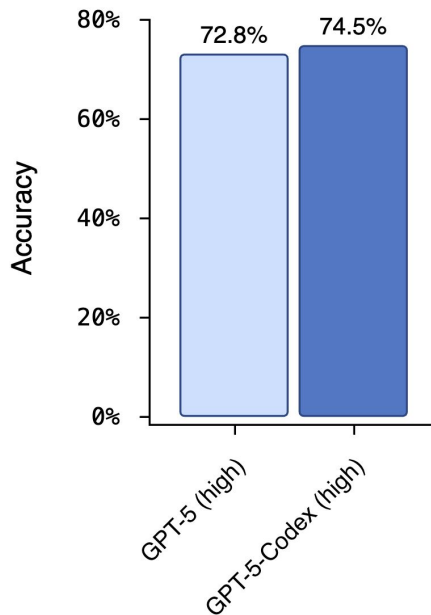
# AI News

- Google's Nano Banana 🍌🍌🍌
  - restore old photos
  - try out different hairstyles
  - imagine yourselves as 3D model figurines

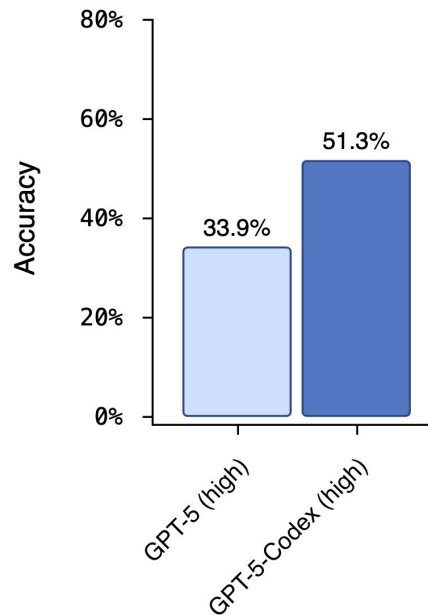
# AI News (cont'd)

- [OpenAI's GPT-5-Codex](#): GPT-5 optimized for agentic coding

**SWE-bench Verified (n=500)** 



**Code refactoring tasks** 



# AI News (cont'd): Thinking Machines

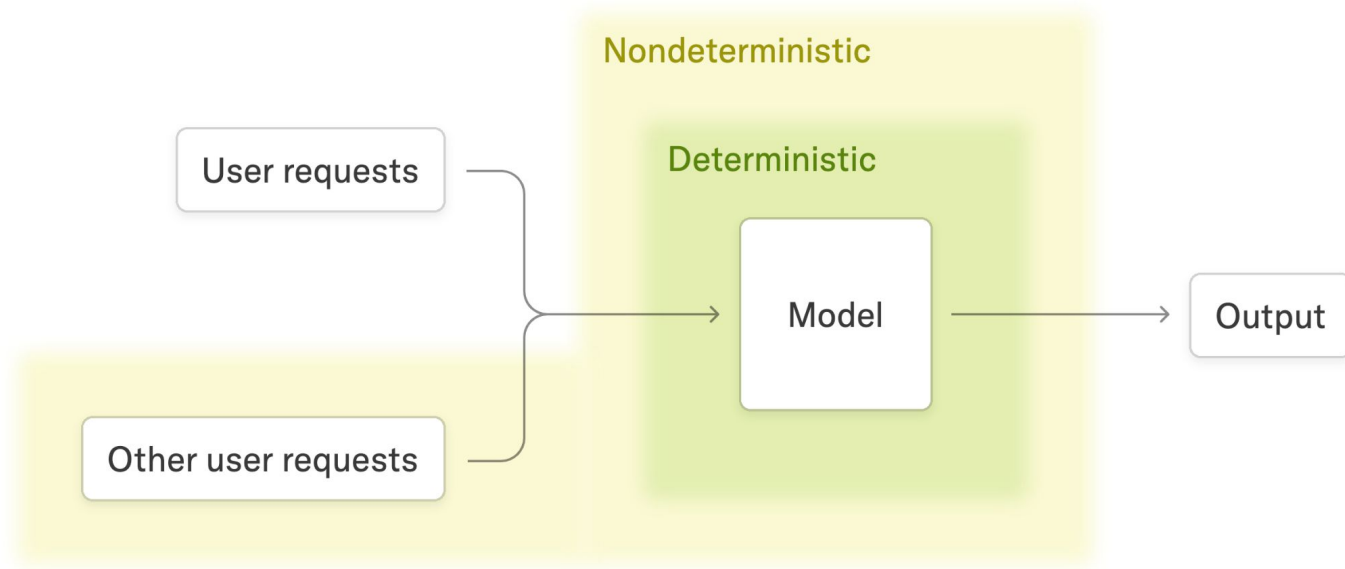
- Defeating Nondeterminism in LLM Inference

Temperature 1.00

Max tokens 2048

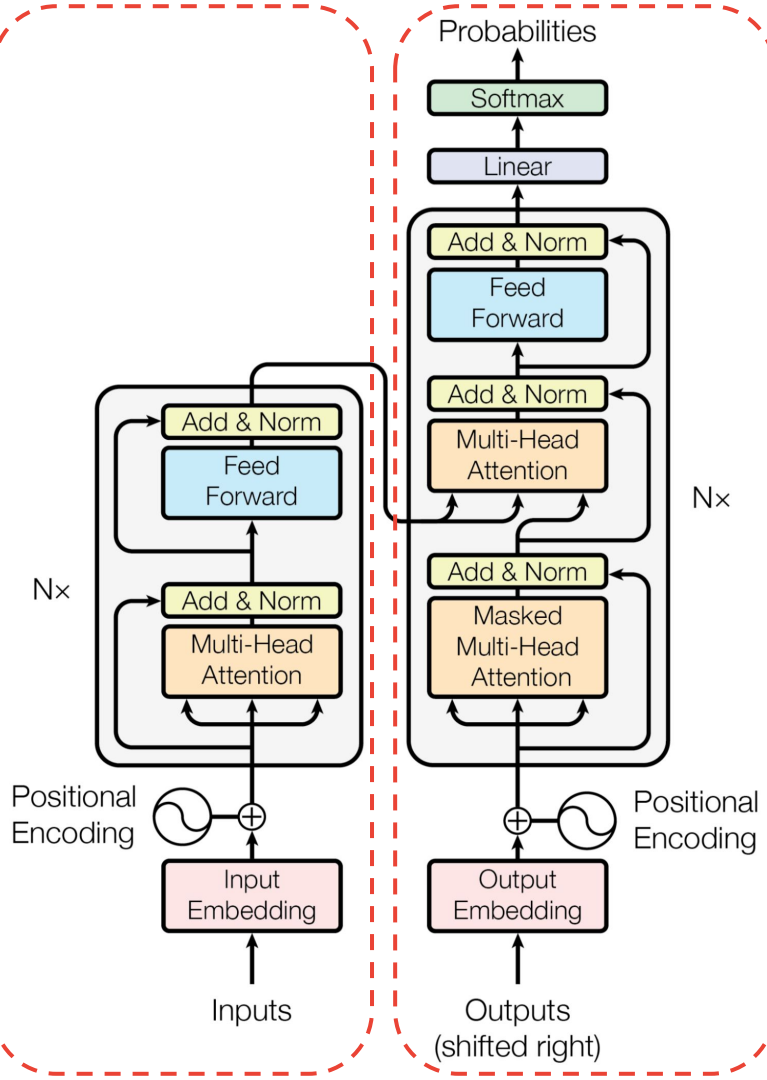
Top P 1.00

Store logs ☒



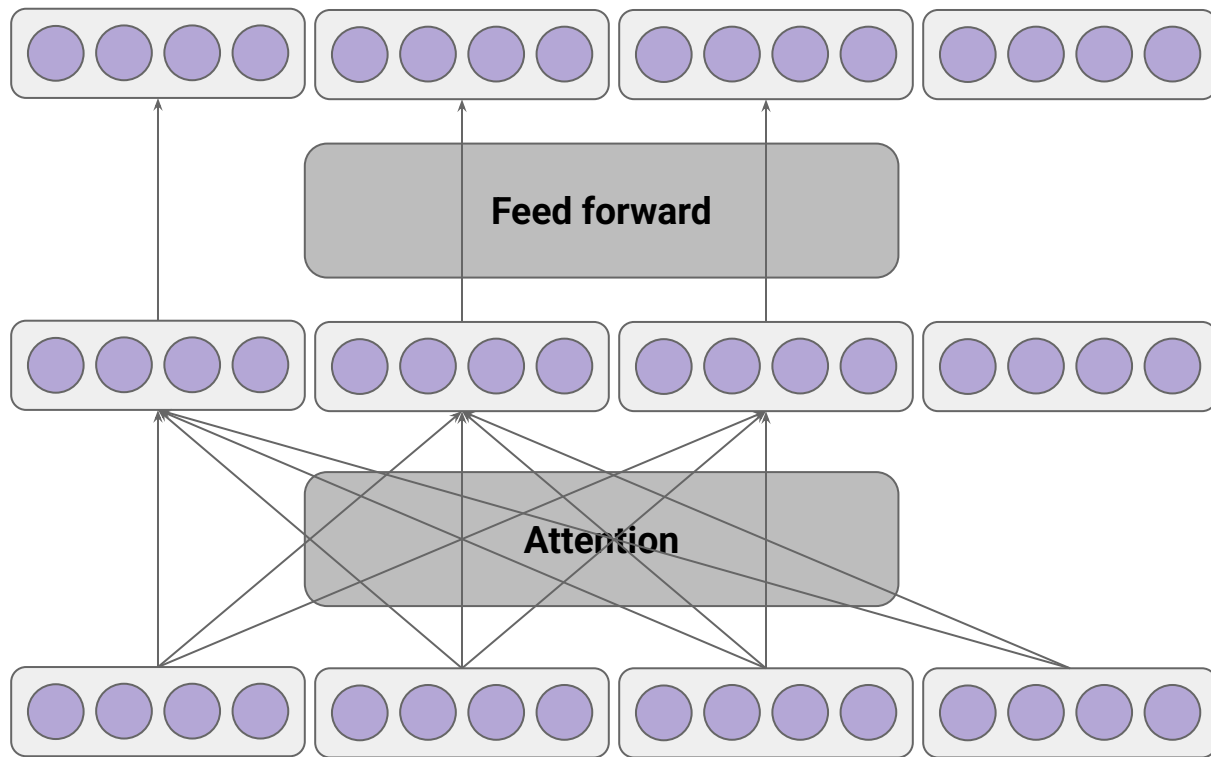
# Transformers (cont'd)

**encoder**



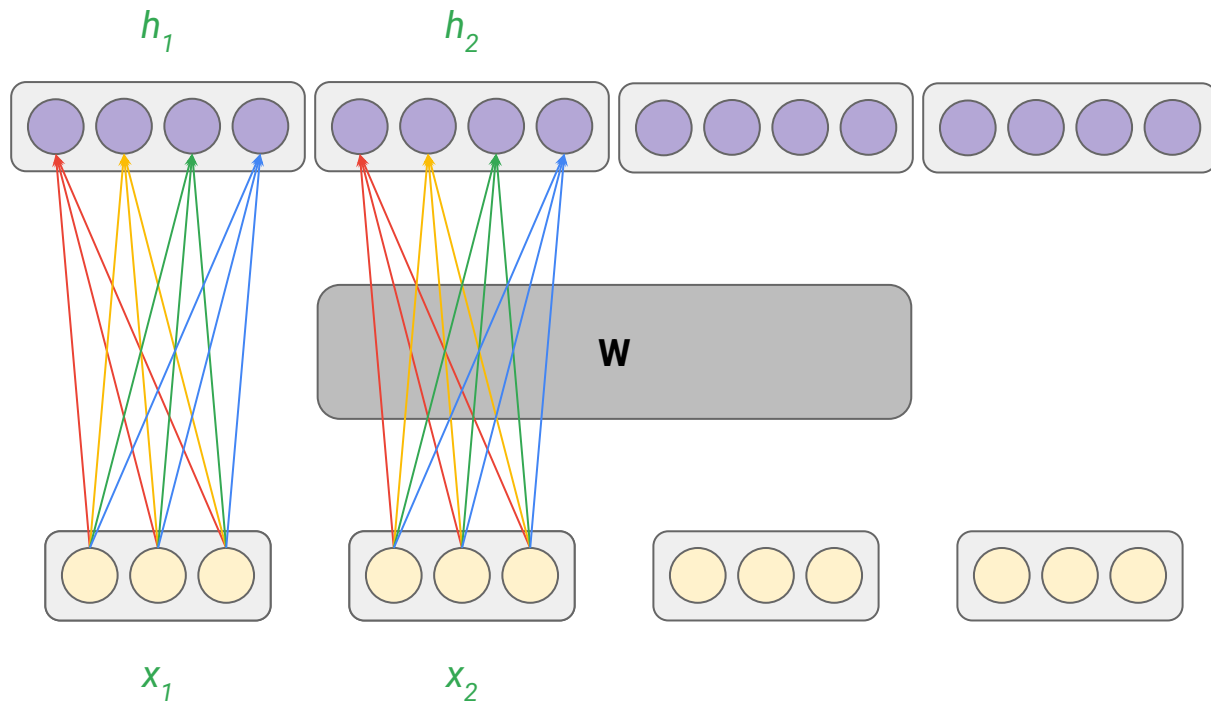
**decoder**

# Encoder





# Position-wise Feed-Forward Networks



We multiply the weight matrix  $W$  (size  $4 \times 3$ ) with the embeddings matrix  $X$  (size  $3 \times 2$ ):

$$H = WX$$

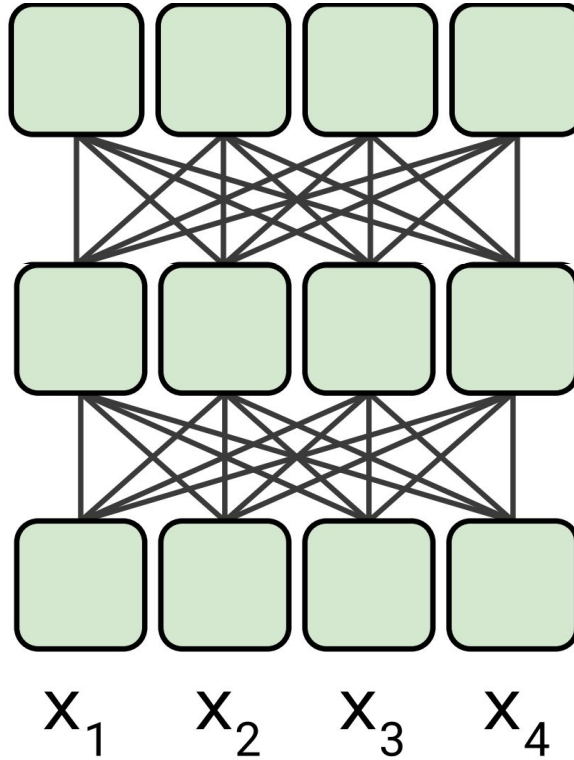
Performing the multiplication:

$$H = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} \begin{bmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \\ x_{13} & x_{23} \end{bmatrix}$$

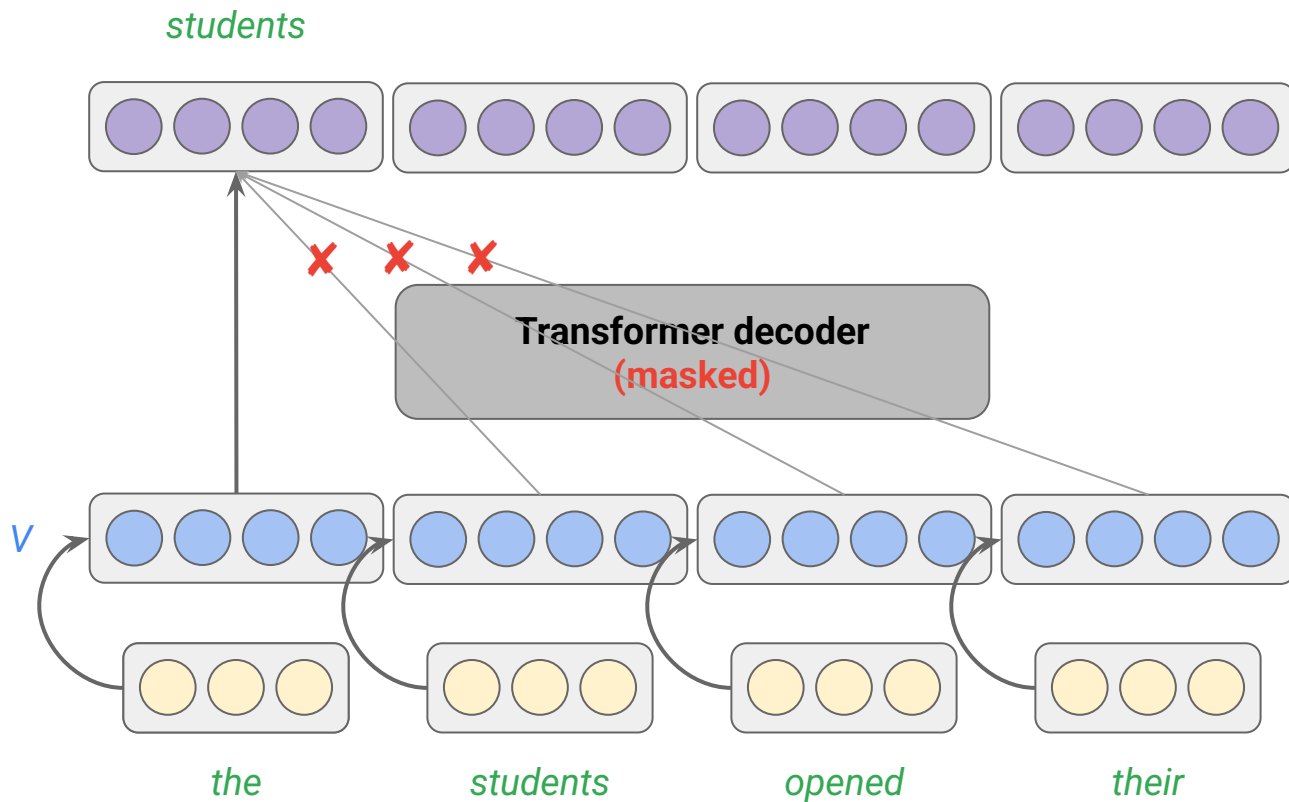
This results in:

$$H = \begin{bmatrix} \overset{h_1}{w_{11}x_{11} + w_{12}x_{12} + w_{13}x_{13}} & \overset{h_2}{w_{11}x_{21} + w_{12}x_{22} + w_{13}x_{23}} \\ w_{21}x_{11} + w_{22}x_{12} + w_{23}x_{13} & w_{21}x_{21} + w_{22}x_{22} + w_{23}x_{23} \\ w_{31}x_{11} + w_{32}x_{12} + w_{33}x_{13} & w_{31}x_{21} + w_{32}x_{22} + w_{33}x_{23} \\ w_{41}x_{11} + w_{42}x_{12} + w_{43}x_{13} & w_{41}x_{21} + w_{42}x_{22} + w_{43}x_{23} \end{bmatrix}$$

