

The Era of BERT

CS 5624: Natural Language Processing



Spring 2025

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

Tu Vu

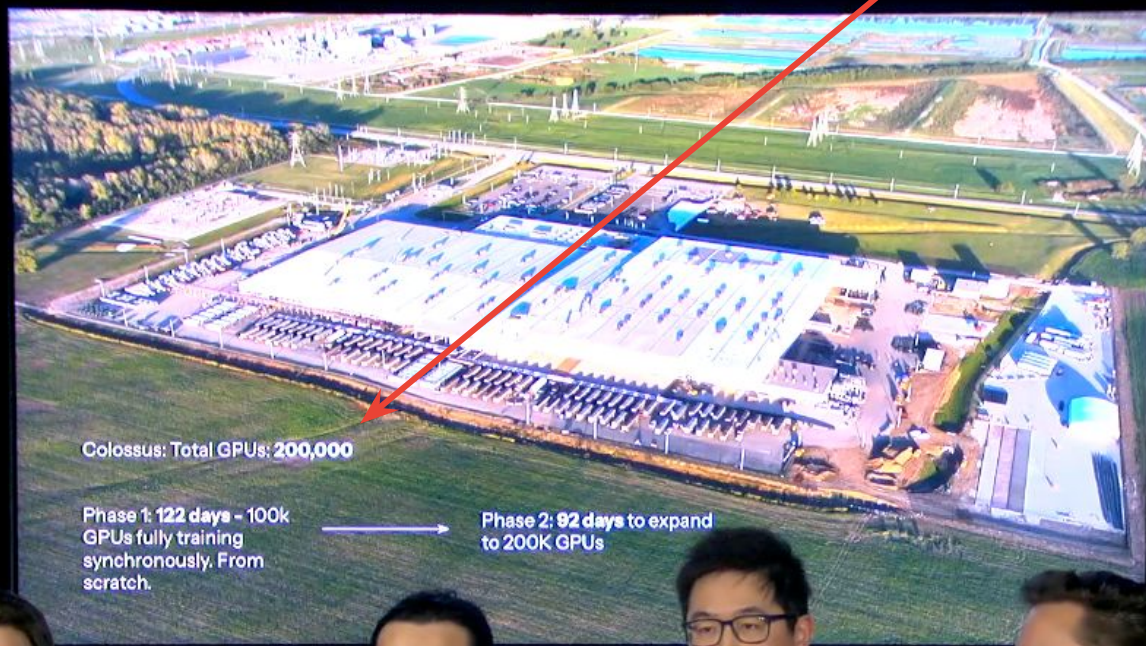


Logistics

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

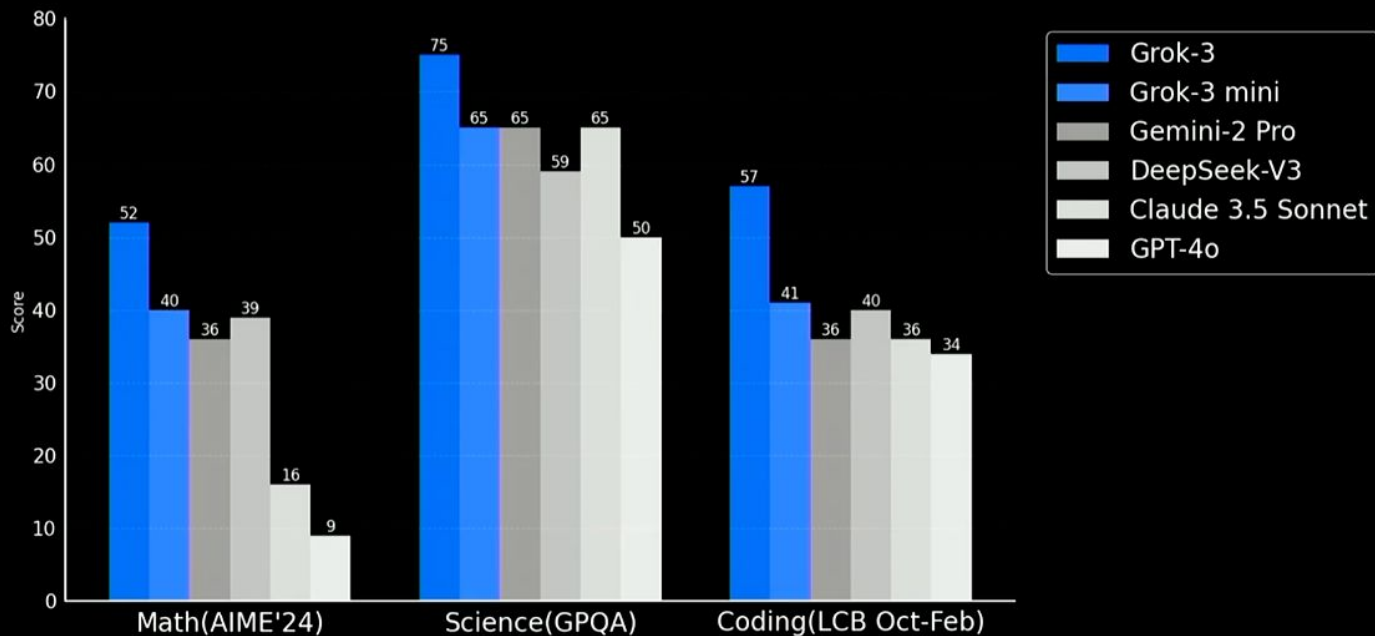
Grok-3 came out last night

200,000 GPUs



Grok-3

Benchmarks



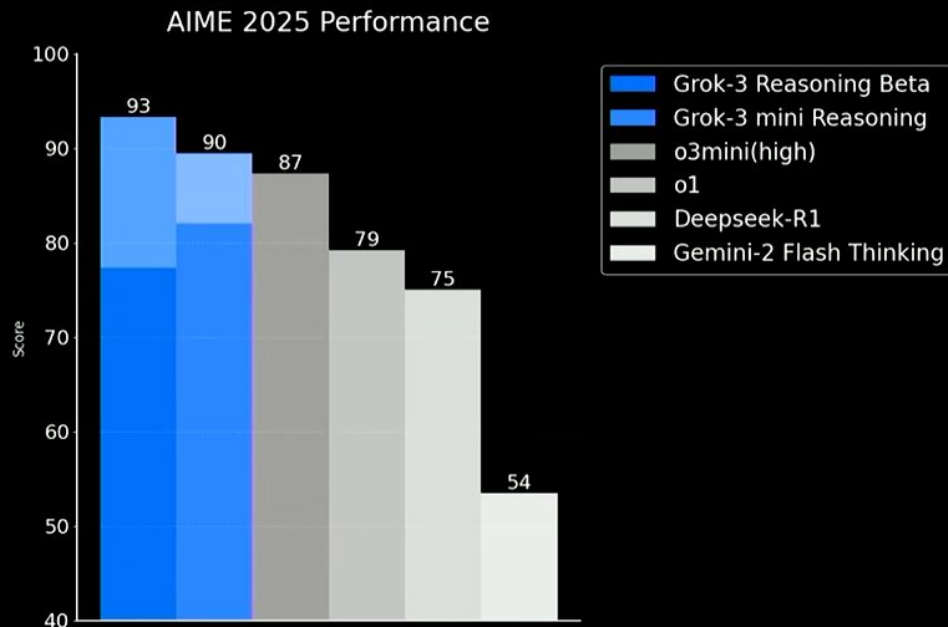
LIVE

1.5M views



Grok-3 (cont'd)

Reasoning + Test-Time Compute

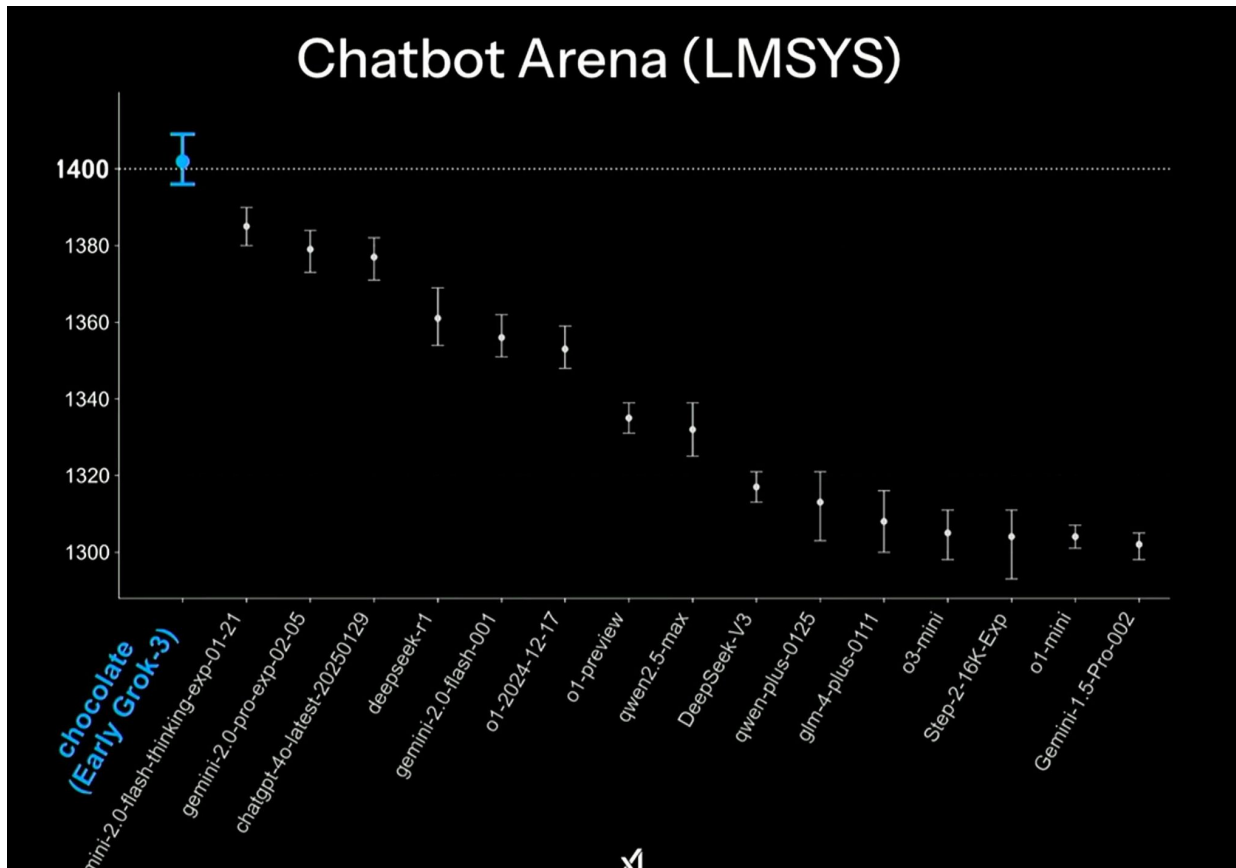


LIVE

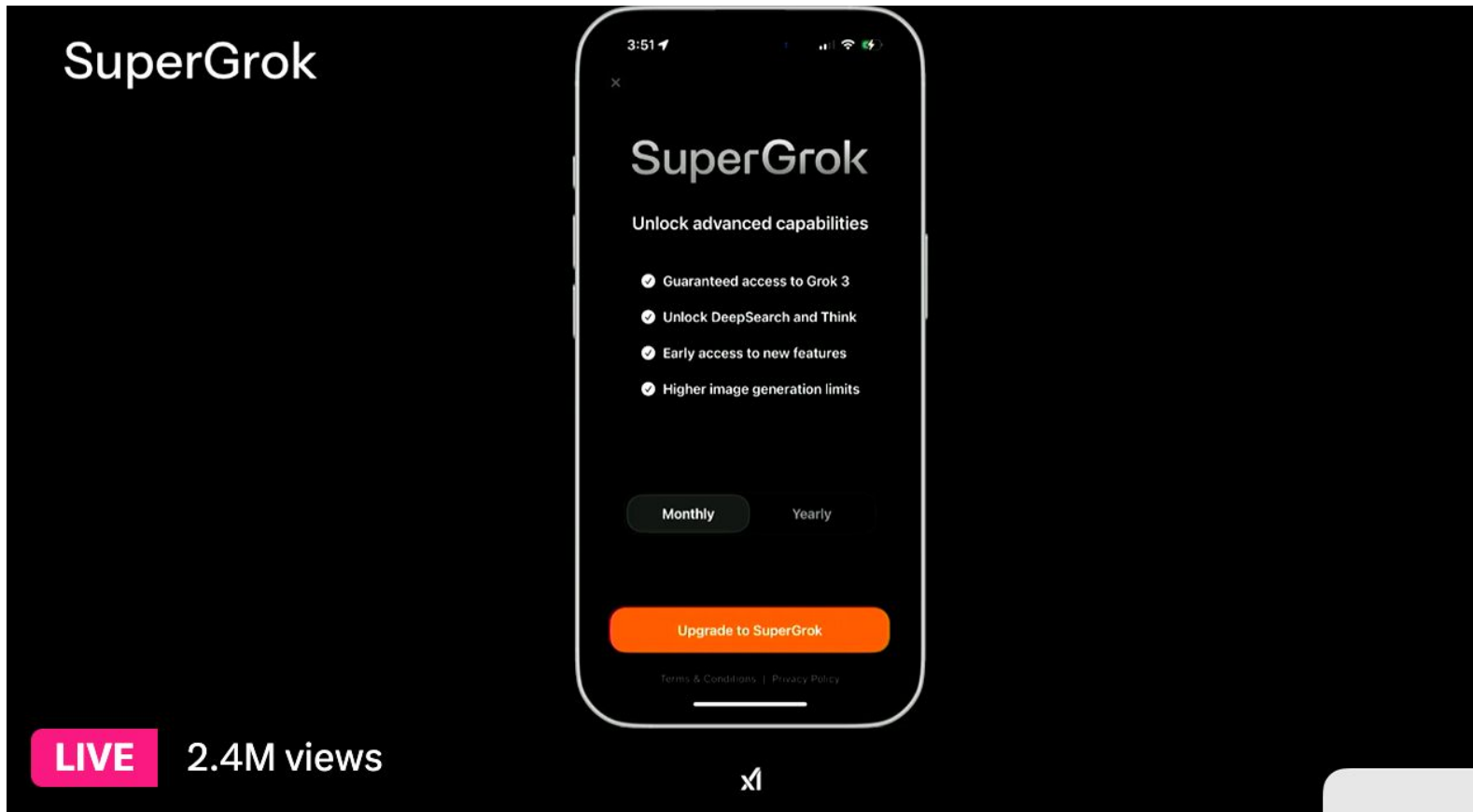
1.8M views

x1

Grok-3 (cont'd)