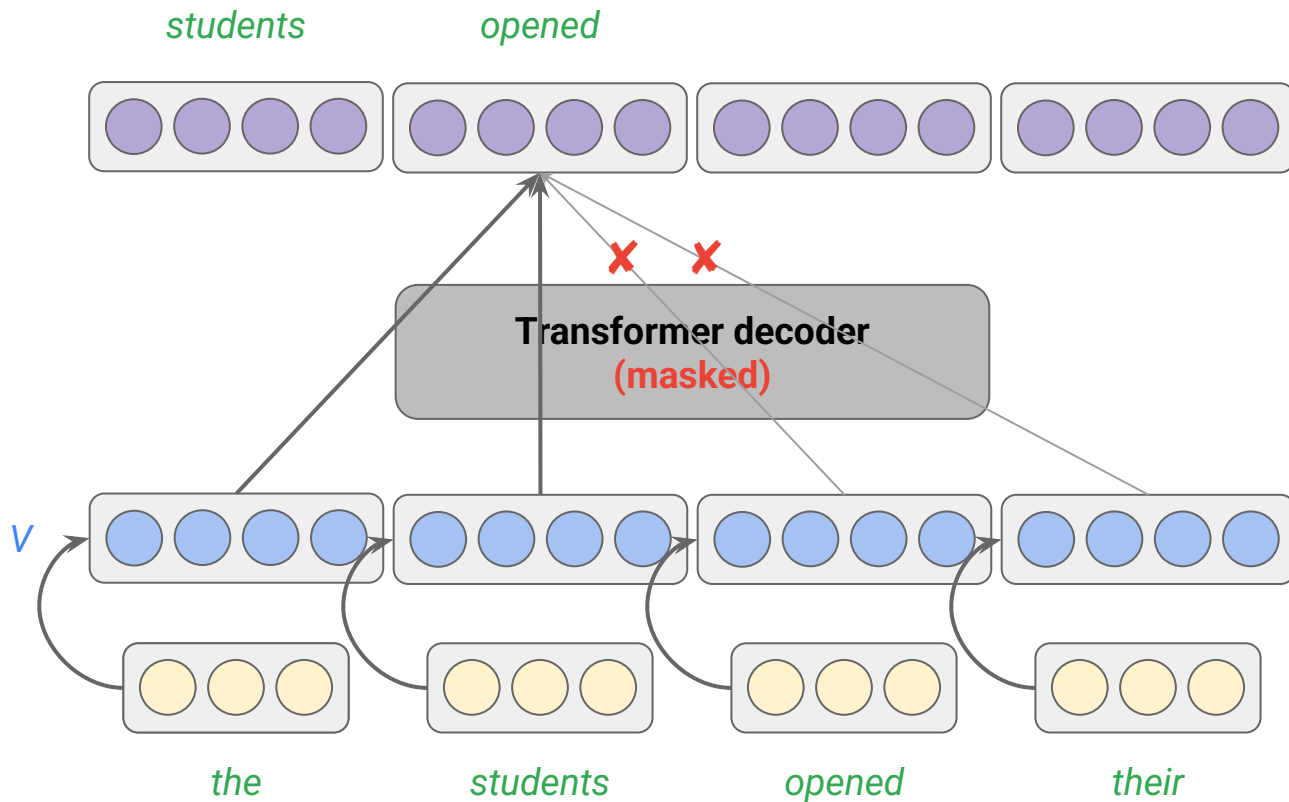
## Encoder (cont'd)



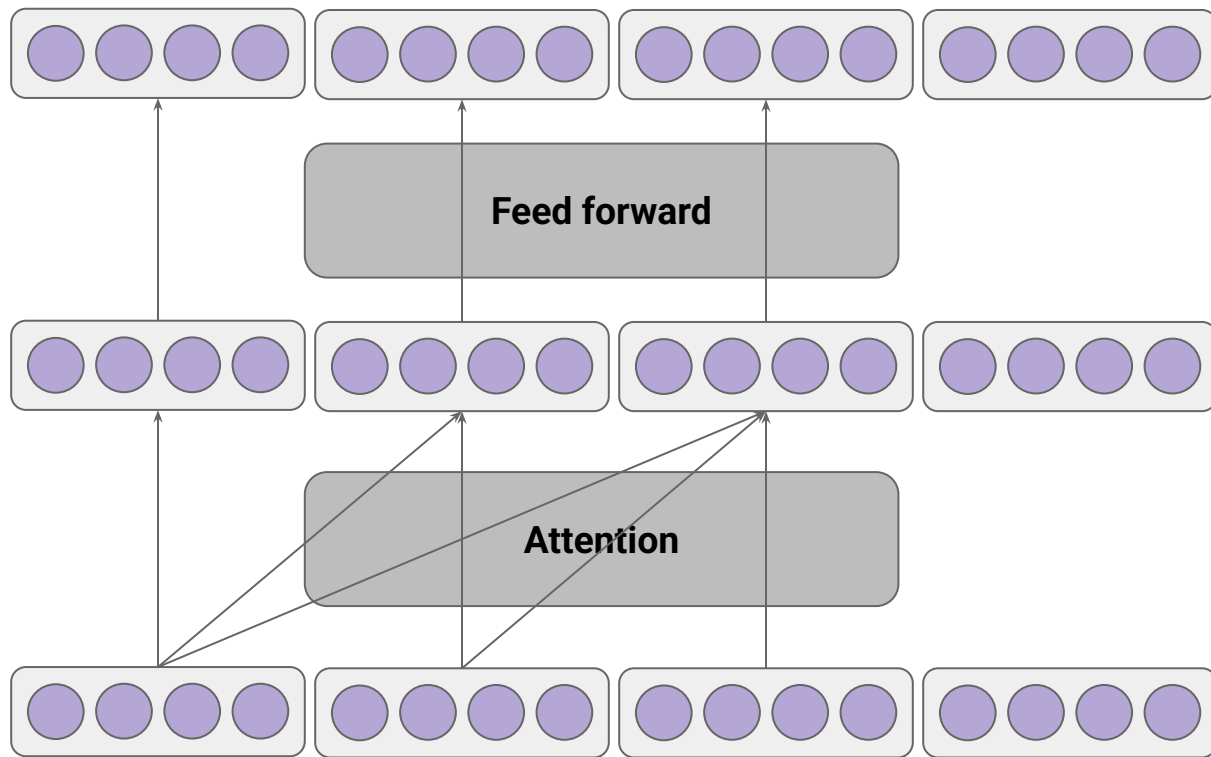
# Decoder



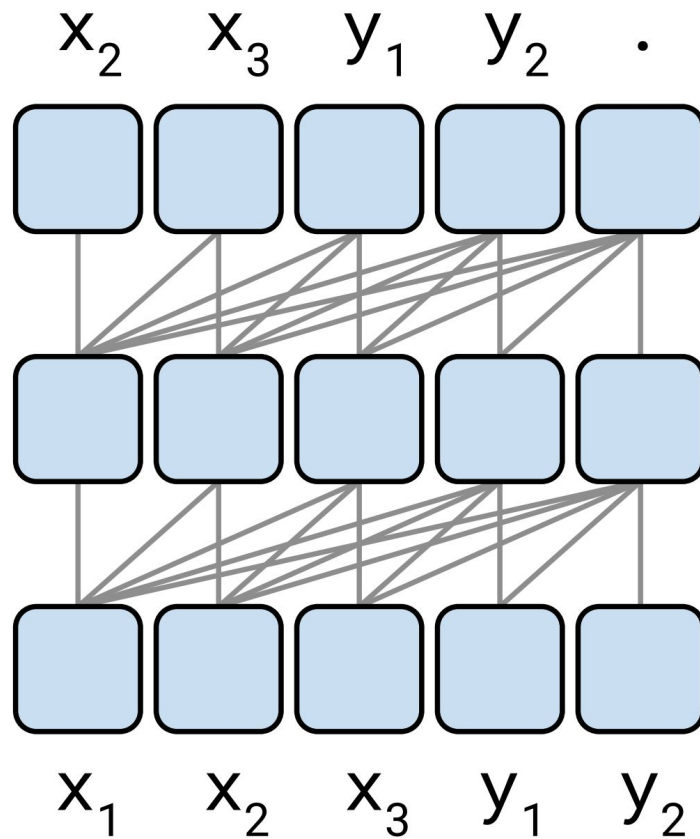
# Decoder (cont'd)



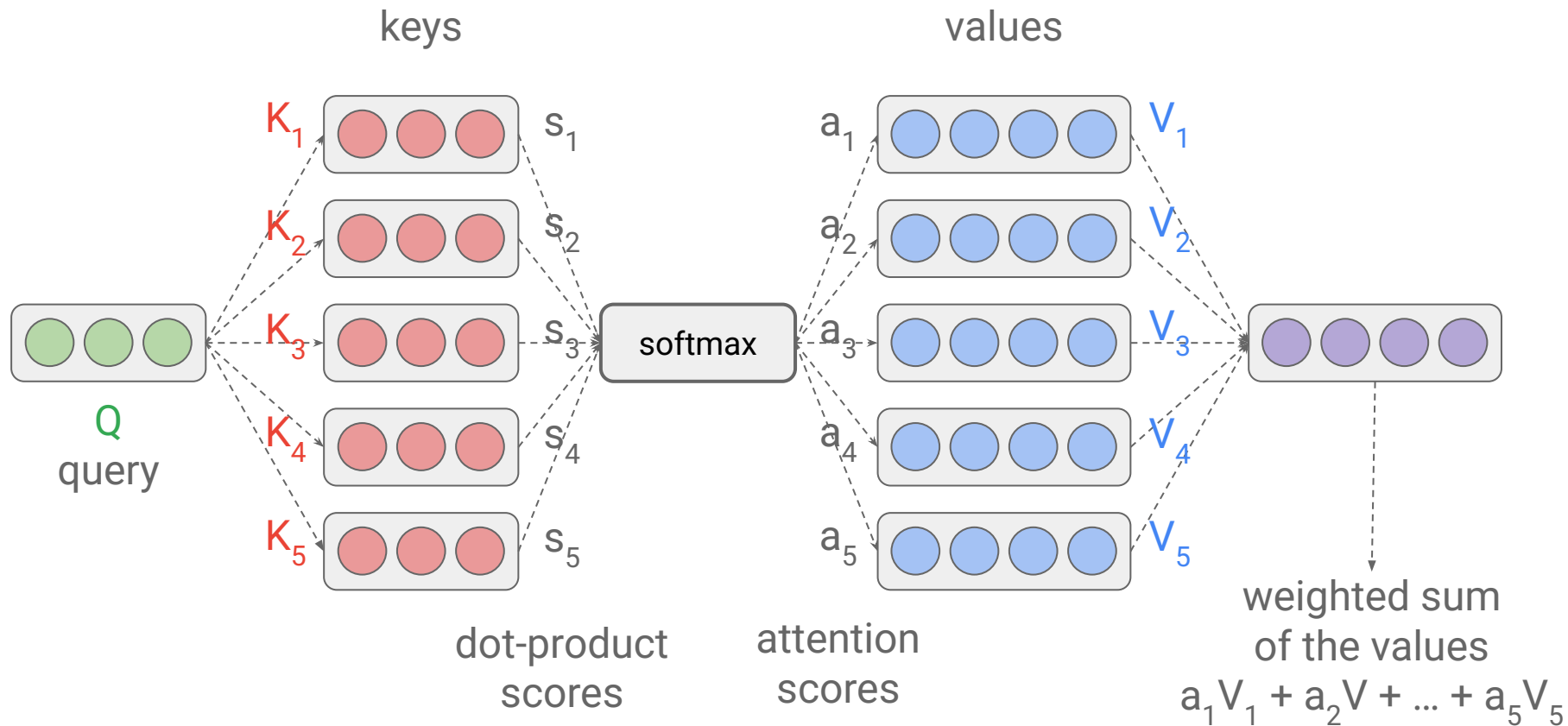
## Decoder (cont'd)