Grok-3 (cont'd)



The image shows a mobile application interface for SuperGrok. On the left, the text "SuperGrok" is displayed in white. In the center, a rounded rectangular modal is shown, representing a mobile phone screen. The modal has a dark background and contains the following elements:

- At the top left of the modal, a small "x" icon for closing the modal.
- The "SuperGrok" logo in white.
- The heading "Unlock advanced capabilities" in white.
- A list of four benefits, each preceded by a white checkmark icon:
 - Guaranteed access to Grok 3
 - Unlock DeepSearch and Think
 - Early access to new features
 - Higher image generation limits
- Two subscription options: "Monthly" (highlighted with a dark grey background) and "Yearly" (in white text).
- A large orange button with the text "Upgrade to SuperGrok".
- At the bottom, small white text for "Terms & Conditions" and "Privacy Policy" separated by a vertical bar.

Below the modal, on the left, there is a pink rectangular badge with the word "LIVE" in white, followed by the text "2.4M views" in white. At the bottom center, there is a small white logo that resembles a stylized "X" or a similar symbol.

Big tech companies are considering open-sourcing older AI models



Sam Altman



@sama



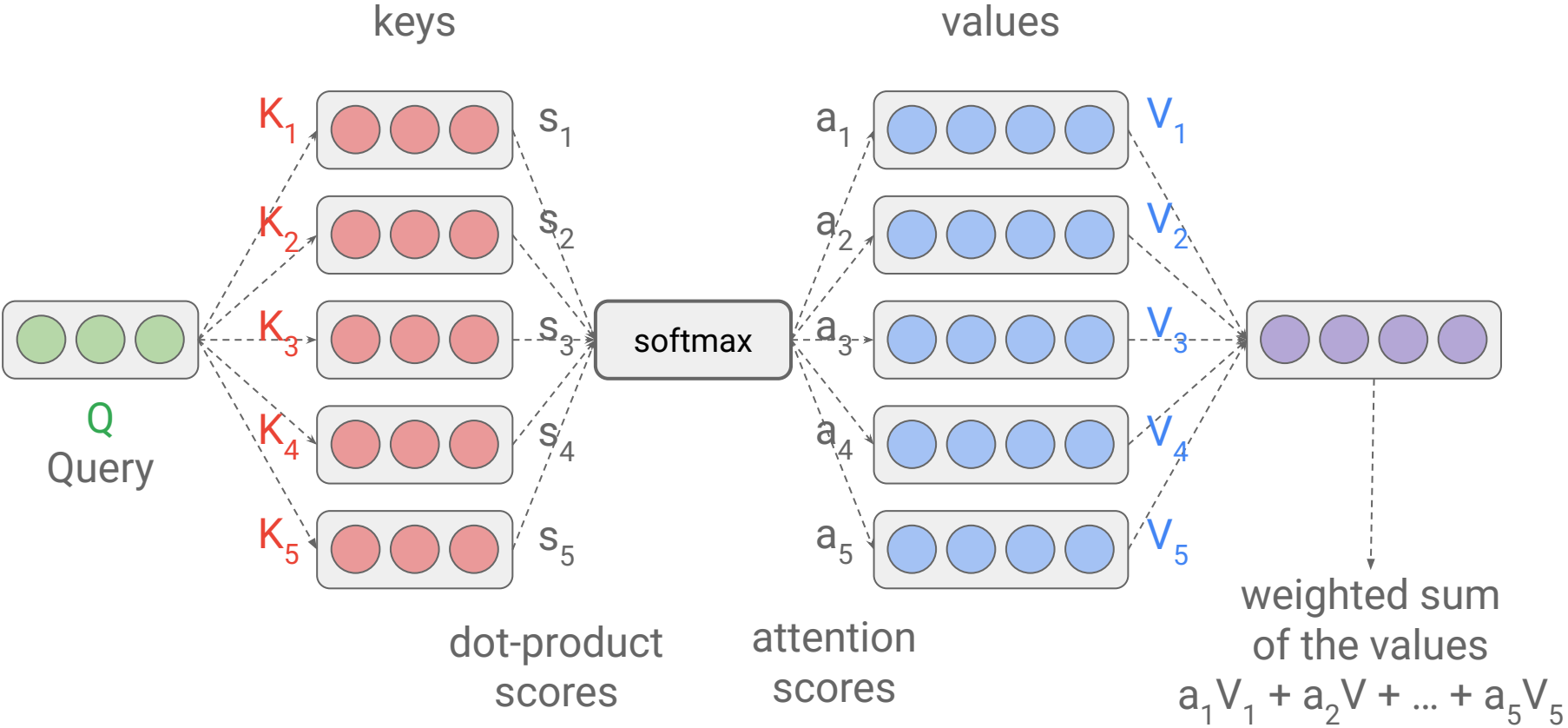
for our next open source project, would it be more useful to do an o3-mini level model that is pretty small but still needs to run on GPUs, or the best phone-sized model we can do?

o3-mini

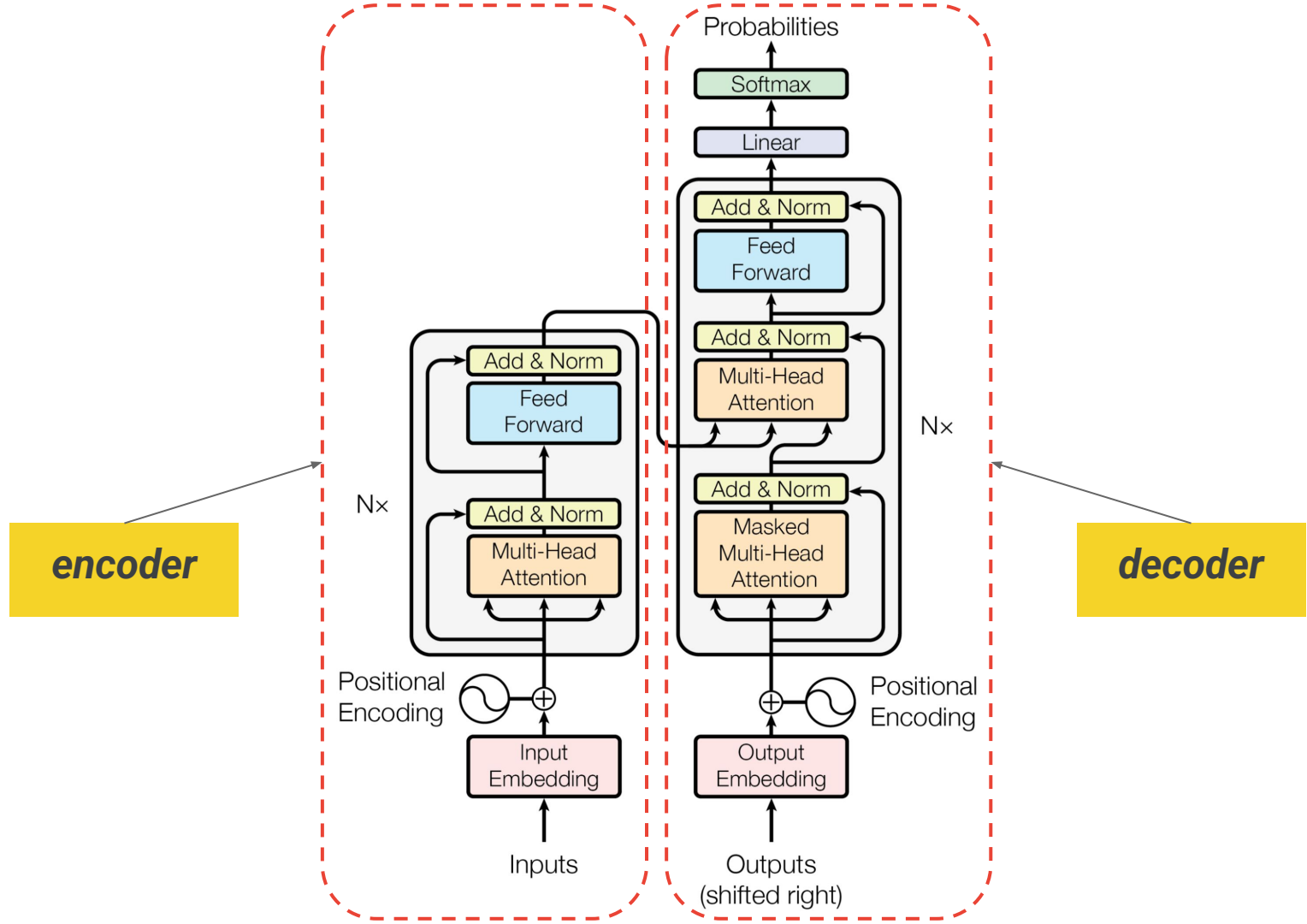
phone-sized model

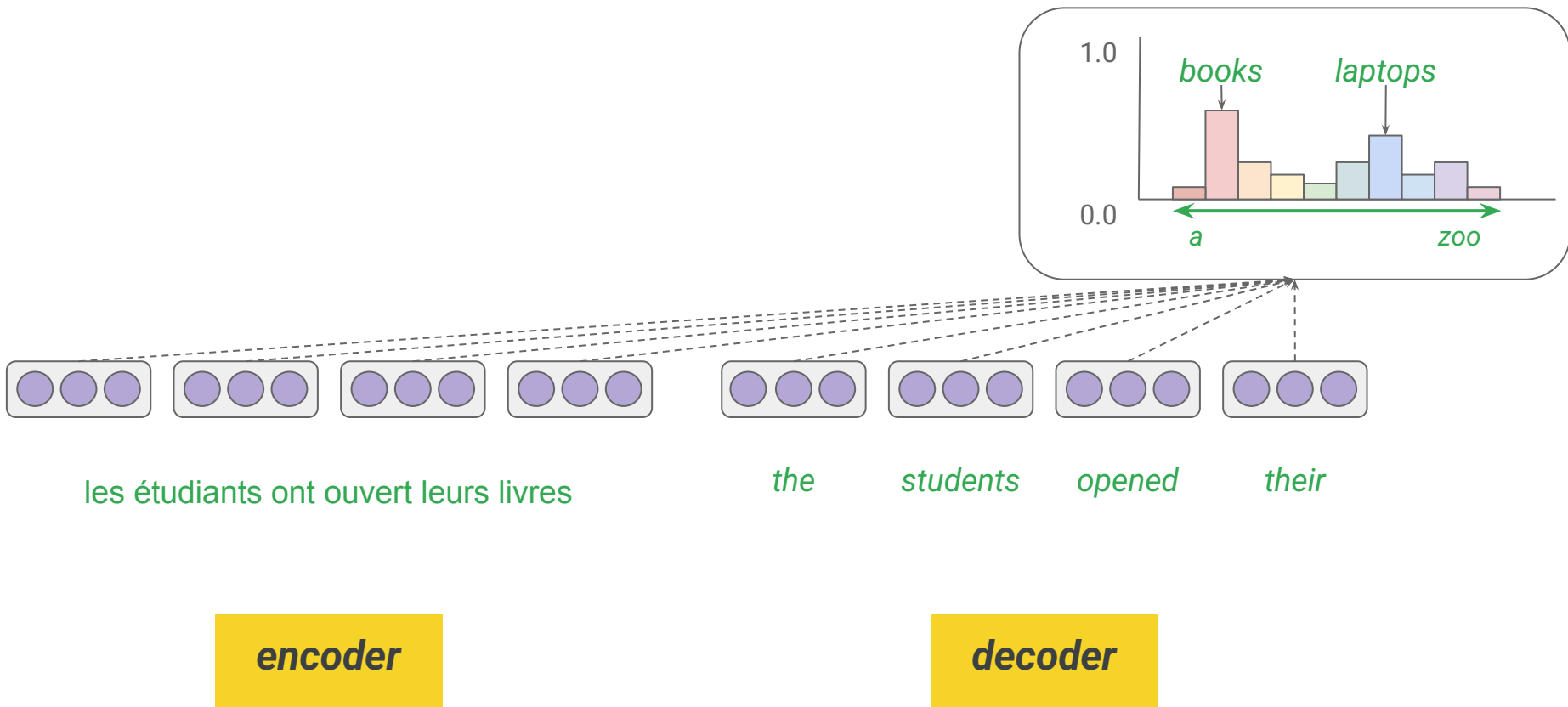
Transformers review

Attention mechanism

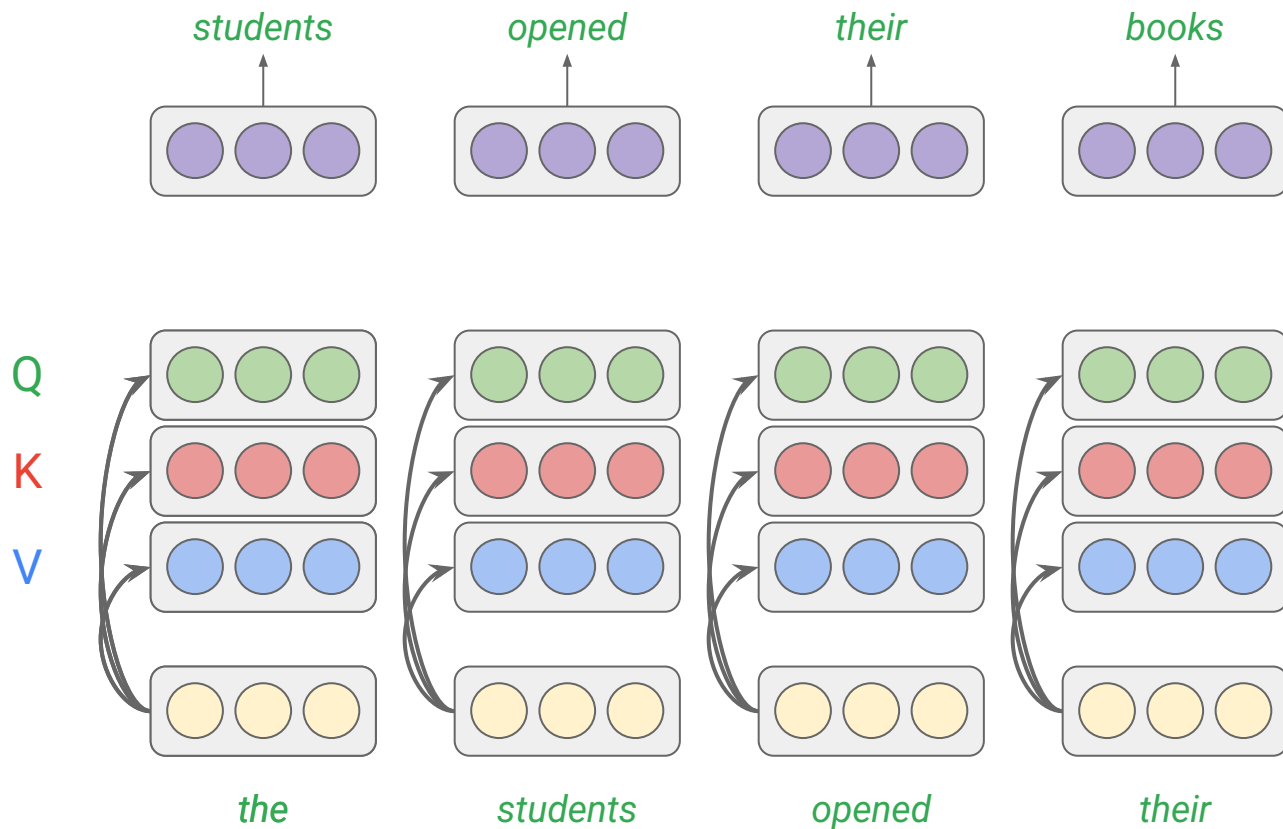


Transformers architecture





Self-attention



$$Q = X \cdot W_Q$$

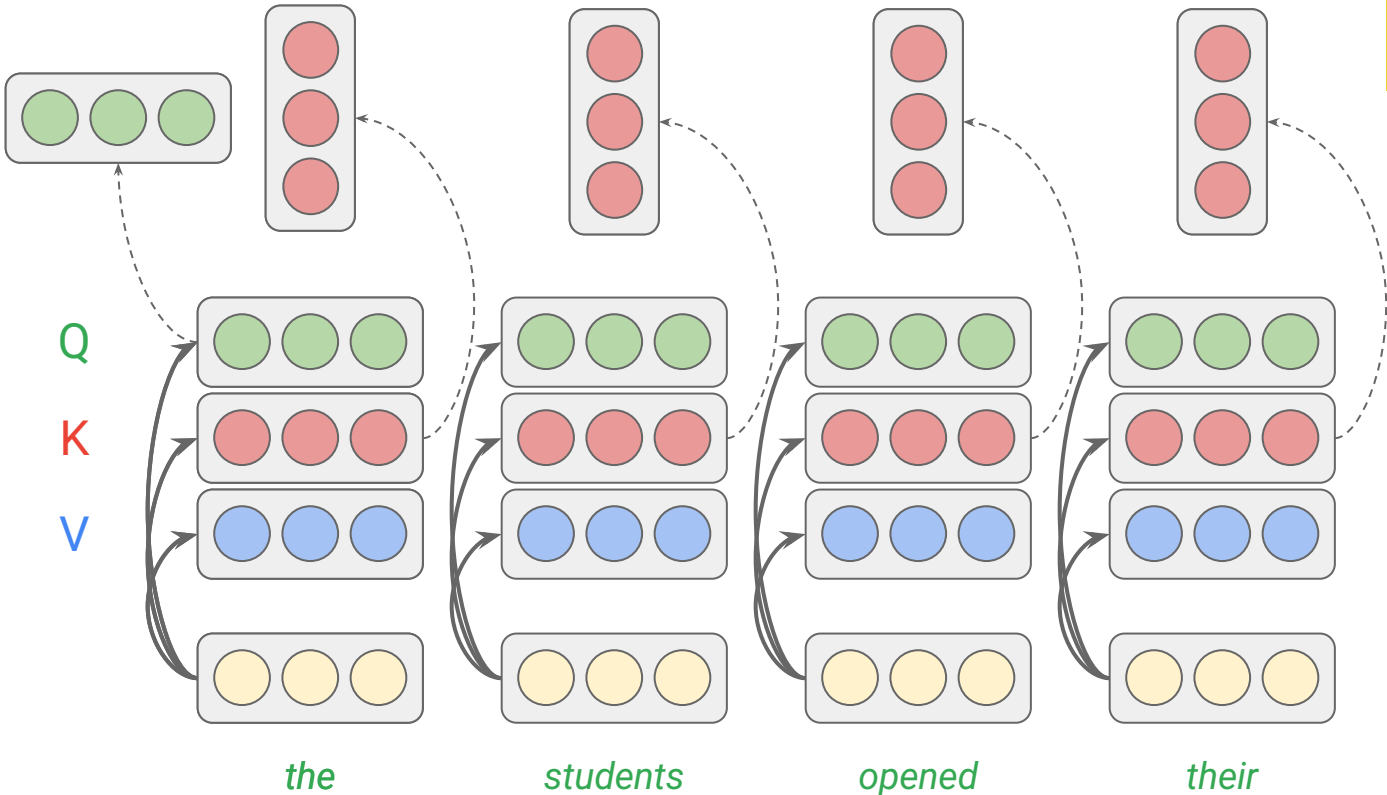
$$K = X \cdot W_K$$

$$V = X \cdot W_V$$

linear
projections

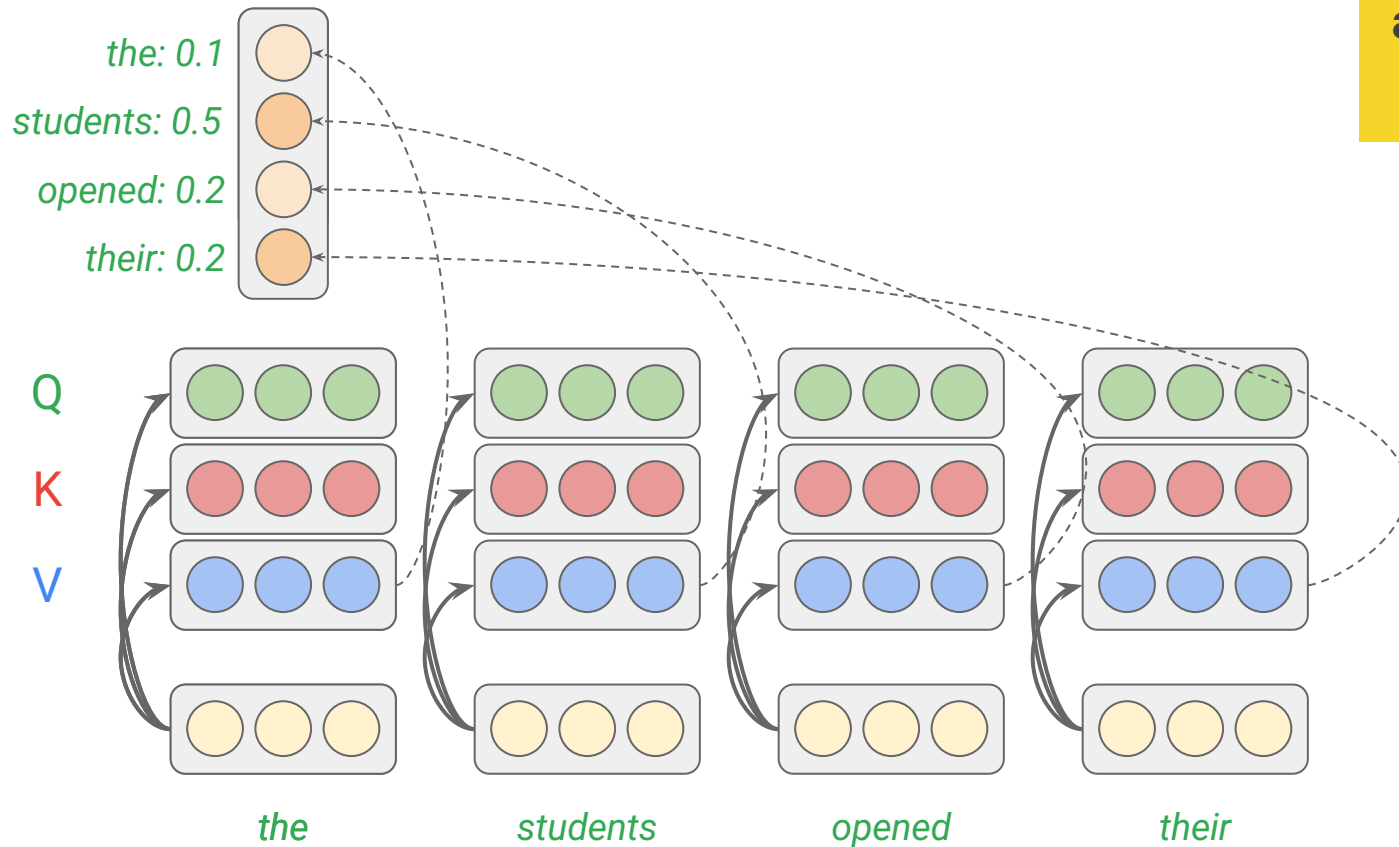
Self-attention (cont'd)

all computations are parallelized



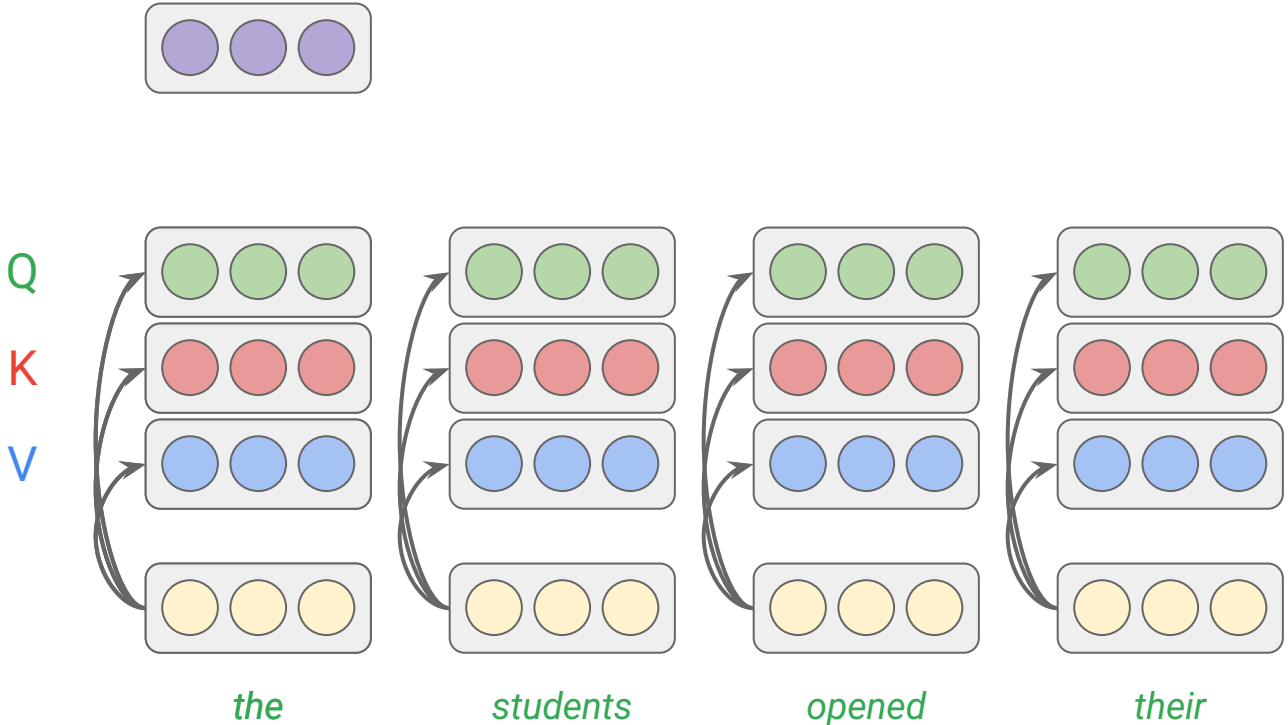
Self-attention (cont'd)