## Decoder (cont'd)

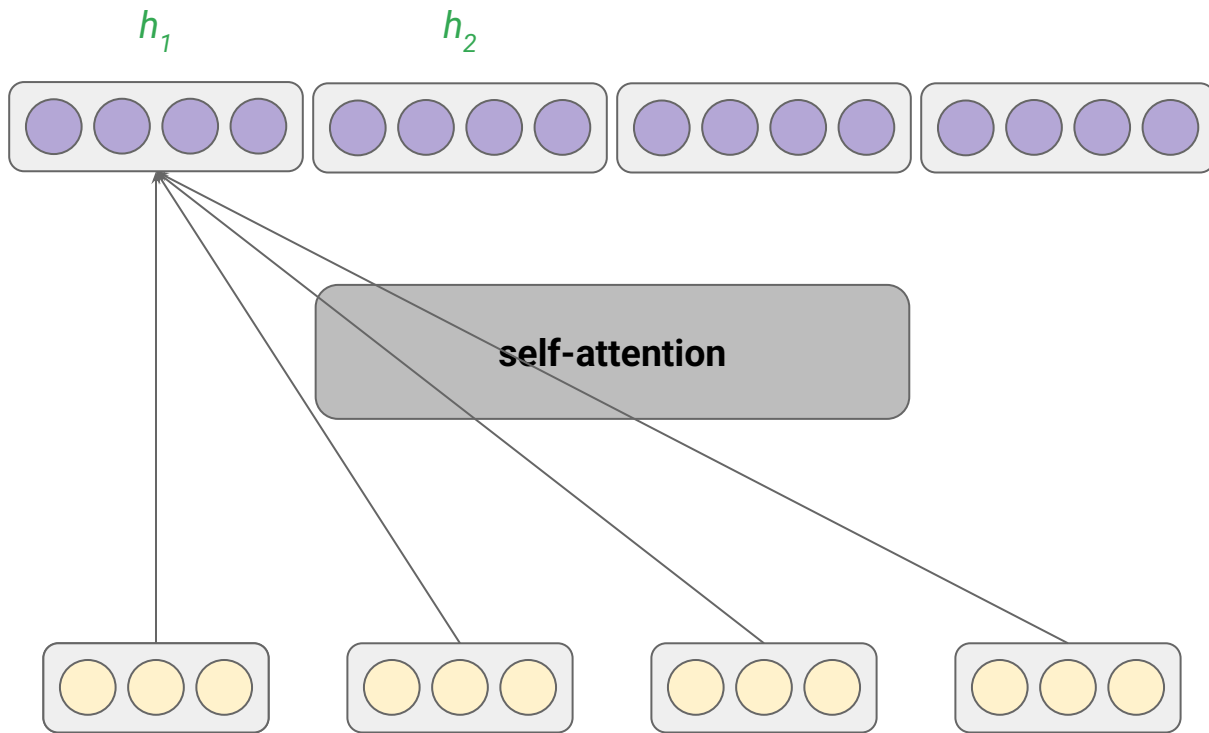


# Attention mechanism

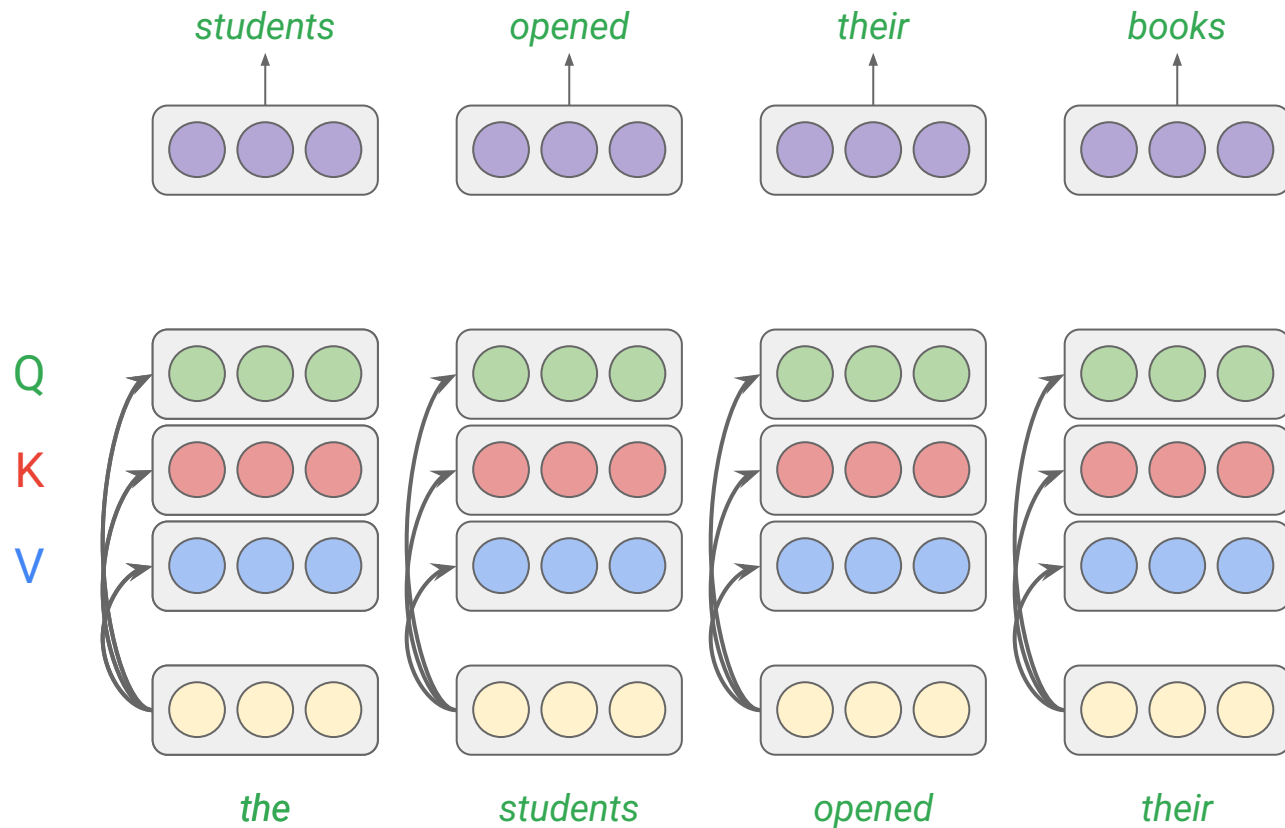




# Self-attention



# Self-attention (cont'd)



$$Q = X \cdot W_Q$$

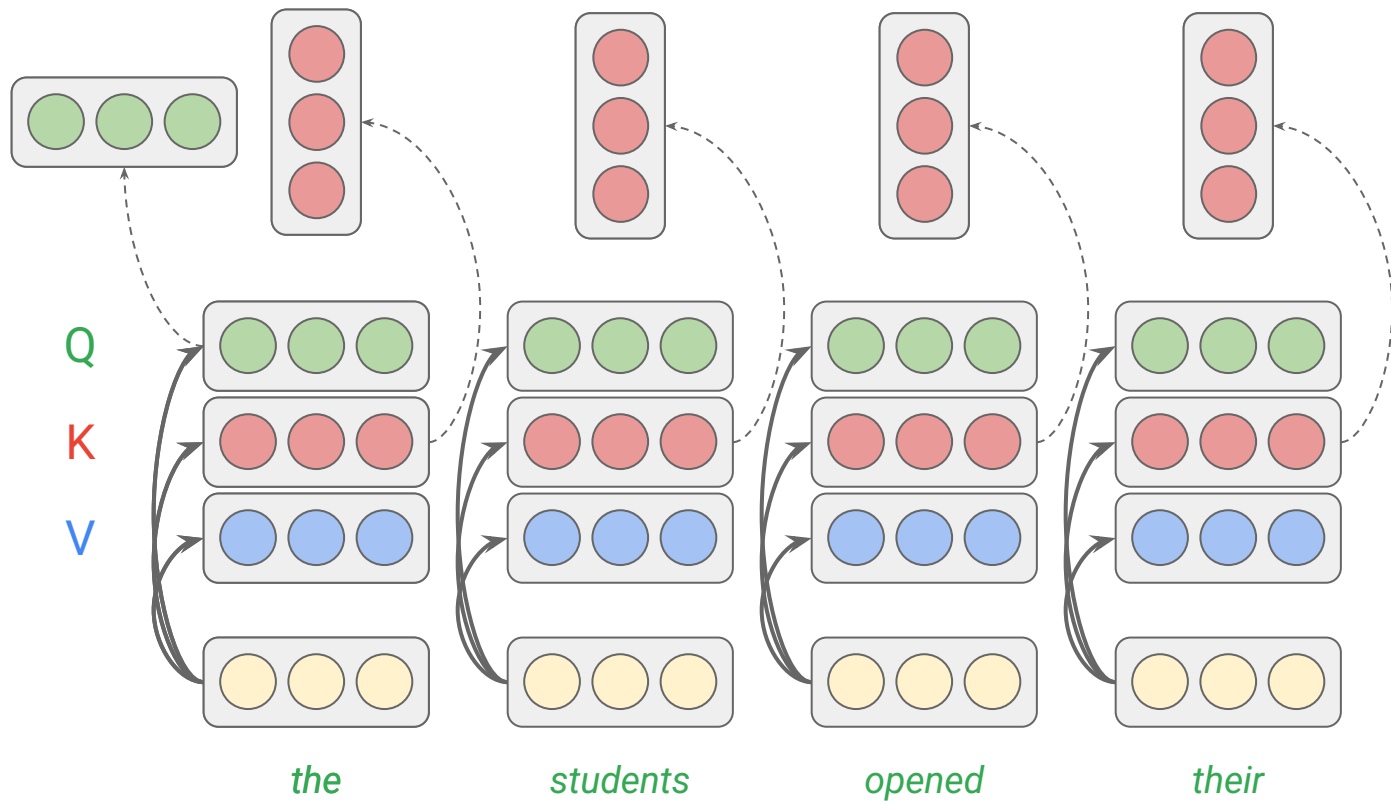
$$K = X \cdot W_K$$

$$V = X \cdot W_V$$

**linear  
projections**

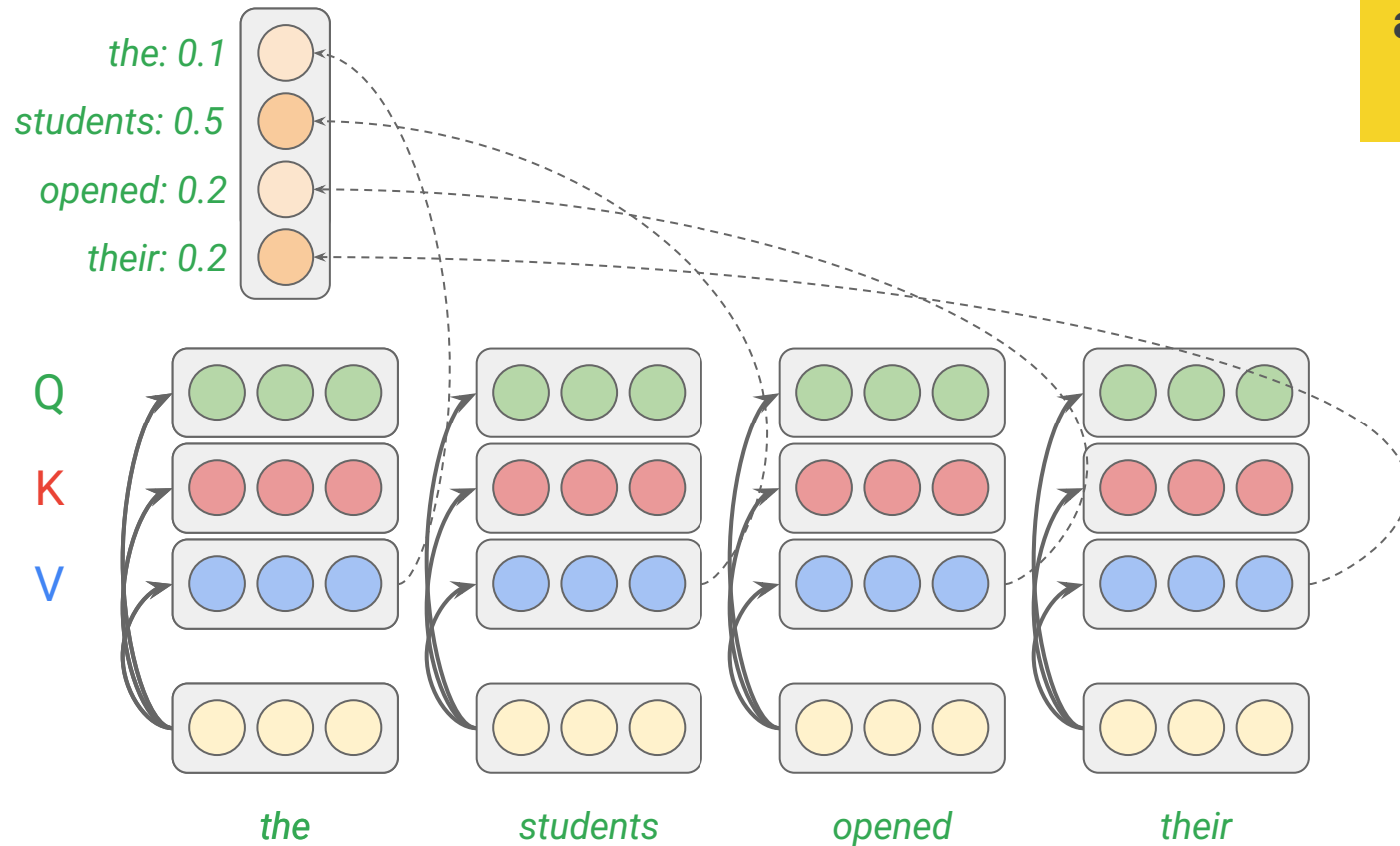
# Self-attention (cont'd)

*all computations  
are parallelized*



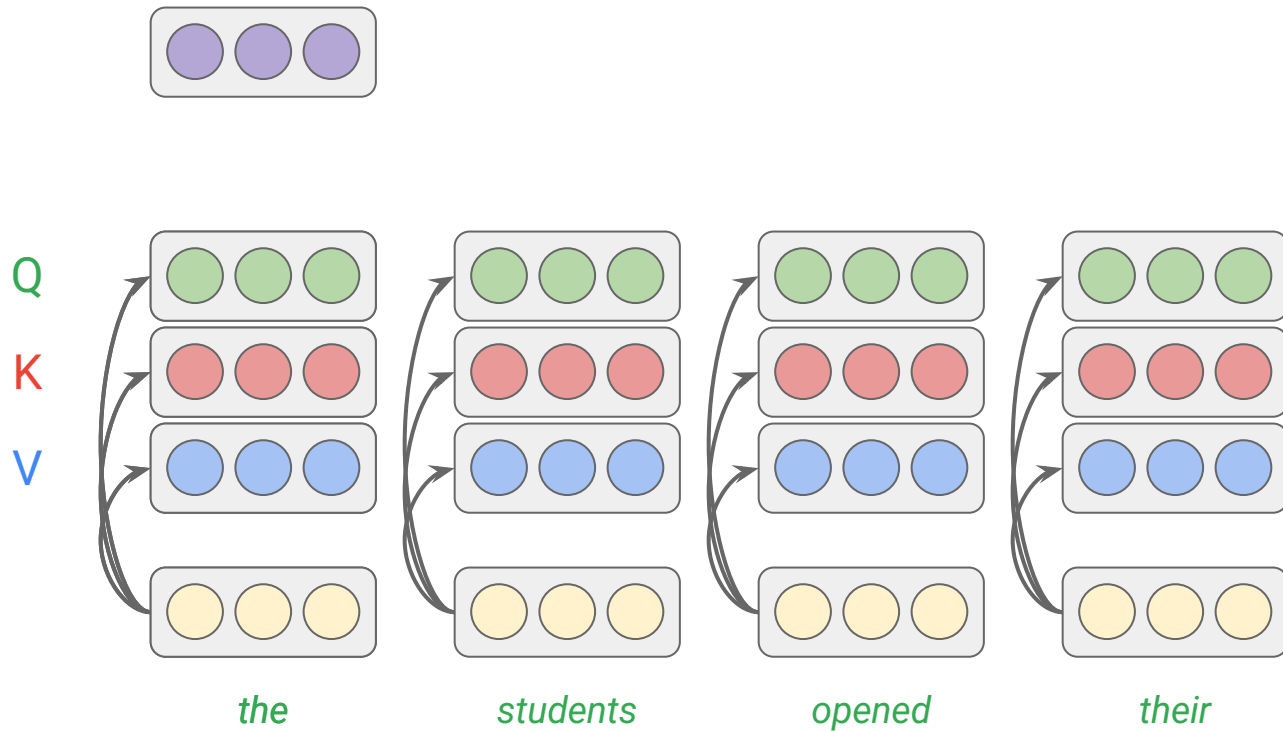
# Self-attention (cont'd)

*all computations  
are parallelized*



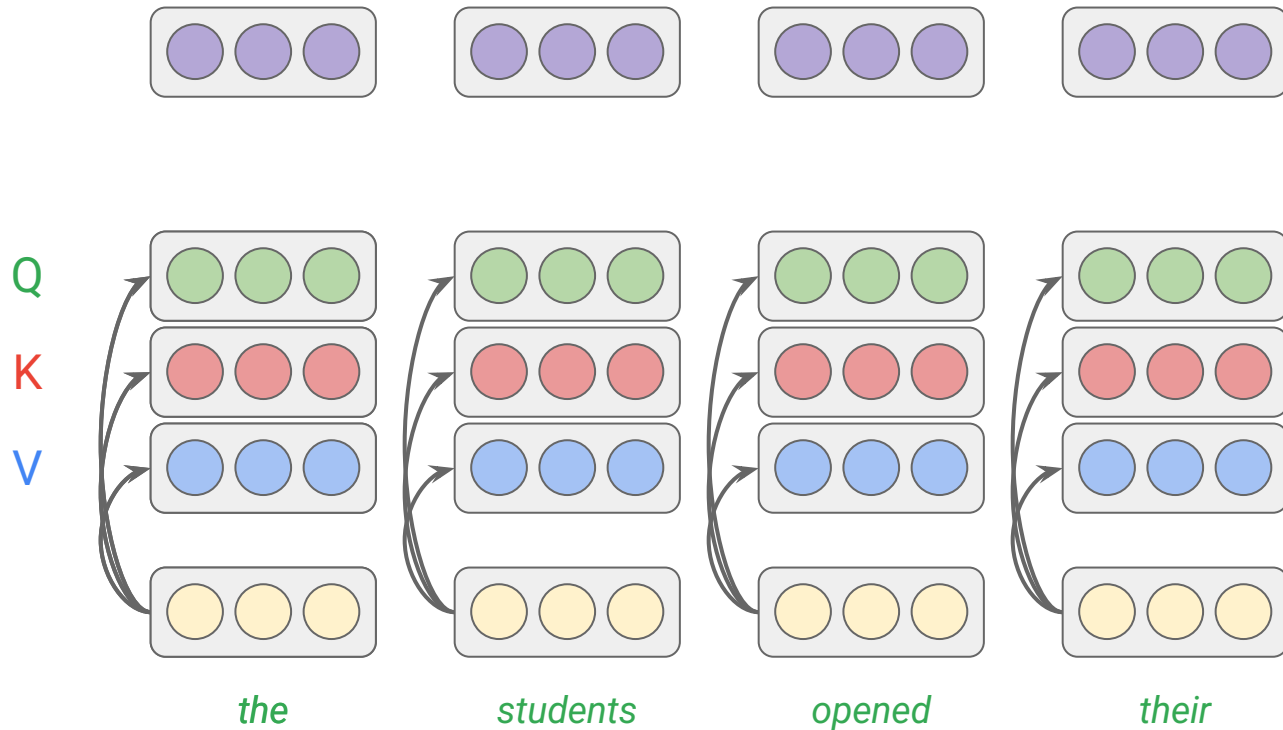
# Self-attention (cont'd)