all computations are parallelized



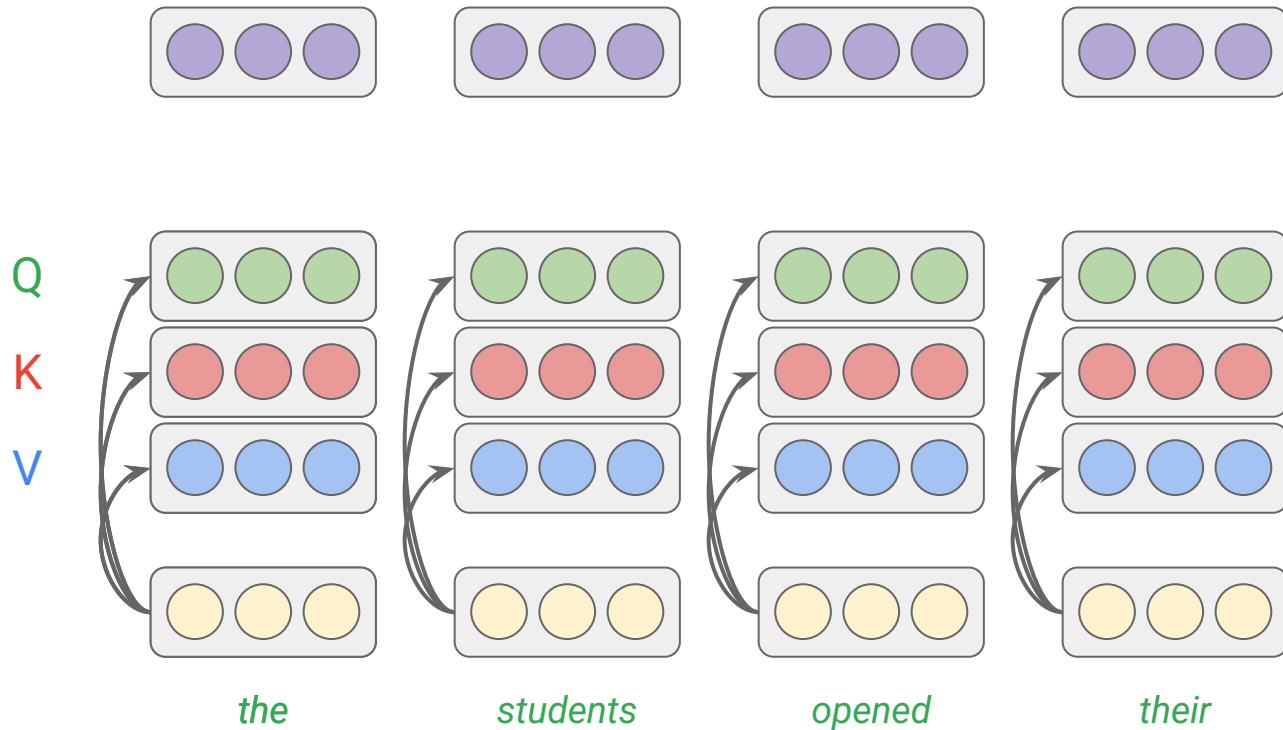
Self-attention (cont'd)

all computations are parallelized



Self-attention (cont'd)

all computations
are parallelized



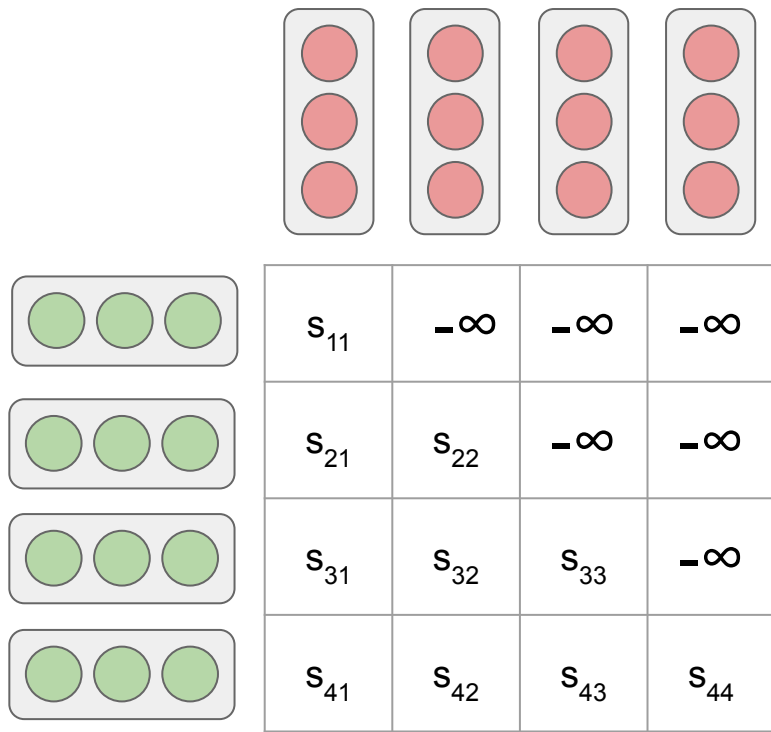
All computations are parallelized

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

d_k : scaling factor

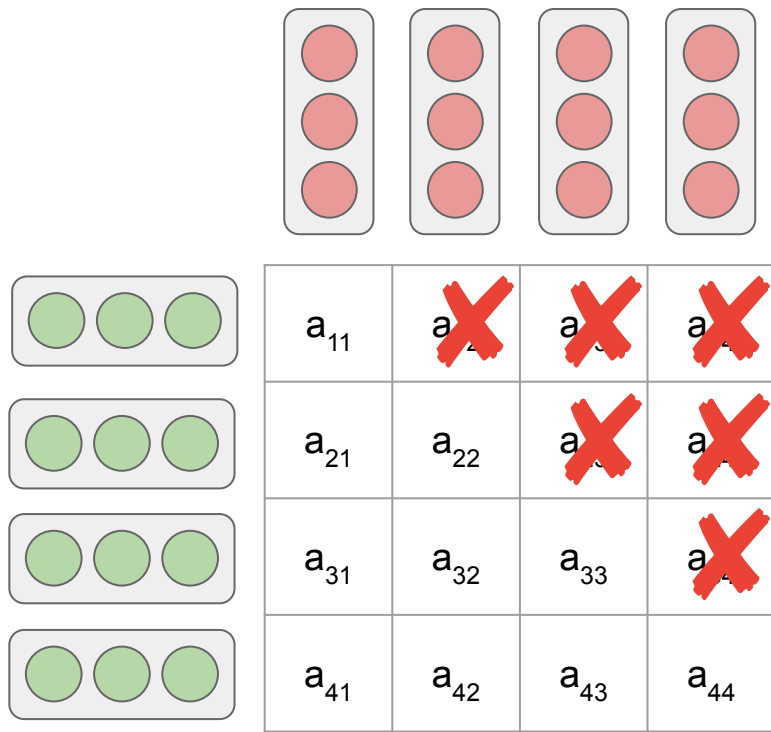
large products push the softmax function into regions where it has extremely small gradients

Self-attention in the decoder



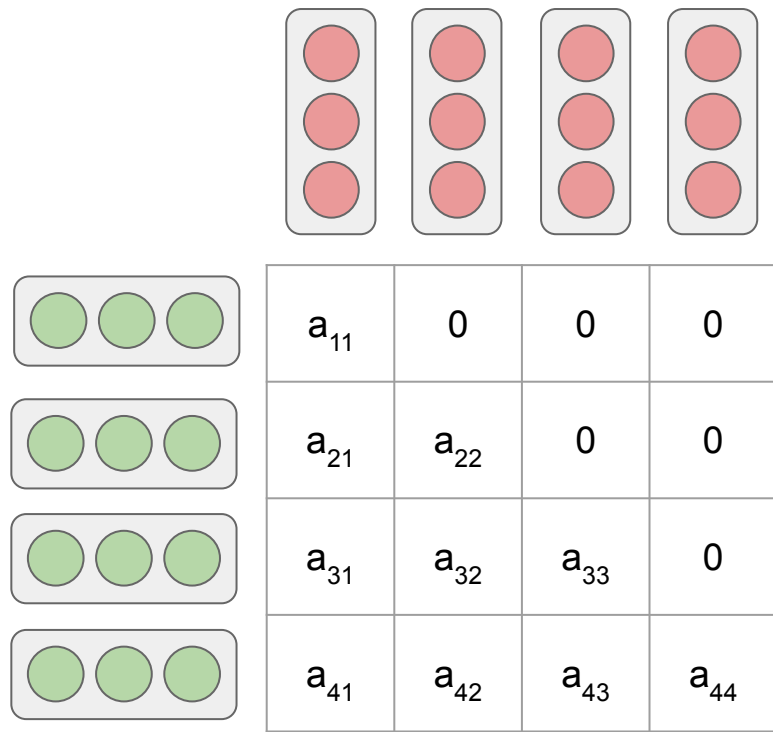
masking out (setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections

Self-attention in the decoder (cont'd)



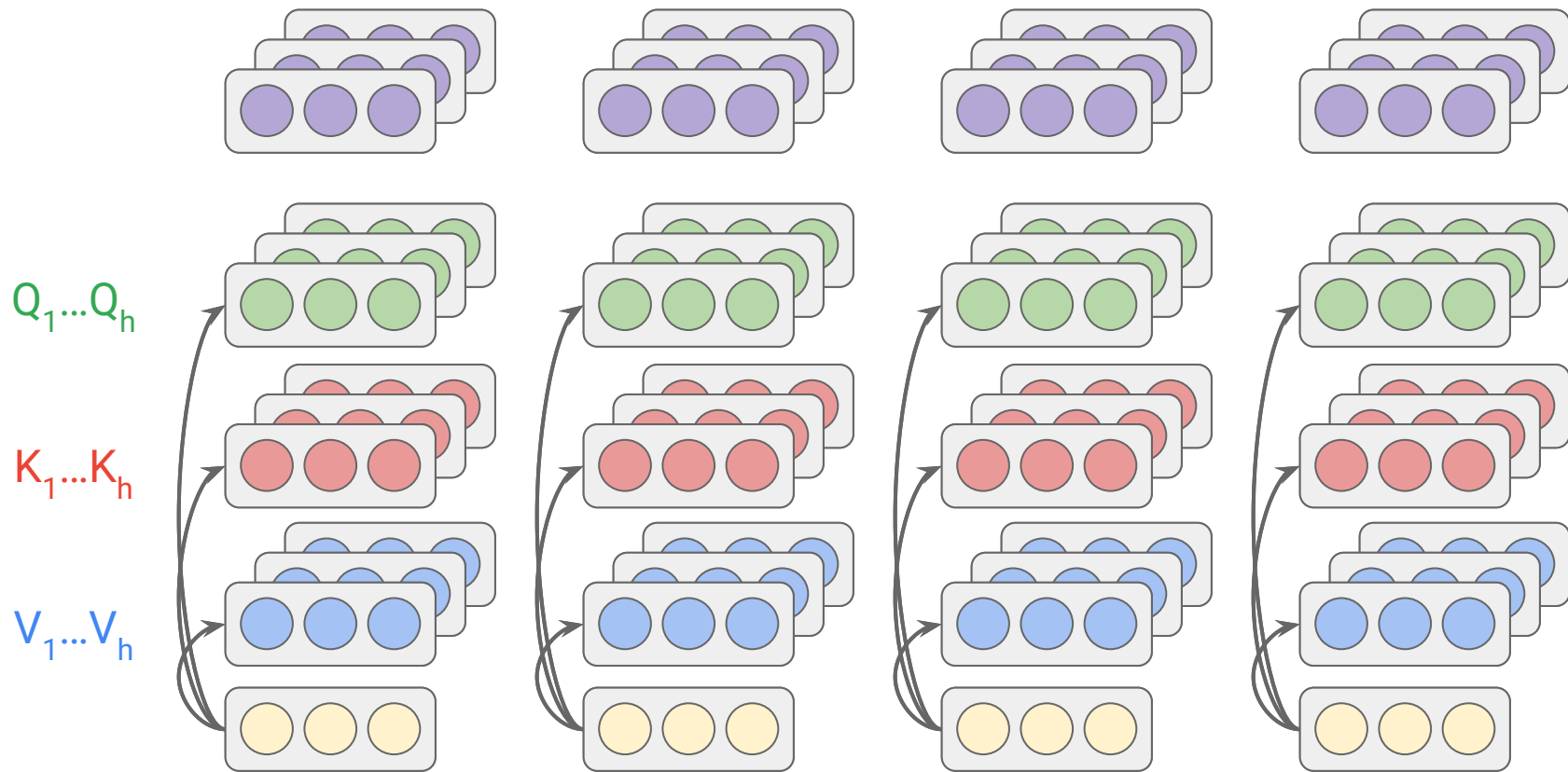
masking out all values in the input of the softmax which correspond to illegal connections

Self-attention in the decoder (cont'd)

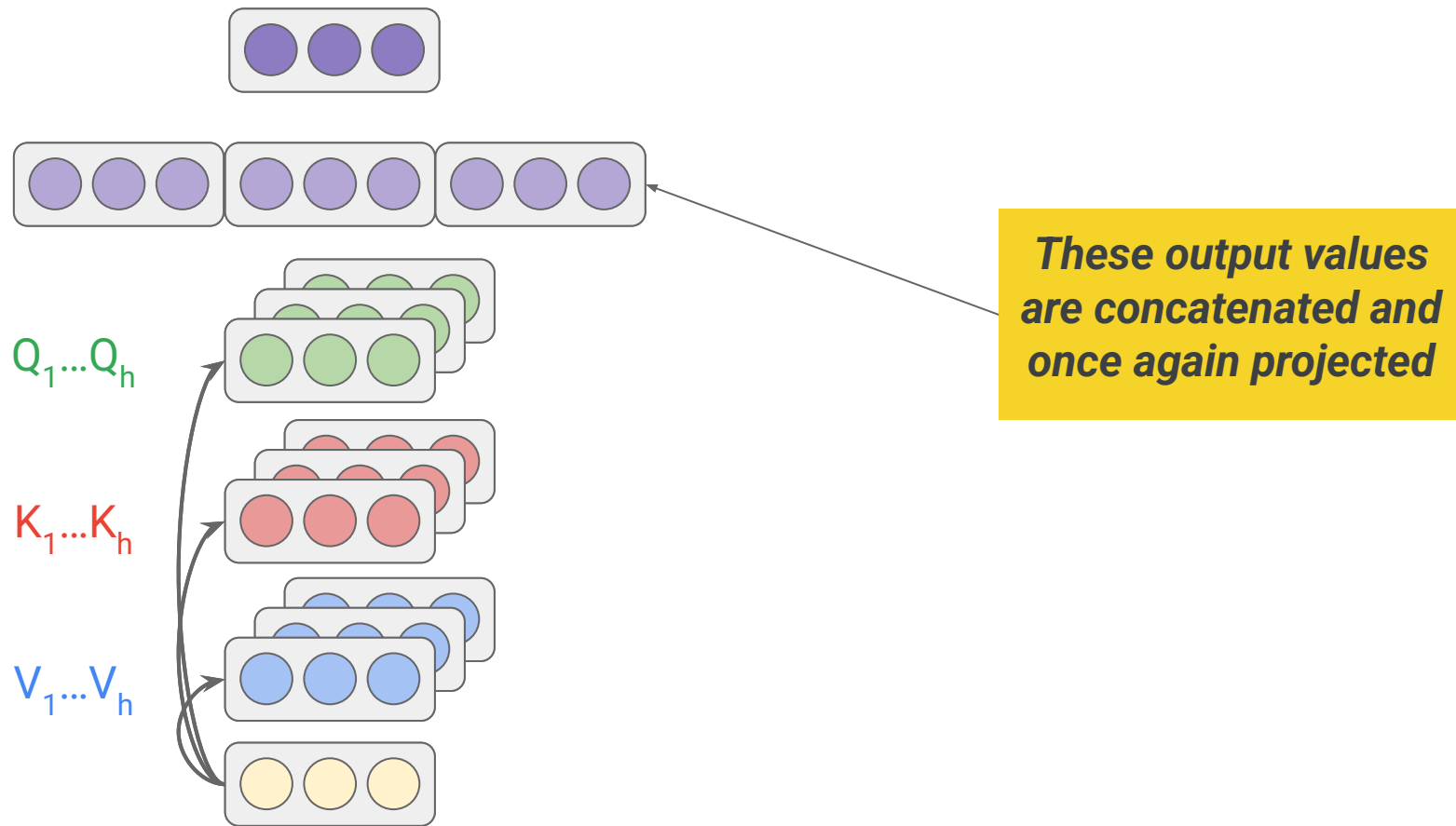


masking out all values in the input of the softmax which correspond to illegal connections

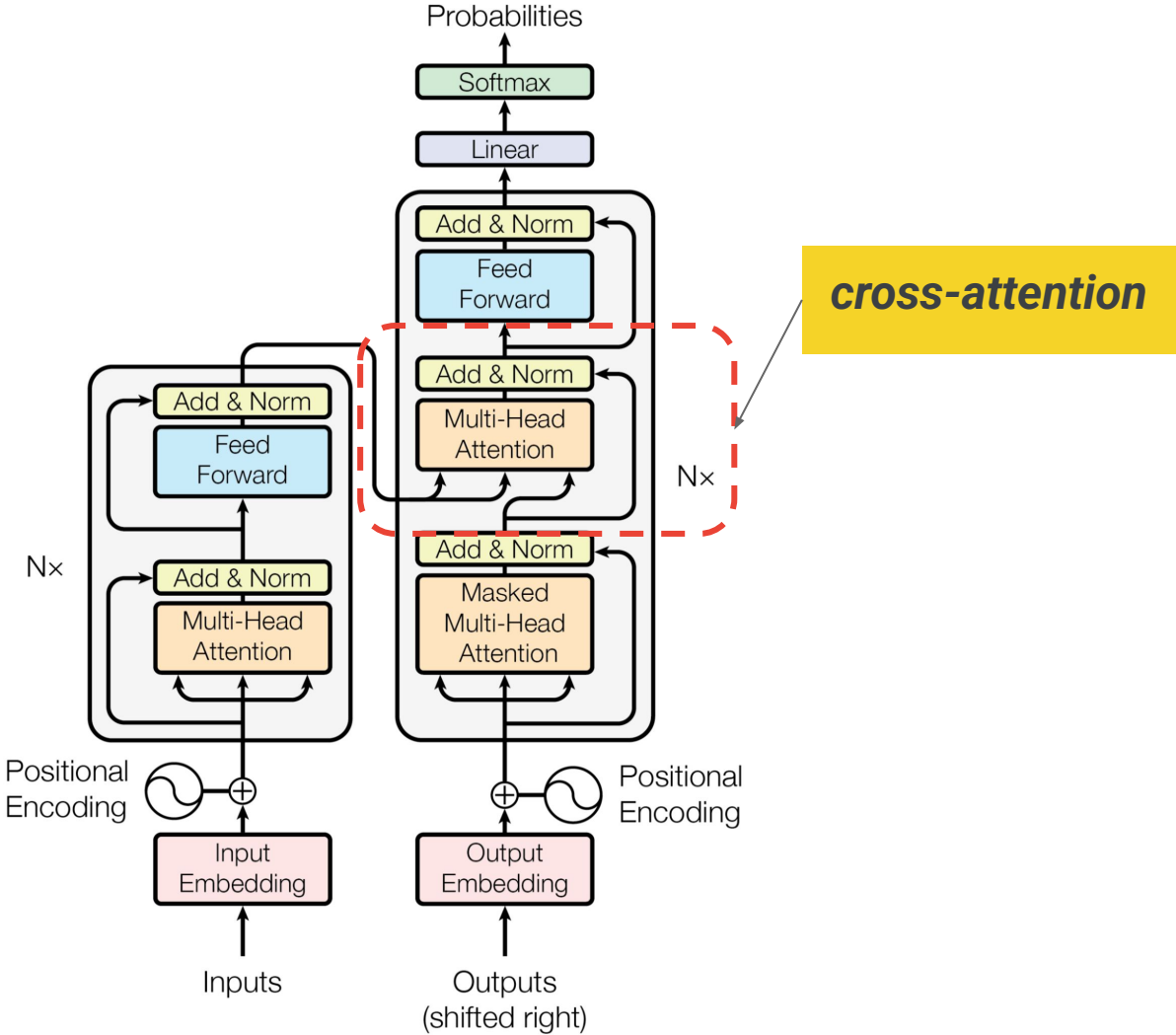
Multi-head attention



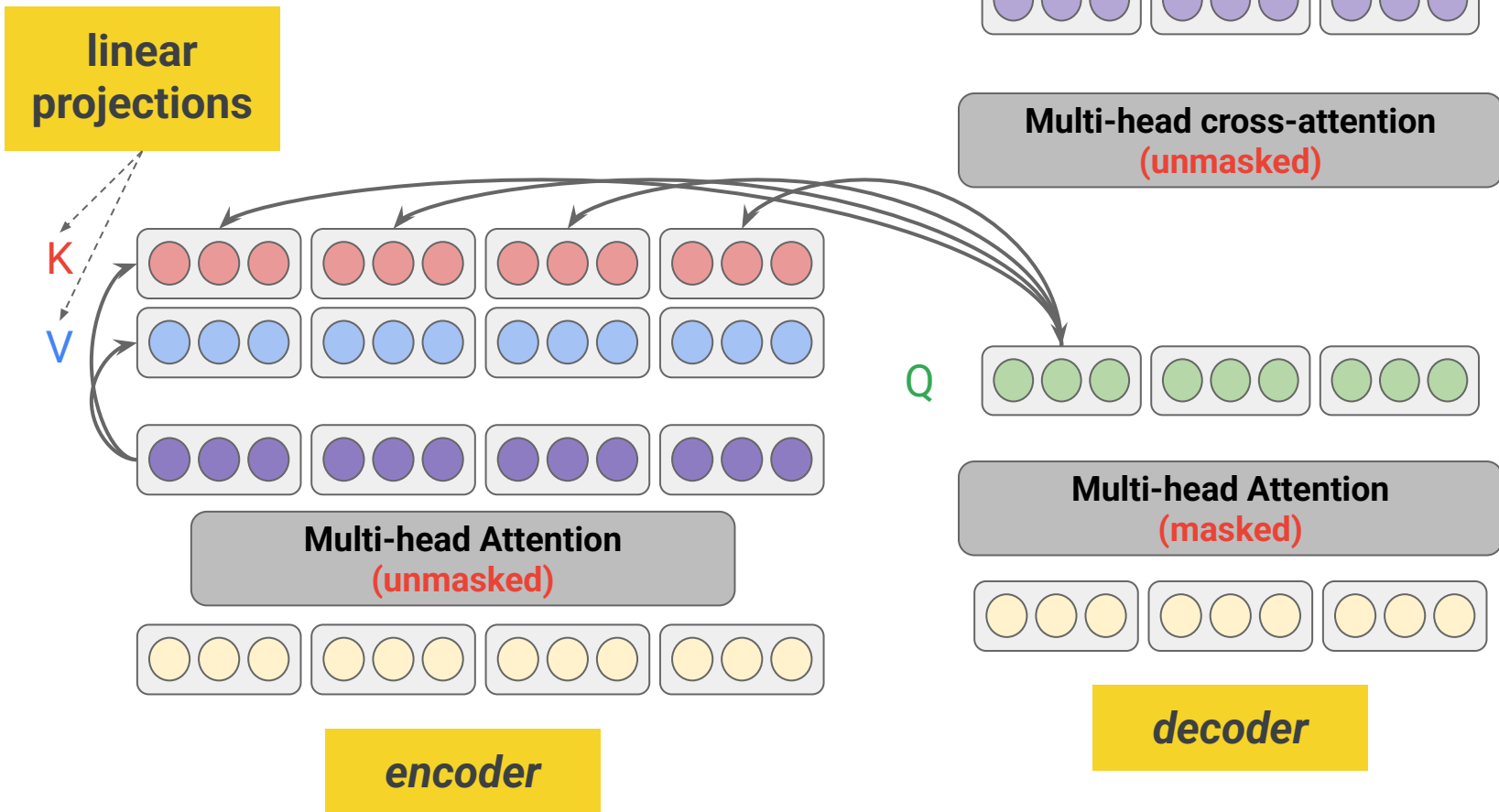
Multi-head attention (cont'd)



Cross-attention in the decoder



Cross-attention in the decoder

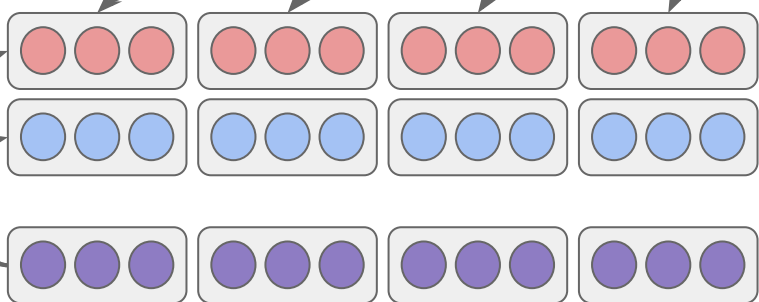


Cross-attention in the decoder (cont'd)

linear projections

K

V



Multi-head Attention
(unmasked)

encoder



Multi-head cross-attention
(unmasked)