*all computations  
are parallelized*

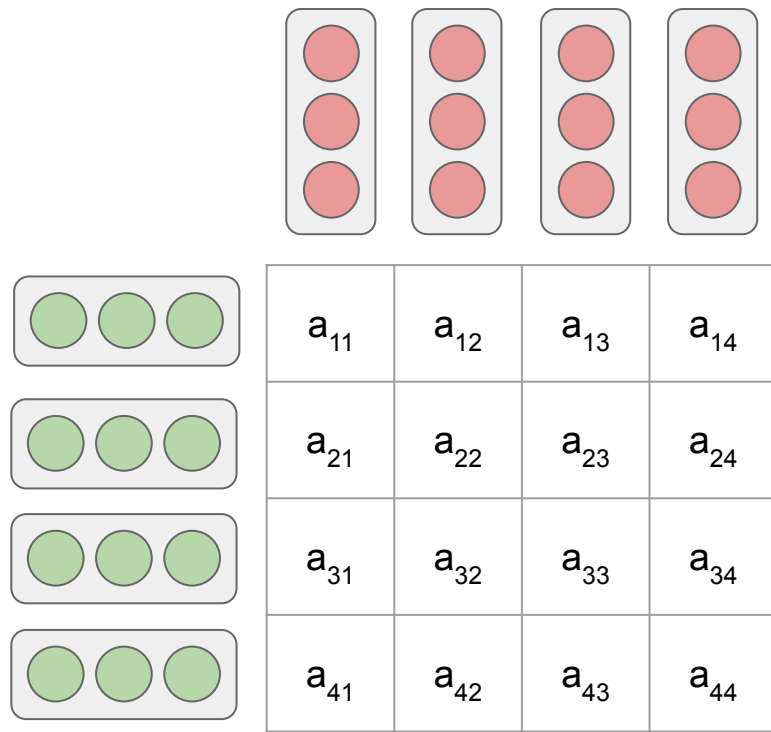


# Self-attention (cont'd)

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

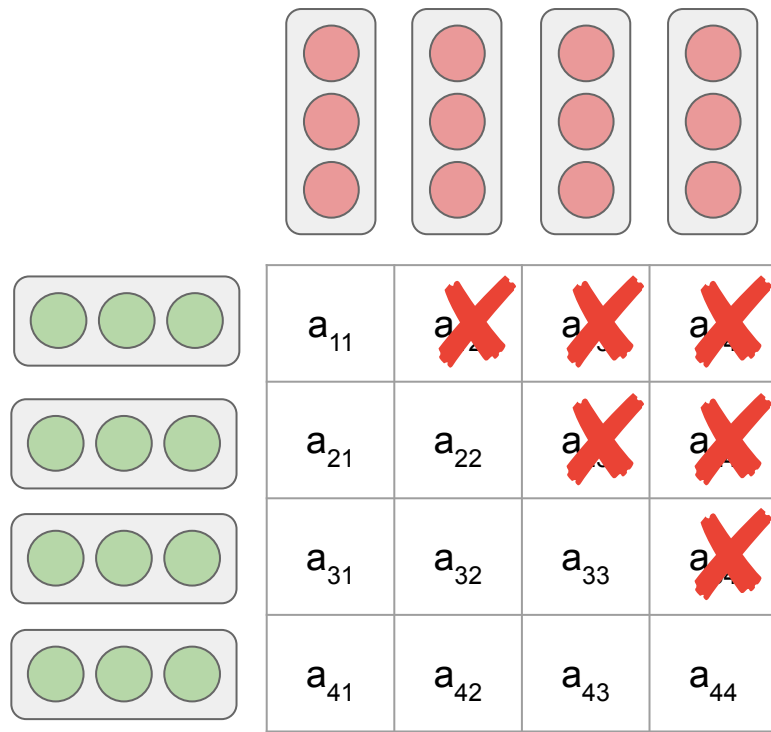


# Quadratic complexity



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

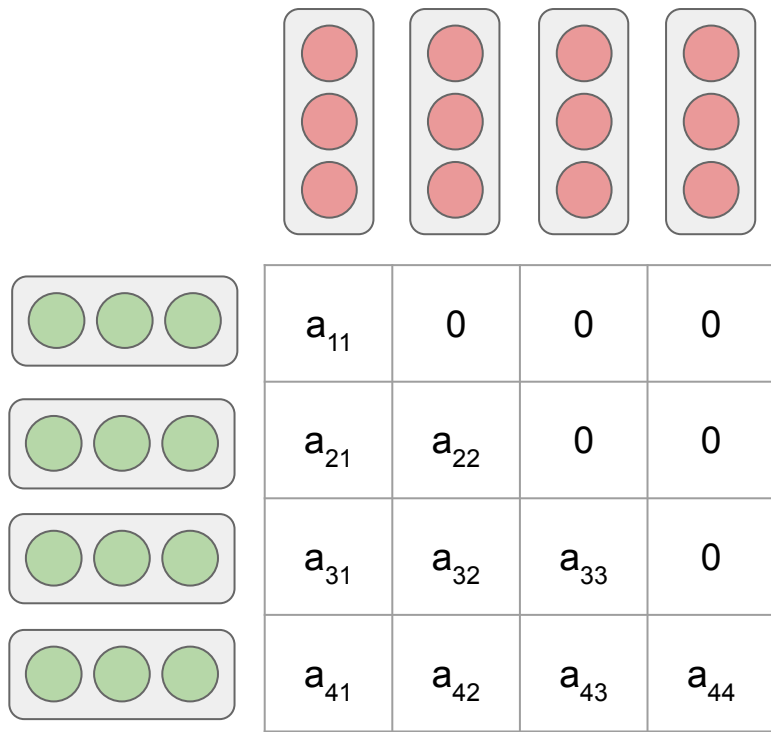
# Self-attention in the decoder



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

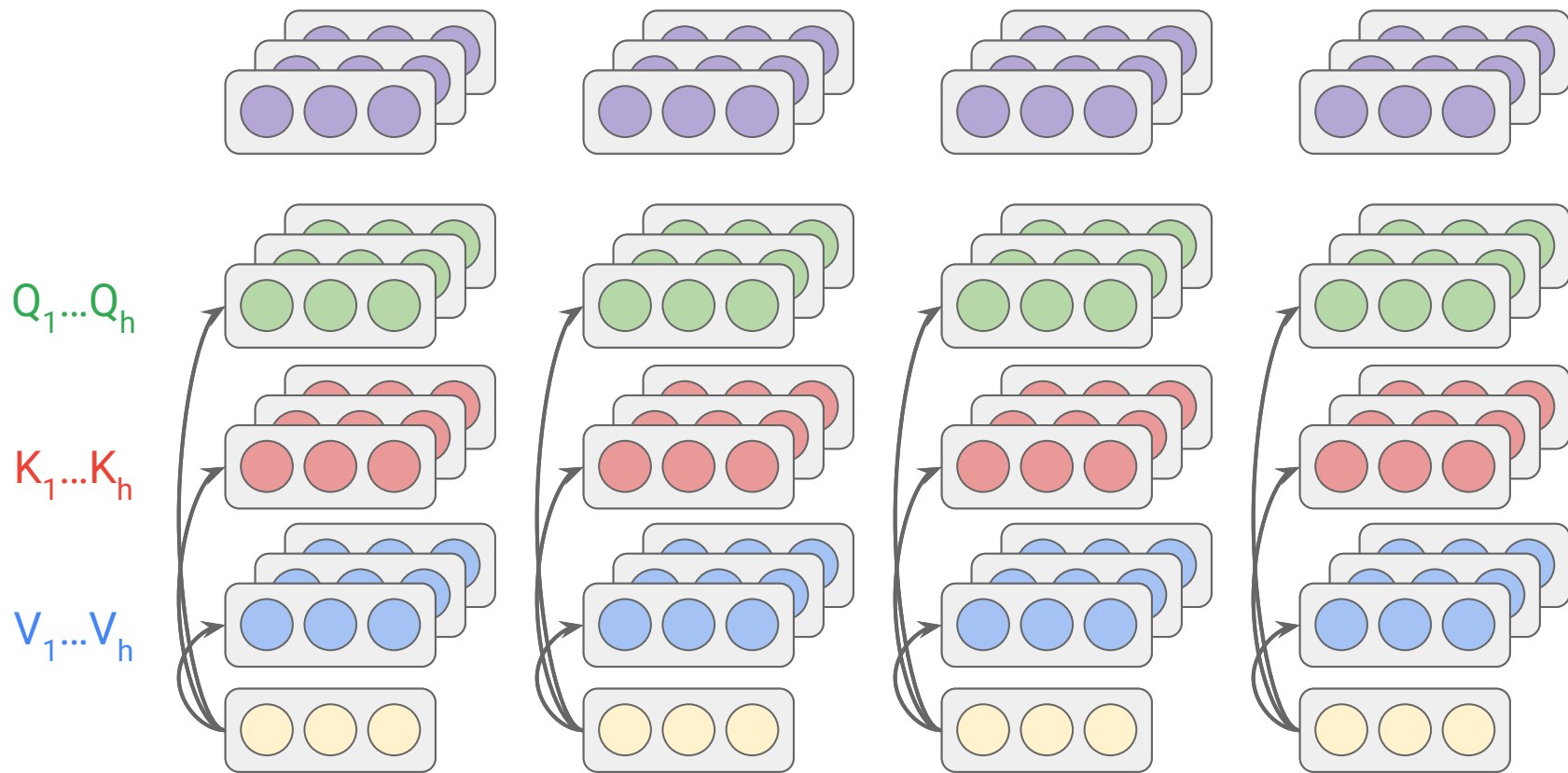


# Self-attention in the decoder (cont'd)

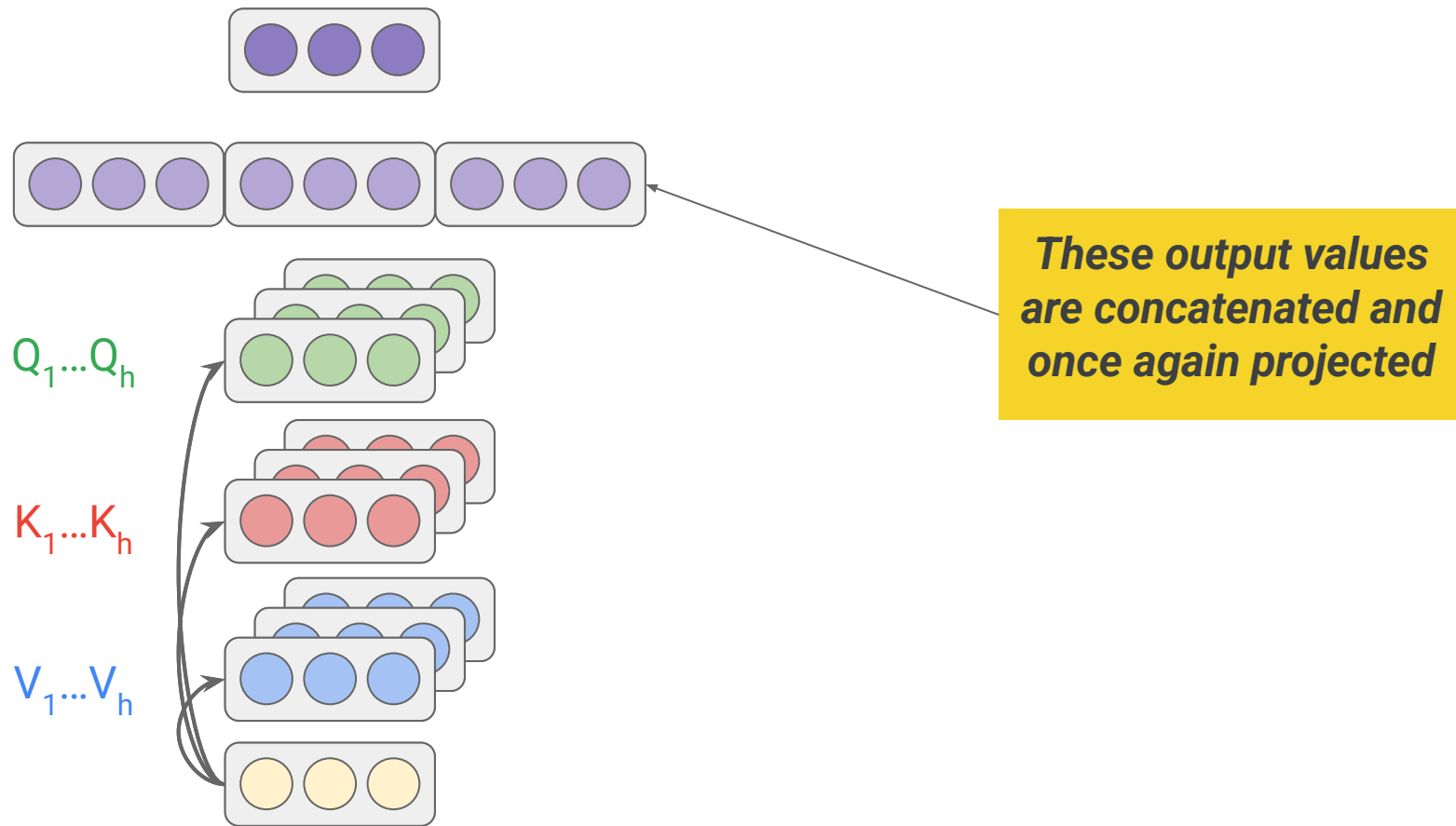


*masking out all values in the input of the softmax which correspond to illegal connections*

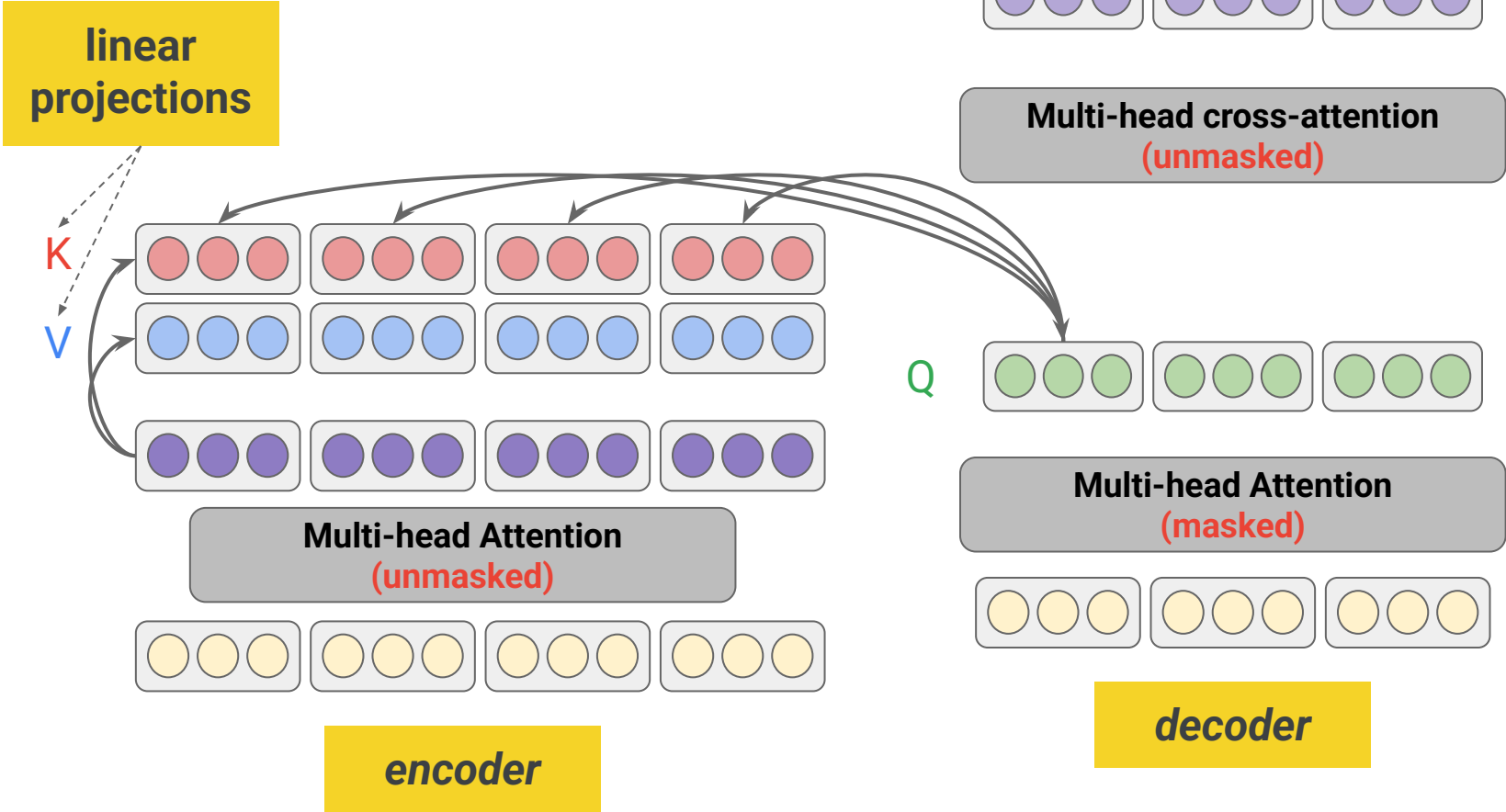
# Multi-head attention



## Multi-head attention (cont'd)

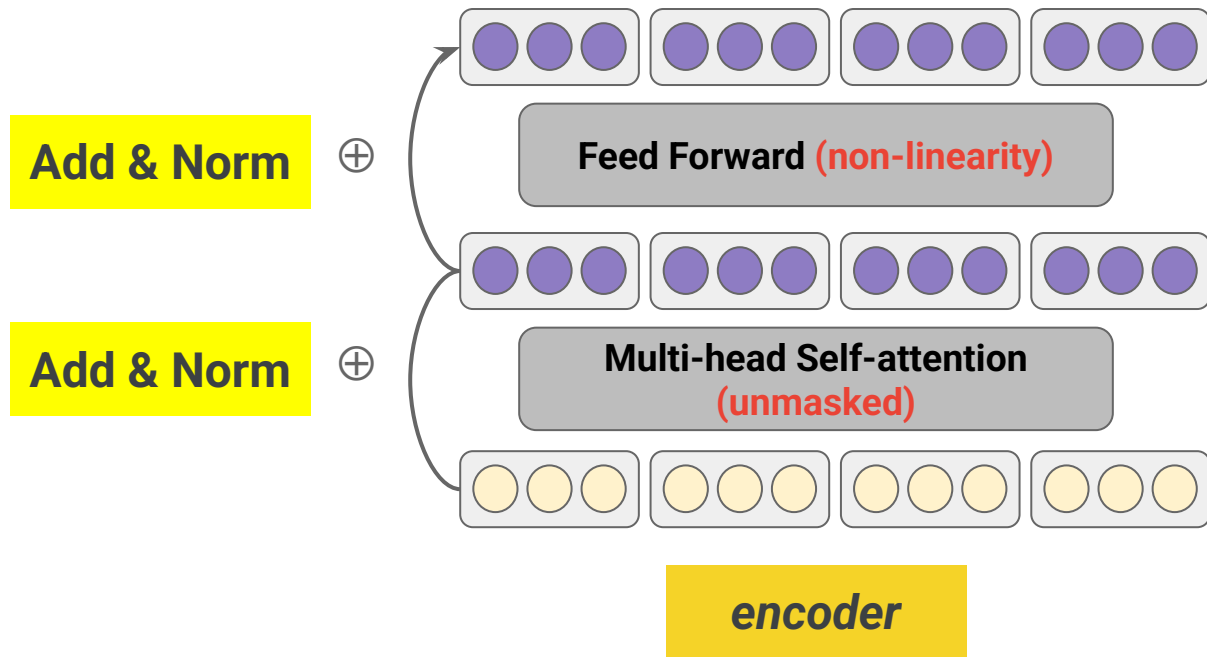


# Cross-attention in the decoder

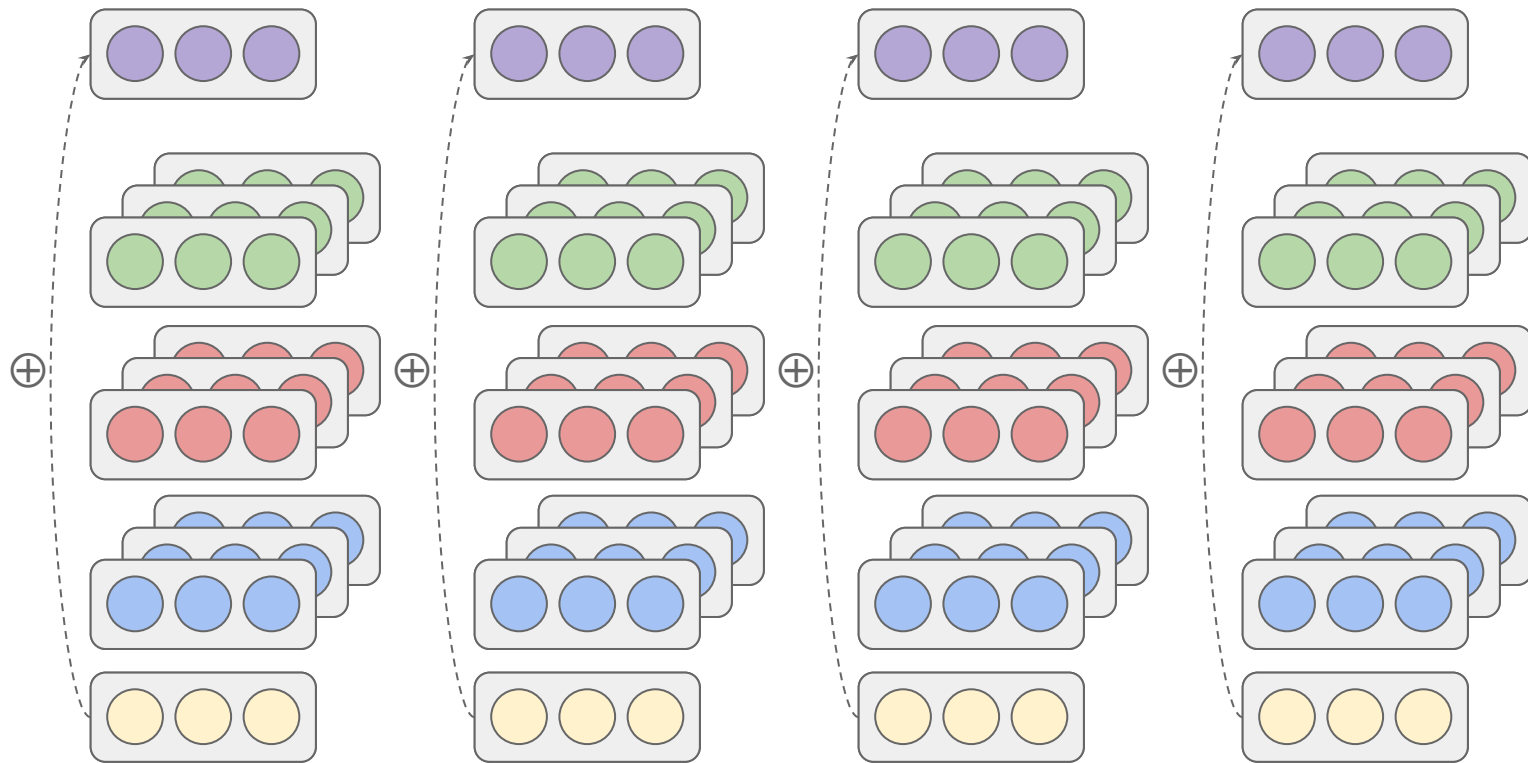




# Encoder (one layer)



# Residual connection



## Residual connection

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



# Layer normalization

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

where  $\mu$  and  $\sigma$  are the mean and standard deviation of the activations, and  $\gamma, \beta$  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

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

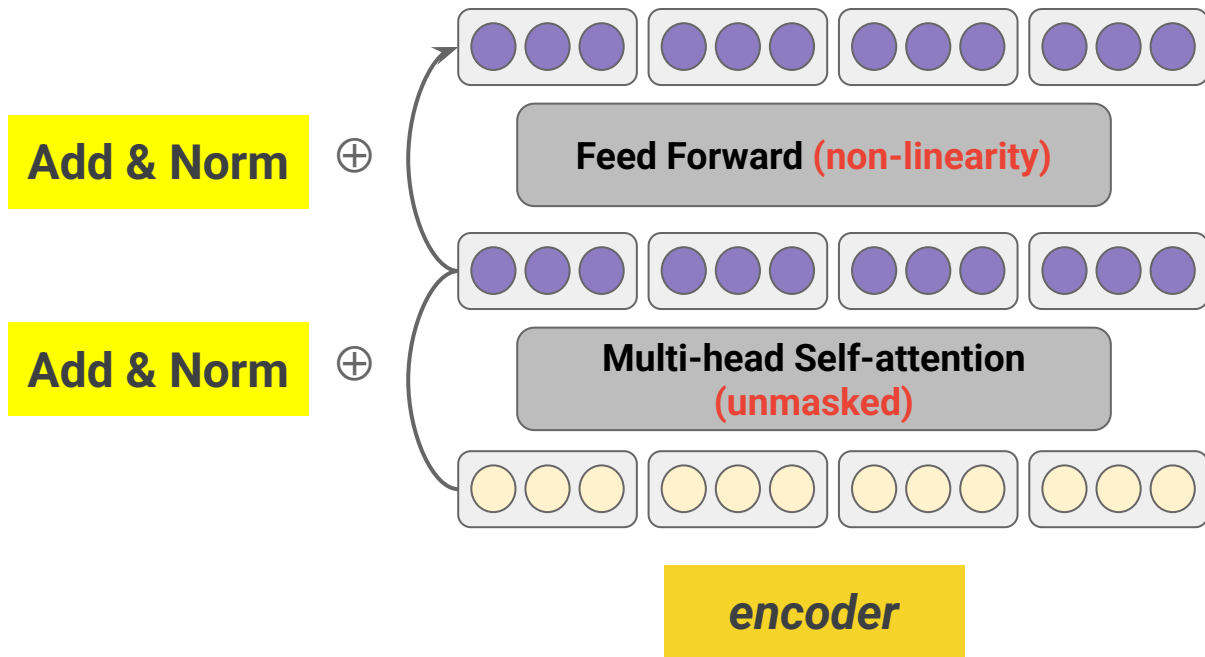
# Position-wise Feed-Forward Networks

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

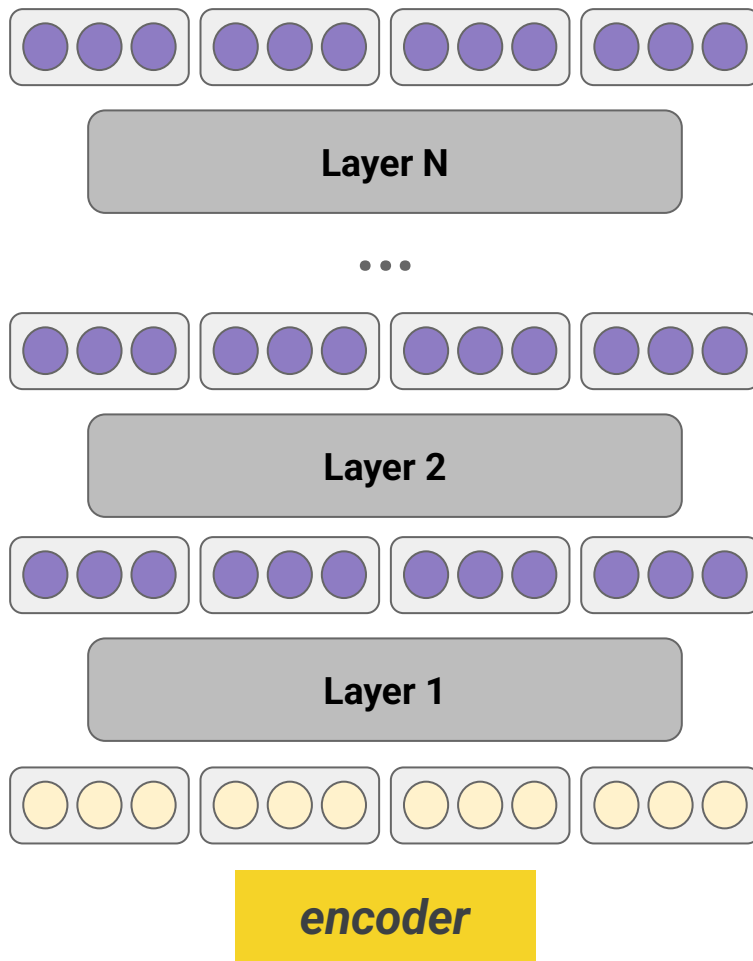


**ReLU (Rectified  
Linear Unit)**

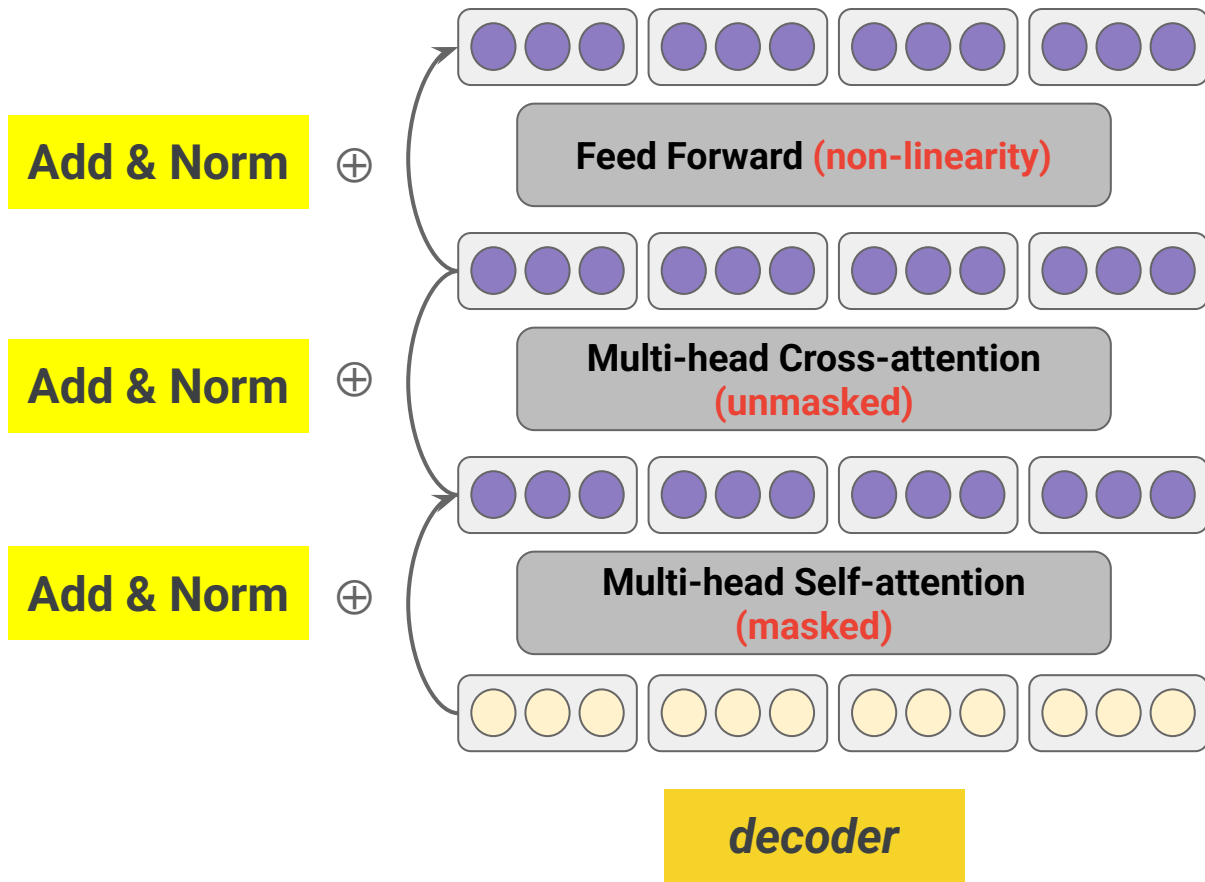
# Encoder (one layer)



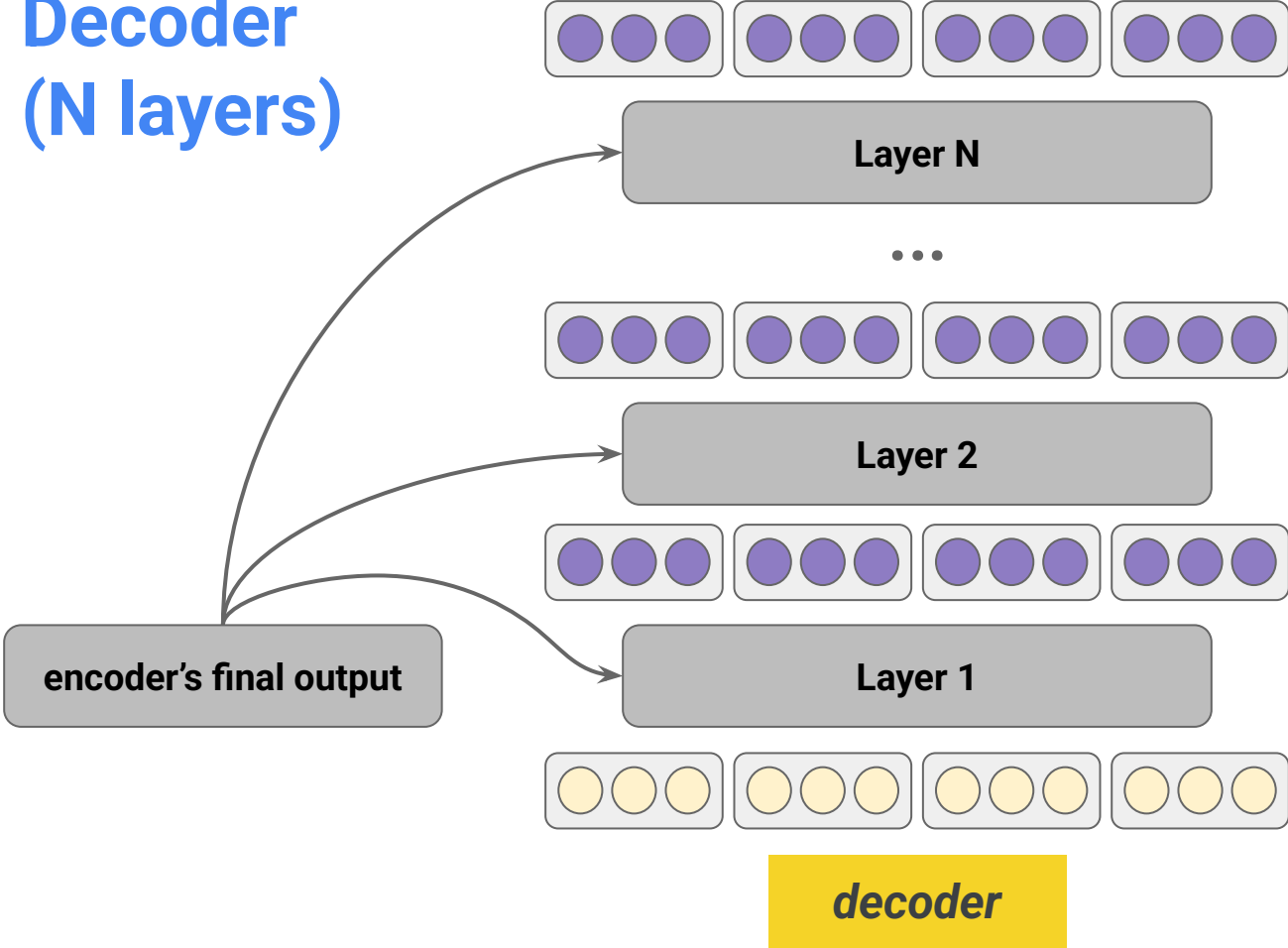
# Encoder (N layers)



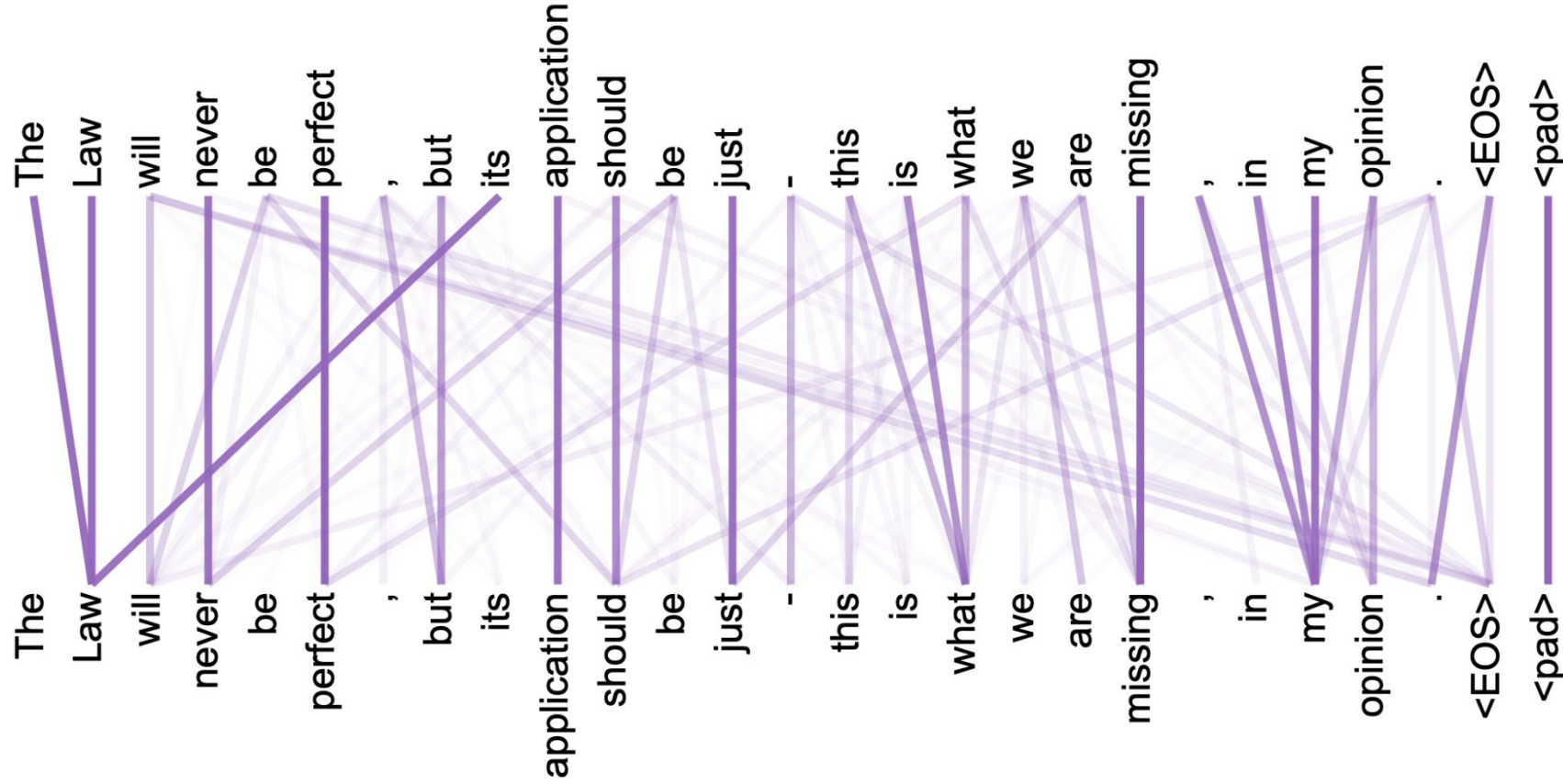
# Decoder (one layer)