Q

Multi-head Attention
(masked)

decoder

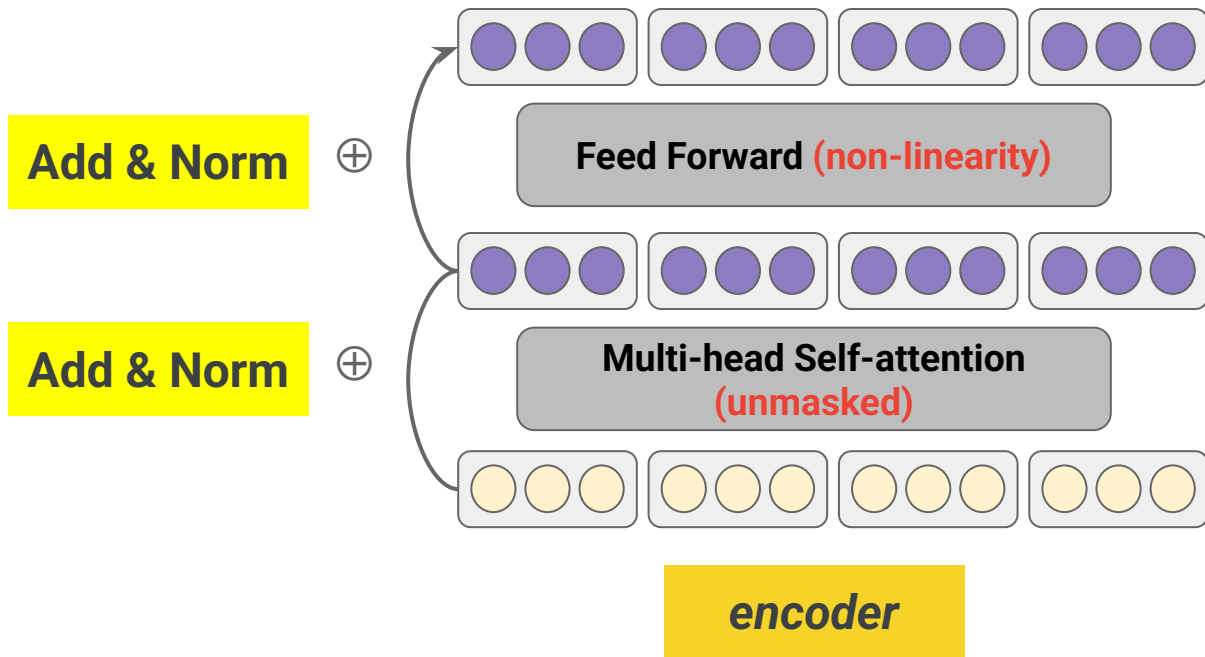
residual connections

+

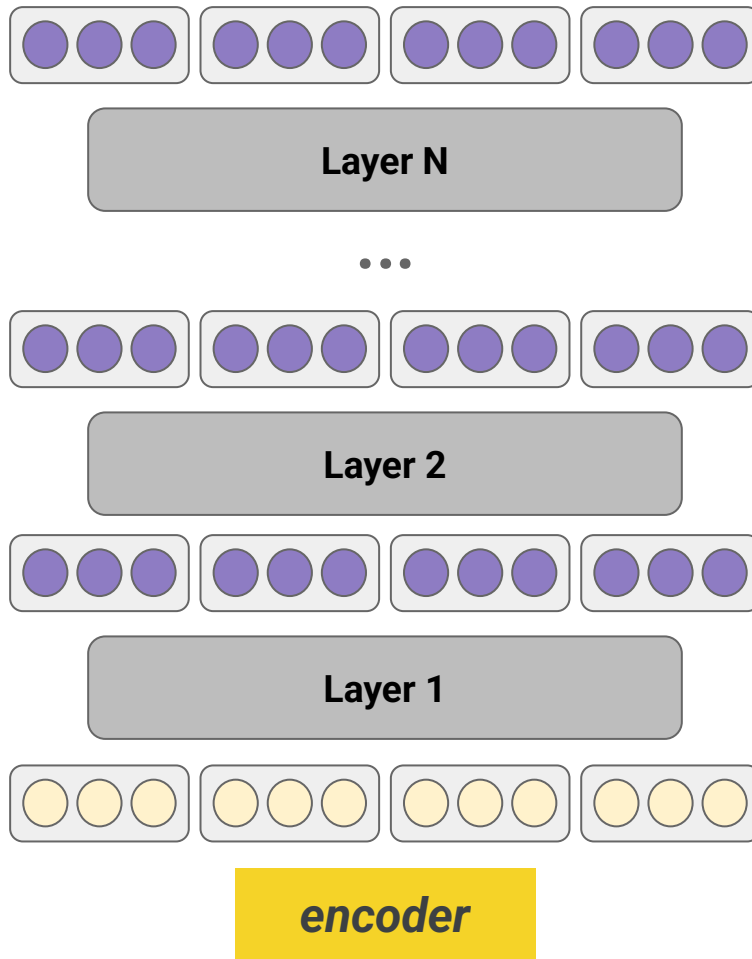
+

residual connections

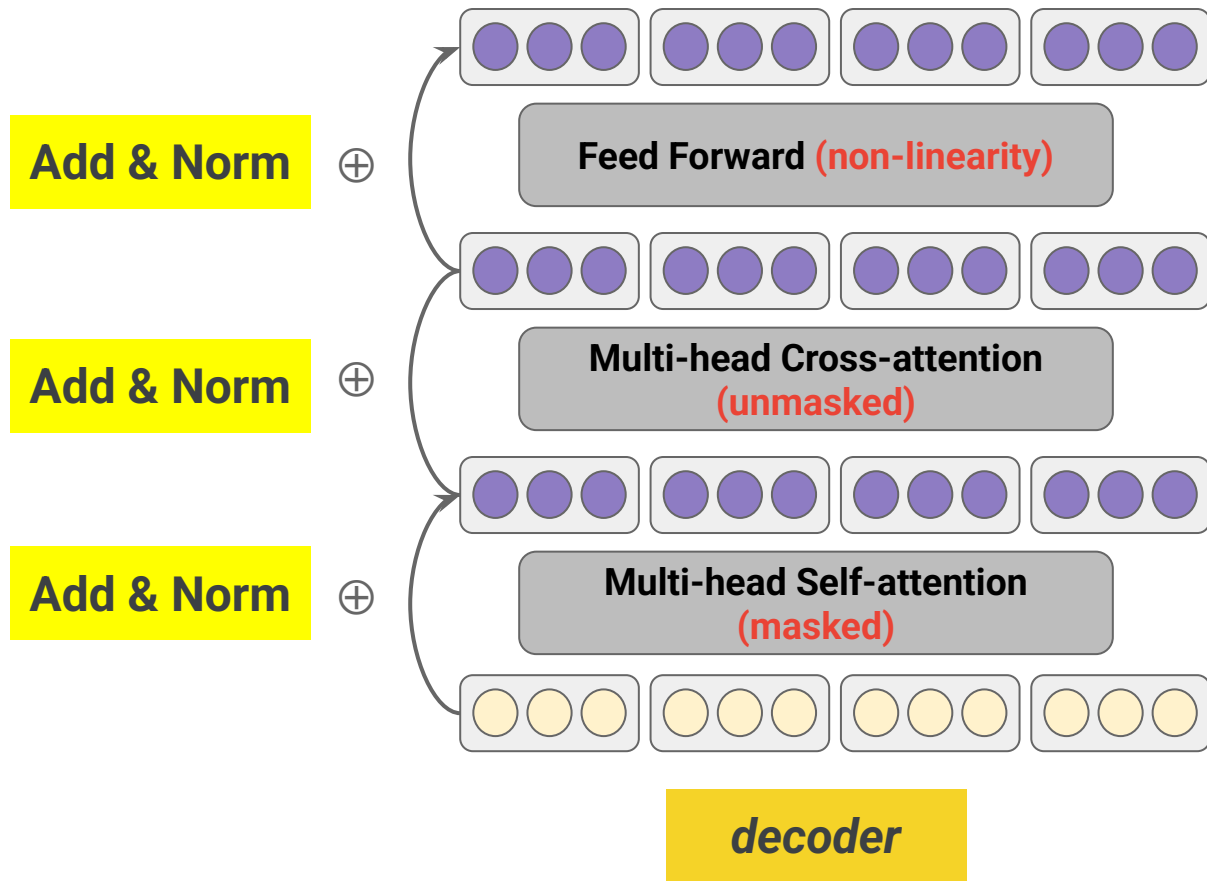
Encoder (one layer)



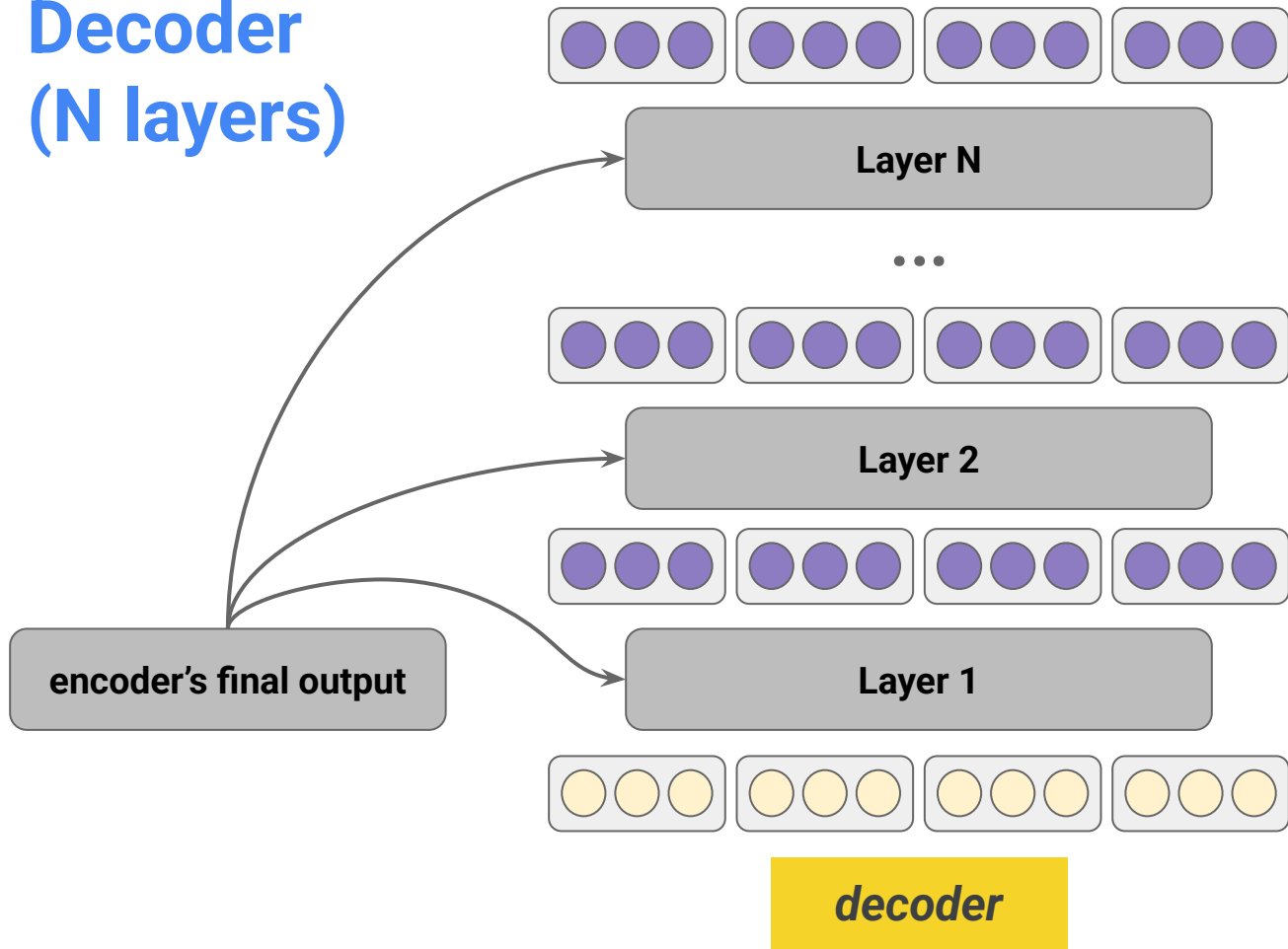
Encoder (N layers)



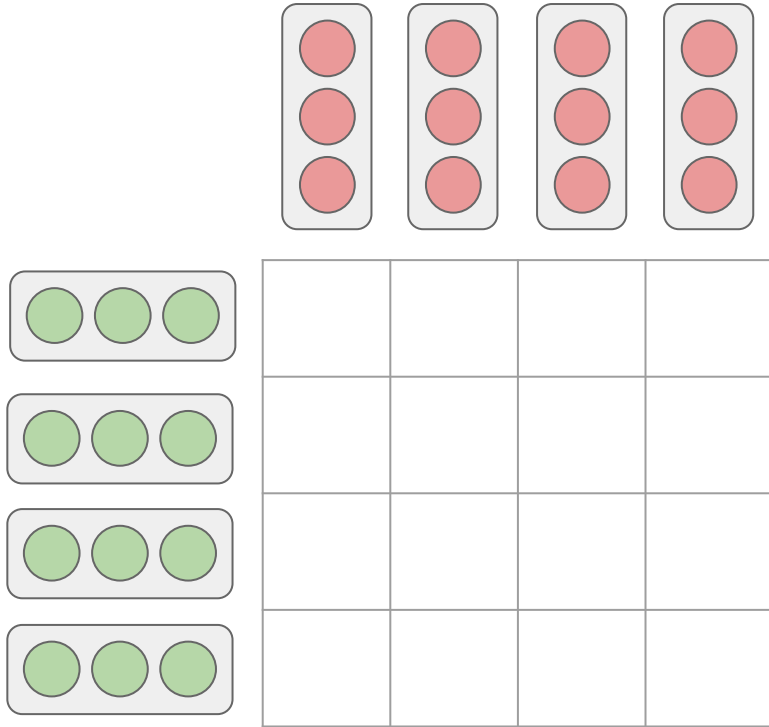
Decoder (one layer)



Decoder (N layers)