**Decoder  
(N layers)**

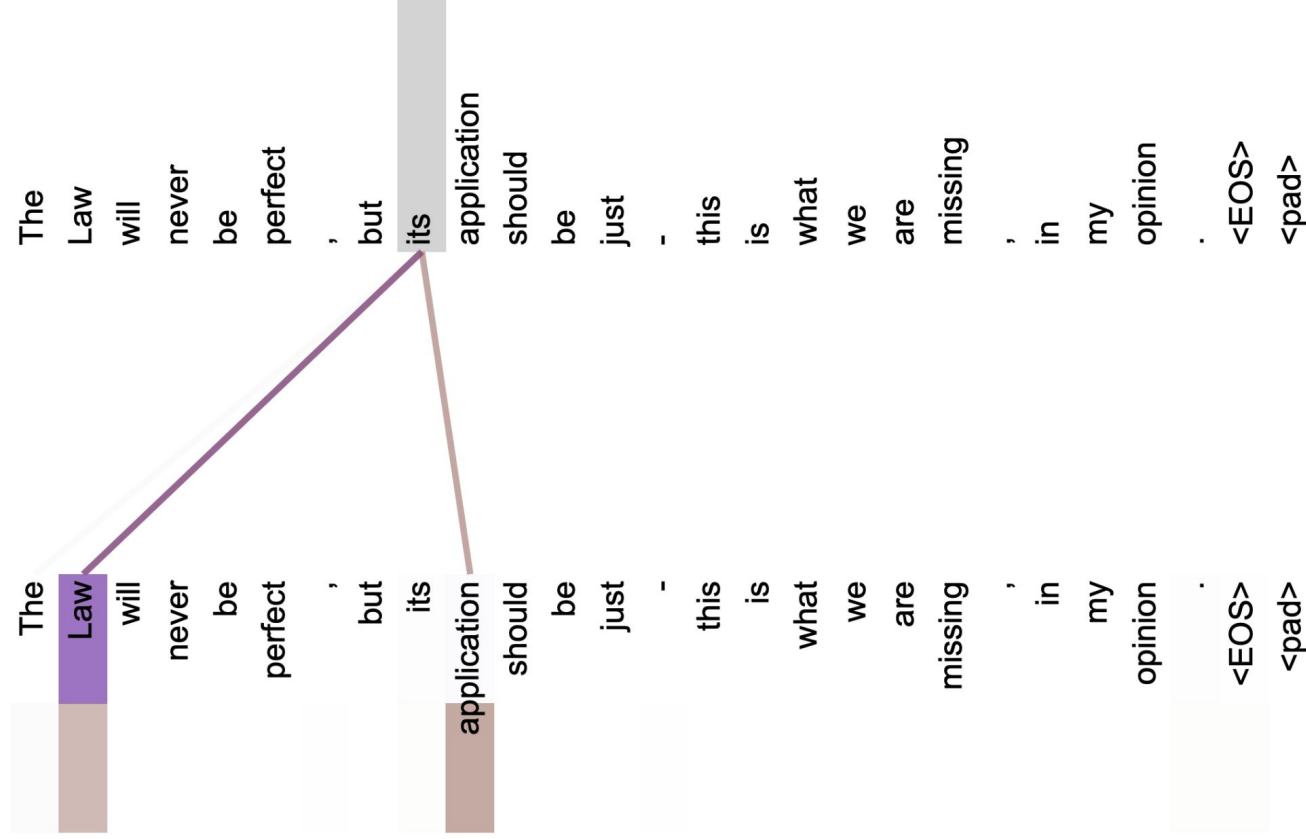


# Attention visualizations





# Attention visualizations (cont'd)

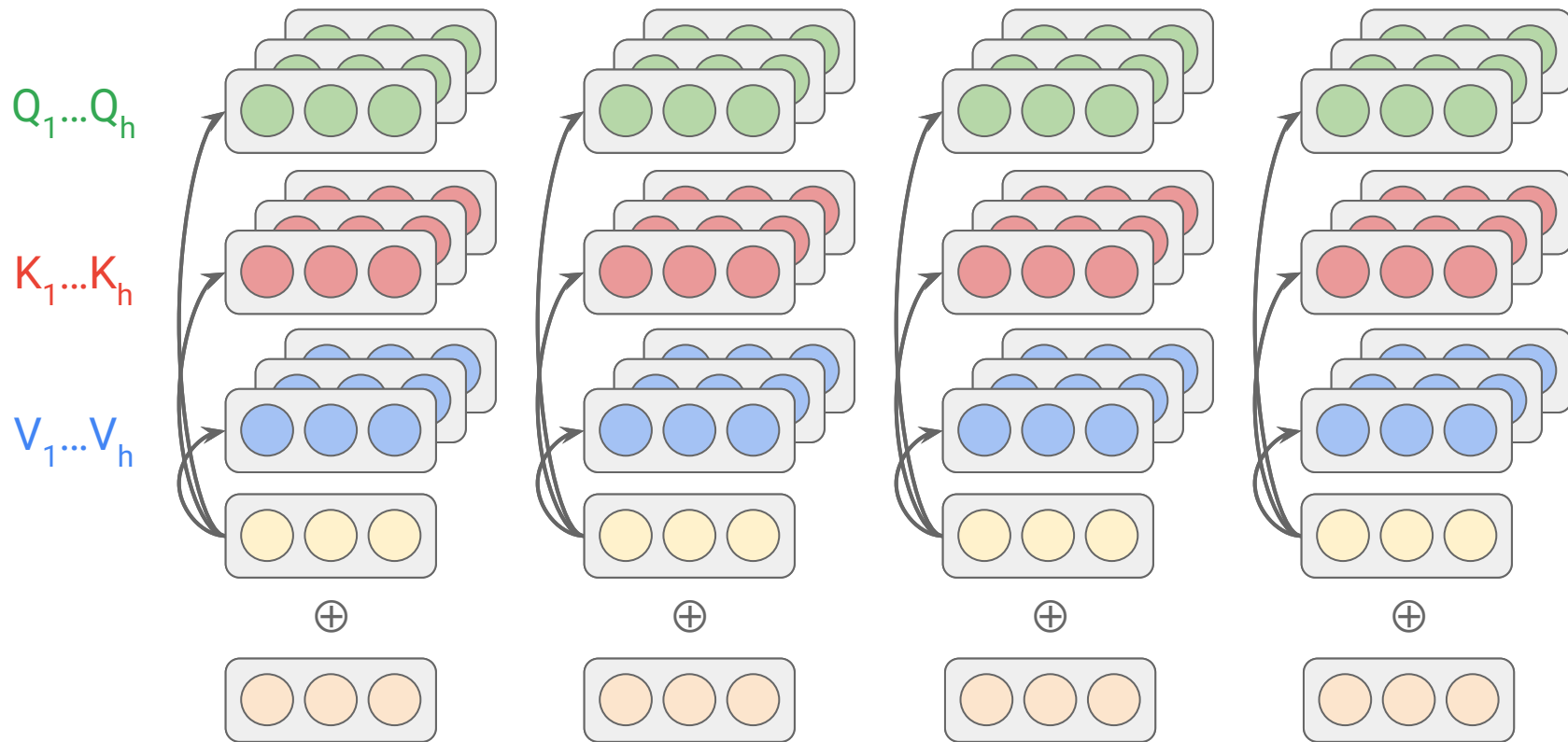


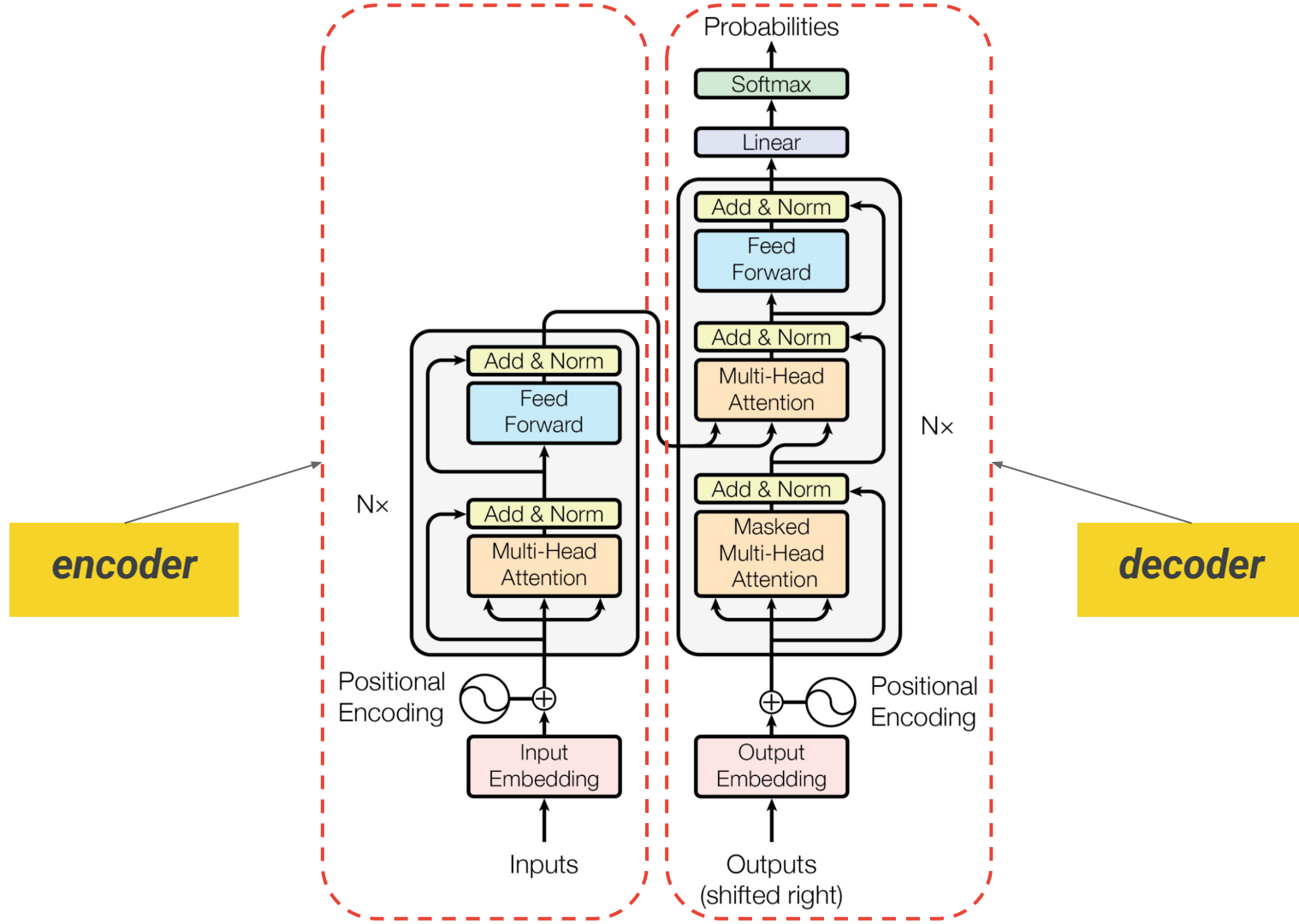
## Sinusoidal positional encoding

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

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

# Positional Encoding (cont'd)





# Different model architectures

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



Image created by Gemini

# **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

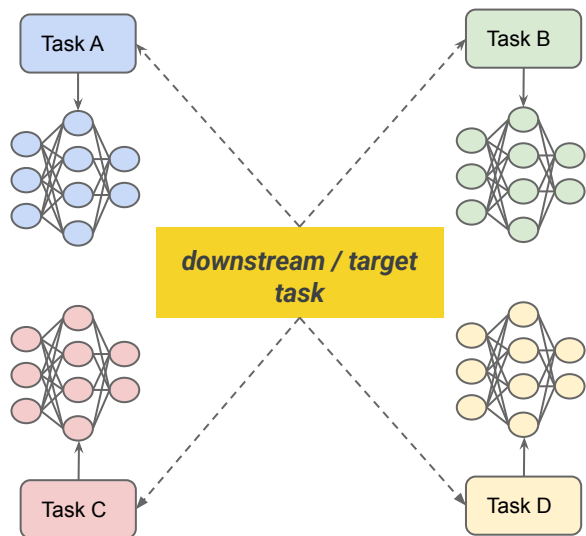
**Jacob Devlin   Ming-Wei Chang   Kenton Lee   Kristina Toutanova**

**Google AI Language**

`{jacobdevlin, mingweichang, kentonl, kristout}@google.com`

# A learning paradigm shift

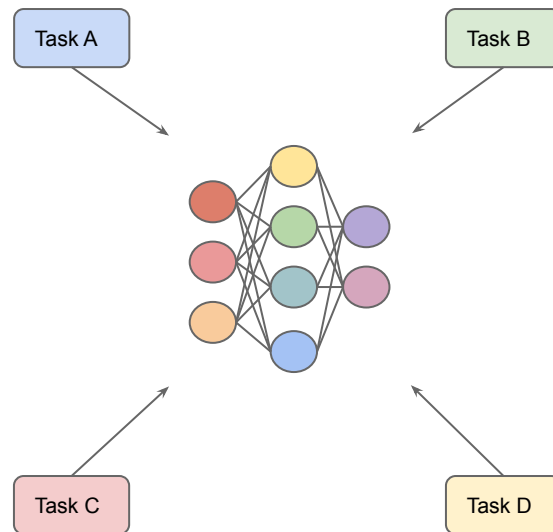
training task-specific models  
from scratch



before  
BERT



pretraining and then adapting



since  
BERT



Image created by Gemini

## Deep contextualized word representations

**Matthew E. Peters<sup>†</sup>, Mark Neumann<sup>†</sup>, Mohit Iyyer<sup>†</sup>, Matt Gardner<sup>†</sup>,**  
`{matthewp, markn, mohiti, mattg}@allenai.org`

**Christopher Clark\*, Kenton Lee\*, Luke Zettlemoyer<sup>†\*</sup>**  
`{csquared, kentonl, lsz}@cs.washington.edu`

<sup>†</sup>Allen Institute for Artificial Intelligence

\*Paul G. Allen School of Computer Science & Engineering, University of Washington

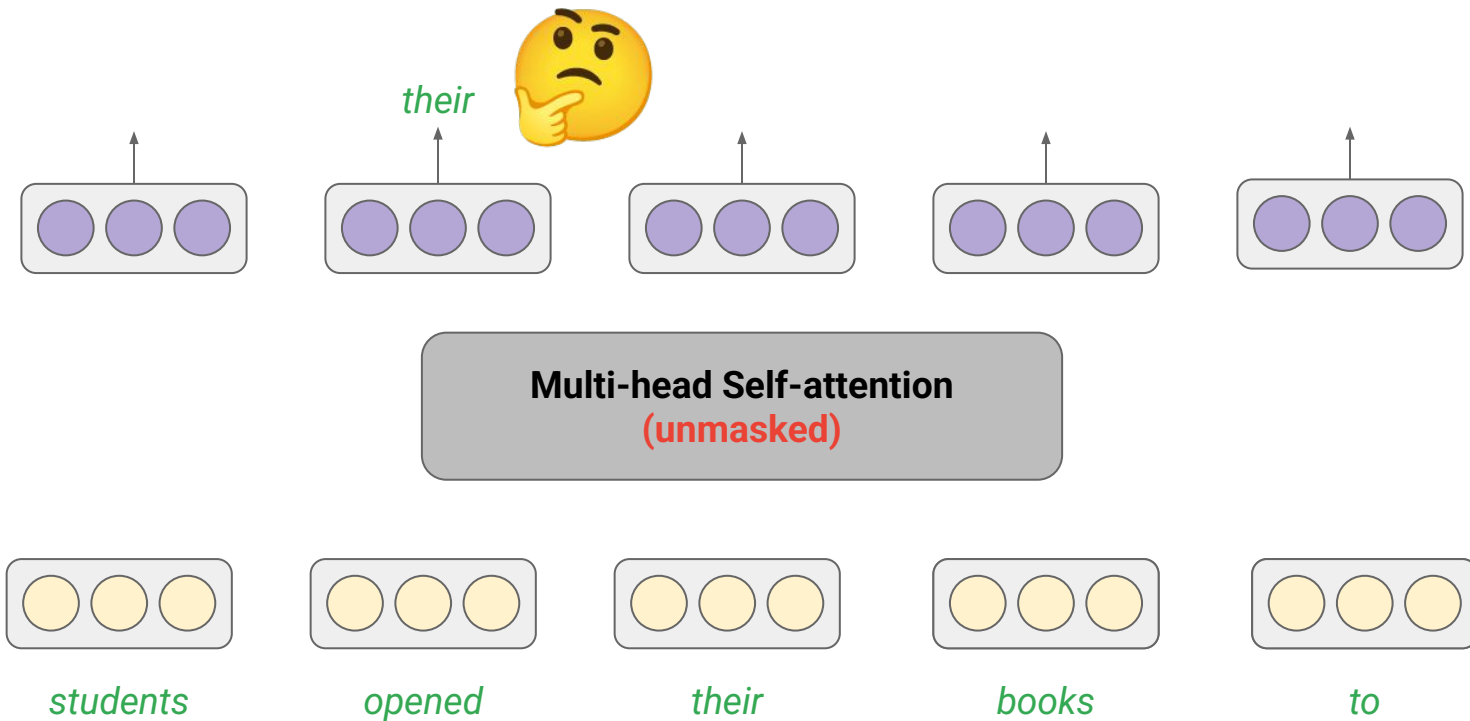


# BERT vs. ELMo

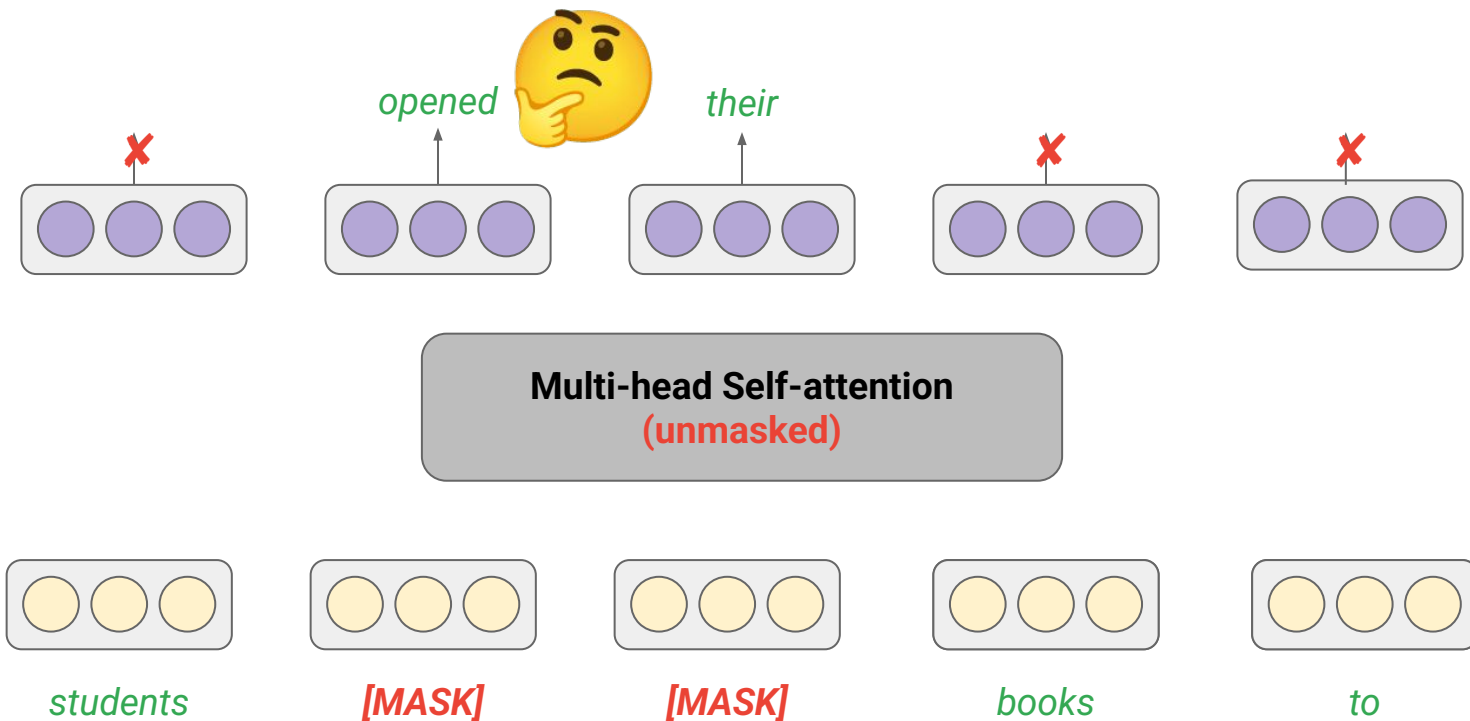
	BERT	ELMo
Model	Transformers	Bidirectional LSTM (Long Short-Term Memory, a variant of RNN)
Pre-training objective(s)	Masked language modeling + next sentence prediction	Left-to-right language modeling
Adaptation method	Fine-tuning	Feature-based (pretrained representations as additional features to task-specific models)

# Pretraining

# Language modeling using a Transformer encoder

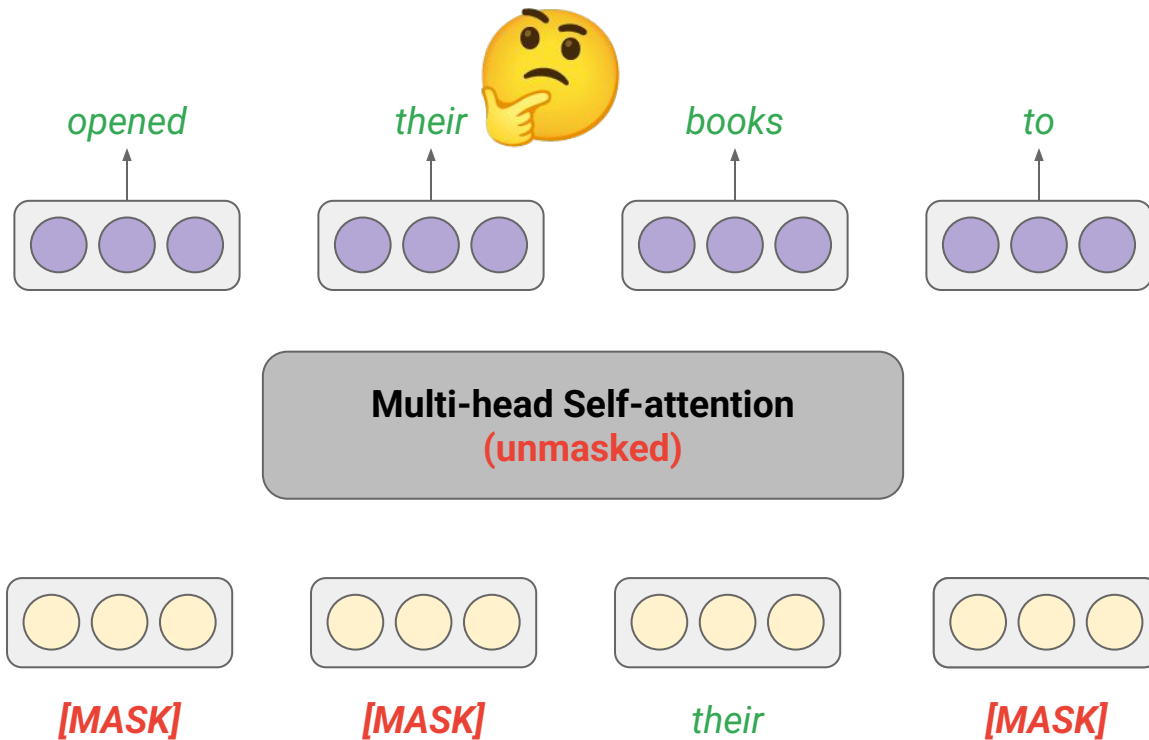


# Masked language modeling



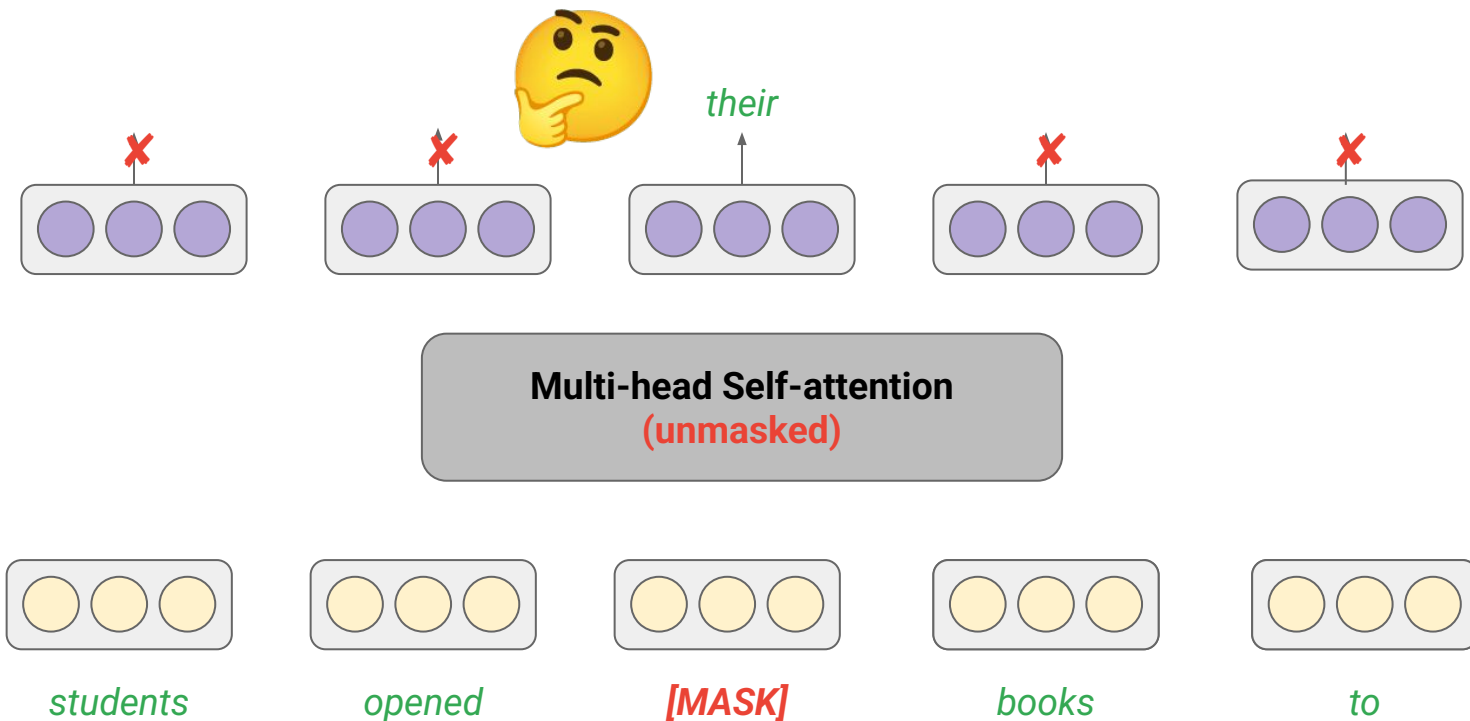
**15% - 30% of all tokens in each sequence are masked at random**

# What if we mask more tokens?



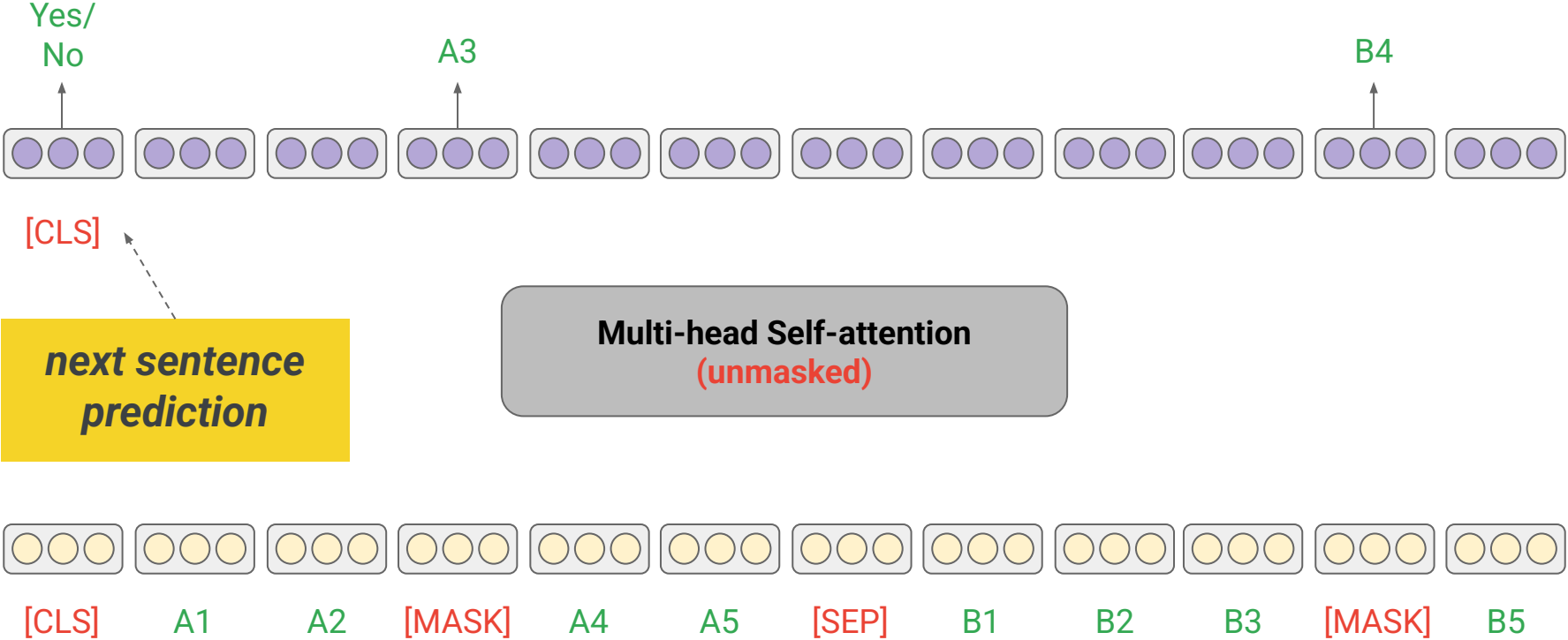
**15% - 30% of all tokens in each sequence are masked at random**

# What if we mask less tokens?

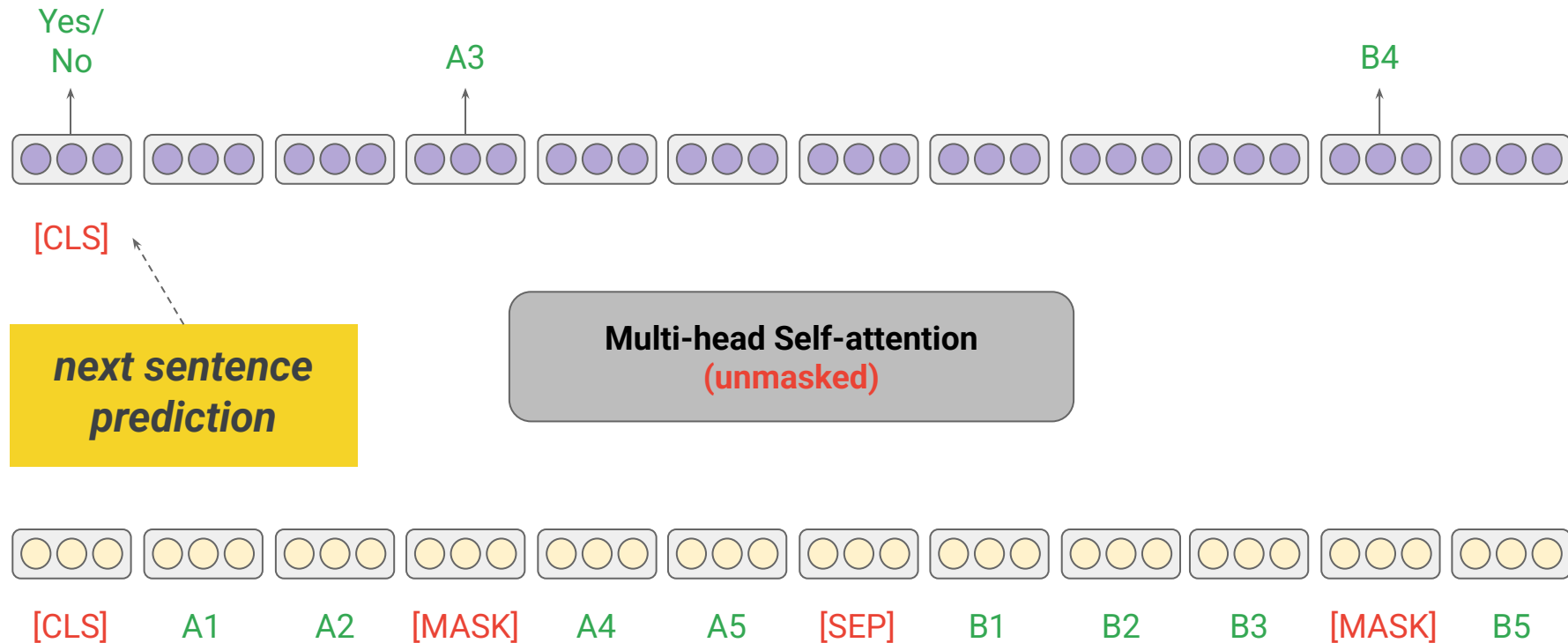


**15% - 30% of all tokens in each sequence are masked at random**

# CLS & SEP tokens

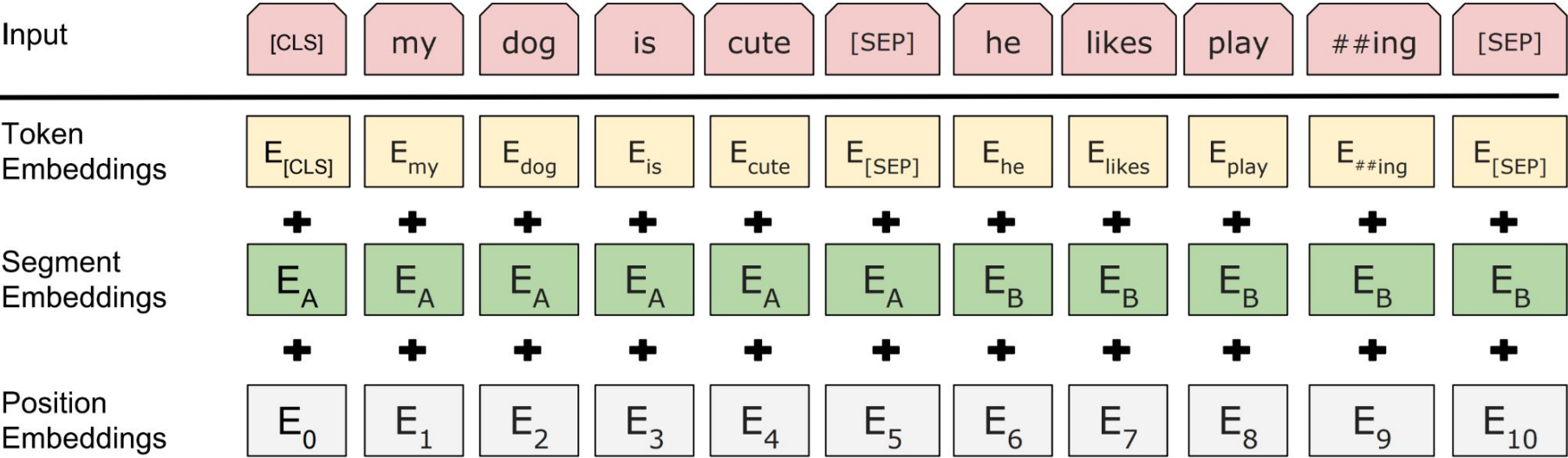


# BERT Pretraining

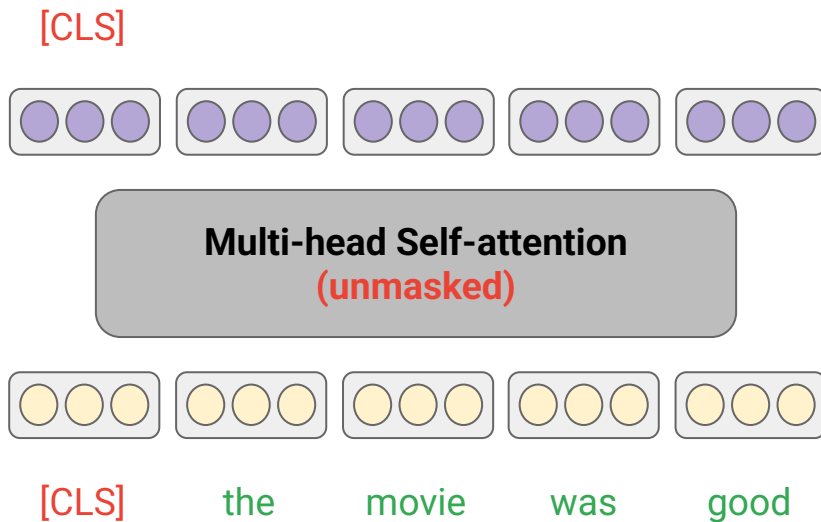
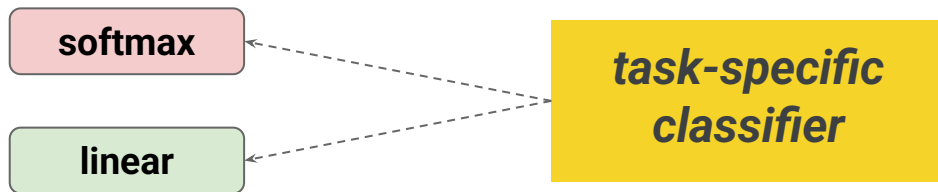