Quadratic complexity



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

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

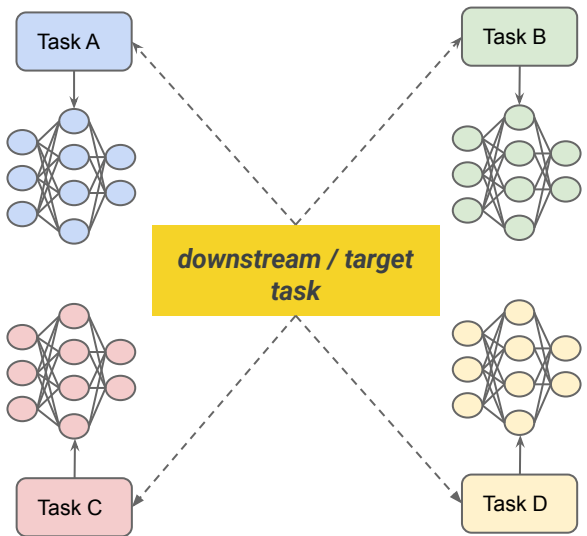
Google AI Language

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

A learning paradigm shift

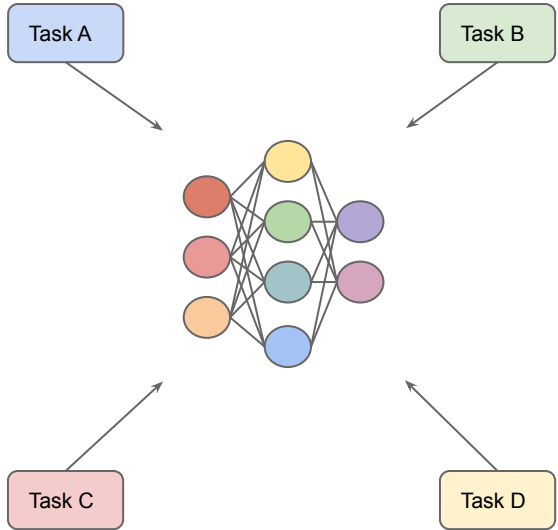


training task-specific models from scratch



before BERT

pretraining and then adapting



since BERT

Neural network diagrams adapted from Colin Raffel's talk at Stanford MLSys Seminars



Image created by Gemini

Deep contextualized word representations

Matthew E. Peters[†], Mark Neumann[†], Mohit Iyyer[†], Matt Gardner[†],
`{matthewp, markn, mohiti, mattg}@allenai.org`

Christopher Clark*, Kenton Lee*, Luke Zettlemoyer^{†*}
`{csquared, kentonl, lsz}@cs.washington.edu`

[†]Allen Institute for Artificial Intelligence

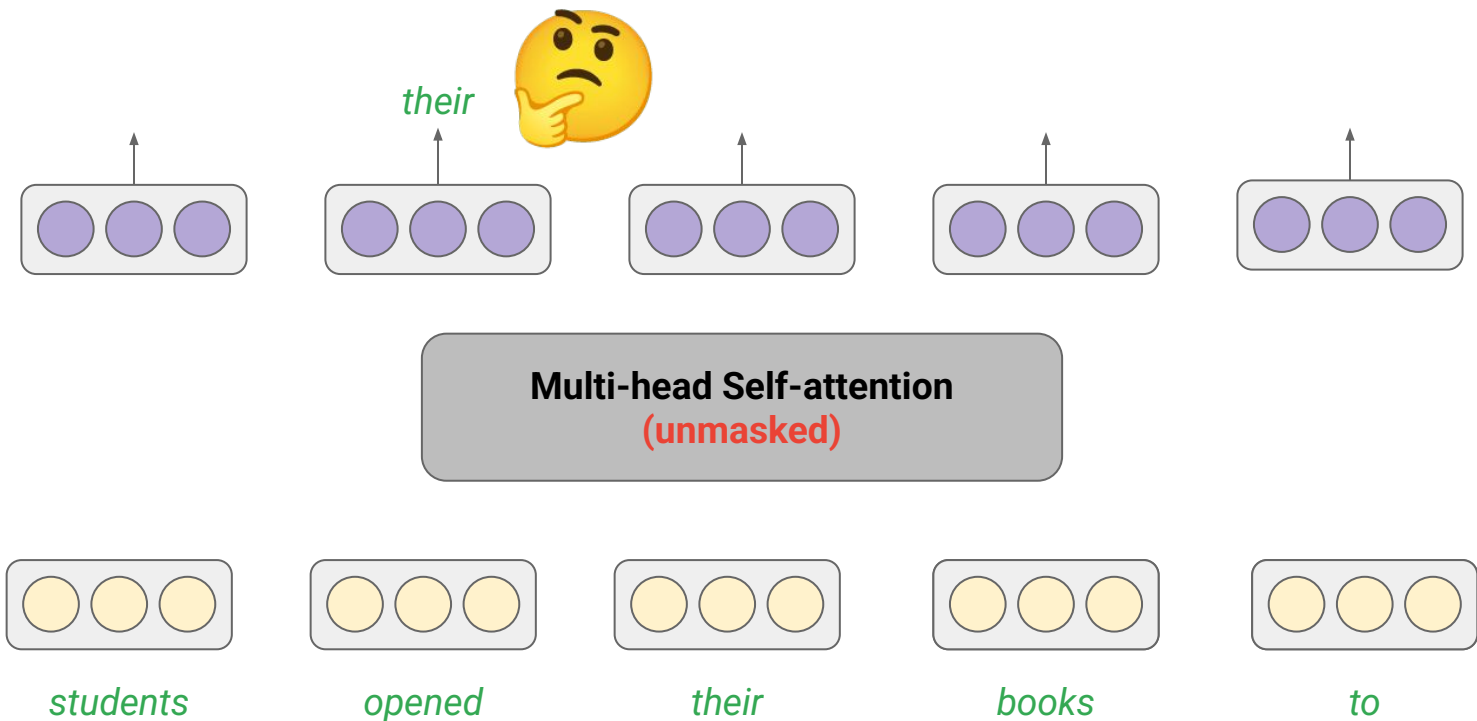
*Paul G. Allen School of Computer Science & Engineering, University of Washington

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)

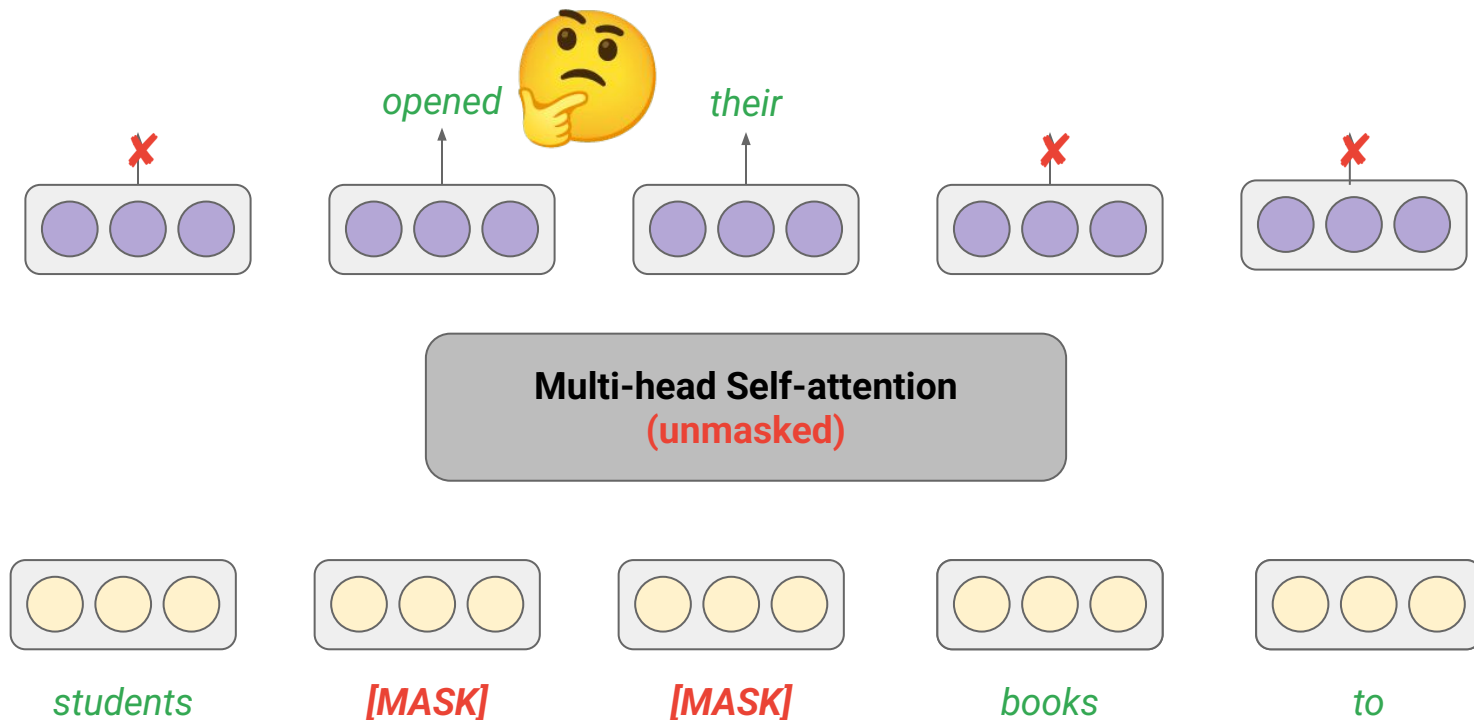
Pretraining

Language modeling using a Transformer encoder



Masked language modeling

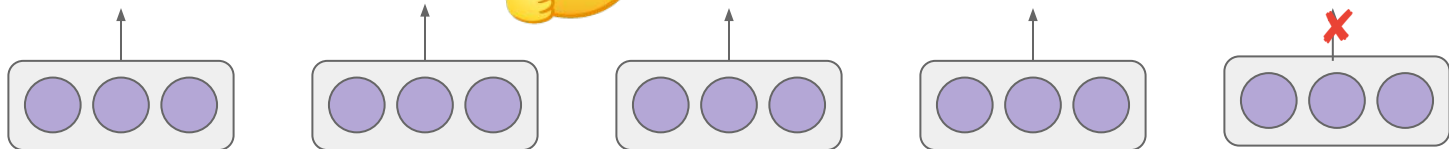
cloze task



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

What if we mask more tokens?

cloze task



Multi-head Self-attention
(unmasked)



[MASK]

[MASK]

[MASK]

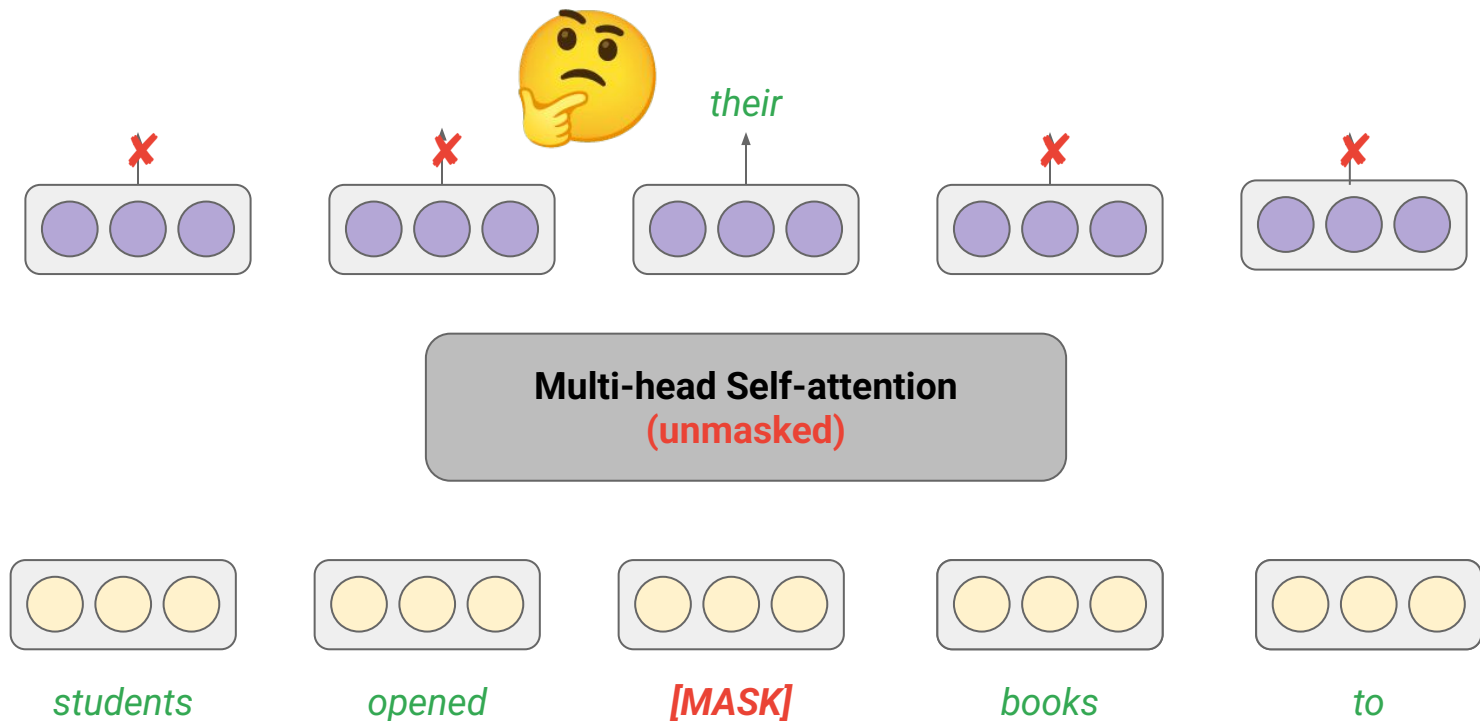
[MASK]

to

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

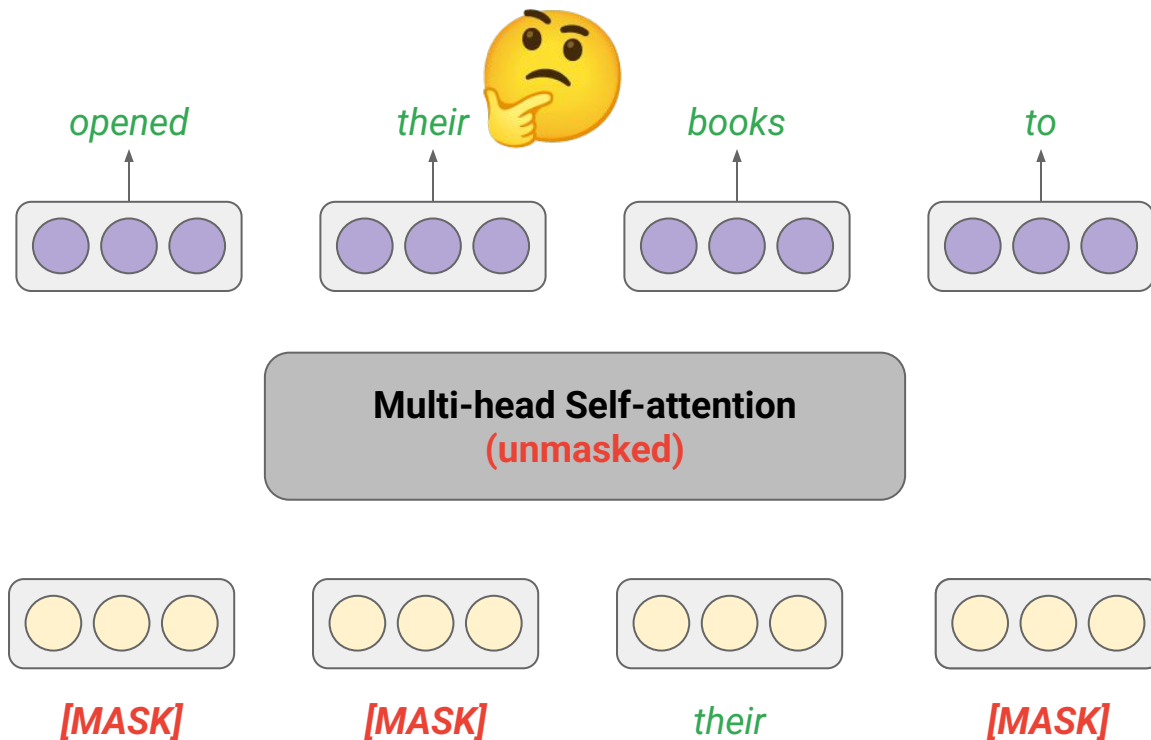
What if we mask less tokens?

cloze task



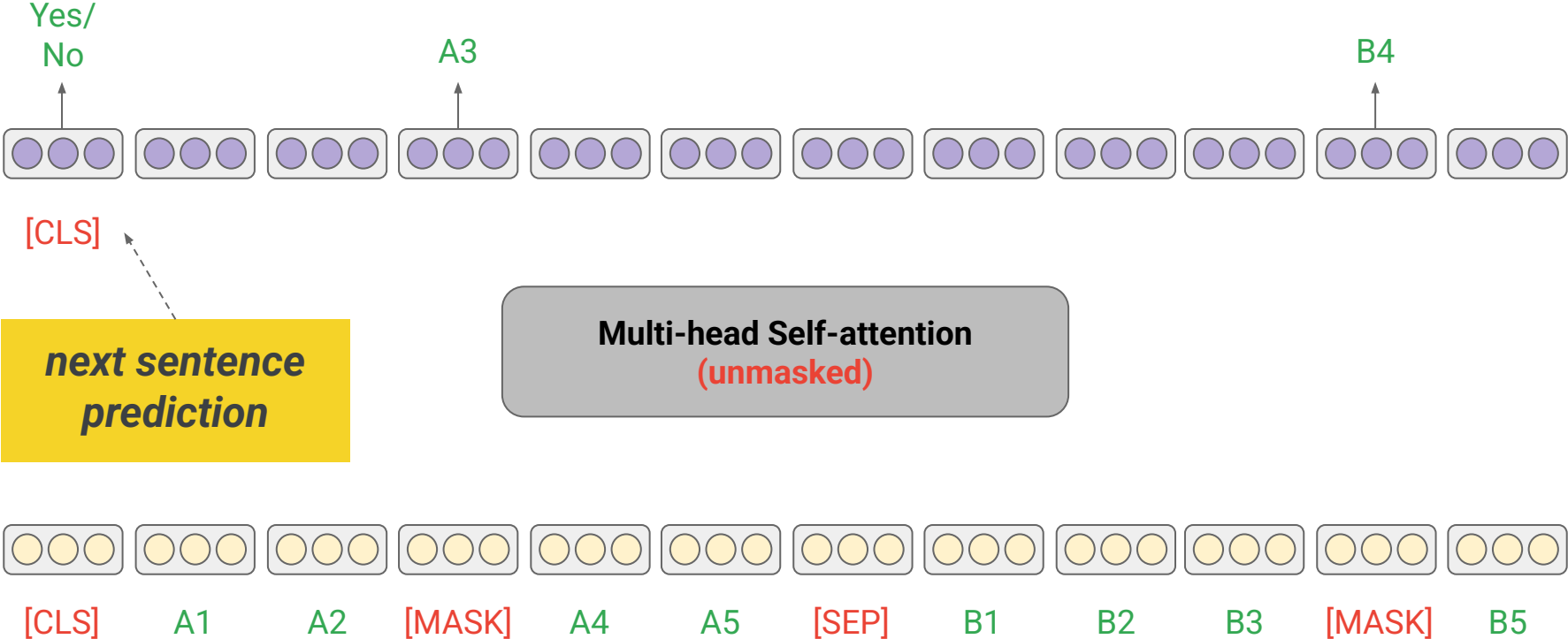
15% - 30% of all tokens in each sequence are masked at random

What if we mask more tokens?

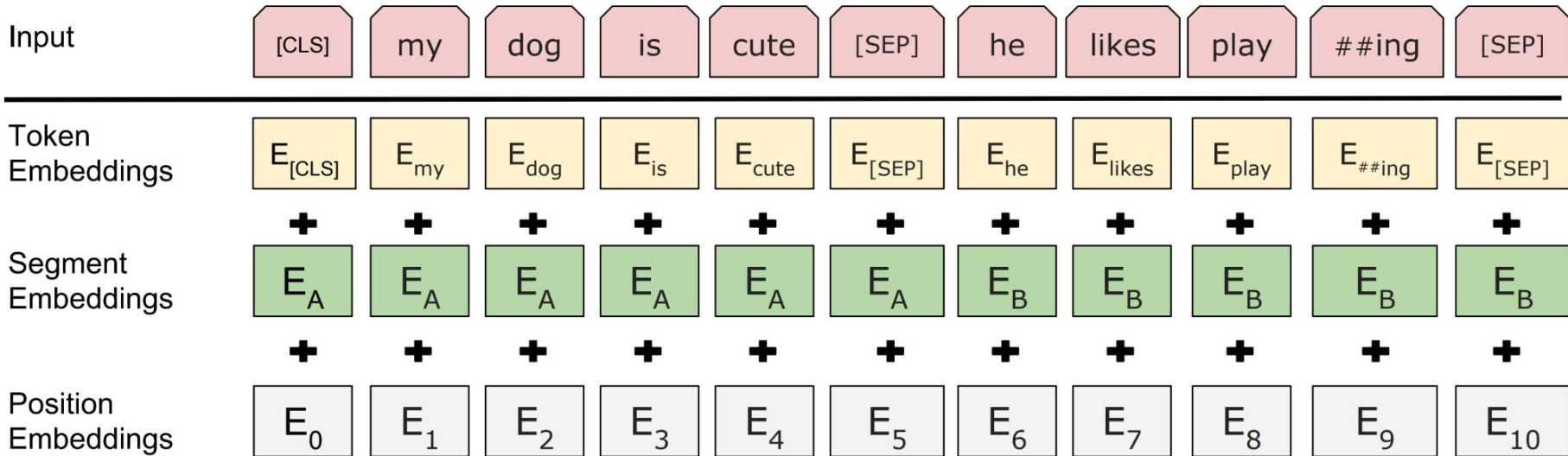


15% - 30% of all tokens in each sequence are masked at random

CLS & SEP tokens



BERT input representation



Fine-tuning

softmax

linear

[CLS]



Multi-head Self-attention
(unmasked)

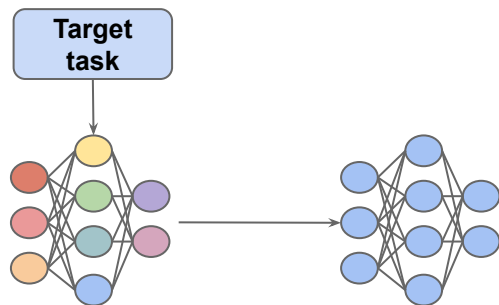


[CLS]

the movie was good

Intermediate-task transfer / fine-tuning

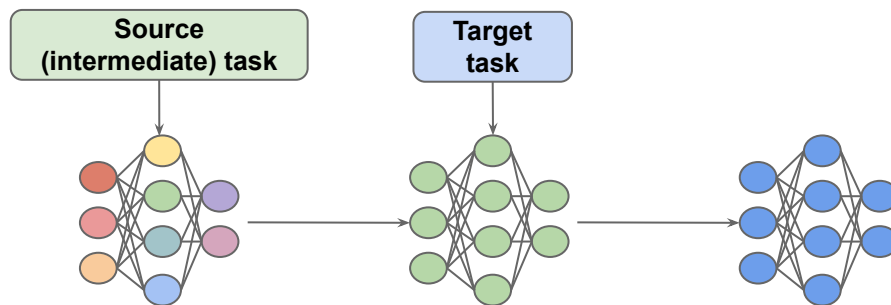
Standard Fine-tuning



BERT_{BASE} → RTE: 63.5 ± 2.3

BERT_{LARGE} → RTE: 68.6 ± 7.2

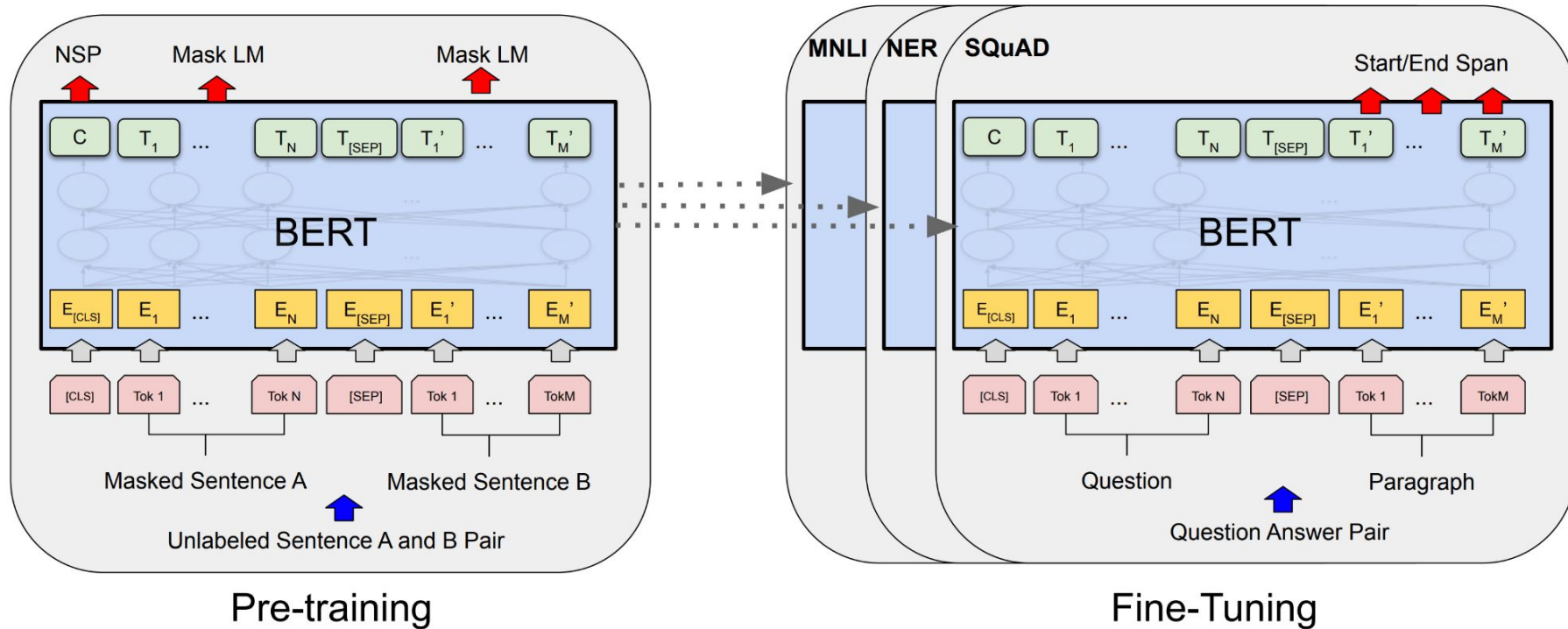
Intermediate-task transfer



BERT_{BASE} → MNLI → RTE: 78.1 ± 1.9

BERT_{LARGE} → MNLI → RTE: 82.3 ± 1.4

BERT Pretraining & Fine-tuning



Thank you!