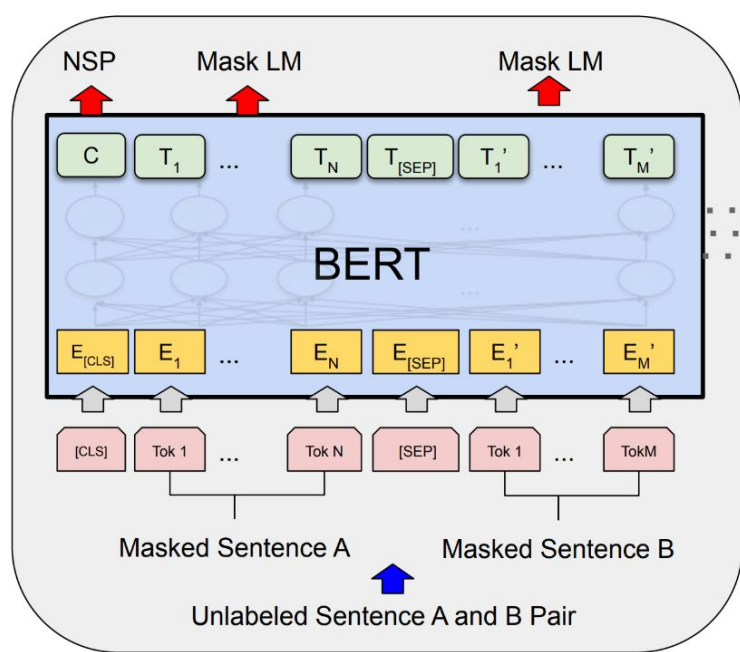
# BERT input representation



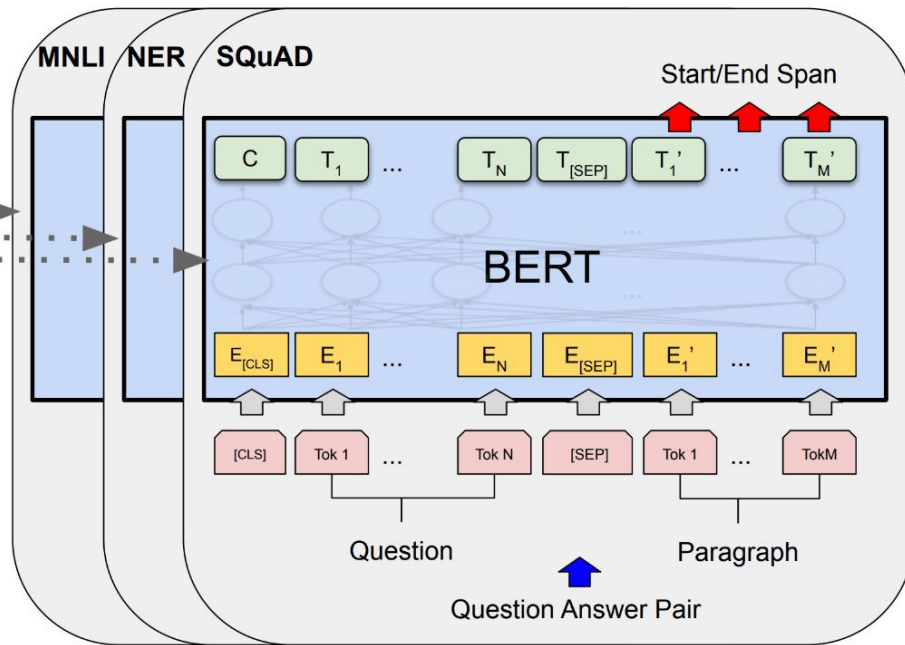
# BERT Fine-tuning



# BERT Pretraining & Fine-tuning



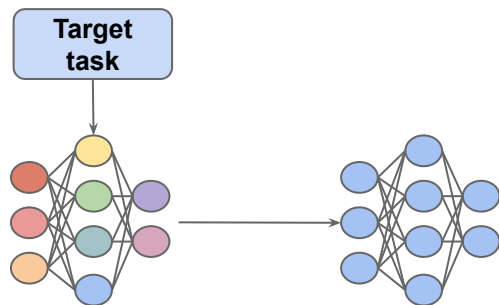
Pre-training



Fine-Tuning

# Intermediate-task transfer / fine-tuning

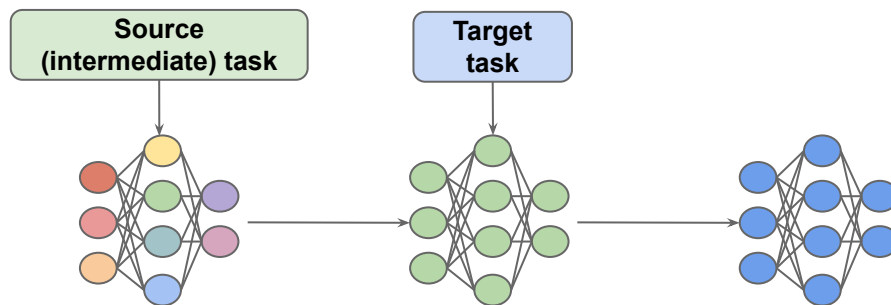
## Standard Fine-tuning



**BERT<sub>BASE</sub>** → **RTE**:  $63.5 \pm 2.3$

**BERT<sub>LARGE</sub>** → **RTE**:  $68.6 \pm 7.2$

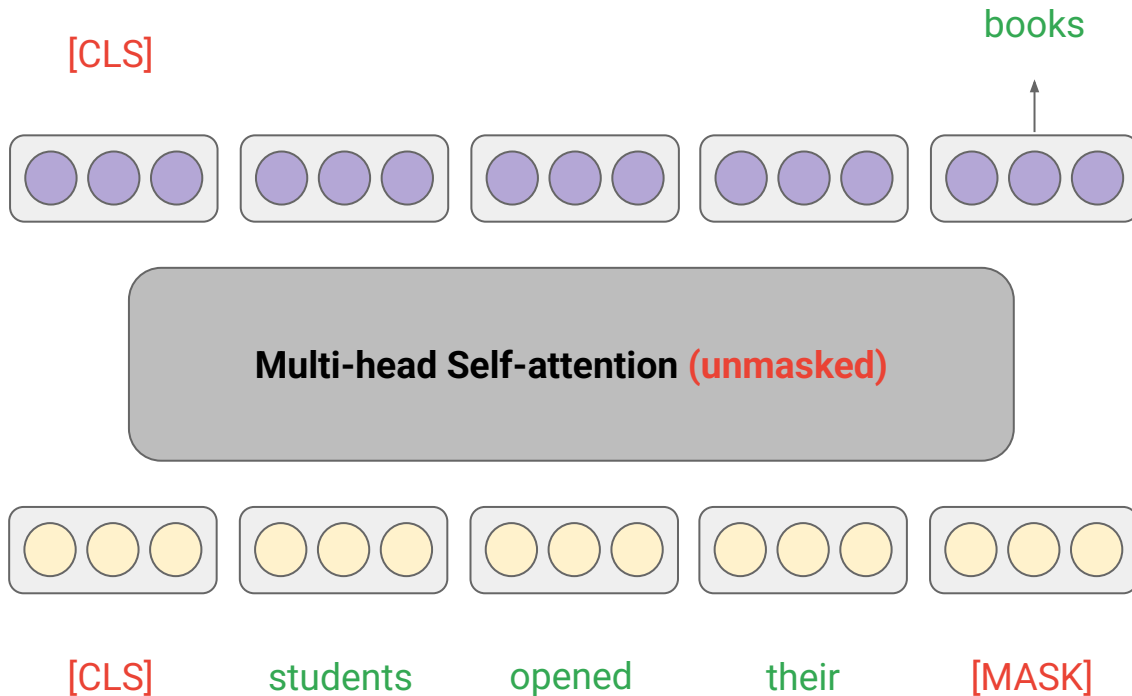
## Intermediate-task transfer



**BERT<sub>BASE</sub>** → **MNLI** → **RTE**:  $78.1 \pm 1.9$

**BERT<sub>LARGE</sub>** → **MNLI** → **RTE**:  $82.3 \pm 1.4$

# Can BERT be used for text generation?



*iterative  
masking and  
unmasking*

# T5: Text-to-Text Transfer Transformer

## Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer

**Colin Raffel\***

CRAFFEL@GMAIL.COM

**Noam Shazeer\***

NOAM@GOOGLE.COM

**Adam Roberts\***

ADAROB@GOOGLE.COM

**Katherine Lee\***

KATHERINELEE@GOOGLE.COM

**Sharan Narang**

SHARANNARANG@GOOGLE.COM

**Michael Matena**

MMATENA@GOOGLE.COM

**Yanqi Zhou**

YANQIZ@GOOGLE.COM

**Wei Li**

MWEILI@GOOGLE.COM

**Peter J. Liu**

PETERJLIU@GOOGLE.COM

*Google, Mountain View, CA 94043, USA*

# T5 Pretraining: Span corruption

**<X>, <Y>: sentinel tokens**

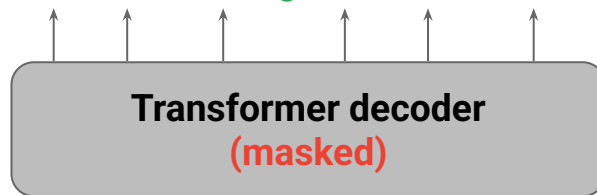


**Transformer encoder**  
(unmasked)

Thank you <X> me to your party <Y> week

*encoder*

<X> for inviting <Y> last <EOS>



<BOS> <X> for inviting <Y> last

*decoder*

Thank you ~~for inviting~~ me to your party ~~last~~ week

# T5 Fine-tuning



**Transformer encoder**  
(unmasked)

sentiment analysis: this movie was good

*encoder*

positive <EOS>

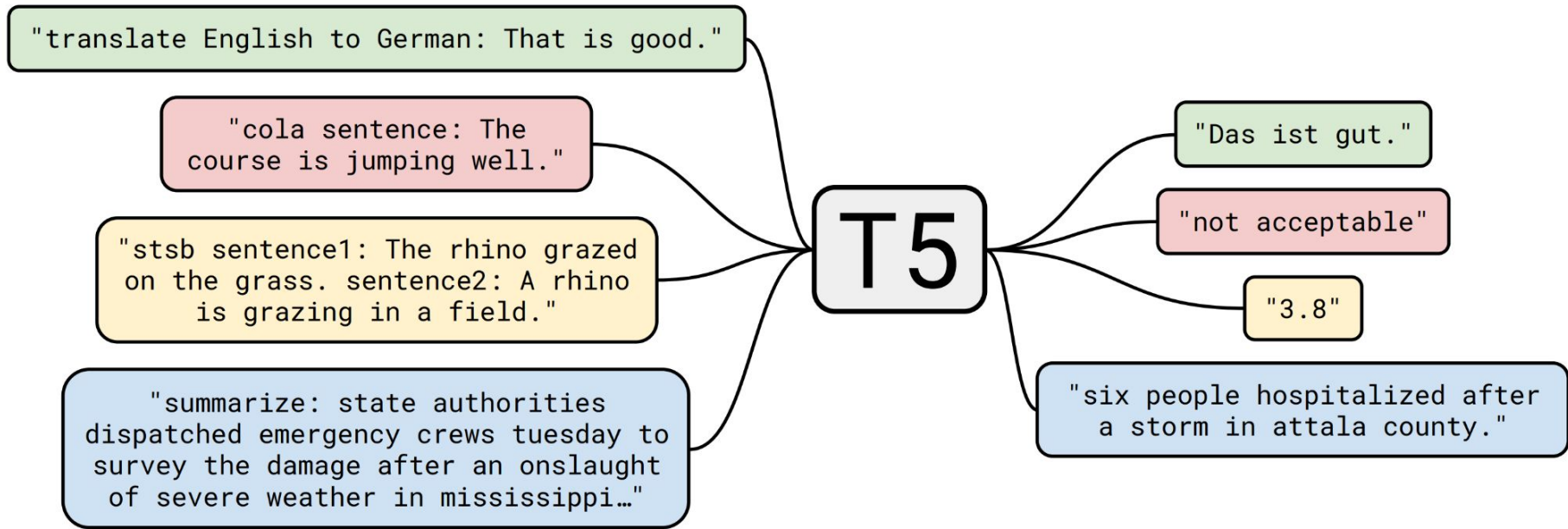
**Transformer decoder**  
(masked)

<BOS> positive

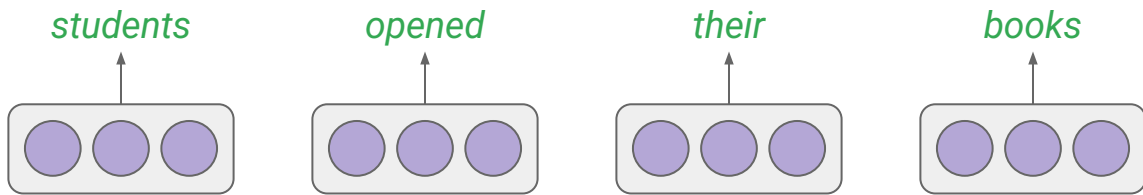
*decoder*



# T5 facilitates multitask learning

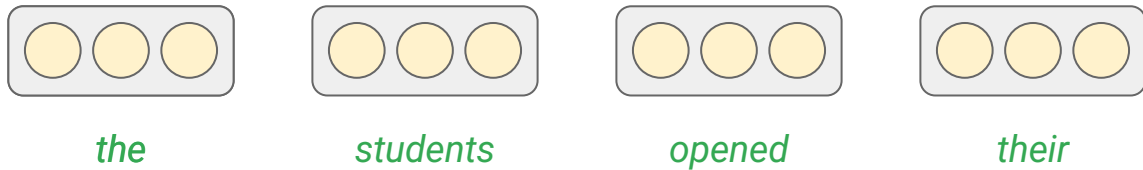


# Decoder-only model



***the architecture  
used in frontier  
LLMs***

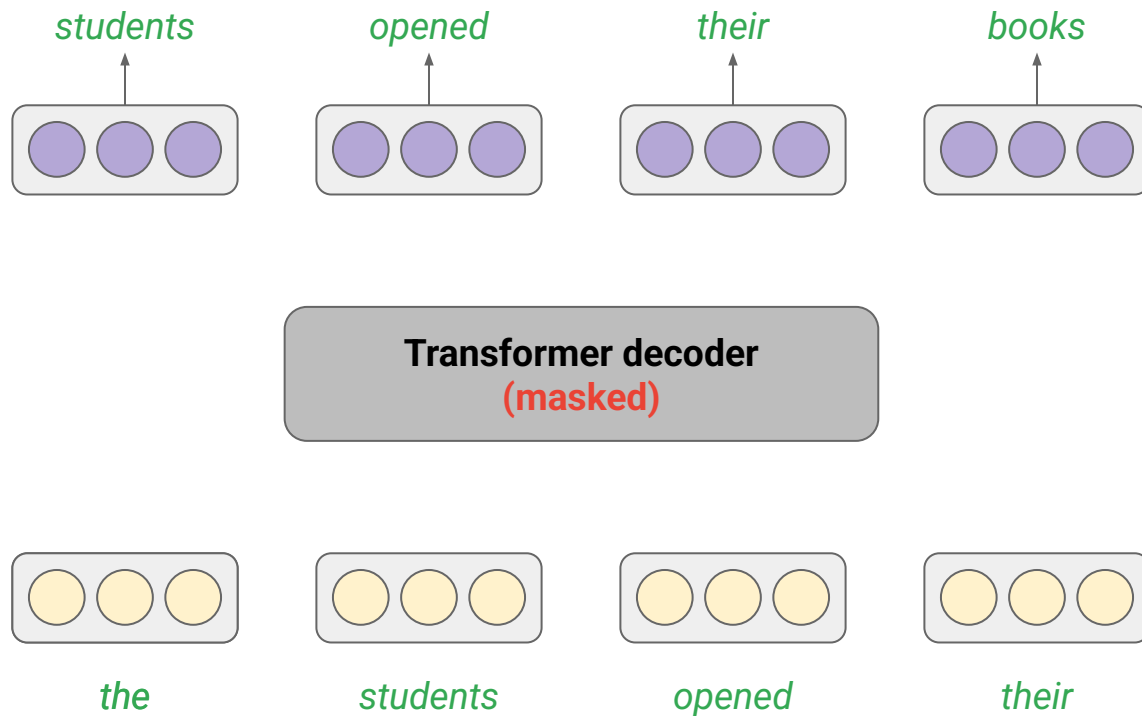
**Transformer decoder  
(masked)**



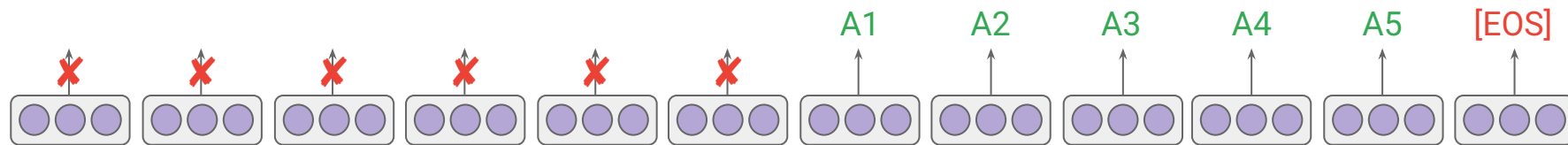
## Note on cross-attention

- Can be used to inject non-text data (e.g., images, structured data, or even sensor readings) into the model

# Pretraining with a causal LM (decoder-only)

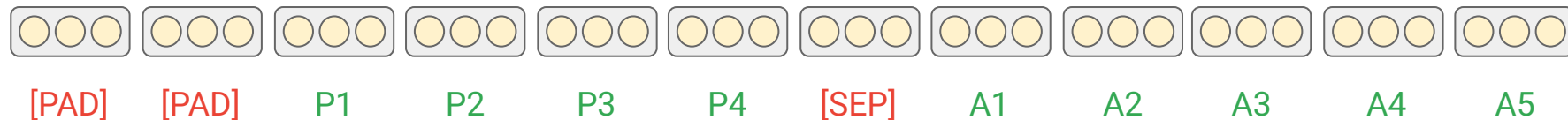


# Training with prefix LM (decoder-only)



*the architecture  
used in frontier  
LLMs*

**Transformer decoder**  
**(partially masked)**



# Different attention mask patterns

K

Q

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$
$x_1$							
$x_2$							
$x_3$							
$x_4$							
$x_5$							
$x_6$							
$x_7$							

***fully-visible***

K

Q

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$
$x_1$							
$x_2$							
$x_3$							
$x_4$							
$x_5$							
$x_6$							
$x_7$							

***causal***

# Different attention mask patterns (cont'd)

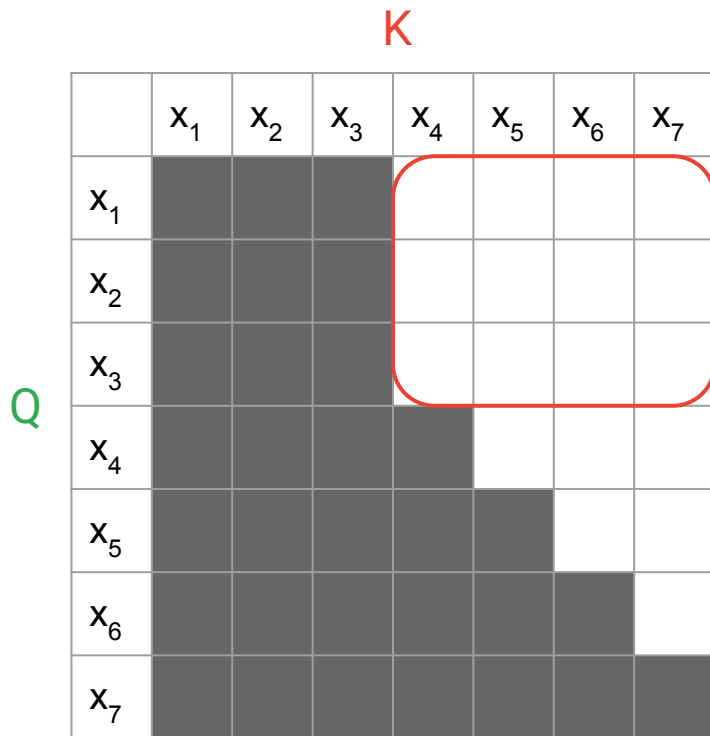
**K**

**Q**

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$
$x_1$							
$x_2$							
$x_3$							
$x_4$							
$x_5$							
$x_6$							
$x_7$							

***Prefix LM***

# Different attention mask patterns (cont'd)



***Why masking  
here?***



**Thank you!**